



САНКТ-ПЕТЕРБУРГСКИЙ ГОСУДАРСТВЕННЫЙ  
УНИВЕРСИТЕТ ЭКОНОМИКИ И ФИНАНСОВ

# Эконометрика

**Учебник для магистров**

Под редакцией члена-корреспондента РАН  
**И. И. Елисеевой**

Допущено Министерством образования и науки  
Российской Федерации в качестве учебника для  
студентов высших учебных заведений, обучающихся  
по экономическим направлениям и специальностям

Москва  
 **ЮРАИТ**  
2012

УДК 33  
ББК 65вб6я73  
Э40

**Рецензенты:**

*Ватник П. А.* — доктор экономических наук, профессор Санкт-Петербургского государственного инженерно-экономического университета;

*Максимова Т. Г.* — доктор экономических наук, профессор Санкт-Петербургского торгово-экономического института.

Э40 **Эконометрика** : учебник для магистров / И. И. Елисева [и др.] ; под ред. И. И. Елисейевой. — М. : Издательство Юрайт, 2012. — 453 с. — Серия : Магистр.

ISBN 978-5-9916-1930-1

Учебник охватывает все основные разделы современного курса эконометрики, отвечающего требованиям подготовки магистров по экономическим направлениям. Рассматриваются этапы возникновения и развития эконометрики, методы построения и оценки качества парной и множественной регрессий. Особое внимание уделяется мультиколлинеарности и гетероскедастичности случайных остатков, а также прогнозированию на основе модели множественной регрессии. Обсуждаются возможности построения регрессии с разнотипными переменными, разные виды регрессии с фиктивными переменными. Освещаются проблемы структурного моделирования. Подробно рассматривается эконометрика временных рядов, начиная с моделирования изолированного временного ряда, моделей по временным рядам, с лаговыми переменными, заканчивая моделями *ARMA*, *ARIMA*, *ARCH* и *GARCH*. Обсуждается проблема коинтеграции. Одна из глав посвящена анализу панельных данных, в рамках которой выделены модель с фиксированными эффектами и модель со случайными эффектами. Обсуждаются проблемы выбора модели и качества подгонки.

Соответствует Федеральному государственному образовательному стандарту высшего профессионального образования третьего поколения.

*Книга предназначена для магистрантов высших учебных заведений и факультетов экономических направлений.*

УДК 33  
ББК 65вб6я73

# Оглавление

Предисловие .....	6
<b>Глава 1. Возникновение и развитие эконометрики. Парная регрессия .....</b>	<b>9</b>
1.1. Возникновение и развитие эконометрики .....	9
1.2. Парная регрессия .....	23
1.3. Свойства остатков .....	30
<i>Контрольные вопросы и задания</i> .....	38
<b>Глава 2. Множественная регрессия .....</b>	<b>39</b>
2.1. Множественная линейная регрессия в скалярной и векторной формах .....	39
2.2. Метод наименьших квадратов и предпосылки его применения для множественной линейной регрессии .....	40
2.3. Следствия выполнения предпосылок Гаусса — Маркова .....	53
2.4. Изучение тесноты связи по множественной регрессии .....	58
2.5. Проверка значимости модели множественной регрессии и ее параметров .....	60
2.6. Множественная линейная регрессия с ограничениями на параметры .....	68
2.7. Нелинейные модели множественной регрессии .....	75
2.8. Выбор наилучшей функции регрессии .....	79
2.9. Метод максимального правдоподобия .....	90
2.10. Прогнозирование по модели множественной регрессии .....	97
2.11. Мультиколлинеарность данных .....	99
2.12. Гетероскедастичность случайных остатков .....	112
2.13. Обобщенный метод наименьших квадратов .....	123
<i>Контрольные вопросы и задания</i> .....	138

<b>Глава 3. Фиктивные переменные .....</b>	<b>141</b>
3.1. Особенности включения в модели регрессии неколичественных показателей.....	141
3.2. Спецификация моделей регрессии с фиктивными независимыми переменными.....	143
3.3. Модели регрессии с фиктивными переменными сдвига .....	144
3.4. Модели регрессии с фиктивными переменными наклона .....	149
3.5. Общий вид модели регрессии с фиктивными переменными.....	151
3.6. Исследование структурных изменений с помощью теста Чоу .....	157
<i>Контрольные вопросы и задания .....</i>	<i>160</i>
<b>Глава 4. Системы эконометрических уравнений.....</b>	<b>161</b>
4.1. Виды систем эконометрических уравнений и методы их оценивания.....	161
4.2. Системы одновременных уравнений.....	164
4.3. Уравнения, кажущиеся несвязанными ....	173
<i>Контрольные вопросы и задания .....</i>	<i>175</i>
<b>Глава 5. Моделирование изолированного динамического ряда.....</b>	<b>176</b>
5.1. Компоненты динамического ряда.....	176
5.2. Автокорреляция уровней динамического ряда и характеристика его структуры .....	184
5.3. Модели тенденции развития .....	191
5.4. Моделирование периодических колебаний .....	227
<i>Контрольные вопросы и задания .....</i>	<i>249</i>
<b>Глава 6. Модели регрессии по временным рядам .....</b>	<b>251</b>
6.1. Специфика изучения взаимосвязей по рядам динамики.....	251
6.2. Учет тенденции при построении модели регрессии.....	253
6.3. Обобщенный метод наименьших квадратов при построении модели регрессии по временным рядам.....	266

6.4.	Учет сезонности при построении модели регрессии.....	277
	<i>Контрольные вопросы и задания</i> .....	282
<b>Глава 7.</b>	<b>Модели с лаговыми переменными.....</b>	<b>283</b>
7.1.	Общая характеристика.....	283
7.2.	Модели с распределенными лагами.....	286
7.3.	Модели авторегрессии.....	302
7.4.	Авторегрессионные процессы и их моделирование (общая характеристика) .....	309
	<i>Контрольные вопросы и задания</i> .....	316
<b>Глава 8.</b>	<b>Модели <i>ARMA</i>, <i>ARIMA</i>, <i>ARCH</i>, <i>GARCH</i>.....</b>	<b>317</b>
8.1.	Стационарный ряд.....	317
8.2.	Базовые модели временных рядов.....	320
8.3.	Теорема декомпозиции Вольда .....	325
8.4.	Частная автокорреляционная функция....	327
8.5.	Модель <i>ARMA</i> .....	329
8.6.	Модель <i>ARIMA</i> .....	342
8.7.	Коинтеграция.....	350
8.8.	Модели <i>ARCH</i> и <i>GARCH</i> .....	355
	<i>Контрольные вопросы и задания</i> .....	368
<b>Глава 9.</b>	<b>Анализ панельных данных .....</b>	<b>370</b>
9.1.	Панельные данные и их преимущества..	370
9.2.	Однонаправленные модели панельных данных.....	376
9.3.	Качество подгонки.....	402
9.4.	Выбор модели.....	405
9.5.	Двухнаправленная модель панельных данных с фиксированными эффектами..	414
	<i>Контрольные вопросы и задания</i> .....	418
	<b>Выдающиеся ученые.....</b>	<b>420</b>
	<b>Литература.....</b>	<b>430</b>
	<b>Предметный указатель .....</b>	<b>433</b>
	<b>Приложения.....</b>	<b>439</b>

---

# Предисловие

Преподавание эконометрики вошло в стандарты третьего поколения для экономических специальностей в качестве федеральной компоненты, т.е. дисциплины, обязательной для изучения обучающимися по экономическим специальностям как на уровне бакалавриата, так и на уровне слушателя магистратуры.

Подготовленный учебник нацелен на обучение в магистратуре. Авторы исходили из того, что современный состав магистрантов довольно неоднороден и включает как продвинутых слушателей, так и тех, кто либо закончил бакалавриат по другому профилю, либо имел трудности в освоении этой дисциплины. Мы стремились излагать как можно проще самые сложные вопросы, причем делать это, привлекая примеры, основанные на российских данных. Почти везде это удалось сделать.

Учебник начинается с краткого исторического обзора возникновения и развития эконометрики и обсуждения парной регрессии с позиций эконометрического анализа, т.е. проблемы спецификации уравнения регрессии и его соответствия исходным данным. Также обсуждаются свойства остатков и показываются возможности применения простейших эконометрических тестов. В последующих главах начинается изложение эконометрического моделирования, более адекватного природе экономических процессов и явлений. Прежде всего, рассматриваются случаи использования пространственных данных и все возникающие при этом сложности: нелинейность

эффектов, мультиколлинеарность переменных, гетероскедастичность случайных остатков. Одна из глав посвящена использованию фиктивных переменных сдвига, наклона, исследованию структурных изменений. Выявленные в ней проблемы получили развитие в главе «Системы эконометрических уравнений», в которой основное внимание уделяется идентификации уравнений, тестированию на экзогенность, а также оценке связности уравнений. При переходе к использованию временных рядов в эконометрическом анализе излагаются, прежде всего, методы исследования изолированного ряда, выделяются проблемы моделирования тренда и периодических колебаний, а также автокорреляции.

От рассмотрения изолированного ряда авторы переходят к построению моделей на основе системы временных рядов. Потребность в такого рода моделях совершенно очевидна, хотя их популярность в последние годы снизилась. Обсуждаются все основные вопросы построения моделей, включая отражение фактора сезонности, а также случай применения обобщенного метода наименьших квадратов для оценки параметров модели.

Более популярны в настоящее время модели с лаговыми переменными, которые подробно описаны в учебнике с учетом метода инструментальных переменных и моделирования авторегрессионных процессов. Здесь же рассмотрены столь популярные модели скользящей средней *ARMA* и *ARIMA*. Дополненные моделями *ARCH* и *GARCH*, они раскрывают особенности проверки ряда на стационарность и его обработки, которая особенно актуальна для финансовых временных рядов с их суперволатильностью. Завершает учебник глава «Анализ панельных данных», охватывающая построение моделей с фиксированными и случайными эффектами, проблемы качества подгонки и выбора модели.

Учебник подготовлен представителями двух эконометрических школ в Санкт-Петербурге: одна — это школа Санкт-Петербургского государственного университета экономики и финансов, другая — это школа Европейского университета. Представители обеих школ смогли реализовать свои знания и методологические новации в процессе разноуровневого преподавания дисциплины «Эконометрика» в СПбГУЭФ. Накопленный опыт совместной работы составил основу подготовки этого

учебника. В учебник включены персональные сведения о тех ученых, чей вклад в математическую статистику и эконометрику оказался решающим.

Распределение авторства по главам учебника: предисловие и гл. 1 — чл.-корр. РАН, д-р экон. наук, зав. кафедрой статистики и эконометрики, проф. И. И. Елисеева, гл. 5—7 — д-р экон. наук, проф. С. В. Курьшева, гл. 2, 3 — канд. экон. наук, доц. Ю. В. Нерадовская, гл. 4 — Л. М. Галиуллина, гл. 8 — Д. В. Беляков, гл. 9 — А. В. Кабачек.

Весь авторский коллектив благодарит Ю. В. Нерадовскую за неоценимую помощь по организации столь сложной коллективной работы. Без ее усилий работа не могла бы завершиться в требуемые сроки.

Руководитель коллектива авторов И. И. Елисеева.

После изучения материалов данного учебника обучающийся должен:

**знать:**

- современные эконометрические методы;
- системы эконометрических уравнений;
- особенности моделей, позволяющих при наличии различной информации решать разные эконометрические задачи;

**уметь:**

- применять методы идентификации и оценивания систем эконометрических уравнений;
- осуществлять эконометрическое прогнозирование на основе различных эконометрических моделей;

**владеть:**

- спецификой эконометрических измерений;
- навыками структурного моделирования для анализа ситуаций.



# **Возникновение и развитие эконометрики. Парная регрессия**

## **1.1. Возникновение и развитие эконометрики**

*Эконометрика* — наука, изучающая количественные и качественные экономические взаимосвязи с помощью математических и статистических методов и моделей. Определение предмета эконометрики было выработано в уставе Эконометрического общества (основано в 1930 г.), который в качестве главных целей выдвигает использование статистики и математики для развития экономической теории. Эконометрика может рассматриваться как наука, состоящая из двух разделов: теоретическая эконометрика и прикладная эконометрика. Теоретическая эконометрика изучает статистические свойства оценок и проверки гипотез, в то время как прикладная эконометрика занимается применением эконометрических методов для оценки тех или иных положений экономической теории. Эконометрика дает инструментарий для экономических измерений, а также методологию оценки параметров моделей микро- и макроэкономики. Кроме того, эконометрика активно используется для прогнозирования экономических процессов как в масштабах экономики в целом или отдельных ее отраслей, так и на уровне предприятий. При этом эконометрика является частью экономической теории, наряду с макро- и микроэкономикой.

Термин «эконометрика» состоит из двух частей: «эконо» — от слова «экономика» и «метрика» — от слова «измерение». Эконометрика входит в обширное семейство дисциплин, посвященных измерениям и применению статистических методов в различных областях науки и практики. К этому семейству относятся, в частности, биометрия, технометрика, наукометрия, психометрия, хемометрия, квалиметрия, клиометрия. Особняком стоит социометрия — этот термин закрепился за статистическими методами анализа взаимоотношений в малых группах (входит в статистический анализ в социологии).

### **Предпосылки возникновения эконометрики**

Первые попытки количественных исследований в экономике относятся ко второй половине XVII в. Они были связаны с одним из новых направлений в экономической теории — политической арифметикой. У. Петти (1623—1687), Ч. Давенант (1656—1714), Г. Кинг (1648—1712) использовали конкретные экономические данные в своих исследованиях и в первую очередь при расчете национального дохода. Это направление стимулировало поиск экономических законов по аналогии с физическими, астрономическими и другими естественнонаучными законами. При этом существование неопределенности в экономике еще не осознавалось.

Важным этапом возникновения эконометрики стало развитие статистической теории в трудах Ф. Гальтона (1822—1911), К. Пирсона (1857—1936) и Ф. Эджворта (1845—1926). Эти ученые предопределили первые применения парной корреляции. Так, ученик К. Пирсона Дж. Э. Юл (1871—1951) определял связь между уровнем бедности и формами помощи бедным. Р. Г. Хукер (1867—1944) измерял связь между уровнем брачности и благосостоянием, в котором использовалось несколько индикаторов благосостояния. Ему также принадлежат первые исследования временных рядов экономических переменных.

С 1830-х гг. наиболее развитые страны стали испытывать необъяснимые с точки зрения экономической науки того времени потрясения — упадок деловой активности, возникновение массовой безработицы. Быстрое промышленное развитие и урбанизация выявили огромный пласт нерешенных социальных

проблем. Уже в конце XIX в. неоклассическая теория стала восприниматься как слишком удаленная от действительности. Теория могла стать убедительной в том случае, если бы она смогла объяснить изменения, происходящие в экономике. Для ее практического применения требовались количественные выражения базовых экономических терминов.

В 1911 г. вышла книга американского экономиста Г. Мура (1869—1958) «Законы заработной платы: эссе по статистической экономике». В своем исследовании Г. Мур провел анализ рынка труда, статистически проверил теорию производительности Дж. Кларка (1847—1938) и изложил основы стратегии объединения пролетариата. Г. Мур показал, что с помощью статистико-математических построений, основанных на фактических данных, можно разработать основу для социальной политики. В это же время итальянский экономист Р. Бенини (1862—1956) впервые использовал множественную регрессию при оценке функции спроса.

Значительный вклад в становление эконометрики внесли исследования цикличности экономики. Первым цикличность экономики обнаружил К. Жюгляр (1819—1905). Он выявил 7—11-летние циклы инвестиций. Сразу после него Дж. Китчин (1861—1932) выявил 3—5-летнюю периодичность обновления оборотных средств, С. Кузнец (1901—1985), лауреат Нобелевской премии по экономике за 1971 г., обнаружил 15—20-летние циклы в строительстве, а Н. Кондратьев (1892—1938) выявил «длинные волны» в экономике продолжительностью 45—60 лет.

Важным этапом формирования эконометрики явилось построение экономических барометров. Оно основано на идее, согласно которой существуют показатели, которые изменяются раньше других и поэтому могут служить сигналами изменений последних. Первым и самым известным стал Гарвардский барометр, который был создан в 1903 г. под руководством У. Персона (1878—1937) и У. Митчелла (1874—1948). Он состоял из кривых, характеризующих фондовый, товарный и денежный рынки. Каждая из кривых представляла собой среднюю арифметическую из входящих в нее нескольких показателей. Эти ряды предварительно обрабатывались путем исключения тенденции сезонности и приведения колебаний отдельных кривых к сравнимому масштабу колеблемости. Успех использования

Гарвардского барометра вызвал появление многих аналогичных барометров в других странах. Однако приблизительно с 1925 г. он потерял свою чувствительность. Его крах объясняется появлением мощного регулирующего фактора в экономике США. В этих условиях основным методом макроэкономического анализа становится метод построения межотраслевого баланса В. В. Леонтьева (1906—1999). В это же время начали строиться экономические модели, использующие методы гармонического анализа. Эти методы были перенесены в экономику из астрономии, метеорологии и физики.

### **Развитие эконометрики**

К 1930 г. сложились все предпосылки для выделения эконометрики в отдельную науку. Стало ясно, что для более глубокого понимания экономических процессов стоит использовать в той или иной степени статистику и математику. Возникла необходимость появления новой науки, объединяющей исследования, проводимые в этом направлении. 29 декабря 1930 г. по инициативе И. Фишера (1867—1947), Р. Фриша (1895—1973), Я. Тинбергена (1903—1994), Й. Шумпетера (1883—1950), О. Андерсона (1887—1960) и других ученых было создано Эконометрическое общество. В 1933 г. Р. Фриш основал журнал «Эконометрика», который и сейчас имеет большое значение для развития эконометрики. А уже в 1941 г. появился первый учебник по новой научной дисциплине, написанный Я. Тинбергенем. В 1969 г. Р. Фриш и Я. Тинберген стали первыми исследователями, получившими Нобелевскую премию по экономике за создание и применение динамических моделей для анализа экономических процессов.

До 1970-х гг. эконометрика понималась как эмпирическая оценка моделей, созданных в рамках экономической теории. По мнению эконометристов того времени, статистические данные должны были защитить теорию от догматизма. При этом подавляющее большинство экономических моделей, построенных в тот период, были кейнсианскими. Но начиная с 1970-х гг. формальные методы стали использоваться при выборе теоретических концепций. При этом эконометрикой стали активно пользоваться и монетаристы.

В 1980 г. вторую Нобелевскую премию по экономике получил американский экономист и эконометрист Л. Клейн (р. 1920)

за создание экономических моделей и их применение для анализа колебаний экономики и экономической политики. Совместно с А. Голдбергом (1930—2009) он создал одну из самых известных моделей американской экономики, известной как «модель Клейна — Голдберга». Модель состояла из взаимосвязанных одновременных и направленных рядов уравнений, решение которых давало картину производства в стране. Говоря об этой модели, Р. Дж. Болл отмечал: «Как эмпирическое представление об основах кейнсианской системы эта модель стала, возможно, самой знаменитой среди моделей крупных национальных хозяйств до появления других моделей в 1960-е гг.». Л. Клейн также организовал широко известный проект «Link» для интеграции статистических моделей разных стран в единую общую систему с целью улучшения понимания международных экономических связей и прогнозирования в области мировой торговли. В это время активно развивались не только макро-, но микроэконометрика. Пионерами в этом направлении выступили Дж. Хекман (р. 1944) и Д. Макфадден (р. 1937). Они разработали теорию и методы, которые широко используются в статистическом анализе поведения индивидуумов и домохозяйств как в экономике, так и в других общественных науках. Так, Дж. Хекман решил проблему смещения выборки из-за селективности данных и самоотбора. Для этого он предложил использовать метод коррекции Хекмана, который благодаря своей эффективности и простоте в использовании стал широко использоваться в эмпирических исследованиях. Основной вклад Д. Макфаддена в науку заключается в развитии методов для анализа дискретного выбора. В 1974 г. он разработал условный логит-анализ, который сразу был признан фундаментальным достижением экономической науки. Также он создал эконометрические методы для оценки производственных технологий и исследования факторов, лежащих в основе спроса фирм на капитал и рабочую силу. Выдающиеся достижения этих ученых были отмечены Нобелевской премией по экономике в 1990 г.

Важным событием для развития эконометрики стало появление компьютеров. Благодаря им мощное развитие получил статистический анализ временных рядов. Дж. Бокс (р. 1919) и Г. Дженкинс (1933—1982) создали модель *ARIMA* в 1970 г., а К. Симс (р. 1942) и другие ученые — модели *VAR* в начале 1980-х гг. Расширение эконометрических исследований

стимулировало и бурное развитие финансовых рынков и производных инструментов. Это привело Дж. Тобина (1918—2002) к разработке моделей с использованием цензурированных данных (Нобелевская премия по экономике 1981 г.).

Большое влияние на современную эконометрику оказал и Т. Хаавельмо (1911—1999). Он показал, как можно использовать методы математической статистики для того, чтобы получать обоснованные заключения о сложных экономических взаимосвязях исходя из случайной выборки эмпирических наблюдений. Разработанные им методы можно использовать и для оценивания соотношений, полученных на основе экономических теорий, и для проверки этих теорий. В 1989 г. ему присудили Нобелевскую премию по экономике за прояснение вероятностных основ эконометрики и анализ одномерных экономических структур.

Т. Хаавельмо рассматривал экономические ряды как реализацию случайных процессов. Главные проблемы, возникающие при работе с такими данными, — это нестационарность и сильная волатильность. Если переменные нестационарны, то возникает риск установить связь там, где ее нет. Вариантом решения данной проблемы является переход от уровней ряда к их разностям. Недостаток данного метода — сложность экономической интерпретации полученных результатов. Для решения этой проблемы К. Грэнджер (1934—2009) предложил концепцию коинтеграции как стационарной комбинации между нестационарными переменными. Им была создана модель корректировки отклонений (*ECM*), для которой он разработал методы оценивания ее параметров, обобщения и тестирования. Коинтеграция применяется в случае, если краткосрочная динамика отражает значительные дестабилизирующие факторы, а долгосрочная динамика стремится к экономическому равновесию. Модели, созданные К. Грэнджером, были обобщены С. Йохансеном (р. 1939) в 1990 г. для многомерного случая. В 2003 г. К. Грэнджер совместно с Р. Энглom (р. 1942) удостоились Нобелевской премии. Р. Энгл известен как создатель моделей с меняющейся во времени волатильностью (модели *ARCH*). Эти модели получили широкое распространение на финансовых рынках.

Одним из основных бурно развивающихся направлений эконометрики является непараметрическая эконометрика. **Непараметрическая эконометрика** — раздел эконометрики,

который не требует спецификации функциональных форм оцениваемых объектов. Вместо этого данные сами формируют модель. Непараметрические методы становятся все более популярными в прикладных исследованиях. Они наилучшим образом подходят для анализа большого объема данных при малом количестве переменных. Также эти методы применяются в тех случаях, когда обычные параметрические спецификации не дают возможности решения поставленной задачи. Непараметрические методы не включают гипотезы о распределении, что иногда является полезным в прикладном исследовании. Основные методы построения гибких моделей — это ядерные методы, сглаживание сплайнами, методы ближайших соседей, нейронные сети, фрактальный анализ и гибкие методы сглаживания с помощью рядов данных.

Иногда к непараметрической эконометрике относят эконометрический анализ нечисловых математических понятий, принадлежащих к тем или иным классам объектов нечисловой природы, таким, как нечеткие множества, интервалы, распределения вероятностей и т.д. Так, в статистике интервальных данных, где элементами выборки являются не числа, а интервалы, изучены практически все задачи классической прикладной математической статистики, в частности задачи регрессионного анализа, планирования эксперимента, сравнения альтернатив и принятия решений в условиях интервальной неопределенности и т.д. Для данной отрасли науки разработана общая схема исследования, включающая расчет двух основных характеристик — максимально возможного отклонения статистики, вызванного интервальностью исходных данных, и рационального объема выборки (превышение которого не дает существенного повышения точности оценивания и статистических выводов, связанных с проверкой гипотез). Также разработаны подходы к учету интервальной неопределенности в основных постановках регрессионного, дискриминантного и кластерного анализа.

### Специфика экономических измерений<sup>1</sup>

Специфические особенности экономических данных можно свести к следующим пяти группам:

<sup>1</sup> Подробнее см.: Эконометрика : учебник / И. И. Елисеева, С. В. Курышева, Н. М. Гордеенко [и др.] ; под ред. И. И. Елисеевой. 2-е изд., перераб. и доп. М. : Финансы и статистика, 2007.

- измеряться могут только операционально определенные данные. При этом экономические измерения подвержены сильному влиянию теоретических представлений о данных величинах;
- неэкспериментальный характер данных и короткие ряды наблюдений, которые ставят под сомнение адекватность полученных результатов;
- экономические данные, как правило, являются косвенными;
- изменчивость единиц измерения;
- влияние инструмента измерения на изучаемый объект.

## Эконометрические методы

### 1. Регрессионный анализ

**Регрессионный анализ** — статистический метод исследования связи между зависимой переменной  $y$  и одной или несколькими независимыми переменными  $x_1, x_2, \dots, x_p$ . При этом терминология зависимых и независимых переменных отражает лишь математическую зависимость переменных, выражающуюся в сопряженности изменений значений переменных, а не причинно-следственные отношения. Для адекватного описания сложных внутренне неоднородных экономических процессов, как правило, применяются системы эконометрических уравнений. В более простых случаях можно использовать изолированные уравнения (уравнения регрессии).

### 2. Анализ временных рядов

**Анализ временных рядов** — совокупность математико-статистических методов, предназначенных для выявления структуры временных рядов и прогноза. Определение структуры временного ряда необходимо для того, чтобы построить математическую модель такого явления, которое служит источником анализируемого временного ряда. Прогноз будущих значений временного ряда используется при принятии решений. Прогнозирование также интересно тем, что оно рационализирует анализ временных рядов отдельно от экономической теории.



Как правило, при прогнозировании исходят из некоторой заданной параметрической модели. При этом используются стандартные методы параметрического оценивания (метод наименьших квадратов (МНК), метод максимального правдоподобия (ММП), метод моментов). Также достаточно разработаны методы непараметрического оценивания для нечетко заданных моделей.

### 3. Панельный анализ

Панельные данные представляют собой прослеженные во времени пространственные микроэкономические выборки, т.е. они состоят из наблюдений одних и тех же экономических единиц в последовательные периоды времени. Панельные данные состоят из тех измерений: признаки — объекты — время. Их использование дает ряд существенных преимуществ при оценке параметров регрессионных зависимостей, так как они позволяют проводить и анализ временных рядов, и анализ пространственных выборок. С помощью подобных данных изучают бедность, безработицу, преступность, а также оценивают результативность государственных программ в области социальной политики.

#### **Критика и апологетика эконометрики**

Во многом определяющим для развития эконометрики стал спор Я. Тинбергена и Дж. М. Кейнса (1883—1946) об эконометрическом методе исследования. В статье «Professor Tinbergen's Method» М. Кейнс написал, что Я. Тинберген «предпочитает лабиринты арифметики лабиринтам логики». Он сказал, что эконометрический анализ становится похож на «детские головоломки, в которых вам нужно написать ваш возраст, умножить на что-то, прибавить еще что-то, вычесть и в конце концов получить число зверя из Откровения св. Иоанна Богослова».

М. Кейнс утверждал, что исследовательский потенциал анализа множественной корреляции во многом зависит от экономиста. По его мнению, данный метод применим, только когда экономист в состоянии заранее представить правильный и безукоризненно полный анализ значимых факторов. При этом возникает проблема использования неполного набора объясняющих переменных (смещенная оценка, вызванная пропуском

переменных); построение моделей, содержащих ненаблюдаемые переменные (такие, как рациональные ожидания), полученные при помощи плохо измеренных данных, основанных на индексах; получение ложной корреляции в результате использования замещающих переменных и одновременности.

На эту критику Я. Тинберген ответил, что «нерелевантные объясняющие переменные можно трактовать как случайные остатки, не коррелирующие систематически с другими объясняющими переменными. Если математическая форма соотношения задана, то можно представить определенные данные о вероятностных распределениях остатков». При этом объясняющие факторы можно измерить, а независимость остатков можно проверить впоследствии, изучая их автокорреляцию. При этом экономист не должен забывать об ограниченности метода и проверке достоверности данных.

М. Кейнс также попытался предъявить методу множественной регрессии, являющемуся прикладным, те требования, которым должен отвечать фундаментальный метод. Он настаивал на необходимости истинности предпосылок, соизмеримости условий, независимости рассматриваемых факторов, характере функций и т.д., при этом он не ответил на вопрос о том, как проверить истинность, что взять в качестве критериев истинности, соизмеримости и независимости. Современная же научная методология отказалась от принципа верификации предпосылок и перешла к верификации выводов или оценке точности прогноза.

Критиковал М. Кейнс и введение фактора времени в уравнение регрессии. Очевидно, что использование линейного тренда означает, что между первым и последним годами временного ряда проводится прямая линия. В результате очень многое зависит от того, какие годы выбраны для исследования. Разбирая пример временного ряда, взятого с 1919 по 1933 г. из книги Я. Тинбергена, он отмечал, что «возникает парадокс, состоящий в том, что экономика США характеризовалась серьезным понижательным трендом за весь период, в том числе и за период, закончившийся в 1929 г.». Суммарно изменения достигают 20%, при этом если бы Я. Тинберген исследовал временной ряд, заканчивающийся на 1929 г., то он использовал бы растущий тренд вместо понижательного для анализа тех же самых лет. Трендовая компонента, по мнению М. Кейнса, очень похожа на метод корректировки неудачных

результатов и затемняет тот факт, что «данное объяснение на самом деле ошибочно».

При этом, по его мнению, непонятно, «в какой степени кривые и уравнения считаются не более чем частью описания и исторического анализа с целью подбора кривых и в какой степени с их помощью делаются индуктивные выводы относительно будущего или прошлого». М. Кейнс выразил сомнения в ценности такого подхода. По его словам, очевидно, что данный метод «представляет собой не самый ясный способ описания прошлого». Самое важное условие при таком анализе состоит в том, что «экономическая среда на протяжении некоторого периода времени должна оставаться неизменной и однородной во всех значимых отношениях, за исключением колебаний тех факторов, которые рассматриваются отдельно. Но быть уверенными, что такие условия сохранятся в будущем, даже если они обнаруживаются в прошлом, нельзя».

На это Я. Тинберген ответил утверждением, что «зачастую сам вид кривых подсказывает, что некоторый фактор, не упомянутый в большинстве учебников по экономике, имеет огромную важность. Представив численное значение одного или нескольких коэффициентов регрессии, можно критиковать одну или несколько использовавшихся ранее теорий». Я. Тинберген привел пример такой ситуации, когда «множество теоретиков соглашались с тем, что ставка процента является существенным фактором спроса на деньги или инвестиционной активности, а полученные результаты после анализа указывают на то, что такое влияние незначительно или, по меньшей мере, было таковым в США в течение данного периода времени».

М. Кейнс считал очень важным вопрос о предполагаемой линейности соотношений. Он утверждал, что не обнаружил какого-либо примера нелинейной корреляции. Он говорил о том, что не понимает, анализ каких эмпирических данных заставляет использовать нелинейную корреляцию. Однако, по словам Я. Тинбергена, «диаграммы рассеяния позволяют понять, является ли некоторая корреляция линейной или нет. Нелинейность ни в коем случае не является произвольной манипуляцией с коэффициентами». Строго говоря, для каждого значения объясняющей переменной возможен только один коэффициент, и, с учетом непрерывности, требуется, чтобы эти коэффициенты не колебались слишком сильно. М. Кейнс

очень плохо относился к линейным соотношениям, он называл их «смехотворными». Однако есть причины, в силу которых степень их смехотворности снижается:

- на малых интервалах неразрывную функцию можно аппроксимировать линейными функциями;
- наблюдение за экономическими данными показывает, что линейные соотношения часто встречаются на практике. При этом логично начинать анализ, опираясь на самую простую предпосылку, которая коррелирована с общей теорией. По словам Я. Тинбергена, «такой подход очень часто встречается в индуктивной части любой исследовательской работы. Также существует теоретическое обоснование линейности, согласно которому для больших масс индивидов совместная реакция будет носить значительно более линейный характер, чем какая-либо индивидуальная реакция».

Критика эконометрики М. Кейнсом главным образом обусловлена отличием его подхода к экономической науке от подхода экономического мейнстрима. Основным пунктом этого расхождения является вопрос: «Следует ли трактовать экономику как точную науку?» Сам М. Кейнс давал отрицательный ответ на этот вопрос. В рамках его традиции экономическая среда изменчива и непредсказуема, а большинство экономических переменных связано между собой множеством сложных нелинейных зависимостей. Из этого следуют нестабильность коэффициентов корреляции и невозможность решения предсказательных задач. Поэтому экономическая наука не может претендовать на точные количественные измерения. Она должна быть основана на реалистичных предпосылках и содержать инструменты, помогающие понять и объяснить экономическую среду. Подход же Я. Тинбергена вполне согласуется с современным мейнстримом: экономический анализ должен быть как можно более формализованным и нацеленным на решение конкретных количественных задач. В рамках данного подхода экономическая наука должна быть точной, а объект ее изучения аналогичен объектам технических и естественнонаучных дисциплин.

### Последующая критика

Несмотря на потенциальные возможности, эконометрика не получила поддержки у многих крупных экономистов. В начале

1970-х гг. Д. Уорсвик (1916—2001) резко критиковал экономистов-математиков за «отсутствие связи с конкретными фактами». Он утверждал, что эконометристы «занимаются не столько изобретением средств систематизации и измерения имеющихся фактов, сколько созданием неисчислимого множества претендующих на это способов». В это же время Ф. Браун (1906—1994) утверждал, что «построение регрессий временных рядов годится только для обмана». В. Леонтьев охарактеризовал эконометрику как «попытку компенсировать бросающийся в глаза недостаток имеющихся данных путем широко использования все более и более изощренных статистических приемов». В подобном же духе высказывался и Дж. Хикс (1904—1989), он говорил о том, что «не следует преувеличивать значение эконометрических методов в экономической теории». А Э. Лимер (р. 1944) писал, что «существует две вещи, процесс изготовления которых лучше не видеть: сосиски и эконометрические оценки».

Резко отрицательно к эконометрике относились и представители австрийской школы экономики. Так, Л. фон Мизес (1881—1973) писал: «Введенные в заблуждение идеей, что науки о человеческой деятельности должны подражать методу естественных наук, великое множество авторов поглощены квантификацией экономики. Они думают, что экономика должна подражать химии, которая развилась от качественного к количественному состоянию. Их девиз — позитивистский принцип: наука — это измерение. Но они не в состоянии понять, что в области человеческой деятельности статистика — это всегда история, и что гипотетические корреляции и функции не описывают ничего, кроме того, что случилось в определенный момент времени в определенной географической области как результат деятельности определенного числа людей. Как метод экономического анализа эконометрика — ребяческая игра с числами, которая не добавляет чего-либо в разъяснение проблем экономической действительности».

К более детализированной критике множественной регрессии со времен М. Кейнса также добавилась невозможность отделения мультиколлинеарности от попытки изолированной оценки объясняющей переменной, неправильная спецификация динамических реакций и длинных лагов, предположение о линейности без точного знания соответствующих значений

регрессии, некорректная предварительная фильтрация данных, необоснованные выводы из корреляции, непостоянство параметров уравнения регрессии, отождествление экономической и статистической значимости и невозможность соотнесения экономической теории с эконометрикой, а также неадекватный объем выборки.

Благодаря этой и некоторой другой критике методология прикладных исследований была пересмотрена. Согласно классической эконометрической методологии полученные результаты считаются более адекватными, если изучаемые переменные более сильно коррелированы, предсказания точнее соответствуют данным и чем более значимыми являются полученные оценки с точки зрения  $t$ - или  $F$ -статистик. Отводится значительное место тому, как наиболее эффективным образом организовать перебор потенциальных объясняющих переменных, чтобы наилучшим образом предсказать объясняемую переменную, при этом, чтобы коэффициент детерминации был как можно большим, а  $F$ -статистика как можно более значимой. Если получены неудовлетворительные результаты в критериях спецификации, то исследователь, пользующийся традиционной методологией вместо того, чтобы пересмотреть модель, начинает применять более продвинутые методы оценивания. В рамках этого подхода характерно стремление получить «наилучший» результат вместо стремления получить результат осмысленный и надежный. Однако на современном этапе развития эконометрики предпочтение отдается тем моделям, которые проходят диагностические критерии, даже если они имеют более низкий коэффициент детерминации.

Сегодня эконометрика занимает достойное место в ряду экономических наук. В мире выпускается ряд научных журналов, полностью посвященных эконометрике, в том числе: *Journal of Econometrics* (Швеция), *Econometric Reviews* (США), *Econometrica* (США), *Sankhya*, *Indian Journal of Statistics. Ser. D. Quantitative Economics* (Индия), *Publications Econometriques* (Франция). Эконометрику изучают в ведущих мировых университетах, пришло понимание, что без эконометрических методов невозможно проводить современный макро- и микроэкономический анализ.

Признание эконометрики в качестве важного исследовательского инструмента подтверждается тем, что подавляющее большинство

лауреатов премии имени Нобеля Банка Швеции являются эконометристами (начиная с Р. Фриша и Я. Тинбергена, премия 1969 г., и заканчивая Д. Мортенсенем (р. 1939), К. Писсаридесом (р. 1948), П. Даймондом (р. 1940), премия 2010 г.).

В России также издаются специализированные журналы. К ним относятся «Прикладная эконометрика» и «Квантиль» (электронный журнал). Отдельные публикации по эконометрике можно найти в журналах «Экономика и математические методы», «Вопросы статистики», «Журнал Новой экономической ассоциации», «Финансы и бизнес» и некоторых других.

## 1.2. Парная регрессия

Парной регрессией называется условное математическое ожидание переменной  $y$  как функции от переменной  $x$ . Модель парной линейной регрессии имеет вид

$$y_i = \alpha_0 + \alpha_1 x_i + \varepsilon_i, \quad i = 1, \dots, n,$$

где  $y$  — зависимая переменная;  $\alpha_0, \alpha_1$  — неизвестные параметры регрессии;  $x$  — независимая переменная;  $\varepsilon$  — случайная составляющая;  $n$  — число наблюдений.

Для нахождения оценок параметров  $\alpha_0$  и  $\alpha_1$  часто используется метод наименьших квадратов (МНК), сводящийся к минимизации по  $\alpha_0$  и  $\alpha_1$  суммы квадратов отклонений:

$$SS_e = \sum_{i=1}^n (y_i - \alpha_0 - \alpha_1 x_i)^2.$$

Оценки МНК параметров имеют вид

$$a_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}; \quad a_0 = \bar{y} - a_1 \bar{x}.$$

Функция выборочной линейной регрессии будет иметь вид

$$\hat{y}_i = a_0 + a_1 x_i, \quad i = 1, \dots, n.$$

**Регрессионными остатками называются разности фактических и оцененных значений зависимой переменной:**

$$e_i = y_i - \hat{y}_i, \quad i = 1, \dots, n. \quad (1.1)$$

Если случайные величины  $\varepsilon_i, \varepsilon_j$  при всех  $i \neq j, i = 1, \dots, n, j = 1, \dots, n$ , являются некоррелированными и  $\varepsilon_i \sim N(0, \sigma_\varepsilon^2)$ , а  $X$  — детерминированный вектор, то случайная величина

$$\frac{\alpha_1 - \alpha_1}{s_{\alpha_1}}$$

имеет распределение Стьюдента с  $(n - 2)$  степенями свободы.

Гипотеза  $H_0: \alpha_1 = c$  о конкретном значении коэффициента  $\alpha_1$  при двусторонней альтернативной гипотезе  $H_1: \alpha_1 \neq c$  проверяется с помощью тестовой статистики:

$$t = \frac{\alpha_1 - c}{s_{\alpha_1}}.$$

Если  $|t| > t_{T(\alpha/2, (n-2))}$ , где  $\alpha$  — выбранный уровень значимости, то основная гипотеза отвергается.

Если же альтернативная гипотеза односторонняя

$$H_1: \alpha_1 < c \quad (H_1: \alpha_1 > c),$$

то основная гипотеза отвергается при  $t < t_{T[\alpha, (n-2)]}$  ( $t > t_{T[\alpha, (n-2)]}$ ).

Аналогично формулируется и проверяется гипотеза о конкретном значении свободного члена  $\alpha_0$ .

Проверка гипотезы о равенстве параметра нулю (обычно при двусторонней альтернативной гипотезе) называется проверкой гипотезы о значимости параметра. Если гипотеза о равенстве параметра нулю не отвергается, то этот параметр называется незначимым.

При оценке параметров регрессии статистическими пакетами *Excel*, *Eviews*, *STATA*, *SPSS* и др. автоматически проводится проверка гипотез о значимости коэффициентов. Для тестовой  $t$ -статистики вычисляется  $p$ -value ( $p$ -значение) — минимальный уровень значимости, при котором основная гипотеза отвергается. Если  $p$ -value превышает выбранный уровень



значимости, то основная гипотеза (о равенстве коэффициента нулю) не отвергается.

Если коэффициент  $\alpha_1$  является незначимым, то между переменными  $x$  и  $y$  не существует статистически значимой линейной связи. Если коэффициент  $\alpha_1$  является значимым, то его оценка интерпретируется следующим образом: при увеличении  $x$  на одну единицу  $y$  изменяется на  $\alpha_1$  единиц (в сторону увеличения при положительном  $\alpha_1$  и в сторону уменьшения при отрицательном  $\alpha_1$ ).

### Пример

Имеются статистические данные о значениях двух показателей в разрезе 48 субъектов РФ, которые приведены в табл. 1.1.

1. Поступление налогов, сборов и иных обязательных платежей в консолидированный бюджет РФ (без поступлений ЕСН) в 2009 г., млн руб.

2. Количество занятых в Российской Федерации в 2009 г., тыс. человек.

Зависимая переменная ( $y$ ) — поступление налогов, сборов и иных обязательных платежей в консолидированный бюджет РФ («поступление налогов»). Независимая переменная — «количество занятых» ( $x_1$ ).

**Таблица 1.1.** Некоторые экономические показатели деятельности субъектов РФ в 2009 г.

Субъект РФ	Поступление налогов, млн руб.	Количество занятых, тыс. человек
	$y$	$x_1$
Республика Ингушетия	1422,20	107,20
Еврейская автономная область	2529,70	82,30
Республика Тыва	2629,10	101,60
Республика Алтай	2764,30	87,60
Карачаево-Черкесская Республика	3347,50	188,30
Республика Калмыкия	3914,20	121,90
Республика Адыгея	4400,80	187,10
Республика Северная Осетия-Алания	5904,00	326,50
Магаданская область	6956,70	97,10
Кабардино-Балкарская Республика	7595,10	352,50
Республика Хакасия	9257,80	254,70
Чукотский автономный округ	9317,10	30,50

Окончание табл. 1.1

Субъект РФ	Поступление налогов, млн руб.	Количество занятых, тыс. человек
	$y$	$x_1$
Республика Марий Эл	9978,80	323,70
Псковская область	10 144,80	323,30
Чеченская Республика	10 215,40	357,00
Республика Карелия	11 349,50	337,80
Курганская область	12 046,90	393,00
Республика Мордовия	12 061,40	439,00
Костромская область	12 104,20	340,90
Камчатский край	13 042,40	190,00
Орловская область	13 104,30	375,00
Ивановская область	13 396,40	491,20
Республика Дагестан	14 170,30	1104,10
Тамбовская область	14 227,00	499,50
Новгородская область	16 868,50	322,50
Республика Бурятия	18 019,40	392,30
Смоленская область	18 950,30	505,40
Курская область	19 995,50	536,50
Забайкальский край	20 445,60	482,00
Липецкая область	21 220,80	575,50
Ульяновская область	21 360,00	619,10
Пензенская область	21 418,80	634,30
Кировская область	21 477,10	684,00
Чувашская Республика	21 816,30	608,40
Астраханская область	22 824,90	475,80
Брянская область	23 579,30	569,80
Амурская область	23 702,60	417,30
Калужская область	24 007,20	530,50
Тульская область	27 581,20	746,60
Вологодская область	28 057,50	617,80
Алтайский край	29 815,50	1125,50
Тверская область	32 236,50	687,40
Белгородская область	32 657,40	754,90
Владимирская область	32 672,70	688,40
Мурманская область	34 351,10	482,20
Воронежская область	36 050,40	1042,40
Рязанская область	36 544,30	522,00
Калининградская область	37 136,90	459,50

Источник: данные Росстата.

По исходным данным (см. табл. 1.1) с помощью *MS Excel* были произведены расчеты и получены значения параметров уравнения парной линейной регрессии (табл. 1.2).

**Таблица 1.2.** Результаты оценки параметров МНК

Параметры	Коэффициенты уравнения парной регрессии	Стандартная ошибка	t-статистика	p-value	Нижние интервальные оценки 95%	Верхние интервальные оценки 95%
у-пересечение	3855,796	2084,617931	1,849641742	0,070796	-340,321851	8051,915
$x_1$	29,80708	4,045027426	7,368820075	2,56E-09	21,66486153	37,9493

Таблица 1.2 включает оценки параметров, их среднеквадратические ошибки, вероятности ошибочного решения (*p-value*), нижние и верхние интервальные оценки параметров с вероятностью 95%. Согласно полученным значениям уравнение парной регрессии запишется в виде

$$y = 3855,80 + 29,81x_1 + e.$$

Для того чтобы определить, на сколько процентов изменится значение  $y$  при изменении  $x$  на 1%, рассчитывается **коэффициент эластичности**

$$\varepsilon = b \frac{\bar{x}}{\bar{y}} = 29,81 \frac{449,83}{17\,263,95} = 0,78\%,$$

где  $b = 29,81$  — коэффициент регрессии, показывающий, что с увеличением количества занятых на 1 тыс. человек, поступление налогов возрастает в среднем на 29,81 млн руб.

Полученное уравнение регрессии статистически значимо. Об этом свидетельствуют результаты дисперсионного анализа (табл. 1.3).

**Таблица 1.3.** Дисперсионный анализ

Источники вариации	df	SS	MS	F-критерий
Регрессия	1	2 696 977 935	2 696 977 935	2,55649E-09
Остаток	46	2 284 753 335	49 668 550,77	
Итого	47	4 981 731 270		

В первой графе табл. 1.3 показаны источники вариации зависимой переменной; во второй — число степеней свободы;

в третьей — суммы квадратов отклонений; в четвертой — суммы квадратов отклонений, приходящиеся на одну степень свободы; в пятой — значение  $F$ -критерия.

Для парной линейной регрессии число степеней свободы равно числу параметров  $p$  минус единица:

$$df = p - 1 = m;$$

$$df = 2 - 1 = 1.$$

Число степеней свободы для остаточной вариации равно

$$df_{\text{ост}} = n - p = n - m - 1,$$

где  $m$  — число независимых переменных.

Сумма квадратов отклонений для регрессии называется **объясненной** (или **факторной**) и определяется по формуле

$$SS_{\text{факт}} = \sum_{(i)} (\hat{y}_i - \bar{y})^2.$$

Остаточная сумма квадратов отклонений имеет вид

$$SS_{\text{ост}} = \sum_{(i)} (y_i - \hat{y}_i)^2.$$

В соответствии с правилом сложения дисперсий сумма квадратов отклонений объясненной и остаточной вариации есть не что иное, как общая вариация зависимой переменной:

$$SS_{\text{общ}} = SS_{\text{факт}} + SS_{\text{ост}}.$$

Тогда

$$MS_{\text{факт}} = \frac{SS_{\text{факт}}}{df_{\text{факт}}};$$

$$MS_{\text{ост}} = \frac{SS_{\text{ост}}}{df_{\text{ост}}},$$

где  $MS$  — сумма квадратов отклонений в расчете на одну степень свободы.

Отсюда имеем формулу

$$F = \frac{MS_{\text{факт}}}{MS_{\text{ост}}}$$

В нашем примере  $F = 54,3$ . Полученное значение  $F$ -критерия необходимо сравнить с табличным значением, соответствующим гипотезе  $H_0$ :  $\sigma_{\text{факт}}^2 = \sigma_{\text{ост}}^2$ . Распределение  $F$ -статистики зависит от числа степеней свободы числителя ( $df_{\text{факт}}$ ) и знаменателя ( $df_{\text{ост}}$ ), а также от уровня значимости, т.е. вероятности ошибочного отклонения  $H_0$ . На 5%-ном уровне значимости  $df_{\text{факт}} = 1$ ,  $df_{\text{ост}} = 46$ ,  $F_{\text{таб}} = 4,05$ . Поскольку  $F > F_{\text{таб}}$ , гипотеза  $H_0$  не принимается.

В табл. 1.4 приведены значения зависимой переменной, рассчитанные по уравнению регрессии  $\hat{y} = 3855,80 + 29,81x_1$ , а также значения остатков  $e_i = y_i - \hat{y}_i$ .

**Таблица 1.4.** Вывод остатка

Наблюдение	Расчетное значение поступления налогов, млн руб.	Остаток, млн руб.
1	7051,115	-5628,915243
2	6308,919	-3779,218968
3	6884,196	-4255,095599
4	6466,896	-3702,596489
5	9468,469	-6120,969374
6	7489,279	-3575,07309
7	9432,701	-5031,90 879
8	13 587,81	-7683,807734
9	6750,064	206,636258
10	14 362,79	-6767,691796
11	11 447,66	-2189,85944
12	4764,912	4552,18774
13	13 504,35	-3525,547912
14	13 492,43	-3347,62508
15	14 496,92	-4281,523653
16	13 924,63	-2575,12773
17	15 569,98	-3523,078508
18	16 941,1	-4879,704156
19	14 017,03	-1912,829676
20	9519,141	3523,258591

Окончание табл. 1.4

Наблюдение	Расчетное значение поступления налогов, млн руб.	Остаток, млн руб.
21	15 033,45	-1929,15108
22	18 497,03	-5100,633695
23	36 765,79	-22 595,4926
24	18 744,43	-4517,432453
25	13 468,58	3399,920583
26	15 549,11	2470,286448
27	18 920,29	30,00 77866
28	19 847,29	148,2056123
29	18 222,81	2222,791434
30	21 009,77	211,0295195
31	22 309,36	-949,3591381
32	22 762,43	-1343,626743
33	24 243,84	-2766,738585
34	21 990,42	-174,1233895
35	18 038	4786,895326
36	20 839,87	2739,429872
37	16 294,29	7408,309465
38	19 668,45	4338,748088
39	26 109,76	1471,438251
40	22 270,61	5786,890065
41	37 403,66	-7588,164097
42	24 345,18	7891,317346
43	26 357,16	6300,239493
44	24 374,99	8297,710266
45	18 228,77	16 122,33002
46	34 926,7	1123,704193
47	19 415,09	17 129,20826
48	17 552,15	19 584,75072

### 1.3. Свойства остатков

*Первое свойство остатков* следует из уравнения (1.1), которое показывает, что

$$\sum_{i=1}^n e_i x_i = 0,$$

т.е. остатки и объясняющая переменная не коррелированы.

**Второе свойство остатков**, которым обладают оценки, полученные МНК, представляется в виде

$$\sum_{i=1}^n y_i = \sum_{i=1}^n \hat{y}_i,$$

или

$$\sum_{i=1}^n e_i \hat{y}_i = 0,$$

т.е. остатки и предсказанные значения  $\hat{y}_i$  не коррелированы.

**Третье свойство остатков** ( $e_i$ ) имеет вид

$$E(e_i) = 0$$

математическое ожидание остатков равно нулю. В выборке  $\bar{e} = 0$ .

**Четвертое свойство остатков**: остатки имеют постоянную дисперсию, т.е.  $\text{var}(e_i) = \sigma^2$  для всех  $i = 1, 2, \dots, N$ . Дисперсия остатков равна

$$\sigma_e^2 = 1 \quad (\text{в выборке } s_e^2 = 1).$$

**Пятое свойство остатков**: остатки не коррелированы между собой

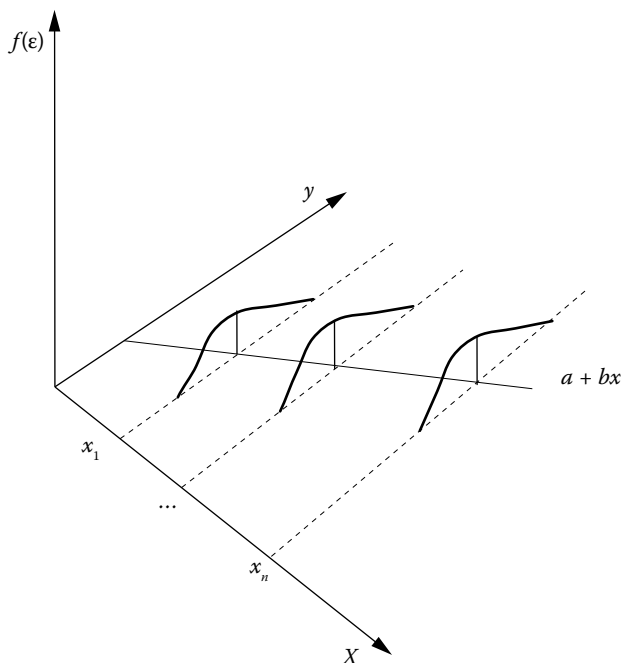
$$E(\varepsilon_i \varepsilon_j) = 0$$

$$\text{для } i \neq j; \quad i, j = 1, 2, \dots, n.$$

Зная остаток для  $i$ -го наблюдения, мы ничего не можем сказать об остатке для  $j$ -го наблюдения ( $i \neq j$ ). Заметим, что это свойство остатков исчезает при построении регрессии по временным рядам, в которых наблюдение каждого последующего года (месяца, квартала) зависит от наблюдения предыдущего года (месяца, квартала).

Четвертое и пятое свойства остатков можно проиллюстрировать графически (рис. 1.1).

Постоянство дисперсии остатков называют **гомоскедастичностью** остатков. Если же дисперсия остатков не постоянна, то имеет место **гетероскедастичность** остатков.



**Рис. 1.1.** Распределение остатков при фиксированных значениях независимой переменной  $X$

Дисперсия остатков регрессии  $\sigma_\varepsilon^2$  неизвестна и должна быть оценена. При этом нужно принять во внимание, что дисперсии оцененных параметров  $s_b^2$  и  $s_a^2$  зависят от дисперсии остатков  $s_\varepsilon^2$ . МНК-оценка дисперсии остатков имеет вид

$$s_\varepsilon^2 = \sum_{i=1}^n e_i^2 / (n-2),$$

или

$$E(s_\varepsilon^2) = E\left[\sum_{i=1}^n e_i^2 / (n-2)\right] = \sigma_\varepsilon^2.$$

Если распределение остатков не нормально, то наилучшим методом их оценки будет не МНК, а метод максимального правдоподобия.



Измерение ошибки аппроксимации определяется величиной коэффициента детерминации  $R^2$ :

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}.$$

Эта величина принимает значения от нуля до единицы.

Коэффициент  $R^2$  может быть представлен также в виде следующей формулы:

$$R^2 = 1 - \frac{\sum_{i=1}^n e_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2},$$

где  $e_i$  — остаток (величина отклонения от линии регрессии, которая минимизируется МНК).

Если  $\sum_{i=1}^n e_i^2$  велика, то регрессия не объясняет вариацию  $y$ .

В этом случае коэффициент детерминации  $R^2$  будет невелик (близок к нулю). Если многие точки фактических наблюдений лежат на линии регрессии или располагаются вблизи от нее,  $\sum_{i=1}^n e_i^2$  мала и, соответственно, коэффициент детерминации  $R^2$

будет близок к единице. Если все наблюдения располагаются на линии регрессии, то  $y_i = \hat{y}_i$ ,  $e_i = 0$ , а значит, и  $\sum_{i=1}^n e_i^2 = 0$ , тогда

$R^2 = 1$ . Если же точки, соответствующие наблюдениям, не лежат на линии регрессии, то  $\sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = 0$ , а значит, регрессия

не объясняет вариацию  $y_i$ . В случае если  $\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n e_i^2$ ,  $R^2 = 0$ .

Тогда наилучшей аппроксимацией данных будет линия  $\hat{y}_i = \bar{y}$  для всех  $i$ , т.е. это горизонтальная линия, параллельная оси абсцисс и проходящая через точку  $y_i = \bar{y}$ , что соответствует случаю полной независимости переменных  $x$  и  $y$ .

Степень аппроксимации данных выборки, полученной регрессией  $\hat{y} = a + bx$ , оценивается с помощью средней ошибки аппроксимации

$$\bar{e} = \frac{\sum |y_i - \bar{y}|}{n}.$$

Большей информативностью обладает средняя относительная ошибка аппроксимации

$$\bar{e}^* = \left( \frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{\sum_{i=1}^n y_i} \right) \cdot 100.$$

Значения средней относительной ошибки аппроксимации, не превышающие 10%, свидетельствуют о хорошем соответствии линии регрессии исходным данным.

Коэффициент детерминации  $R^2$  имеет две трактовки:

- 1) это квадрат коэффициента парной корреляции между фактическими и расчетными значениями зависимой переменной, т.е.  $R^2 = r_{\hat{y}y}^2$ ;
- 2) это квадрат коэффициента парной корреляции между  $y$  и  $x$ , т.е.  $R^2 = r_{yx}^2$ .

$$r_{yx}^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}.$$

Коэффициент парной корреляции — это мера тесноты линейной связи:

$$r_{yx} = \frac{\sum_{(i)} (y_i - \bar{y})(x_i - \bar{x})}{\sqrt{\sum (y_i - \bar{y})^2 \sum (x_i - \bar{x})^2}}; \quad -1 < r_{yx} < 1.$$

Можно представить коэффициент парной корреляции  $r_{yx}$  через коэффициент регрессии  $b$ , определяющий наклон регрессии к оси абсцисс:

$$r_{yx} = b \frac{s_x}{s_y}.$$

По данным рассматриваемого примера 1.1, коэффициент парной корреляции равен  $r_{yx} = 0,73$ , а коэффициент детерминации  $r_{yx}^2 = 0,54$ , т.е. на 54% вариация объема собираемых

налогов зависит от количества занятых, а на 46% — от других факторов.

Подчеркнем, что коэффициент парной корреляции представляет собой меру **линейной** связи между  $y$  и  $x$ . Если, например, существует полная квадратическая связь между  $y$  и  $x$ , которая описывается уравнением регрессии  $\hat{y} = a + bx + cx^2$ , то  $r_{yx}$ , а значит и  $r_{yx}^2$ , может иметь значения, далекие от единицы. В этом случае нужно провести преобразование переменных с тем, чтобы линеаризовать уравнение регрессии. Примем  $x^2 = z$ , тогда параболическое уравнение регрессии будет иметь вид  $\hat{y} = a + bx + cz$ , т.е. это уравнение адекватно линейному уравнению множественной регрессии (см. гл. 2).

Сформулированные свойства остатков проверяются после нахождения параметров уравнения регрессии. По уравнению регрессии находятся расчетные (предсказанные) значения зависимой переменной ( $\hat{y}$ ). После этого производится расчет остатков  $e_i = y_i - \hat{y}_i$ . Затем строится график остатков (*residual plot*). В случае гомоскедастичности остатков положительные и отрицательные значения остатков чередуются и находятся в области, параллельной оси абсцисс (рис. 1.2).

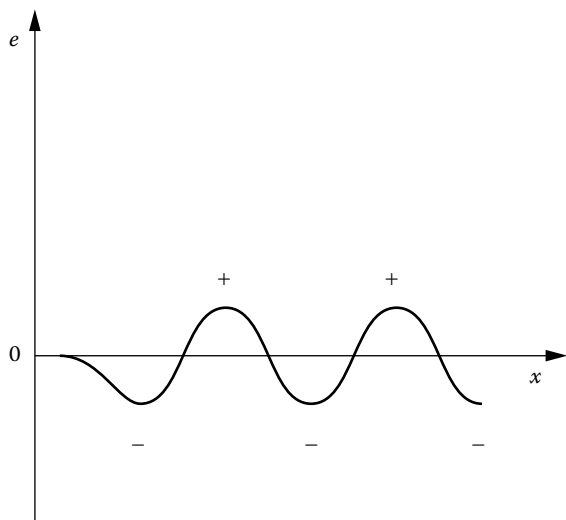
График остатков по данным нашего примера о зависимости выпуска продукции от стоимости основных фондов представлен на рис. 1.3. Расположение «облака» остатков позволяет предположить наличие гетероскедастичности.

Кроме визуального анализа остатков существует ряд специальных тестов, позволяющих выявить гетероскедастичность остатков: тесты Гольдфелда — Квандта, Парка, Глейзера, Уайта, ранговой корреляции Спирмена и др. Названные тесты рассмотрены в гл. 2.

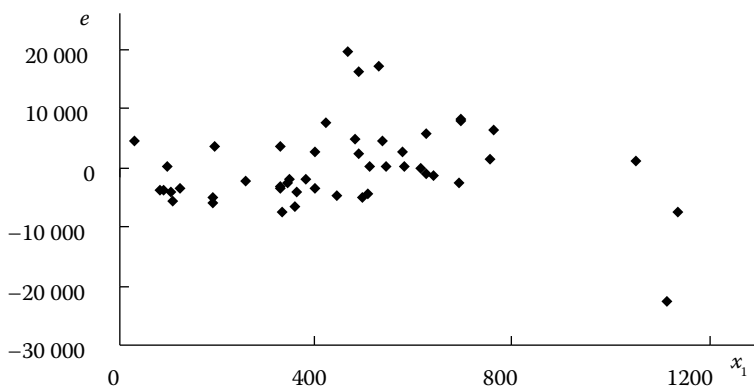
Наличие гетероскедастичности сказывается на точности предсказания значения зависимой переменной на основе регрессии. Ошибка предсказания может быть представлена формулой

$$y_0 - \hat{y}_0 = [y_0 - E(y_0)] + [E(y_0) - \hat{y}_0] = e_0 + [E(y_0) - \hat{y}_0],$$

где  $y_0$  — фактическое значение;  $\hat{y}_0$  — предсказанное значение для объекта 0;  $[y_0 - E(y_0)]$  — отклонение фактического значения от «истинного» в генеральной совокупности;  $[E(y_0) - \hat{y}_0]$  — отклонение «истинного» значения от предсказанного по регрессии.



**Рис. 1.2.** График остатков (случай гомоскедастичности)



**Рис. 1.3.** График остатков по данным примера

Тогда дисперсия ошибки предсказания имеет вид

$$\begin{aligned} \text{var}(e_0) + \text{var}[E(y_0) - \hat{y}_0] + 2\text{cov}[e_0 E(y_0) - \hat{y}_0] = \\ = \sigma_{\text{ост}}^2 [1 + (1/n)] + (x_0 - \bar{x})^2 / \sum_{i=1}^n (x_i - \bar{x})^2. \end{aligned}$$

Таким образом, можно утверждать, что ошибка предсказания значения  $y$  на основе уравнения регрессии зависит от остаточной дисперсии  $y$  (дисперсии остатков), от объема выборки  $n$ , от того, насколько значение  $x_0$  (для объекта предсказания) отличается от среднего значения по наблюдаемым данным. Чем меньше остаточная дисперсия, чем больше объем выборки  $n$  и чем меньше вариация  $x$ , т.е.  $\sum (x_i - \bar{x})^2$ , а также чем ближе  $x_0$  к  $\bar{x}$ , тем меньше ошибка предсказания.

Как уже отмечалось, на основе уравнения регрессии могут быть получены точечные прогнозные значения,  $\hat{y}_i$ , и интервальные оценки. Например, 95%-ный доверительный интервал значения  $\hat{y}_0$  при  $x_0$  имеет вид

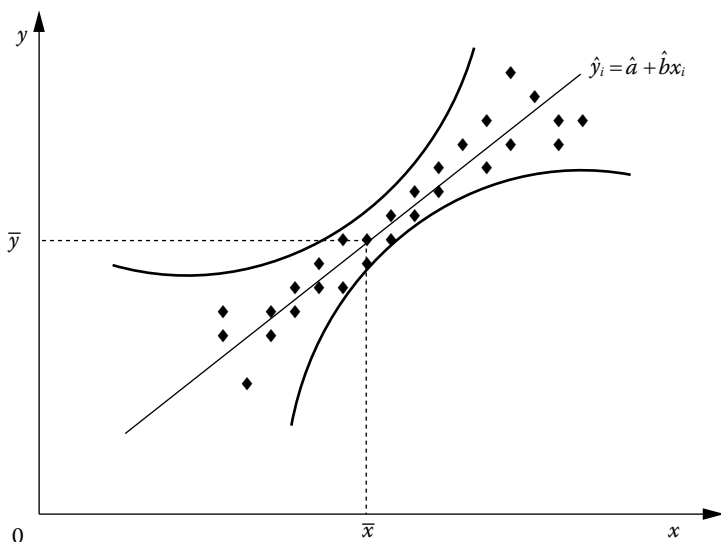
$$(a + bx_0) \pm t_{0,025; n-2} \left\{ s_{\text{ост}} \left[ 1 + (1/n) + (x_0 - \bar{x})^2 / \sum_{i=1}^n (x_i - \bar{x})^2 \right]^{1/2} \right\},$$

где  $s_{\text{ост}}$  — выборочная оценка остаточной дисперсии  $y$  (заменяет  $\sigma_{\text{ост}}$ );  $t_{0,025; n-2}$  — критическое 2,5%-ное значение, полученное по таблице  $t$ -распределения с  $(n - 2)$  степенями свободы.

Как показано на рис. 1.4, границы доверительного интервала представляют собой гиперболу. Самое «узкое» значение интервала — в точке  $\bar{x}$  и ее окрестностях; чем больше удалено значение  $x_i$  от  $\bar{x}$ , тем шире становится доверительный интервал.

Устранение из выборки регионов, для которых значение остатка оказалось большим (в нашем примере это Калининградская, Рязанская, Воронежская области), позволило бы существенно улучшить свойства регрессии и повысить точность предсказания объема выпуска.

В заключение отметим, что парная регрессия довольно редко выступает в качестве эконометрической модели, поскольку исследуемые экономические явления формируются под влиянием не одного, а нескольких факторов. Этим обстоятельством



**Рис. 1.4.** 95%-ный доверительный интервал

определяется гораздо бóльшая распространенность множественной регрессии в экономическом моделировании.

## Контрольные вопросы и задания

1. Дайте определение парной регрессии.
2. Что такое «линия регрессии»?
3. Поясните экономическую сущность параметров уравнения парной линейной регрессии.
4. Как производится оценка параметров уравнения парной регрессии?
5. Как производится оценка качества уравнения в целом?
6. Каковы свойства остатков модели парной регрессии?
7. Как измеряется ошибка аппроксимации?
8. Какие трактовки коэффициента детерминации вам известны?
9. Как осуществляется линейаризация модели?
10. Как анализируется график остатков модели регрессии?
11. От чего зависит точность предсказания значения зависимой переменной по уравнению парной регрессии?

**2.1. Множественная линейная  
регрессия в скалярной  
и векторной формах**

Уравнение множественной регрессии отражает корреляционную связь результативной (зависимой) переменной  $y$  и нескольких объясняющих (независимых)  $x_1, x_2, \dots, x_p$ :

$$y = f(x_1, x_2, \dots, x_p, \varepsilon), \quad (2.1)$$

где  $y$  — результативная переменная (зависимая, объясняемая);  $x_1, x_2, \dots, x_p$  — объясняющие переменные (независимые);  $\varepsilon$  — случайный остаток;  $f$  — некая математическая функция.

Если в качестве функции в формуле (2.1) выбрана линейная, уравнение регрессии называется *уравнением множественной линейной регрессии* и имеет вид

$$y = \alpha_0 + \alpha_1 x_1 + \dots + \alpha_p x_p + \varepsilon,$$

где  $\alpha_0, \alpha_1, \dots, \alpha_p$  — параметры.

В ряде случаев удобнее пользоваться матричной записью этого уравнения

$$Y = X\alpha + \varepsilon,$$

где  $X$  — матрица значений независимых переменных

$$X = \begin{pmatrix} 1 & x_{11} & \dots & x_{p1} \\ 1 & x_{12} & \dots & x_{p2} \\ \dots & \dots & \dots & \dots \\ 1 & x_{1n} & \dots & x_{pn} \end{pmatrix}.$$

Первый столбец этой матрицы состоит из единиц, которые рассматриваются как значения дополнительной переменной, на которую умножается свободный член. В матрицу  $X$  входит  $p$  независимых переменных, принимающих  $n$  значений.

$Y$ ,  $\alpha$ ,  $\varepsilon$  — матрицы-столбцы значений зависимой переменной (матрица  $Y$ ), параметров регрессии (матрица  $\alpha$ ) и случайных остатков (матрица  $\varepsilon$ ):

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}; \alpha = \begin{pmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_p \end{pmatrix}; \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}.$$

Векторы  $Y$  и  $\varepsilon$  включают в себя по  $n$  значений зависимой переменной и случайных остатков, вектор  $\alpha$  состоит из значений  $(p + 1)$  параметра — свободного члена  $\alpha_0$  и  $p$  коэффициентов регрессии  $\alpha_j$ .

## 2.2. Метод наименьших квадратов и предпосылки его применения для множественной линейной регрессии

Для определения значений параметров уравнения множественной регрессии используют числовую информацию, которая рассматривается как *выборочная*. Поэтому рассчитанные на ее основе величины называют *оценками параметров*,



подчеркивая тем самым их возможную неточность из-за неполноты информации. Оценки параметров могут меняться от выборки к выборке, поэтому они рассматриваются как **случайные** величины.

Далее будем различать обозначения параметров и случайных остатков, полученных по выборке (т.е. по имеющимся у исследователя данным), и значения параметров и случайных остатков по генеральной совокупности. Для обозначения выборочных значений будут использоваться латинские буквы:  $a_0, a_1, a_2, \dots, a_p, e$ ; в матричной форме:  $\mathbf{a}, \mathbf{e}$ ; для генеральных значений — греческие буквы:  $\alpha, \alpha_1, \alpha_2, \dots, \alpha_p, \varepsilon$ ; в матричной форме:  $\boldsymbol{\alpha}, \boldsymbol{\varepsilon}$ .

Так как найденные параметры являются лишь выборочными оценками неизвестных параметров по генеральной совокупности, то возникает вопрос об их качестве. Характеристиками качества полученных оценок параметров регрессии являются их несмещенность, эффективность и состоятельность.

Оценка параметра является **несмещенной**, если ее математическое ожидание равно оцениваемому параметру. Например, математическое ожидание оценки коэффициента регрессии  $a_j$  равно его значению в генеральной совокупности  $\alpha_j$ :

$$Ma_j = \alpha_j.$$

Оценка параметра является **эффективной**, если она имеет наименьшую дисперсию среди всех возможных оценок данного параметра по выборкам одного и того же объема:

$$M(a_{j,eff} - \alpha_j)^2 = \sigma_{a_j,eff}^2 = \min_k \sigma_{\alpha_j,k}^2$$

где  $a_{j,eff}$  — эффективная оценка параметра  $\alpha_j$ ;  $\sigma_{a_j,eff}^2$  — дисперсия эффективной оценки параметра  $\alpha_j$ ;  $\sigma_{\alpha_j,k}^2$  — дисперсия оценки параметра  $\alpha_j$ , полученной при применении  $k$ -го метода.

Оценка параметра является **состоятельной**, если с увеличением числа наблюдений она стремится к значению параметра в генеральной совокупности:

$$a_j \xrightarrow{n \rightarrow \infty} \alpha_j.$$

Простейшим методом оценки параметров множественной регрессии является МНК. **МНК-оценки** будут несмещенными, эффективными и состоятельными при выполнении определенных требований, называемых **предпосылками МНК**. Эти требования касаются статистических свойств исходных данных:

- независимые переменные являются неслучайными величинами, не связанными между собой;
- зависимая переменная является случайной величиной, не ограниченной сверху или снизу;
- для каждого набора значений независимых переменных результирующая переменная рассматривается как отдельная случайная величина  $Y_i$ . Ее распределение описывается нормальным законом с математическим ожиданием, равным выровненному значению зависимой переменной:

$$MY_i = \alpha_0 + \alpha_1 x_{1i} + \alpha_2 x_{2i} + \dots + \alpha_p x_{pi} = \hat{Y}_i.$$

Например, в модели с тремя независимыми переменными имеется наблюдение номер десять с  $x_1 = 3$ ;  $x_2 = 5$ ;  $x_3 = 2$ . Для этих значений объясняющих переменных результирующая переменная является случайной величиной  $Y_{10}$  с математическим ожиданием  $MY_{10} = \hat{Y}_{10}$ .

Вторая характеристика нормального распределения — среднее квадратическое отклонение  $\sigma_Y$  может быть любым, однако оно должно быть одинаковым для всех случайных величин  $Y_i$ :

$$\sigma_{Y_i} = \sigma_{Y_j} = \sigma_Y \quad \forall i, j;$$

- разные случайные величины  $Y_i$  и  $Y_j$  должны быть независимы друг от друга:

$$r_{Y_i Y_j} = 0, \quad i \neq j.$$

Линейная модель регрессии, для которой выполняются эти условия, называется **классической нормальной линейной моделью**.

Так как случайная величина  $Y_i$  в линейной регрессии представляет собой сумму неслучайной величины  $\hat{Y}_i$  и случайного остатка  $\varepsilon_i$ , то все перечисленные выше требования к случайной величине  $Y_i$  можно сформулировать в форме требований к случайным остаткам модели регрессии  $\varepsilon_i$ , которые (кроме последнего требования) называются *условиями Гаусса — Маркова*.

1. Математическое ожидание случайного остатка равно нулю:

$$M\varepsilon_i = 0. \quad (2.2)$$

2. Дисперсия случайных остатков одинакова для различных  $i$  и  $j$ :

$$\sigma_{\varepsilon_i}^2 = \sigma_{\varepsilon_j}^2 = \sigma_{\varepsilon}^2. \quad (2.3)$$

3. Случайные остатки не зависят друг от друга (не автокоррелированы):

$$r_{\varepsilon_i, \varepsilon_j} = 0, \quad i \neq j. \quad (2.4)$$

4. Случайные остатки не зависят от значений независимых переменных, входящих в модель регрессии:

$$r_{\varepsilon_i, X_j} = 0. \quad (2.5)$$

5. Случайные остатки распределены по нормальному закону распределения.

Метод наименьших квадратов основывается на принципе минимизации квадратов отклонений фактических значений результативного признака  $y$  от его выровненных значений  $\hat{y}$ , рассчитанных по уравнению регрессии

$$SS_e = \sum (y - \hat{y})^2 \rightarrow \min. \quad (2.6)$$

Для множественной линейной регрессии выражение (2.6) будет иметь вид

$$SS_e = \sum (y - a_0 - a_1x_1 - a_2x_2 - \dots - a_px_p)^2 \rightarrow \min.$$

Для нахождения экстремума по каждому из неизвестных параметров  $a_0, a_1, \dots, a_p$  рассчитывается производная функции  $SS_e$  и полученное выражение приравняется к нулю:

$$\begin{cases} \frac{dSS_e}{da_0} = \sum(-2)(y - a_0 - a_1x_1 - a_2x_2 - \dots - a_px_p) = 0, \\ \frac{dSS_e}{da_1} = \sum(-2x_1)(y - a_0 - a_1x_1 - a_2x_2 - \dots - a_px_p) = 0, \\ \frac{dSS_e}{da_2} = \sum(-2x_2)(y - a_0 - a_1x_1 - a_2x_2 - \dots - a_px_p) = 0, \\ \dots \\ \frac{dSS_e}{da_p} = \sum(-2x_p)(y - a_0 - a_1x_1 - a_2x_2 - \dots - a_px_p) = 0. \end{cases} \quad (2.7)$$

После преобразований система уравнений (2.7) имеет вид

$$\begin{cases} \sum y = na_0 + a_1 \sum x_1 + a_2 \sum x_2 + \dots + a_p \sum x_p, \\ \sum x_1 y = a_0 \sum x_1 + a_1 \sum x_1^2 + a_2 \sum x_1 x_2 + \dots + a_p \sum x_1 x_p, \\ \sum x_2 y = a_0 \sum x_2 + a_1 \sum x_1 x_2 + a_2 \sum x_2^2 + \dots + a_p \sum x_2 x_p, \\ \dots \\ \sum x_p y = a_0 \sum x_p + a_1 \sum x_1 x_p + a_2 \sum x_2 x_p + \dots + a_p \sum x_p^2. \end{cases} \quad (2.8)$$

Система (2.8) называется *системой нормальных уравнений*.

В матричной форме применение МНК приводит к следующему результату:

$$\mathbf{a} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}. \quad (2.9)$$

### Пример

Имеются статистические данные о значениях четырех показателей в разрезе 48 субъектов РФ, которые приведены в табл. 2.1.

1. Поступление налогов, сборов и иных обязательных платежей в консолидированный бюджет РФ (без поступлений ЕСН) в 2009 г., млн. руб.
2. Количество занятых в Российской Федерации в 2009 г., тыс. человек.
3. Объем отгруженных товаров собственного производства, выполненных работ и оказанных услуг собственными силами

по виду экономической деятельности «Обрабатывающие производства» в Российской Федерации в 2009 г., млн руб.

4. Объем отгруженных товаров собственного производства, выполненных работ и услуг собственными силами по виду экономической деятельности «Производство и распределение электроэнергии, газа и воды» в Российской Федерации в 2009 г., млн руб.

Предположим, что зависимой переменной  $y$  в данном случае является «поступление налогов, сборов и иных обязательных платежей в консолидированный бюджет РФ». Для краткости будем называть эту переменную «поступление налогов». Независимыми переменными являются остальные три переменные, которые мы будем называть и обозначать как «количество занятых»  $x_1$ , «отгрузка в обрабатывающих производствах»  $x_2$ , «производство энергии»  $x_3$ . Допустим также, что связь между зависимой и независимыми переменными может быть выражена в виде линейной функции регрессии

$$y = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_3 + \varepsilon,$$

а случайные остатки удовлетворяют условиям Гаусса — Маркова.

**Таблица 2.1.** Некоторые экономические показатели деятельности субъектов РФ в 2009 г.

Субъект РФ	Поступление налогов, млн руб.	Количество занятых, тыс. человек	Отгрузка в обрабатывающих производствах, млн руб.	Производство энергии, млн руб.
	$y$	$x_1$	$x_2$	$x_3$
Республика Ингушетия	1422,20	107,20	266	733
Еврейская автономная область	2529,70	82,30	2865	2040
Республика Тыва	2629,10	101,60	431	2023
Республика Алтай	2764,30	87,60	1228	1176
Карачаево-Черкесская Республика	3347,50	188,30	10 921	4275
Республика Калмыкия	3914,20	121,90	928	1495
Республика Адыгея	4400,80	187,10	12 565	3000
Республика Северная Осетия — Алания	5904,00	326,50	11 088	3337
Магаданская область	6956,70	97,10	2486	8305
Кабардино-Балкарская Республика	7595,10	352,50	17 609	5790

Продолжение табл. 2.1

Субъект РФ	Поступление налогов, млн руб.	Количество занятых, тыс. человек	Отгрузка в обрабатывающих производствах, млн руб.	Производство энергии, млн руб.
	$y$	$x_1$	$x_2$	$x_3$
Республика Хакасия	9257,80	254,70	39 640	17 634
Чукотский автономный округ	9317,10	30,50	531	6226
Республика Марий Эл	9978,80	323,70	46 180	7489
Псковская область	10 144,80	323,30	32 074	5640
Чеченская Республика	10 215,40	357,00	579	6170
Республика Карелия	11 349,50	337,80	39 962	14 684
Курганская область	12 046,90	393,00	38 308	12 093
Республика Мордовия	12 061,40	439,00	65 507	7312
Костромская область	12 104,20	340,90	50 532	20 922
Камчатский край	13 042,40	190,00	11 245	12 721
Орловская область	13 104,30	375,00	38 089	9612
Ивановская область	13 396,40	491,20	42 865	18 506
Республика Дагестан	14 170,30	1104,10	21 031	12 573
Тамбовская область	14 227,00	499,50	47 738	9522
Новгородская область	16 868,50	322,50	80 915	8989
Республика Бурятия	18 019,40	392,30	29 660	12 532
Смоленская область	18 950,30	505,40	78 278	43 604
Курская область	19 995,50	536,50	67 241	43 733
Забайкальский край	20 445,60	482,00	7910	13 687
Липецкая область	21 220,80	575,50	228 812	17 311
Ульяновская область	21 360,00	619,10	76 523	16 471
Пензенская область	21 418,80	634,30	71 307	12 061
Кировская область	21 477,10	684,00	76 151	20 857
Чувашская Республика	21 816,30	608,40	85 926	17 071
Астраханская область	22 824,90	475,80	34 576	10 532
Брянская область	23 579,30	569,80	57 187	10 519
Амурская область	23 702,60	417,30	16 412	16 512
Калужская область	24 007,20	530,50	161 769	10 369
Тульская область	27 581,20	746,60	182 031	24 376
Вологодская область	28 057,50	617,80	236 267	23 180
Алтайский край	29 815,50	1125,50	115 197	24 804
Тверская область	32 236,50	687,40	103 158	44 961
Белгородская область	32 657,40	754,90	233 608	18 773
Владимирская область	32 672,70	688,40	142 867	20 093
Мурманская область	34 351,10	482,20	49 081	34 395

Окончание табл. 2.1

Субъект РФ	Поступление налогов, млн руб.	Количество занятых, тыс. человек	Отгрузка в обрабатывающих производствах, млн руб.	Производство энергии, млн руб.
	$y$	$x_1$	$x_2$	$x_3$
Воронежская область	36 050,40	1042,40	125 343	39 170
Рязанская область	36 544,30	522,00	95 522	23 932
Калининградская область	37 136,90	459,50	147 573	15 429

Источник: данные Росстата.

Применяя к исходным данным (см. табл. 2.1) МНК, оценим параметры регрессии. Система нормальных уравнений для нашего примера имеет вид

$$\begin{cases} \sum y = na_0 + a_1 \sum x_1 + a_2 \sum x_2 + a_3 \sum x_3, \\ \sum x_1 y = a_0 \sum x_1 + a_1 \sum x_1^2 + a_2 \sum x_1 x_2 + a_3 \sum x_1 x_3, \\ \sum x_2 y = a_0 \sum x_2 + a_1 \sum x_1 x_2 + a_2 \sum x_2^2 + a_3 \sum x_2 x_3, \\ \sum x_3 y = a_0 \sum x_3 + a_1 \sum x_1 x_3 + a_2 \sum x_2 x_3 + a_3 \sum x_3^2. \end{cases}$$

После вычислений получаем уравнение

$$y = 3226,44 + 12,45x_1 + 0,06x_2 + 0,31x_3 + e.$$

Из уравнения регрессии следует, что между сбором налогов и независимыми переменными, входящими в модель, наблюдается прямая связь. Напомним, что коэффициенты при независимых переменных называются **коэффициентами регрессии**. Они являются абсолютными показателями силы связи и характеризуют среднее изменение зависимой переменной при единичном изменении независимой переменной — множителя данного коэффициента при условии неизменности остальных независимых переменных, включенных в уравнение (модель) регрессии.

В частности, можно сделать вывод, что с изменением количества занятых на 1 тыс. человек поступление налогов в среднем изменится в ту же сторону на 12,45 млн руб. при неизменном объеме отгрузки в обрабатывающих производствах и производстве энергии.

Изменение объема отгрузки в обрабатывающих производствах на 1 млн руб. приведет к изменению поступления налогов

в среднем на 0,06 млн руб. при неизменных значениях количества занятых и производства энергии.

При изменении производства энергии на 1 млн руб. поступление налогов в среднем изменится на 0,31 млн руб. при неизменных значениях количества занятых и объема отгрузки в обрабатывающих производствах.

Величины коэффициентов регрессии определяются не только силой связи между показателями, но и масштабом их измерения, и поэтому несравнимы между собой. Для сопоставления независимых переменных по силе их влияния на результативную используют относительные показатели силы связи — **коэффициенты эластичности**.

Общая формула коэффициента эластичности по переменной  $x_j$  имеет вид

$$\varepsilon_{x_j} = \frac{dy}{dx_j} \cdot \frac{x_j}{\hat{y}},$$

где  $dy/dx_j$  — частная производная функции регрессии по переменной  $x_j$ ;  $\hat{y}$  — выровненное значение зависимой переменной  $y$  при заданном значении переменной  $x_j$ .

Отметим, что коэффициенты эластичности могут быть рассчитаны для любой функции — как линейной, так и нелинейной. Например, для множественной линейной регрессии коэффициент эластичности по переменной  $x_j$  будет равен

$$\varepsilon_{x_j} = a_j \cdot \frac{x_j}{a_0 + a_1 x_1 + a_2 x_2 + \dots + a_j x_j + \dots + a_p x_p}. \quad (2.10)$$

Из формулы (2.10) следует, что коэффициент эластичности для линейной функции зависит от конкретных значений независимых переменных, включенных в модель. Так как коэффициент эластичности измеряет влияние переменной  $x_j$  на переменную  $y$ , то значения остальных независимых переменных принято фиксировать на их среднем уровне:

$$\varepsilon_{x_j} = a_j \cdot \frac{x_j}{a_0 + a_1 \bar{x}_1 + a_2 \bar{x}_2 + \dots + a_j \bar{x}_j + \dots + a_p \bar{x}_p}.$$



Таким образом, при фиксированных значениях других переменных существует целый ряд коэффициентов эластичности по переменной  $x_j$ , определяемый областью ее значений. Они называются *частными* коэффициентами эластичности. Если зафиксировать значение  $x_j$  на среднем уровне, получим *средний* коэффициент эластичности (или коэффициент эластичности для среднего значения  $x_j$ ):

$$\bar{\Theta}_{x_j} = a_j \cdot \frac{\bar{x}_j}{a_0 + a_1 \bar{x}_1 + a_2 \bar{x}_2 + \dots + a_j \bar{x}_j + \dots + a_p \bar{x}_p}. \quad (2.11)$$

Как следует из МНК для линейной регрессии, выражение в знаменателе (2.11) равно среднему значению зависимой переменной, что позволяет упростить формулу:

$$\bar{\Theta}_{x_j} = a_j \cdot \frac{\bar{x}_j}{\bar{y}}.$$

Коэффициенты эластичности показывают, на сколько процентов в среднем изменится зависимая переменная при изменении переменной  $x_j$  на 1% и значениях других независимых переменных, фиксированных на средних уровнях. Так как величина коэффициента эластичности зависит от выбранного значения независимой переменной, при формулировке вывода следует указывать это значение переменной  $x_j$  и соответствующее ему выровненное значение зависимой переменной.

Рассмотрим средние коэффициенты эластичности для регрессии, полученной в примере 2.1.

Для характеристики силы связи количества занятых и поступления налогов рассчитаем коэффициент эластичности

$$\bar{\Theta}_{x_1} = 12,45 \cdot \frac{449,83}{17\,263,95} = 0,32.$$

Следовательно, с изменением количества занятых на 1% от среднего уровня поступление налогов в среднем изменится в ту же сторону на 0,32% своего среднего уровня при неизменном объеме отгрузки в обрабатывающих производствах и производстве энергии.

Коэффициенты эластичности для переменных «объем отгрузки в обрабатывающих производствах» и «производство энергии» рассчитываются и интерпретируются аналогично. Они равны

$$\bar{\varepsilon}_{x_2} = 0,22; \quad \bar{\varepsilon}_{x_3} = 0,27.$$

Сопоставляя полученные коэффициенты эластичности, можно сделать вывод о том, что наиболее сильное влияние на сбор налогов оказывает число занятых.

Другими относительными показателями силы связи являются **стандартизованные коэффициенты регрессии**. В отличие от коэффициентов эластичности, они рассчитываются только для случая множественной линейной регрессии по формуле

$$b_j = a_j \frac{\sigma_{x_j}}{\sigma_y}.$$

Для примера 2.1 стандартизованные коэффициенты регрессии равны

$$b_1 = 12,45 \frac{251,48}{10182,54} = 0,31; \quad b_2 = 0,37; \quad b_3 = 0,34.$$

Стандартизованные коэффициенты регрессии могут быть проинтерпретированы. Величина  $b_j$  показывает, на сколько своих средних квадратических отклонений  $\sigma_y$  в среднем изменится зависимая переменная при изменении переменной  $x_j$  на одно свое среднее квадратическое отклонение  $\sigma_{x_j}$  при фиксированных значениях других переменных, включенных в модель регрессии. Например, при изменении количества занятых на  $1\sigma_{x_1}$  поступление налогов в среднем изменится в ту же сторону на  $0,32\sigma_y$  при неизменном объеме отгрузки в обрабатывающих производствах и производстве энергии.

Сопоставление стандартизованных коэффициентов регрессии для примера 2.1 приводит к другим, чем по коэффициентам эластичности, выводам о том, какая переменная оказывает более сильное влияние на результативную переменную. В данном случае это  $x_2$ , т.е. объем отгрузки в обрабатывающих производствах. Такое изменение в выводах объясняется разным уровнем вариации независимых переменных. Можно

показать, что отношение стандартизованных коэффициентов регрессии равно

$$\frac{b_i}{b_j} = \frac{\partial_{x_i} v_{x_i}}{\partial_{x_j} v_{x_j}},$$

где  $v_{x_i}, v_{x_j}$  — коэффициенты вариации переменных  $x_i$  и  $x_j$ .

Стандартизованные коэффициенты регрессии являются коэффициентами в функции регрессии, построенной на основе **стандартизованных переменных** — переменных, полученных в результате преобразования исходных данных по формулам

$$t_{yi} = \frac{y_i - \bar{y}}{\sigma_y}; \quad t_{ji} = \frac{x_{ji} - \bar{x}_j}{\sigma_{x_j}},$$

где  $t_{yi}$  —  $i$ -е значение стандартизованной зависимой переменной;  $t_{ji}$  —  $i$ -е значение стандартизованной  $j$ -й независимой переменной;  $\bar{y}, \bar{x}_j$  — средние значения переменных  $y, x_j$ ;  $\sigma_y, \sigma_{x_j}$  — средние квадратические отклонения переменных  $y, x_j$ .

Уравнение регрессии, включающее стандартизованные переменные, называют **уравнением регрессии в стандартизованном масштабе** (форме, виде):

$$t_y = b_1 t_1 + b_2 t_2 + \dots + b_p t_p + \delta, \quad (2.12)$$

где  $\delta$  — случайный остаток.

Так как уравнение  $y = \alpha_0 + \alpha_1 x_1 + \dots + \alpha_p x_p + \varepsilon$  построено на основе непреобразованных переменных, его еще называют **уравнением регрессии в натуральном масштабе**.

Особенностью уравнения регрессии в стандартизованном масштабе является отсутствие в нем свободного члена. Это свойство является следствием двух утверждений:

- среднее значение стандартизованной переменной равно нулю;
- свободный член уравнения множественной линейной регрессии равен

$$\alpha_0 = \bar{y} - \alpha_1 \bar{x}_1 - \alpha_2 \bar{x}_2 - \dots - \alpha_p \bar{x}_p.$$

Таким образом, прямая, построенная на основе уравнения регрессии в стандартизованном масштабе, всегда проходит через начало координат. Этим же свойством обладает уравнение регрессии, построенное на основе **центрированных переменных**, которые рассчитываются по формулам

$$t_{yi}^* = y_i - \bar{y}; \quad t_{ji}^* = x_{ji} - \bar{x}_j. \quad (2.13)$$

При использовании центрированных переменных коэффициенты регрессии и случайные остатки равны коэффициентам регрессии и случайным остаткам уравнения в натуральном масштабе:

$$t_y^* = a_1 t_1^* + a_2 t_2^* + \dots + a_p t_p^* + e. \quad (2.14)$$

На рис. 2.1 дана геометрическая интерпретация МНК для множественной линейной регрессии. Каждую переменную, образующую столбец матрицы  $X$ , можно рассматривать как вектор-столбец  $X_j$ , где  $j = 0 \dots p$ . Так как векторы  $X_j$  образуют многомерное пространство, изобразить их на рисунке невозможно. Можно, однако, воспользоваться тем, что вектор  $\hat{Y}$  (содержащий выровненные значения зависимой переменной) является линейной комбинацией векторов  $X_j$ . На рис. 2.1 это вектор  $\hat{Y} = Xa$ , лежащий в подпространстве, образованном независимыми переменными, которое называют **пространством оценок**. Вектор  $Y$  лежит вне этого подпространства. Для удобства будем считать, что все векторы выходят из точки 0. Тогда начало вектора  $\hat{Y}$  также будет исходить из точки 0. При формировании проекции вектора  $Y$  на пространство оценок наименьшую длину будет иметь перпендикуляр к этому пространству — вектор случайных остатков  $e$ . Следовательно, параметры вектора выровненных значений  $\hat{Y} = Xa$  должны быть подобраны таким образом, чтобы векторы  $Y$ ,  $\hat{Y}$ ,  $e$  сформировали прямоугольный треугольник.

Так как вектор  $e$  ортогонален вектору  $\hat{Y}$ , то их произведение равно нулю:

$$\hat{Y}^T e = 0;$$

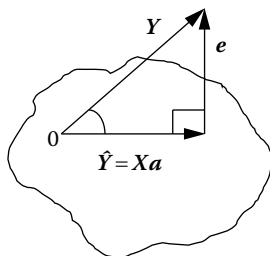


Рис. 2.1. Геометрическая интерпретация МНК

$$\mathbf{a}^T \mathbf{X}^T (\mathbf{Y} - \mathbf{X}\mathbf{a}) = 0;$$

$$\mathbf{a}^T (\mathbf{X}^T \mathbf{Y} - \mathbf{X}^T \mathbf{X}\mathbf{a}) = 0.$$

Из последнего выражения следует, что вектор параметров  $\mathbf{a}$  или равен нулю (нулевая проекция, вектор  $\mathbf{Y}$  ортогонален векторам  $\mathbf{X}$ , переменная  $y$  не зависит от переменных  $x$ ), или  $\mathbf{X}^T \mathbf{Y} = \mathbf{X}^T \mathbf{X}\mathbf{a}$ .

Откуда

$$\mathbf{a} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}.$$

Таким образом, на основе геометрической интерпретации мы получили ту же формулу для оценки параметров регрессии, что и аналитическим путем.

### 2.3. Следствия выполнения предпосылок Гаусса — Маркова

**Теорема Гаусса — Маркова** гласит, что при выполнении предпосылок (2.2) — (2.5) оценка параметров множественной регрессии, полученная при применении метода наименьших квадратов,  $\mathbf{a} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ , является наиболее эффективной, т. е. обладает наименьшей дисперсией в классе линейных несмещенных оценок (*Best Linear Unbiased Estimator — BLUE*).

Докажем несмещенность МНК-оценок.

Найдем математическое ожидание оценок параметров множественной линейной регрессии. Используем формулу (2.9), разложив величину  $Y$  на неслучайную и случайную составляющие:

$$Ma = M[(X^T X)^{-1} X^T Y] = M[(X^T X)^{-1} X^T (X\alpha + \varepsilon)].$$

Раскроем скобки внутри выражения под знаком математического ожидания. Математическое ожидание суммы переменных равно сумме математических ожиданий каждой переменной:

$$Ma = M[(X^T X)^{-1} X^T X\alpha] + M[(X^T X)^{-1} X^T \varepsilon].$$

В первом слагаемом произведение матриц  $(X^T X)^{-1} X^T X$  дает единичную матрицу  $I$ , во втором слагаемом выражение  $(X^T X)^{-1} X^T$  можно вынести за скобки как неслучайную величину, а математическое ожидание случайных остатков равно нулю (условие 1). Таким образом, имеем выражение

$$Ma = M(I\alpha) + (X^T X)^{-1} X^T M\varepsilon = \alpha, \quad (2.15)$$

где  $I$  — единичная матрица.

Несмещенность МНК-оценок доказана. Отметим, что из выражения (2.15) следует, что

$$a = \alpha + (X^T X)^{-1} X^T \varepsilon. \quad (2.16)$$

Так как оценки параметров уравнения множественной регрессии могут варьировать, можно оценить их дисперсию и ковариацию, обобщив полученные данные в *ковариационной матрице оценок параметров уравнения регрессии*  $\sum_a$ :

$$\sum_a = \begin{pmatrix} \text{cov}(a_0, a_0) & \text{cov}(a_0, a_1) & \text{cov}(a_0, a_2) & \dots & \text{cov}(a_0, a_p) \\ \text{cov}(a_1, a_0) & \text{cov}(a_1, a_1) & \text{cov}(a_1, a_2) & \dots & \text{cov}(a_1, a_p) \\ \text{cov}(a_2, a_0) & \text{cov}(a_2, a_1) & \text{cov}(a_2, a_2) & \dots & \text{cov}(a_2, a_p) \\ \dots & \dots & \dots & \dots & \dots \\ \text{cov}(a_p, a_0) & \text{cov}(a_p, a_1) & \text{cov}(a_p, a_2) & \dots & \text{cov}(a_p, a_p) \end{pmatrix}. \quad (2.17)$$

Заметим, что в матрице (2.17) нумерация строк и столбцов начинается с нуля. Нулевые строка и столбец введены для учета свободного члена уравнения регрессии и соблюдения нумерации коэффициентов регрессии.

Ковариация двух оценок параметров  $a_i$  и  $a_j$  рассчитывается по формуле

$$\text{cov}(a_i, a_j) = M\{[a_i - M(a_i)][a_j - M(a_j)]\}. \quad (2.18)$$

Из формулы (2.18) следует, что ковариация оценки параметра с самой собой равна ее дисперсии:

$$\text{cov}(a_j, a_j) = M\{[a_j - M(a_j)][a_j - M(a_j)]\} = M[a_j - M(a_j)]^2 = \sigma_{a_j}^2.$$

В матричной форме ковариационную матрицу оценок параметров уравнения регрессии можно записать в виде

$$\Sigma_a = M\{[\mathbf{a} - M(\mathbf{a})][\mathbf{a} - M(\mathbf{a})]^T\} = M[(\mathbf{a} - \boldsymbol{\alpha})(\mathbf{a} - \boldsymbol{\alpha})^T]. \quad (2.19)$$

Преобразуем выражение (2.19) с учетом выражения (2.16):

$$\begin{aligned} M[(\mathbf{a} - \boldsymbol{\alpha})(\mathbf{a} - \boldsymbol{\alpha})^T] &= M[(\boldsymbol{\alpha} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\varepsilon} - \boldsymbol{\alpha})[\boldsymbol{\alpha} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\varepsilon} - \boldsymbol{\alpha}]^T] = \\ &= M[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\varepsilon} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\varepsilon}]^T = M[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}]. \end{aligned}$$

В полученном выражении случайным является только произведение  $\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^T$ , математическое ожидание остальных множителей как детерминированных величин равно им самим. Таким образом, имеем выражение

$$\Sigma_a = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} M(\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^T). \quad (2.20)$$

В выражении (2.20) сомножители, стоящие до математического ожидания, можно представить в виде

$$\mathbf{Z}^{-1} \mathbf{Z} \mathbf{Z}^{-1} = \mathbf{Z}^{-1},$$

где  $\mathbf{Z} = \mathbf{X}^T \mathbf{X}$ .

Математическое ожидание  $M(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^T)$  представляет собой ковариационную матрицу случайных остатков вида

$$M(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^T) = \boldsymbol{\Omega} = \begin{pmatrix} M(\varepsilon_1\varepsilon_1) & M(\varepsilon_1\varepsilon_2) & M(\varepsilon_1\varepsilon_3) & \dots & M(\varepsilon_1\varepsilon_n) \\ M(\varepsilon_2\varepsilon_1) & M(\varepsilon_2\varepsilon_2) & M(\varepsilon_2\varepsilon_3) & \dots & M(\varepsilon_2\varepsilon_n) \\ M(\varepsilon_3\varepsilon_1) & M(\varepsilon_3\varepsilon_2) & M(\varepsilon_3\varepsilon_3) & \dots & M(\varepsilon_3\varepsilon_n) \\ \dots & \dots & \dots & \dots & \dots \\ M(\varepsilon_n\varepsilon_1) & M(\varepsilon_n\varepsilon_2) & M(\varepsilon_n\varepsilon_3) & \dots & M(\varepsilon_n\varepsilon_n) \end{pmatrix},$$

$$\text{или } \boldsymbol{\Omega} = \begin{pmatrix} \sigma_\varepsilon^2 & \text{cov}(\varepsilon_1\varepsilon_2) & \text{cov}(\varepsilon_1\varepsilon_3) & \dots & \text{cov}(\varepsilon_1\varepsilon_n) \\ \text{cov}(\varepsilon_2\varepsilon_1) & \sigma_\varepsilon^2 & \text{cov}(\varepsilon_2\varepsilon_3) & \dots & \text{cov}(\varepsilon_2\varepsilon_n) \\ \text{cov}(\varepsilon_3\varepsilon_1) & \text{cov}(\varepsilon_3\varepsilon_2) & \sigma_\varepsilon^2 & \dots & \text{cov}(\varepsilon_3\varepsilon_n) \\ \dots & \dots & \dots & \dots & \dots \\ \text{cov}(\varepsilon_n\varepsilon_1) & \text{cov}(\varepsilon_n\varepsilon_2) & \text{cov}(\varepsilon_n\varepsilon_3) & \dots & \sigma_\varepsilon^2 \end{pmatrix}. \quad (2.21)$$

В силу условия Гаусса — Маркова о равенстве математического ожидания случайных остатков нулю (условие 1), а также постоянстве дисперсии случайных остатков (условие 2), получаем выражения

$$M(\varepsilon_i\varepsilon_j) = M[(\varepsilon_i - M\varepsilon_i)(\varepsilon_j - M\varepsilon_j)] = \text{cov}(\varepsilon_i\varepsilon_j);$$

$$M(\varepsilon_i\varepsilon_i) = M(\varepsilon_i - M\varepsilon_i)^2 = \sigma_{\varepsilon_i}^2 = \sigma_\varepsilon^2.$$

Согласно условию Гаусса — Маркова о независимости случайных остатков (условие 3) элементы матрицы (2.21), не стоящие на главной диагонали, равны нулю, т.е. матрица  $\boldsymbol{\Omega}$  является скалярной:

$$\boldsymbol{\Omega} = \begin{pmatrix} \sigma_\varepsilon^2 & 0 & 0 & \dots & 0 \\ 0 & \sigma_\varepsilon^2 & 0 & \dots & 0 \\ 0 & 0 & \sigma_\varepsilon^2 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & \sigma_\varepsilon^2 \end{pmatrix} = \sigma_\varepsilon^2 \mathbf{I}_n,$$

где  $\mathbf{I}_n$  — единичная матрица порядка  $n$ .

Вернувшись к рассмотрению ковариационной матрицы оценок параметров уравнения регрессии, получим выражение



$$\Sigma_a = (X^T X)^{-1} \sigma_\varepsilon^2 I_n = \sigma_\varepsilon^2 (X^T X)^{-1}. \quad (2.22)$$

На главной диагонали матрицы  $\Sigma_a$  находятся дисперсии параметров уравнения множественной регрессии. Их величины используются для оценки значимости указанных параметров. Отметим, что в выражении (2.22) дисперсия случайных остатков неизвестна и должна быть оценена по имеющимся у исследователя данным. Можно показать, что несмещенная оценка дисперсии случайных остатков  $\sigma_\varepsilon^2$ , которая обозначается как  $s_\varepsilon^2$ , равна

$$s_\varepsilon^2 = \frac{\sum_{i=1}^n e_i^2}{n - m - 1}, \quad (2.23)$$

где  $n$  — количество наблюдений;  $m$  — количество параметров в уравнении регрессии без учета свободного члена.

Таким образом, ковариационная матрица оценок параметров уравнения множественной регрессии будет иметь вид

$$S_a^2 = s_\varepsilon^2 (X^T X)^{-1}, \quad (2.24)$$

а дисперсия оценки параметра  $j$  ( $a_0$  при  $j = 0$ ,  $a_1, a_2, \dots, a_p$  при  $j = 1, 2, \dots, p$ ), являющаяся диагональным элементом матрицы  $S_a^2$ , может быть оценена по формуле

$$S_{a_j}^2 = s_\varepsilon^2 [(X^T X)^{-1}]_{jj}, \quad (2.25)$$

где  $[(X^T X)^{-1}]_{jj}$  — элемент  $jj$  матрицы  $S_a^2$ .

Можно показать, что оценки параметров уравнения множественной регрессии  $a$  и их дисперсии  $S_a^2$  при выполнении условия о распределении остатков по нормальному закону (условие 5) являются **независимыми**.

## 2.4. Изучение тесноты связи по множественной регрессии

В множественной линейной регрессии показателями тесноты связи являются *коэффициенты множественной корреляции и детерминации*. Коэффициент множественной детерминации равен

$$R^2 = 1 - \frac{\sum (y - \hat{y})^2}{\sum (y - \bar{y})^2} = 1 - \frac{SS_e}{SS_{\text{общ}}},$$

или, с учетом правила сложения дисперсий,

$$R^2 = \frac{\sum (\hat{y} - \bar{y})^2}{\sum (y - \bar{y})^2} = \frac{SS_{\text{факт}}}{SS_{\text{общ}}},$$

где  $SS_e$  — остаточная сумма квадратов,  $SS_e = \sum (y - \hat{y})^2$ ;  $SS_{\text{факт}}$  — факторная сумма квадратов,  $SS_{\text{факт}} = \sum (\hat{y} - \bar{y})^2$ ;  $SS_{\text{общ}}$  — общая сумма квадратов,  $SS_{\text{общ}} = \sum (y - \bar{y})^2$ .

Коэффициент множественной корреляции равен корню из коэффициента множественной детерминации:

$$R = \sqrt{R^2}.$$

Оба показателя — коэффициенты множественной детерминации и корреляции изменяются в границах от нуля до единицы. Коэффициент множественной детерминации показывает долю дисперсии результативной переменной в ее общей дисперсии, объясненную вариацией независимых переменных, входящих в модель регрессии. При формулировке выводов, как правило, полученное значение  $R^2$  умножают на 100%. Например, для модели

$$y = 3226,44 + 12,45x_1 + 0,06x_2 + 0,31x_3 + e,$$

полученной по данным примера 2.1,  $R^2 = 0,73$ , следовательно, вариация объясняющих переменных, входящих в уравнение регрессии, на 73% обусловила вариацию результативной переменной. Коэффициент множественной корреляции

интерпретируют по степени его близости к единице: тесная связь, умеренная, слабая. Для примера 2.1 он составил

$$R = \sqrt{0,73} \approx 0,85,$$

что говорит о тесной связи между зависимой и независимыми переменными.

Показатели тесноты связи можно рассматривать как характеристики качества аппроксимации конкретной моделью (уравнением регрессии) исследуемой зависимости. При прочих равных условиях, чем выше  $R^2$  или  $R$ , тем выше качество аппроксимации. При интерпретации этих показателей следует учитывать, что на их величину оказывают влияние как наличие связи между показателями, так и функция, выбранная для описания связи (функция регрессии), а также соотношение количества наблюдений и количества параметров уравнения регрессии.

Очевидно, что если показатели связаны нелинейно, то чем больше кривизна функции регрессии, тем меньше будет показатель детерминации (и корреляции), рассчитанный для линейной регрессии по тем же данным. Но усложнение функции, увеличение ее параметров при неизменном объеме наблюдений автоматически приводит к увеличению показателей тесноты связи. Это объясняется сокращением степеней свободы у остаточной суммы квадратов. Для учета этой особенности рассчитывают скорректированный коэффициент детерминации

$$R_{\text{скорр}}^2 = 1 - \frac{\sum (y - \hat{y})^2}{\sum (y - \bar{y})^2} \cdot \frac{n-1}{n-m-1}. \quad (2.26)$$

Разделив числитель и знаменатель формулы (2.26) на общую сумму квадратов и проведя преобразования, можно получить выражение, связывающее скорректированный коэффициент детерминации с исходным  $R^2$ :

$$R_{\text{скорр}}^2 = 1 - (1 - R^2) \cdot \frac{n-1}{n-m-1}.$$

Для примера 2.1 скорректированный коэффициент детерминации равен

$$R_{\text{скорр}}^2 = 1 - (1 - 0,73) \cdot \frac{48 - 1}{48 - 3 - 1} \approx 0,71.$$

Скорректированный коэффициент детерминации применяется для оценки реальной тесноты связи между результатом и факторами и сравнения моделей с разным числом параметров. В первом случае обращают внимание на близость скорректированного и нескорректированного коэффициентов детерминации. Если эти показатели близки к единице и различаются незначительно, модель считается хорошей. Этот вывод можно отнести, в частности, к нашему примеру.

При сравнении разных моделей предпочтение, при прочих равных условиях, отдается той, у которой больше скорректированный коэффициент детерминации.

Следует отметить, что область применения скорректированного коэффициента детерминации ограничивается только этими задачами. Его нельзя использовать в формулах, где применяется обычный коэффициент детерминации. Скорректированный коэффициент детерминации нельзя интерпретировать как долю вариации зависимой переменной, объясненную вариацией независимых переменных, включенных в модель регрессии.

## 2.5. Проверка значимости модели множественной регрессии и ее параметров

Для оценки значимости параметров уравнения множественной регрессии используют критерий Стьюдента. Напомним, что *значимость параметров означает их отличие от нуля* с высокой долей вероятности. Нулевой гипотезой в данном случае является утверждение

$$H_0: \alpha_j = 0,$$

альтернативной

$$H_1: \alpha_j \neq 0.$$

Фактическое значение  $t$ -критерия определяется по формуле

$$t_{\phi} = \frac{a_j - \alpha_j}{s_{a_j}} = \frac{a_j}{s_{a_j}}. \quad (2.27)$$

В формуле (2.27) под оценкой параметра  $a_j$  понимается как коэффициент регрессии, так и свободный член (при  $j = 0$ ). Величина среднего квадратического отклонения оцениваемого параметра  $a_j$  определяется как корень из дисперсии  $s_{a_j}^2$ , рас-

считанной по формуле (2.25). Величину  $s_{a_j}$  называют **стандартной ошибкой** параметра  $a_j$ .

Формулу  $s_{a_j}$  для оценки коэффициента регрессии  $a_j$  (т.е. для  $j = \overline{1, p}$ ) можно привести к виду

$$s_{a_j} = \frac{\sigma_y}{\sigma_{x_j}} \sqrt{\frac{1 - R_{yx_1 \dots x_p}^2}{(1 - R_{x_j x_1 \dots x_{j-1} x_{j+1} \dots x_p}^2)(n - m - 1)}}, \quad (2.28)$$

где  $\sigma_y$  — среднее квадратическое отклонение результативной переменной  $y$ ;  $\sigma_{x_j}$  — среднее квадратическое отклонение объясняющей переменной  $x_j$ , являющейся сомножителем коэффициента  $a_j$ ;  $R_{yx_1 \dots x_p}^2$  — коэффициент детерминации, найденный для уравнения зависимости переменной  $y$  от переменных  $x_1 \dots x_p$ , включая  $x_j$ ;  $R_{x_j x_1 \dots x_{j-1} x_{j+1} \dots x_p}^2$  —

коэффициент детерминации, найденный для уравнения зависимости переменной  $x_j$  от других переменных  $x_1 \dots x_p$ , входящих в рассматриваемую модель множественной регрессии.

Теоретическое значение  $t$ -критерия находят по таблице значений критерия Стьюдента для уровня значимости  $\alpha$  и числа степеней свободы  $df = n - m - 1$ . Уровень значимости  $\alpha$  представляет собой вероятность ошибки первого рода, т.е. вероятность отвергнуть гипотезу  $H_0$ , когда она верна. Как правило,  $\alpha$  выбирают равным 0,1; 0,05 или 0,01.

Нулевая гипотеза о незначимости параметра  $a_j$  отвергается, если выполняется неравенство

$$\left| t_{\phi} \right| = \left| \frac{a_j}{s_{a_j}} \right| > t_{T(\alpha, df)} \quad (2.29)$$

где  $t_{T(\alpha, df)}$  — теоретическое значение критерия Стьюдента.

На основе выражения (2.29) можно построить также доверительный интервал для оцениваемого параметра  $\alpha_j$ :

$$a_j - t_{T(\alpha, df)} \cdot s_{a_j} \leq \alpha_j \leq a_j + t_{T(\alpha, df)} \cdot s_{a_j}. \quad (2.30)$$

Выражение (2.30) позволяет как оценить значимость параметра, так и дать его экономическую интерпретацию (если оценивается коэффициент регрессии). Очевидно, что параметр  $\alpha_j$  будет значим, если в доверительный интервал (2.30) не входит нуль, т.е. с большой долей вероятности оцениваемый параметр не равен нулю.

Так как коэффициент регрессии является абсолютным показателем силы связи, границы доверительного интервала  $a_{j, \min}$  и  $a_{j, \max}$  для него также можно интерпретировать аналогичным образом: с вероятностью  $(1 - \alpha)$  при единичном изменении независимой переменной  $x_j$  зависимая переменная  $y$  изменится не меньше, чем на  $a_{j, \min}$ , и не больше, чем на  $a_{j, \max}$ .

Рассмотрим результаты оценки значимости параметров для примера 2.1. Стандартные ошибки параметров равны

$$s_{a_0} = \frac{\sum(y - \hat{y})^2}{n - m - 1} \sqrt{[(X^T X)^{-1}]_{00}} \approx 1651,80;$$

$$s_{a_1} = \frac{\sum(y - \hat{y})^2}{n - m - 1} \sqrt{[(X^T X)^{-1}]_{11}} \approx 4,41;$$

$$s_{a_2} = \frac{\sum(y - \hat{y})^2}{n - m - 1} \sqrt{[(X^T X)^{-1}]_{22}} \approx 0,016;$$

$$s_{a_3} = \frac{\sum(y - \hat{y})^2}{n - m - 1} \sqrt{[(X^T X)^{-1}]_{33}} \approx 0,09.$$

Напомним, что под знаком корня в квадратных скобках стоит элемент матрицы  $(X^T X)^{-1}$ , который находится на пересече-

нии  $j$ -й строки и  $j$ -го столбца, номер  $j$  равен номеру оцениваемого параметра.

Фактическое значение критерия Стьюдента равно

$$t_{a_0} = \frac{a_0}{s_{a_0}} = \frac{3226,44}{1651,8} \approx 1,95;$$

$$t_{a_1} = \frac{a_1}{s_{a_1}} = \frac{12,45}{4,41} \approx 2,82;$$

$$t_{a_2} = \frac{a_2}{s_{a_2}} = \frac{0,06}{0,016} \approx 3,75;$$

$$t_{a_3} = \frac{a_3}{s_{a_3}} = \frac{0,31}{0,09} \approx 3,44.$$

Табличное значение  $t$ -критерия для  $df = 48 - 3 - 1 = 44$  и уровне значимости  $\alpha = 0,05$  составляет 2,0153, следовательно, все параметры, кроме свободного члена, значимы ( $|t_{a_j}| > t_T$ ).

Найдем границы доверительных интервалов для коэффициентов регрессии.

$$\alpha_1: 12,45 \pm 2,0153 \cdot 4,41; \alpha_1 \in [3,56; 21,34];$$

$$\alpha_2: 0,06 \pm 2,0153 \cdot 0,016; \alpha_2 \in [0,028; 0,092];$$

$$\alpha_3: 0,31 \pm 2,0153 \cdot 0,09; \alpha_3 \in [0,13; 0,18].$$

Отметим, что, руководствуясь значениями границ доверительных интервалов, можно сделать те же выводы о значимости коэффициентов регрессии (так как нуль не попадает в доверительный интервал). Выводы в данном случае и не могли быть иными, чем при сравнении фактического и табличного

значений критерия Стьюдента, так как формула (2.30) является следствием формулы (2.29). Дадим экономическую интерпретацию границ доверительных интервалов для коэффициентов регрессии.

Коэффициент  $\alpha_1$  является характеристикой силы связи между объемом поступления налогов и количеством занятых. С учетом значений границ доверительного интервала для  $\alpha_1$  можно сказать, что изменение количества занятых на 1 тыс. человек приведет к изменению (с вероятностью 0,95 ( $= 1 - \alpha$ )) поступления налогов не менее чем на 3,56 млн руб. и не более чем на 21,34 млн руб. при неизменном объеме отгрузки в обрабатывающих производствах и производстве энергии. Для двух других коэффициентов регрессии выводы будут следующими.

Изменение объема отгрузки в обрабатывающих производствах на 1 млн руб. приведет к изменению (с вероятностью 0,95 ( $= 1 - \alpha$ )) поступления налогов не менее чем на 0,028 млн руб. и не более чем на 0,092 млн руб. при неизменных значениях количества занятых и производства энергии.

При изменении производства энергии на 1 млн руб. поступление налогов изменится (с вероятностью 0,95 ( $= 1 - \alpha$ )) не менее чем на 0,13 млн руб. и не более чем на 0,18 млн руб. при неизменных значениях количества занятых и объема отгрузки в обрабатывающих производствах.

Как было отмечено в параграфе 2.2, при построении модели регрессии с использованием *центрированных переменных* коэффициенты регрессии не отличаются от коэффициентов регрессии в натуральной форме. Это утверждение относится также к величине стандартных ошибок коэффициентов регрессии и, следовательно, к фактическим значениям критерия Стьюдента.

При использовании *стандартизованных переменных* меняется масштаб их измерения, что приводит к другим, чем в исходной регрессии, значениям параметров (стандартизованных коэффициентов регрессии) и их стандартных ошибок. Однако фактические значения критерия Стьюдента для параметров уравнения в стандартизованном масштабе совпадают с теми значениями, которые были получены по уравнению в натуральном масштабе.



Для оценки значимости всего уравнения регрессии в целом используется **критерий Фишера (F-критерий)**, который в данном случае называют также **общим F-критерием**. Под незначимостью уравнения регрессии понимается одновременное равенство нулю (с высокой долей вероятности) всех коэффициентов регрессии в генеральной совокупности:

$$H_0: \alpha_1 = \alpha_2 = \dots = \alpha_p = 0.$$

Фактическое значение F-критерия определяется как соотношение факторной и остаточной сумм квадратов, рассчитанных по уравнению регрессии и скорректированных на число степеней свободы:

$$F = \frac{SS_{\text{факт}}}{m} : \frac{SS_e}{n-m-1} = \frac{SS_{\text{факт}}}{SS_e} \cdot \frac{n-m-1}{m}, \quad (2.31)$$

где  $SS_{\text{факт}} = \sum (\hat{y} - \bar{y})^2$  — факторная сумма квадратов;  $SS_e = \sum (y - \hat{y})^2$  — остаточная сумма квадратов.

Теоретическое значение F-критерия находят по таблице значений критерия Фишера для уровня значимости  $\alpha$ , числа степеней свободы  $df_1 = m$  и  $df_2 = n - m - 1$ . Нулевая гипотеза отвергается, если

$$F > F_{T(\alpha, df_1, df_2)},$$

где  $F_{T(\alpha, df_1, df_2)}$  — теоретическое значение критерия Фишера.

Отметим, что если модель незначима, то незначимы и показатели корреляции, рассчитанные по ней. Действительно, если

$$a_1 = a_2 = \dots = a_p = 0,$$

$$\text{то } \hat{y} = a_0$$

и линия регрессии параллельна оси абсцисс. Кроме того, из системы нормальных уравнений, полученной по методу наименьших квадратов (2.8), следует, что  $\bar{y} = a_0 + a_1 \bar{x}_1 + a_2 \bar{x}_2 + \dots + a_p \bar{x}_p$ .

При нулевых значениях всех коэффициентов регрессии имеем выражение

$$\bar{y} = a_0.$$

Тогда

$$R^2 = \frac{\sum(\hat{y} - \bar{y})^2}{\sum(y - \bar{y})^2} = \frac{\sum(a_0 - a_0)^2}{\sum(y - \bar{y})^2} = 0,$$

т.е. при равенстве всех коэффициентов регрессии нулю (их статистической незначимости) коэффициент детерминации также будет равен нулю (статистически незначим).

Формулу (2.31) расчета  $F$ -критерия можно преобразовать, разделив факторную и остаточную суммы квадратов на общую сумму квадратов:

$$SS_{\text{общ}} = \sum(y - \bar{y})^2.$$

После простых преобразований получаем выражение

$$F = \frac{R^2}{1 - R^2} \cdot \frac{n - m - 1}{m}.$$

Расчет общего  $F$ -критерия можно оформить в виде таблицы дисперсионного анализа (табл. 2.2).

**Таблица 2.2.** Анализ статистической значимости модели множественной регрессии

Источники вариации	Число степеней свободы $df$	Сумма квадратов $SS$	Дисперсия на одну степень свободы $MS = SS / df$	$F$ -критерий Фишера	
				фактическое значение	табличное значение для $\alpha = 0,05$
Общая	$n - 1$	$SS_{\text{общ}} = \sum(y - \bar{y})^2$			
Регрессия	$m$	$SS_{\text{факт}} = \sum(\hat{y} - \bar{y})^2$	$MS_{\text{факт}}$	$F_{\text{факт}} = MS_{\text{факт}} / MS_e$	$F_{\text{табл}} (df_1 = m; df_2 = n - m - 1)$
Остаточная	$n - m - 1$	$SS_e = \sum(y - \hat{y})^2$	$MS_e$		

Аналогичную таблицу дисперсионного анализа можно увидеть в результатах компьютерной обработки данных. Ее отличие

от приведенной выше таблицы заключается в содержании последнего столбца. В нашем случае это теоретическое значение критерия Фишера. В компьютерных вариантах в последнем столбце приводится значение вероятности допустить ошибку первого рода (отвергнуть верную нулевую гипотезу), которая соответствует фактическому значению  $F$ -критерия. В *Excel* эта величина называется «значимость  $F$ ». Обозначим величину, выдаваемую компьютером в таблице дисперсионного анализа, как  $\alpha_{\text{факт}}$ . Ее значение можно проинтерпретировать следующим образом: если теоретическое значение  $F$ -критерия равно его фактическому значению, то вероятность ошибки первого рода (уровень значимости) равна  $\alpha_{\text{факт}}$ .

Выбирая для определения табличного значения критерия некий уровень значимости  $\alpha$ , мы соглашаемся на величину ошибки, равную  $\alpha$ . Следовательно, если  $\alpha_{\text{факт}} < \alpha$ , то фактическая ошибка будет меньше запланированной и можно говорить о значимости уравнения регрессии при заданном уровне значимости  $\alpha$ .

Проверим на статистическую значимость уравнение регрессии, полученное в примере 2.1. Фактическое значение  $F$ -критерия равно

$$F = \frac{0,73}{1-0,73} \cdot \frac{48-3-1}{3} \approx 40.$$

Табличное значение критерия Фишера для  $\alpha = 0,05$ , числа степеней свободы  $df_1 = 3$  и  $df_2 = 44$  равно 2,82. Так как фактическое значение  $F$ -критерия больше табличного, уравнение регрессии значимо с вероятностью  $1 - \alpha = 1 - 0,05 = 0,95$ . Следовательно, значим также коэффициент детерминации, т.е. он с большой долей вероятности отличен от нуля.

При использовании опции «Регрессия» в ППП *Excel* для данного примера получена следующая таблица дисперсионного анализа (табл. 2.3).

**Таблица 2.3.** Таблица дисперсионного анализа, полученная при применении опции «Регрессия» в ППП *Excel*

Дисперсионный анализ					
	$df$	$SS$	$MS$	$F$	Значимость $F$
Регрессия	3	3 652 714 368	1 217 571 456	40,31035571	1,10224E - 12
Остаток	44	1 329 016 902	30 204 929,59		
Итого	47	4 981 731 270			

Фактическое значение  $F$ -критерия содержится в предпоследнем столбце данной таблицы. Отметим, что его значение отличается от приведенного выше из-за ошибок округления. В последнем столбце табл. 2.3 приведена вероятность допустить ошибку первого рода. Она равна  $1,10224E-12$ , т.е.  $0,00000000000110224$ . Нами задана максимальная величина этой вероятности, равная  $0,05$ . Так как фактическое значение вероятности допустить ошибку первого рода меньше (значительно меньше) установленного нами максимального, нулевая гипотеза о незначимости уравнения регрессии должна быть отвергнута.

## 2.6. Множественная линейная регрессия с ограничениями на параметры

В эконометрических моделях на значения параметров уравнения множественной регрессии могут накладываться ограничения, связанные с величиной этих параметров, взаимосвязи их друг с другом. Примеры таких ограничений:

$$\alpha_i = 0;$$

$$\alpha_i + \alpha_j = 1;$$

$$\alpha_i = \alpha_j.$$

При наличии ограничений на параметры уравнения регрессии возникают два вопроса. Первый: если ограничения справедливы, то каким образом обеспечить их выполнение в формируемой модели? Второй: насколько значимо модель без ограничений на параметры отличается от модели с ограничениями?

В простейшем случае при наличии ограничений на параметры уравнение регрессии может быть преобразовано таким образом, чтобы учесть имеющиеся ограничения в самой структуре модели. Например, в линейной модели регрессии с тремя независимыми переменными

$$y = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_3 + \varepsilon \quad (2.32)$$

имеется ограничение вида

$$\alpha_1 + \alpha_3 = 1. \quad (2.33)$$

Тогда можно провести следующие преобразования:

$$\alpha_1 = 1 - \alpha_3; \quad (2.34)$$

$$y = \alpha_0 + (1 - \alpha_3)x_1 + \alpha_2 x_2 + \alpha_3 x_3 + \varepsilon;$$

$$(y - x_1) = \alpha_0 + \alpha_2 x_2 + \alpha_3 (x_3 - x_1) + \varepsilon.$$

Рассчитав новые переменные

$$y_2 = y - x_1;$$

$$x_4 = x_3 - x_1,$$

получим новое уравнение регрессии

$$y_2 = \alpha_0 + \alpha_2 x_2 + \alpha_3 x_4 + \varepsilon. \quad (2.35)$$

После применения к уравнению (2.35) метода наименьших квадратов будут получены оценки неизвестных параметров  $\alpha_2$  и  $\alpha_3$ , а по формуле (2.34) можно оценить параметр  $\alpha_1$ .

В общем случае необходимо сформулировать общую линейную гипотезу, содержащую ограничения на параметры. Если представить одно линейное ограничение на параметры в виде выражения

$$c_{10}\alpha_0 + c_{11}\alpha_1 + c_{12}\alpha_2 + \dots + c_{1p}\alpha_p = q_1, \quad (2.36)$$

то при  $k$  линейных ограничениях выражение (2.36) можно записать в матричной форме:

$$C\alpha = Q, \quad (2.37)$$

где  $\mathbf{C}$  — матрица размерностью  $k \times (p + 1)$ ;  $\mathbf{Q}$  — вектор размерностью  $k \times 1$ .

Очевидно, что если на параметр  $\alpha_j$  не наложено ограничений, то коэффициент  $c_{1j}$  в выражении (2.36) и, в общем виде, коэффициенты  $c_{ij}$  в выражении (2.37) будут равны нулю.

Например, для выражения (2.32) с ограничением (2.33) матрицы  $\mathbf{C}$  и  $\mathbf{Q}$  будут иметь вид

$$\mathbf{C} = (0101); \quad \mathbf{Q} = (1).$$

Тогда

$$\mathbf{C}\boldsymbol{\alpha} = (0101) \cdot \begin{pmatrix} \alpha_0 \\ \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{pmatrix} = (\alpha_1 + \alpha_3) = (1) = \mathbf{Q}.$$

Если для выражения (2.32) добавить еще одно ограничение:  $\alpha_2 = 0$ , то матрицы  $\mathbf{C}$  и  $\mathbf{Q}$  будут иметь вид

$$\mathbf{C} = \begin{pmatrix} 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{pmatrix}, \quad \mathbf{Q} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}.$$

Для решения задачи минимизации целевой функции квадратов случайных остатков (метод наименьших квадратов) с учетом линейных ограничений можно воспользоваться **методом множителей Лагранжа**.

Метод множителей Лагранжа заключается в нахождении минимума функции вида

$$\Phi = (\mathbf{Y} - \mathbf{X}\boldsymbol{\alpha}_0)^T (\mathbf{Y} - \mathbf{X}\boldsymbol{\alpha}_0) - \mathbf{M}^T (\mathbf{C}\boldsymbol{\alpha}_0 - \mathbf{Q}),$$

где  $\boldsymbol{\alpha}_0$  — вектор-столбец параметров регрессии с ограничениями;  $\mathbf{M}$  — вектор множителей Лагранжа размерностью  $k \times 1$ .

Функция  $\Phi$  дифференцируется по  $\boldsymbol{\alpha}_0$  и  $\mathbf{M}$ , полученные выражения приравниваются нулю и находится их решение.

Применение модели с линейными ограничениями целесообразно в том случае, если сокращение остаточной дисперсии

в модели без ограничений по сравнению с остаточной дисперсией в модели с ограничениями статистически незначимо.

Ограничение, связанное с равенством коэффициента  $\alpha_j$  регрессии нулю, можно проверить с помощью частного критерия Фишера (частного  $F$ -критерия).

Нулевой гипотезой в данном случае является незначимость величины сокращения остаточной дисперсии при включении в уравнение независимой переменной  $x_j$  (т.е. переходе от модели с ограничениями к модели без ограничений). Фактическое значение частного  $F$ -критерия рассчитывается по формуле

$$F_{x_j} = \frac{SS_e^{(1)} - SS_e^{(2)}}{1} : \frac{SS_e^{(2)}}{(n-m-1)}, \quad (2.38)$$

где  $SS_e^{(1)}$  — остаточная сумма квадратов для модели без независимой переменной  $x_j$ ;  $SS_e^{(2)}$  — остаточная сумма квадратов для модели, включающей независимую переменную  $x_j$ ;  $m$  — количество параметров (без учета свободного члена) в модели с независимой переменной  $x_j$ .

Разделив каждую из дробей на общую сумму квадратов  $SS_{\text{общ}}$  и произведя преобразования, получим выражение

$$F_{x_j} = \frac{R_{yx_1 \dots x_j \dots x_p}^2 - R_{yx_1 \dots x_{j-1} x_{j+1} \dots x_p}^2}{1} : \frac{1 - R_{yx_1 \dots x_j \dots x_p}^2}{(n-m-1)},$$

или  $F_{x_j} = \frac{R_{yx_1 \dots x_j \dots x_p}^2 - R_{yx_1 \dots x_{j-1} x_{j+1} \dots x_p}^2}{1 - R_{yx_1 \dots x_j \dots x_p}^2} \cdot \frac{n-m-1}{1}, \quad (2.39)$

где  $R_{yx_1 \dots x_j \dots x_p}^2$  — коэффициент детерминации для уравнения зависимости переменной  $y$  от переменных  $x_1 \dots x_p$ , включая переменную  $x_j$ ;  $R_{yx_1 \dots x_{j-1} x_{j+1} \dots x_p}^2$  — коэффициент детерминации для уравнения зависимости переменной  $y$  от всех переменных  $x_1 \dots x_p$ , кроме переменной  $x_j$ .

Полученное значение  $F$ -критерия сравнивают с табличным, найденным для  $df_1 = 1$  и  $df_2 = (n - m - 1)$  степеней свободы. Если фактическое значение  $F$ -критерия больше табличного,

нулевая гипотеза отклоняется, т.е. включение рассматриваемой переменной  $x_j$  существенно уменьшает остаточную дисперсию и является статистически значимым и ограничение  $\alpha_j = 0$  не подтверждается.

Можно показать, что для линейной множественной регрессии верно соотношение

$$F_{x_j} = t_{a_j}^2, \quad (2.40)$$

т.е. значимость параметра  $a_j$  при переменной  $x_j$  означает также значимость сокращения остаточной дисперсии при включении данной переменной в модель регрессии.

Если ограничения на равенство нулю касается  $k$  коэффициентов регрессии, то для оценки этого предположения проводят проверку, аналогичную проверке по частному  $F$ -критерию. Отличие состоит в том, что под величиной  $SS_e^{(1)}$  понимается остаточная сумма квадратов по уравнению регрессии без переменных, на коэффициенты при которых накладывается ограничение в виде их равенства нулю. Так как число добавляемых переменных составит  $k$ , то число степеней свободы для разности дисперсий в формуле (2.38) будет также равно  $k$  (а не единице):

$$F_{x_j \dots x_{j+k-1}} = \frac{SS_e^{(1)} - SS_e^{(2)}}{k} : \frac{SS_e^{(2)}}{(n-m-1)}, \quad (2.41)$$

где  $SS_e^{(1)}$  — остаточная сумма квадратов для модели без независимых переменных  $x_j \dots x_{j+k-1}$ ;  $SS_e^{(2)}$  — остаточная сумма квадратов для модели без ограничений (с независимыми переменными  $x_j \dots x_{j+k-1}$ );  $m$  — количество параметров (без учета свободного члена) в модели с независимыми переменными  $x_j \dots x_{j+k-1}$ .

Аналогично формула (2.39) примет вид

$$F_{x_j \dots x_{j+k-1}} = \frac{R_{yx_1 \dots x_j \dots x_p}^2 - R_{yx_1 \dots x_{j-1} x_{j+k} \dots x_p}^2}{1 - R_{yx_1 \dots x_j \dots x_p}^2} \cdot \frac{n-m-1}{k}, \quad (2.42)$$



где  $R^2_{yx_1 \dots x_j \dots x_p}$  — коэффициент детерминации для уравнения зависимости переменной  $y$  от переменных  $x_1 \dots x_p$ , включая  $x_j, \dots, x_{j+k-1}$ ;

$R^2_{yx_1 \dots x_{j-1} x_{j+k} \dots x_p}$  — коэффициент детерминации для уравнения зависимости переменной  $y$  от всех переменных  $x_1 \dots x_p$ , кроме  $x_j, \dots, x_{j+k-1}$ .

В общем случае линейных ограничений вида (2.37) фактическое значение  $F$ -критерия рассчитывается по формулам (2.43) или (2.44):

$$F_0 = \frac{SS_e^{(1)} - SS_e^{(2)}}{k} : \frac{SS_e^{(2)}}{(n-m-1)}, \quad (2.43)$$

где  $SS_e^{(1)}$  — остаточная сумма квадратов для модели с ограничениями;  
 $SS_e^{(2)}$  — остаточная сумма квадратов для модели без ограничений;  
 $k$  — количество ограничений, наложенных на уравнение регрессии;

$$F_0 = \frac{R^2 - R_0^2}{1 - R^2} \cdot \frac{n-m-1}{k}, \quad (2.44)$$

где  $R^2$  — коэффициент детерминации для уравнения регрессии без ограничений;  $R_0^2$  — коэффициент детерминации для уравнения регрессии с ограничениями.

Для случая  $k$  ограничений нулевая гипотеза о незначимости различий остаточных дисперсий в модели с ограничениями и модели без ограничений принимается, если фактическое значение  $F$ -критерия ((2.41)—(2.44)) меньше или равно табличному, взятому для уровня значимости  $\alpha$ , числа степеней свободы  $df_1 = k$ ,  $df_2 = n - m - 1$ , где  $m$  — количество параметров (без учета свободного члена) в модели без ограничений.

Рассмотрим наложение ограничений на модель регрессии из примера 2.1. Простейший случай ограничения предполагает равенство нулю одного из коэффициентов регрессии. Пусть, например, мы предполагаем, что число занятых не оказывает влияния на собираемость налогов. Тогда ограничение будет иметь вид

$$\alpha_1 = 0,$$

а уравнение регрессии с ограничениями запишется в виде

$$y = \alpha_0 + \alpha_2 x_2 + \alpha_3 x_3 + \varepsilon. \quad (2.45)$$

После оценки параметров уравнения (2.45) с помощью МНК получаем выражение

$$y = 5847,83 + 0,080x_2 + 0,423x_3 + e, \quad R^2 = 0,68.$$

Частный  $F$ -критерий для переменной  $x_1$  равен

$$F_{x_1} = \frac{0,73 - 0,68}{1 - 0,73} \cdot \frac{48 - 3 - 1}{1} \approx 8,15.$$

Табличное значение критерия Фишера, найденное для  $df_1 = 1$  и  $df_2 = 44$  степеней свободы, равно 4,06, следовательно, гипотеза о незначимости включения переменной  $x_1$  в уравнение регрессии отклоняется. Иначе говоря, отклоняется гипотеза о наложении ограничения на коэффициент  $\alpha_1$  в виде его равенства нулю.

Проверим выполнение равенства (2.40). Действительно,

$$t_{\alpha_1}^2 = 2,82^2 = 7,95 \approx 8,15 = F_{x_1}.$$

Равенство получилось неточным из-за ошибок округления. При использовании ППП *Excel* имеем равенство

$$t_{\alpha_1}^2 = 2,820130662^2 = 7,953136951 \approx 7,953136948 = F_{x_1}.$$

Как видим, ошибка округления присутствует и в данном случае, но она крайне незначительна и проявляется в восьмом знаке после запятой.

Рассмотрим другой вариант ограничений на модель из примера 2.1. Пусть наложены ограничения в виде равенства нулю двух коэффициентов:

$$\alpha_2 = \alpha_3 = 0,$$

т.е. мы считаем, что объемы отгрузки в обрабатывающих производствах  $x_2$  и производства энергии  $x_3$  не оказывают влияния на сбор налогов. Модель с ограничениями в данном случае представляет парную линейную регрессию вида

$$y = \alpha_0 + \alpha_1 x_1 + \varepsilon.$$

После применения МНК имеем выражение

$$y = 3855,80 + 29,81x_1 + e, \quad R^2 = 0,54.$$

Используя формулу  $F$ -критерия для модели с ограничениями (2.42), получим выражение

$$F_{x_2, x_3} = \frac{0,73 - 0,54}{1 - 0,73} \cdot \frac{48 - 3 - 1}{2} \approx 15,48.$$

Табличное значение критерия Фишера, найденное для  $df_1 = 2$  и  $df_2 = 44$  степеней свободы, равно 3,21, следовательно, гипотеза о необходимости наложения ограничений на коэффициенты  $\alpha_2$  и  $\alpha_3$  в виде их равенства нулю отклоняется.

## 2.7. Нелинейные модели множественной регрессии

В качестве функций множественной регрессии кроме линейной часто используют нелинейные: полиномы разных степеней, обратную, степенную, показательную, полулогарифмическую функции:

$$y = \alpha_0 + \alpha_{11}x_1 + \alpha_{12}x_1^2 + \dots + \alpha_{1k}x_1^k + \dots + \alpha_{p1}x_p + \alpha_{p2}x_{p2}^2 + \dots + \alpha_{pk}x_p^k + \varepsilon$$
 —

полином  $k$ -й степени;

$$y = \frac{1}{\alpha_0 + \alpha_1 x_1 + \dots + \alpha_p x_p + \varepsilon}$$
 — обратная функция;

$$y = \alpha_0 x_1^{\alpha_1} \cdot \dots \cdot x_p^{\alpha_p} \varepsilon$$
 — степенная функция;

$y = \alpha_0 \alpha_1^{x_1} \cdot \dots \cdot \alpha_p^{x_p} \varepsilon$  — показательная функция;

$y = \alpha_0 + \alpha_1 \ln x_1 + \dots + \alpha_p \ln x_p + \varepsilon$  — полулогарифмическая функция,

где  $\alpha_0, \alpha_1, \dots, \alpha_p$  — параметры функций.

Нелинейные функции могут представлять собой также «смешанные» модели. Например, можно построить уравнение множественной регрессии с тремя независимыми переменными

$$y = \alpha_0 \alpha_1^{x_1} x_2^{\alpha_2} \alpha_3^{x_3} \varepsilon.$$

В данном уравнении использованы две функции: показательная (для учета влияния переменных  $x_1$  и  $x_3$ ) и степенная (для учета влияния переменной  $x_2$ ).

При выборе формы регрессии необходимо учитывать ряд обстоятельств.

Во-первых, необходимо принимать во внимание теоретические предпосылки построения модели, выводы, сформированные в экономической теории о характере взаимосвязи показателей, ограничениях, налагаемых на параметры функций. Например, форму степенной функции имеет **производственная функция Кобба — Дугласа**

$$P = \alpha_0 \cdot L^{\alpha_1} \cdot K^{\alpha_2} \cdot \varepsilon,$$

где  $P$  — объем продукции;  $L$  — затраты труда;  $K$  — величина капитала;  $\alpha_0, \alpha_1, \alpha_2$  — неизвестные параметры;  $\varepsilon$  — случайный остаток.

На коэффициенты производственной функции может быть наложено линейное ограничение вида

$$\alpha_1 + \alpha_2 = 1.$$

Тогда

$$P = \alpha_0 \cdot L^{\alpha_1} \cdot K^{1-\alpha_1} \cdot \varepsilon,$$

или

$$\frac{P}{K} = \alpha_0 \cdot \left(\frac{L}{K}\right)^{\alpha_1} \cdot \varepsilon.$$

Во-вторых, при выборе функции регрессии следует учитывать простоту оценки ее параметров, возможность их интерпретации. Все перечисленные выше нелинейные функции являются внутренне линейными, т.е. их можно преобразовать в линейную форму. Оценка параметров внутренне линейных функций производится путем применения МНК к линейризованной форме нелинейной функции.

Для функций, линейных по параметрам (многочлены второй и более высоких степеней), линейризация заключается в том, что независимые переменные, взятые во второй и более высоких степенях, рассматриваются как самостоятельные переменные. Например, функция

$$y = \alpha_0 + \alpha_1 x_1 + \alpha_{21} x_2 + \alpha_{22} x_2^2 + \varepsilon$$

представляет собой множественную регрессию, в которой переменная  $x_1$  представлена в линейной форме, а переменная  $x_2$  — в нелинейной в виде многочлена второго порядка (параболы). Эта функция линейна по параметрам. Переменная  $x_2^2$  рассматривается как третья независимая переменная. Формально можно обозначить переменную  $x_2^2$  как новую переменную  $z$ :

$$y = \alpha_0 + \alpha_1 x_1 + \alpha_{21} x_2 + \alpha_{22} z + \varepsilon,$$

т.е. линейризация функции заключается в замене переменных.

Линейризация функций, в которых переменные связаны мультипликативно, заключается в логарифмировании правой и левой части функции по любому основанию, наиболее часто — по натуральному. Например, функция

$$y = \alpha_0 x_1^{\alpha_1} \alpha_2^{x_2} \varepsilon$$

после линейризации путем логарифмирования принимает вид

$$\ln y = \ln \alpha_0 + \alpha_1 \ln x_1 + \alpha_2 \ln x_2 + \ln \varepsilon.$$

Введя новые обозначения для логарифмированных величин

$$\ln y \equiv Y; \quad \ln \alpha_0 \equiv A_0; \quad \ln x_1 \equiv X_1; \quad \ln \alpha_2 \equiv A_2; \quad \ln \varepsilon \equiv \varepsilon_1,$$

получим выражение

$$Y = A_0 + \alpha_1 X_1 + A_2 X_2 + \varepsilon_1.$$

Для обратной функции вида

$$y = \frac{1}{\alpha_0 + \alpha_1 x_1 + \dots + \alpha_p x_p + \varepsilon}$$

линеаризация заключается в следующих преобразованиях:

$$\frac{1}{y} = \alpha_0 + \alpha_1 x_1 + \dots + \alpha_p x_p + \varepsilon;$$

$$\frac{1}{y} \equiv Y;$$

$$Y = \alpha_0 + \alpha_1 x_1 + \dots + \alpha_p x_p + \varepsilon.$$

В некоторых нелинейных функциях параметры имеют экономическую интерпретацию. Например, в степенной функции  $y = \alpha_0 x_1^{\alpha_1} \dots x_p^{\alpha_p} \varepsilon$  коэффициенты  $\alpha_1 \dots \alpha_p$  являются коэффициентами эластичности зависимой переменной по соответствующим независимым переменным  $x_1 \dots x_p$ .

В показательной функции  $y = \alpha_0 \alpha_1^{x_1} \dots \alpha_p^{x_p} \varepsilon$  коэффициенты  $\alpha_1 \dots \alpha_p$  показывают, во сколько раз в среднем изменится зависимая переменная при изменении соответствующей независимой переменной  $x_1 \dots x_p$  на единицу и при неизменном значении других независимых переменных, включенных в уравнение регрессии.

Эконометрический анализ нелинейных функций: оценка параметров и их значимости, значимости функции в целом, прогнозирование, расчет показателей тесноты связи ведется по линеаризованным формам.

## 2.8. Выбор наилучшей функции регрессии

Если функциональная форма связи между показателями неизвестна, необходимо построить несколько функций регрессии и выбрать из них наилучшую. Перечень функций, которые будут рассмотрены, определяется знаниями и предположениями, экспертными оценками относительно характера связи между показателями.

При прочих равных условиях критериями выбора наилучшей функции являются следующие.

1. Статистическая значимость всех параметров при независимых переменных.
2. Значимость всей функции в целом.
3. Выполнение требований Гаусса — Маркова, предъявляемых к случайным остаткам модели, в первую очередь, постоянство дисперсии и независимость друг от друга.
4. Минимальная доля остаточной дисперсии в общей дисперсии, т.е. максимальная величина коэффициента детерминации  $R^2$ . Если модели регрессии содержат разное количество параметров, вместо  $R^2$  следует сравнивать скорректированные коэффициенты детерминации  $R^2_{\text{скорр}}$ .

Сопоставление функций по коэффициентам детерминации можно проводить только в том случае, если зависимая переменная представлена в моделях в одной и той же форме. Например, в линейной, полулогарифмической, полиномиальной функция зависимая переменная у представлена в исходной форме:

$$y = \alpha_0 + \alpha_1 x_1 + \dots + \alpha_p x_p + \varepsilon \text{ — линейная функция;}$$

$y = \alpha_0 + \alpha_1 \ln x_1 + \dots + \alpha_p \ln x_p + \varepsilon$  — полулогарифмическая функция;

$y = \alpha_0 + \alpha_{11}x_1 + \alpha_{12}x_1^2 + \dots + \alpha_{p1}x_p + \alpha_{p2}x_p^2 + \varepsilon$  — полином второй степени.

В нелинейных функциях, линеаризация которых связана с преобразованием зависимой переменной, коэффициент детерминации рассчитывается для значений преобразованной переменной. Например, при логарифмировании зависимой переменной

$$R^2 = 1 - \frac{\sum (\ln y - \ln \hat{y})^2}{\sum (\ln y - \ln \bar{y})^2},$$

при преобразовании зависимой переменной в обратной функции

$$R^2 = 1 - \frac{\sum (1/y - 1/\hat{y})^2}{\sum (1/y - 1/\bar{y})^2}.$$

Очевидно, что показатели детерминации, рассчитанные по исходным и преобразованным значениям зависимой переменной, сравнивать нельзя. Не является корректным решением и расчет коэффициента детерминации по исходным (нелинейным) функциям. Это связано с тем, что при мультипликативной связи переменных (степенная, показательная функции) параметры могут быть найдены с помощью МНК, только если случайные остатки умножаются на выровненные значения зависимой переменной, а не прибавляются к ней:

$$y = \alpha_0 x_1^{\alpha_1} \dots x_p^{\alpha_p} \varepsilon;$$

$$y = \alpha_0 \alpha_1^{x_1} \dots \alpha_p^{x_p} \varepsilon.$$

Если рассчитать выровненное значение зависимой переменной по исходной функции



$$\hat{y} = \alpha_0 x_1^{\alpha_1} \cdot \dots \cdot x_p^{\alpha_p};$$

$$\hat{y} = a_0 a_1^{x_1} \cdot \dots \cdot a_p^{x_p},$$

то случайные остатки  $e$  не могут быть вычислены как разность фактического и выровненного значений зависимой переменной, что предполагается при расчете коэффициента детерминации:

$$e \neq y - a_0 x_1^{\alpha_1} \cdot \dots \cdot x_p^{\alpha_p};$$

$$e \neq y - a_0 a_1^{x_1} \cdot \dots \cdot a_p^{x_p}.$$

Аналогичный вывод можно сделать для регрессии, построенной на основе обратной функции. Случайные остатки для этой функции могут быть вычислены как

$$e = \frac{1}{y} - a_0 - \sum_{i=1}^p a_i x_i = \frac{1}{y} - \frac{1}{\hat{y}} \neq y - \hat{y}.$$

Иногда коэффициент детерминации по исходным данным можно вычислить в качестве меры приближения выровненных значений зависимой переменной к исходным. В этом случае используют формулу

$$R^2 = 1 - \frac{\sum (y - \hat{y})^2}{\sum (y - \bar{y})^2}.$$

Полученный показатель называют *квази- $R^2$* . Для функций, линеаризуемых путем логарифмирования правой и левой части, значения коэффициента детерминации, рассчитанные по линеаризованной регрессии и по исходной, часто имеют близкие значения.

На основе данных примера 2.1 найдем параметры степенной функции

$$y = \alpha_0 x_1^{\alpha_1} x_2^{\alpha_2} x_3^{\alpha_3} \varepsilon.$$

В линеаризованной форме она будет иметь вид

$$\ln y = \ln \alpha_0 + \alpha_1 \ln x_1 + \alpha_2 \ln x_2 + \alpha_3 \ln x_3 + \ln \varepsilon.$$

После применения МНК к данным примера 2.1, получим функцию

$$\ln y = 2,58 + 0,23 \ln x_1 + 0,07 \ln x_2 + 0,52 \ln x_3 + \ln e.$$

Коэффициент детерминации для этой функции равен  $R^2 = 0,84$ . Случайными остатками в данном случае являются величины  $\ln e$ . Именно к ним предъявляются требования Гаусса — Маркова.

Перейдем к исходной функции:

$$\alpha_0 = e^{2,58} = \exp(2,58) = 13,2;$$

$$y = 13,2 x_1^{0,23} x_2^{0,07} x_3^{0,52} e.$$

Если рассчитать случайные остатки

$$y - \hat{y},$$

то коэффициент детерминации будет равен  $R^2 = 0,68$ .

Различия между двумя показателями детерминации получились заметными (более 0,1), однако показатель детерминации по исходной функции регрессии находится в допустимых границах [0; 1] и может быть соответствующим образом проинтерпретирован.

Рассмотрим по тем же данным регрессию на основе обратной функции

$$y = \frac{1}{\alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_3 + \varepsilon}.$$

После линеаризации и применения МНК получим регрессию в линеаризованном виде:

$$\frac{1}{y} = 0,00025 - 0,00000018x_1 - 0,0000000027x_2 - 0,0000000029x_3 + e;$$

$$R^2 = 0,42.$$

В исходной форме эта регрессия будет иметь вид

$$y = \frac{1}{2,5 \cdot 10^{-4} - 1,8 \cdot 10^{-6} x_1 - 2,7 \cdot 10^{-10} x_2 - 2,9 \cdot 10^{-9} x_3 + e}.$$

Принимая за выровненное значение зависимой переменной значение дроби без учета случайного остатка, рассчитаем  $(y - \hat{y})$  и найдем коэффициент детерминации по исходной модели. Он оказался равным  $R^2 = -40,04$ . Очевидно, что это значение не имеет никакого смысла. Таким образом, мы убедились, что расчет показателя детерминации по исходным значениям преобразуемой в ходе линеаризации зависимой переменной не только некорректен, но и может привести к бессмысленному результату. Следовательно, сопоставление коэффициентов детерминации по различным нелинейным функциям не приведет к желаемому результату — выбору наилучшей функции.

Процедура выбора наилучшего преобразования зависимой переменной была предложена Дж. Боксом и Д. Коксом<sup>1</sup>. В их честь эту процедуру называют **тестом Бокса — Кокса**. В ее основе лежит утверждение о том, что как исходные (наблюдаемые), так и преобразованные значения какой-либо переменной являются частными случаями реализации функции (2.46) при разных значениях  $\lambda$ :

$$y^* = \frac{y^\lambda - 1}{\lambda}. \quad (2.46)$$

В частности,

<sup>1</sup> Box G. E. P., and Cox D. R. An analysis of transformations. Journal of the Royal Statistical Society. 1964. Series B, 26. pp. 211—43.

$$\text{если } \lambda = 1, \text{ то } y^* = \frac{y^1 - 1}{1} = y - 1;$$

$$\text{если } \lambda = -1, \text{ то } y^* = \frac{y^{-1} - 1}{-1} = 1 - \frac{1}{y};$$

$$\text{если } \lambda \rightarrow 0, \text{ то } \lim_{\lambda \rightarrow 0} y^* = \lim_{\lambda \rightarrow 0} \frac{y^\lambda - 1}{\lambda} = \lim_{\lambda \rightarrow 0} \frac{y^\lambda \ln y}{1} = \ln y.$$

В работе Н. Дрейпера, Г. Смита<sup>1</sup> указывается, что при преобразовании по формуле (2.46) величина  $y^*$  может сильно меняться, что приводит к проблемам в анализе и требует специальной программы для нахождения наилучшего значения  $\lambda$ . Поэтому предпочтительнее рассчитывать  $y^*$  по формуле

$$y^* = \frac{y^\lambda - 1}{\lambda \dot{y}^{\lambda-1}}, \quad (2.47)$$

где  $\dot{y}$  — среднее геометрическое из значений зависимой переменной;

$$\dot{y} = \sqrt[\lambda]{y_1 y_2 \cdots y_n}.$$

При применении формулы (2.47) преобразованное значение  $y^*$  будет равно

$$\text{для } \lambda = 1: y^* = \frac{y^1 - 1}{1 \dot{y}^{1-1}} = y - 1;$$

$$\text{для } \lambda = -1: y^* = \frac{y^{-1} - 1}{-1 \dot{y}^{-1-1}} = \dot{y}^2 \left(1 - \frac{1}{y}\right);$$

$$\text{для } \lambda \rightarrow 0: \lim_{\lambda \rightarrow 0} y^* = \lim_{\lambda \rightarrow 0} \frac{y^\lambda - 1}{\lambda \dot{y}^{0-1}} = \dot{y} \lim_{\lambda \rightarrow 0} \frac{y^\lambda \ln y}{1} = \dot{y} \ln y.$$

Тест Бокса — Кокса заключается в поиске такого значения параметра  $\lambda$ , при котором остаточная сумма квадратов  $SS_e$  для модели с преобразованной переменной  $y^*$  будет минимальной. Для этого необходимо выполнить следующие действия.

1. Выбрать конкретные значения  $\lambda$ . Например,  $\lambda = -1; -0,5; 0; 0,5; 1$ . Границы интервала рассматриваемых значений

<sup>1</sup> Дрейпер Н., Смит Г. Прикладной регрессионный анализ : пер. с англ. 3-е изд. М. : ИД «Вильямс», 2007. С. 367—369.

$\lambda$  можно при необходимости расширить, а шаг перебора уменьшить. Например, можно использовать интервал для  $\lambda$  от  $-2$  до  $2$ , а шаг уменьшить до  $0,25$ . Для  $\lambda = 0$  значение  $y^*$  вычисляется как  $\ln y$  (при использовании формулы (2.46)) или  $\dot{y} \ln y$  (при использовании формулы (2.47)).

2. Для каждого  $\lambda$  вычислить:

— параметры регрессии

$$y^* = f(x_1, x_2, \dots, x_p) + \varepsilon,$$

где  $f$  — линейная по параметрам функция.

Отметим, что функция  $f$  должна быть не только линейна по параметрам, но и количество этих параметров должно быть одинаковым для всех сравниваемых функций. В противном случае дальнейшие действия по сопоставлению остаточных сумм квадратов (см. ниже) будут некорректными из-за разного числа степеней свободы этих показателей;

— остаточную сумму квадратов  $SS_e^\lambda$ :

$$SS_e^\lambda = \sum (y^* - \hat{y}^*)^2.$$

3. Выбрать  $\lambda$  с наименьшим  $SS_e^\lambda$ . Соответствующее преобразование  $y^*$  и будет наилучшим. Отметим, что модель, выбранная в качестве наилучшей, должна отвечать требованиям, перечисленным выше (параметры и уравнение регрессии должны быть значимыми, должны соблюдаться требования, предъявляемые к случайным остаткам).

Так как преобразования зависимой переменной предполагают, в частности, ее логарифмирование, тест Бокса — Кокса может быть применен только если все значения этой переменной положительны. Кроме того, как отмечают в своей работе Н. Дрейпер и Г. Смит<sup>1</sup>, тест Бокса — Кокса целесообразно

<sup>1</sup> Дрейпер Н., Смит Г. Указ. соч. С. 369.

применять, если соотношение максимального и минимального значений зависимой переменной больше 10.

Предположим, например, что для показателей примера 2.1 рассматриваются следующие модели регрессии:

$$\text{линейная } y = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_3 + \varepsilon;$$

полиномиальная второй степени

$$y = \alpha + \alpha_{11} x_1 + \alpha_{21} x_2 + \alpha_{31} x_3 + \alpha_{12} x_1^2 + \alpha_{22} x_2^2 + \alpha_{32} x_3^2 + \varepsilon;$$

$$\text{полулогарифмическая } y = \alpha_0 + \alpha_1 \ln x_1 + \alpha_2 \ln x_2 + \alpha_3 \ln x_3 + \varepsilon;$$

$$\text{обратная } y = \frac{1}{\alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_3 + \varepsilon};$$

$$\text{степенная } y = \alpha_0 x_1^{\alpha_1} x_2^{\alpha_2} x_3^{\alpha_3} \varepsilon;$$

$$\text{показательная } y = \alpha_0 \alpha_1^{x_1} \alpha_2^{x_2} \alpha_3^{x_3} \varepsilon.$$

Мы не имеем теоретических оснований предпочесть ту или иную функцию и должны выбрать наилучшую опытным путем. Оценим параметры каждой из функций с помощью МНК, предварительно произведя линеаризацию нелинейных функций. Результаты расчетов параметров, оценки значимости, величины показателей детерминации приведены в табл. 2.4.

**Таблица 2.4.** Результаты оценки параметров линейной и нелинейных функций, представленных в линеаризованной форме

Параметры функции и оценка их значимости по <i>t</i> -критерию	Показатель детерминации $R^2, R^2_{\text{скорр}}$	Общий <i>F</i> -критерий
$y = 3226,44 + 12,45x_1 + 0,06x_2 + 0,31x_3 + \varepsilon$ $(t_e) \quad (1,95) \quad (2,82) \quad (3,76) \quad (3,43)$	$R^2 = 0,73$ $R^2_{\text{скорр}} = 0,71$	$F = 40,3$
$y = -973,24 + 14,88x_1 + 0,11x_2 + 0,83x_3 -$ $(t_e) \quad (-0,38) \quad (1,07) \quad (1,99) \quad (2,64)$ $-6,2 \cdot 10^{-3} x_1^2 - 2,7 \cdot 10^{-7} x_2^2 - 1,2 \cdot 10^{-6} x_3^2 + \varepsilon$ $(-0,63) \quad (-1,27) \quad (-1,98)$	$R^2 = 0,78$ $R^2_{\text{скорр}} = 0,75$	$F = 24,7$

Окончание табл. 2.4

Параметры функции и оценка их значимости по $t$ -критерию	Показатель детерминации $R^2, R^2_{\text{скорр}}$	Общий $F$ -критерий
$y = -60419,90 + 3537,85 \ln x_1 + 836,54 \ln x_2 +$ $(t_a) \quad (-6,20) \quad (1,63) \quad (0,85)$ $+ 5209,06 \ln x_3 + e$ $(3,32)$	$R^2 = 0,67$ $R^2_{\text{скорр}} = 0,65$	$F = 29,6$
$\frac{1}{y} = 2,5 \cdot 10^{-4} - 1,8 \cdot 10^{-6} x_1 - 2,7 \cdot 10^{-10} x_2 - 2,9 \cdot 10^{-9} x_3 + e$ $(t_a) \quad (8,41) \quad (-2,30) \quad (-0,91) \quad (-1,78)$	$R^2 = 0,42$ $R^2_{\text{скорр}} = 0,39$	$F = 10,9$
$\ln y = 2,58 + 0,23 \ln x_1 + 0,07 \ln x_2 + 0,52 \ln x_3 + \ln e$ $(t_a) \quad (4,91) \quad (2,01) \quad (1,27) \quad (6,19)$	$R^2 = 0,84$ $R^2_{\text{скорр}} = 0,83$	$F = 78,7$
$\ln y = 8,39 + 0,001x_1 + 3,4 \cdot 10^{-6} x_2 + 2,3 \cdot 10^{-5} x_3 + \ln e$ $(t_a) \quad (60,17) \quad (3,31) \quad (2,53) \quad (3,03)$	$R^2 = 0,69$ $R^2_{\text{скорр}} = 0,67$	$F = 32,7$

Из рассмотренных функций незначимые параметры при независимых переменных имеют:

- полиномиальная (второй степени);
- обратная;
- полулогарифмическая;
- степенная.

Эти функции следует исключить из процедуры выбора наилучшей функции. Таким образом, необходимо выбрать функцию из двух оставшихся — линейной и показательной. Так как зависимые переменные в линеаризованной форме этих функций различны, а соотношение максимального и минимального значений зависимой переменной больше 10 (равно 26,11), применим тест Бокса — Кокса. Найдем уравнения регрессии и остаточные суммы квадратов для преобразованной переменной  $y^*$ , вычисленной по формуле (2.47). Нас интересуют только две функции, поэтому  $y^*$  необходимо вычислить только для двух значений  $\lambda$ :

- $\lambda = 1$  (для непробранного значения зависимой переменной);
- $\lambda \rightarrow 0$  (для  $\ln y$ ).

Имеем следующие результаты:

- при  $\lambda = 1$

- $$y^* = 3225,44 + 12,45x_1 + 0,06x_2 + 0,31x_3 + e; \quad S\hat{S}_e^{\lambda=1} = 1329016902;$$

$$(t_a) \quad (1,95) \quad (2,82) \quad (3,76) \quad (3,43)$$

— при  $\lambda \rightarrow 0$

$$y^* = 113141,80 + 16,64x_1 + 0,046x_2 + 0,31x_3 + e; \quad SS_e^{\lambda=0} = 1722319309.$$

$(t_a)$     $(60,17)$     $(3,31)$     $(2,53)$     $(3,03)$

Остаточная сумма квадратов при  $\lambda = 1$  меньше, чем остаточная сумма квадратов при  $\lambda \rightarrow 0$ . Следовательно, линейная функция лучше описывает связь между переменными, чем показательная.

Еще раз подчеркнем, что тест Бокса — Кокса предполагает перебор функций при большем количестве возможных значений  $\lambda$ . Мы ограничились только двумя, так как наша задача заключалась в выборе наилучшей функции из перечня предложенных выше.

Модификацией теста Бокса — Кокса является **преобразование Зарембки**, которое применяется для сопоставления только двух форм зависимой переменной — непреобразованной и прологарифмированной. Для применения этого преобразования необходимо выполнить следующие шаги.

1. Рассчитать значение  $y^*$  по формуле

$$y^* = \frac{y}{\hat{y}}. \quad (2.48)$$

2. Оценить параметры двух регрессий:

$$y^* = f(x_1, x_2, \dots, x_p) + \varepsilon \text{ и} \quad (2.49)$$

$$\ln y^* = f(x_1, x_2, \dots, x_p) + \varepsilon, \quad (2.50)$$

где  $f$  — та же функция, которая была в исходных сравниваемых моделях.

3. Рассчитать и сравнить остаточные суммы квадратов  $SS_e$  по регрессиям (2.49):

$$SS_e = \sum (y^* - \hat{y}^*)^2$$

и (2.50):



$$SS_e = \sum (\ln y^* - \ln \hat{y}^*)^2.$$

Регрессия, имеющая меньшее значение  $SS_e$ , будет лучшей.

4. Для оценки значимости различий остаточных сумм квадратов можно рассчитать величину статистики

$$\chi^2 = \left| \frac{n}{2} \ln \frac{SS_e^{\min}}{SS_e^{\max}} \right|,$$

где  $SS_e^{\min}$  — меньшее значение остаточной суммы квадратов;  $SS_e^{\max}$  — большее значение остаточной суммы квадратов.

Полученная величина имеет распределение  $\chi^2$  с числом степеней свободы  $df = 1$ . Если она превышает критическое значение  $\chi_\alpha^2$  ( $\alpha$  — уровень значимости), то различия между остаточными суммами квадратов признаются существенными.

Рассмотрим применение преобразования Зарембки на нашем примере. После преобразования зависимой переменной по формуле (2.48) применим МНК к функциям

$$y^* = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_3 + \varepsilon;$$

$$\ln y^* = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_3 + \ln \varepsilon.$$

Получим следующие результаты:

$$y^* = 0,24 + 9,2 \cdot 10^{-4} x_1 + 4,5 \cdot 10^{-6} x_2 + 2,3 \cdot 10^{-5} x_3 + \varepsilon;$$

$$\ln y^* = -1,12 + 1,2 \cdot 10^{-3} x_1 + 3,42 \cdot 10^{-6} x_2 + 2,31 \cdot 10^{-5} x_3 + \ln \varepsilon.$$

Остаточные суммы квадратов составили

— для регрессии  $y^*$ :  $SS_e = 7,312$ ;

— для регрессии  $\ln y^*$ :  $SS_e = 9,472$ .

Очевидно, что первая модель лучше, т.е. следует отдать предпочтение линейной модели регрессии. Проверим значимость отличий остаточных сумм квадратов с помощью теста  $\chi^2$ .

$$\chi^2 = \left| \frac{48}{2} \ln \left( \frac{7,312}{9,475} \right) \right| = |-6,22| = 6,22.$$

Табличное значение  $\chi^2$  для  $df = 1$  и  $\alpha = 0,05$  составляет 3,84, следовательно, различия между остаточными суммами квадратов для двух сравниваемых моделей статистически значимы. Вывод о том, что линейная модель регрессии лучше аппроксимирует связь между показателями, подтверждается.

## 2.9. Метод максимального правдоподобия

*Метод максимального правдоподобия (ММП)* является одним из наиболее широко используемых методов в статистике и эконометрике. Для его применения необходимо знание закона распределения исследуемой случайной величины.

Пусть имеется некоторая случайная величина  $Y$  с заданным законом распределения  $f(Y)$ . Параметры этого закона неизвестны и их нужно найти. В общем случае величину  $Y$  рассматривают как многомерную, т.е. состоящую из нескольких одномерных величин  $Y^1, Y^2, Y^3 \dots, Y^m$ .

Предположим, что  $Y$  — одномерная случайная величина и ее отдельные значения являются числами. Каждое из них ( $y_1, y_2, y_3, \dots, y_n$ ) рассматривается как реализация не одной случайной величины  $Y$ , а  $n$  случайных величин  $Y_1, Y_2, Y_3 \dots, Y_n$ . То есть:

$y_1$  — реализация случайной величины  $Y_1$ ;

$y_2$  — реализация случайной величины  $Y_2$ ;

$y_3$  — реализация случайной величины  $Y_3$ ;

...

$y_n$  — реализация случайной величины  $Y_n$ .

Параметры закона распределения вектора  $Y$ , состоящего из случайных величин  $Y_1, Y_2, Y_3, \dots, Y_n$ , представляют как вектор  $\Theta$ , состоящий из  $k$  параметров:  $\theta_1, \theta_2, \dots, \theta_k$ . Величины  $Y_1, Y_2, Y_3, \dots, Y_n$  могут быть распределены как с одинаковыми параметрами, так и с различными; некоторые параметры могут совпадать, а другие различаться. Конкретный ответ на этот вопрос зависит от той задачи, которую решает исследователь.

Например, если стоит задача определения параметров закона распределения случайной величины  $Y$ , реализацией которой являются величины  $Y_1, Y_2, Y_3, \dots, Y_n$ , то предполагают, что каждая из этих величин распределена так же, как величина  $Y$ . Иначе говоря, любая величина  $Y_i$  описывается одним и тем же законом распределения  $f(Y_i)$ , причем с одними и теми же параметрами  $\Theta: \theta_1, \theta_2, \dots, \theta_k$ .

Другой пример — нахождение параметров уравнения регрессии. В этом случае каждая величина  $Y_i$  рассматривается как случайная величина, имеющая «собственные» параметры распределения, которые могут частично совпадать с параметрами распределения других случайных величин, а могут и полностью различаться. Более подробно применение ММП для нахождения параметров уравнения регрессии будет рассмотрено ниже.

В рамках метода максимального правдоподобия совокупность имеющихся значений  $u_1, u_2, u_3, \dots, u_n$  рассматривается как некоторая фиксированная, неизменная. То есть закон  $f(Y_i)$  есть функция от заданной величины  $u_i$  и неизвестных параметров  $\Theta$ . Следовательно, для  $n$  наблюдений случайной величины  $Y$  имеется  $n$  законов  $f(Y_i)$ .

Неизвестные параметры этих законов распределения рассматриваются как случайные величины. Они могут меняться, однако при данном наборе значений  $u_1, u_2, u_3, \dots, u_n$  наиболее вероятны конкретные значения параметров. Иначе говоря, вопрос ставится таким образом: каковы должны быть параметры  $\Theta$ , чтобы значения  $u_1, u_2, u_3, \dots, u_n$  были наиболее вероятны?

Для ответа на него нужно найти закон совместного распределения случайных величин  $Y_1, Y_2, Y_3, \dots, Y_n$  —  $f(Y_1, Y_2, Y_3, \dots, Y_n)$ . Если предположить, что наблюдаемые нами величины  $u_1, u_2, u_3, \dots, u_n$  независимы, то он равен произведению  $n$  законов  $f$

( $Y_i$ ) (произведению вероятностей появления данных значений для дискретных случайных величин или произведению плотностей распределения для непрерывных случайных величин):

$$f(Y_1, Y_2, \dots, Y_n) = f(Y_1)f(Y_2) \cdot \dots \cdot f(Y_n),$$

или

$$f(Y_1, Y_2, \dots, Y_n) = \prod_{i=1}^n f(Y_i).$$

Чтобы подчеркнуть тот факт, что в качестве переменных рассматриваются искомые параметры  $\Theta$ , введем в обозначение закона распределения еще один аргумент — вектор параметров  $\Theta$ :

$$f(Y) \equiv f(Y, \Theta);$$

$$f(Y_i) \equiv f(Y_i, \Theta);$$

$$f(Y_1, Y_2, \dots, Y_n) \equiv f(Y_1, Y_2, \dots, Y_n, \Theta).$$

С учетом введенных обозначений закон совместного распределения **независимых** величин  $Y_1, Y_2, Y_3, \dots, Y_n$  с параметрами  $\Theta$  будет записан в виде

$$f(Y_1, Y_2, \dots, Y_n, \Theta) = \prod_{i=1}^n f(Y_i, \Theta). \quad (2.51)$$

Полученную функцию (2.51) называют **функцией максимального правдоподобия** и обозначают  $L(Y, \Theta)$ :

$$L(Y, \Theta) = \prod_{i=1}^n f(Y_i, \Theta).$$

Еще раз подчеркнем тот факт, что в функции максимального правдоподобия значения  $Y$  считаются фиксированными, а переменными являются параметры вектора  $\Theta$  (в частном случае — один параметр). Часто для упрощения процесса нахождения неизвестных параметров  $\Theta$  функцию правдоподобия логарифмируют, получая **логарифмическую функцию правдоподобия**

$$l(Y, \Theta) = \ln[L(Y, \Theta)] = \ln \prod_{i=1}^n f(Y_i, \Theta).$$

Дальнейшее решение по ММП предполагает нахождение таких значений  $\Theta$ , при которых функция правдоподобия (или ее логарифм) достигает максимума. Найденные значения  $\Theta_L$  называют *оценкой максимального правдоподобия*.

Методы нахождения оценки максимального правдоподобия достаточно разнообразны. В простейшем случае функция правдоподобия является непрерывно дифференцируемой и имеет максимум в точке, для которой

$$\frac{dL(Y, \Theta)}{d\Theta} = 0,$$

или

$$\frac{dl(Y, \Theta)}{d\Theta} = 0.$$

В более сложных случаях максимум функции максимального правдоподобия не может быть найден путем дифференцирования и решения уравнения правдоподобия, что требует поиска других алгоритмов его нахождения, в том числе итеративных.

Оценки параметров, полученные с использованием ММП, являются:

- *состоятельными*, т.е. с увеличением объема наблюдений разница между оценкой и фактическим значением параметра приближается к нулю;
- *инвариантными*: если получена оценка параметра  $\theta$ , равная  $\theta_L$ , и имеется непрерывная функция  $q(\theta)$ , то оценкой значения этой функции будет величина  $q(\theta_L)$ . В частности, если с помощью ММП мы оценили величину дисперсии какого-либо показателя ( $\sigma_L^2$ ), то корень из полученной оценки будет оценкой среднего квадратического отклонения ( $\sigma_L$ ), полученной по ММП.
- *асимптотически эффективными*;
- *асимптотически нормально распределенными*.

Последние два утверждения означают, что оценки параметров, полученные по ММП, проявляют свойства эффективности и нормальности при бесконечно большом увеличении объема выборки.

Для нахождения параметров множественной линейной регрессии вида

$$y = \alpha_0 + \alpha_1 x_1 + \dots + \alpha_p x_p + \varepsilon,$$

необходимо знать законы распределения зависимых переменных  $Y_i$  или случайных остатков  $\varepsilon_i$ . Пусть переменная  $Y_i$  распределена по нормальному закону с параметрами  $\mu_i$ ,  $\sigma_i$ . Каждое наблюдаемое значение  $y_i$  имеет, в соответствии с определением регрессии, математическое ожидание  $\mu_i = MY_i$ , равное его теоретическому значению при условии, что известны значения параметров регрессии в генеральной совокупности

$$\mu_i = \hat{Y}_i = \alpha_0 + \alpha_1 x_{i1} + \dots + \alpha_p x_{ip},$$

где  $x_{i1}, \dots, x_{ip}$  — значения независимых переменных в  $i$ -м наблюдении. При выполнении предпосылок применения МНК (предпосылок построения классической нормальной линейной модели), случайные величины  $Y_i$  имеют одинаковую дисперсию

$$\sigma_{Y_i}^2 = \sigma_{Y_j}^2 = \sigma_Y^2, \quad \forall i, j.$$

Дисперсия величины  $Y_i$  определяется по формуле

$$\sigma_{Y_i}^2 = M(Y_i - \mu_i)^2.$$

Преобразуем эту формулу:

$$\sigma_{Y_i}^2 = M(Y_i - \hat{Y}_i)^2 = M\varepsilon_i^2. \quad (2.52)$$

При выполнении условий Гаусса — Маркова о равенстве нулю математического ожидания случайных остатков и постоянстве их дисперсий можно перейти от формулы (2.52) к формуле

$$\sigma_{Y_i}^2 = M(\varepsilon_i - M\varepsilon_i)^2 = \sigma_{\varepsilon_i}^2 = \sigma_{\varepsilon}^2.$$

Иначе говоря, дисперсии случайной величины  $Y_i$  и соответствующих ей случайных остатков совпадают.

Выборочную оценку математического ожидания случайной величины  $Y_i$  будем обозначать  $\hat{y}_i$ :

$$\hat{y}_i = a_0 + a_1 x_{i1} + \dots + a_p x_{ip},$$

а оценку ее дисперсии (постоянной для разных наблюдений) как  $s_y^2$ .

Если предположить независимость отдельных наблюдений  $y_i$ , то получим функцию максимального правдоподобия

$$L(Y, \hat{y}, s_y) = \prod_{i=1}^n \frac{e^{-\frac{(y_i - \hat{y}_i)^2}{2s_y^2}}}{s_y \sqrt{2\pi}} = \prod_{i=1}^n \frac{e^{-\frac{(y_i - a_0 - a_1 x_{i1} - \dots - a_p x_{pi})^2}{2s_y^2}}}{s_y \sqrt{2\pi}}. \quad (2.53)$$

В приведенной функции делитель  $\sqrt{2\pi}$  является константой и не оказывает влияния на нахождение ее максимума. Поэтому для упрощения расчетов он может быть опущен. С учетом этого замечания и после логарифмирования функция (2.53) примет вид

$$l(Y, \hat{y}, s_y) = \sum_{i=1}^n \left( \frac{-(y_i - a_0 - a_1 x_{i1} - \dots - a_p x_{pi})^2}{2s_y^2} - \ln s_y \right).$$

В соответствии с ММП найдем производные логарифмической функции правдоподобия по неизвестным параметрам  $a_0, a_1, \dots, a_p, s_y$ .

$$\left\{ \begin{array}{l} \frac{dl(Y, \hat{y}, s_y)}{da_0} = \frac{\sum_{i=1}^n (y_i - a_0 - a_1 x_{1i} - \dots - a_p x_{pi})}{s_y^2}, \\ \frac{dl(Y, \hat{y}, s_y)}{da_1} = \frac{\sum_{i=1}^n x_{1i} (y_i - a_0 - a_1 x_{1i} - \dots - a_p x_{pi})}{s_y^2}, \\ \dots \\ \frac{dl(Y, \hat{y}, s_y)}{da_p} = \frac{\sum_{i=1}^n x_{pi} (y_i - a_0 - a_1 x_{1i} - \dots - a_p x_{pi})}{s_y^2}, \\ \frac{dl(Y, \hat{y}, s_y)}{ds_y} = \sum_{i=1}^n \left( \frac{(y_i - a_0 - a_1 x_{1i} - \dots - a_p x_{pi})^2}{s_y^3} - \frac{1}{s_y} \right). \end{array} \right.$$

Для нахождения экстремума приравняем полученные выражения к нулю. После преобразований получим систему

$$\left\{ \begin{array}{l} \sum_{i=1}^n y_i = n a_0 + a_1 \sum_{i=1}^n x_{1i} + \dots + a_p \sum_{i=1}^n x_{pi}, \\ \sum_{i=1}^n x_{1i} y_i = a_0 \sum_{i=1}^n x_{1i} + a_1 \sum_{i=1}^n x_{1i}^2 + \dots + a_p \sum_{i=1}^n x_{1i} x_{pi}, \\ \dots \\ \sum_{i=1}^n x_{pi} y_i = a_0 \sum_{i=1}^n x_{pi} + a_1 \sum_{i=1}^n x_{1i} x_{pi} + \dots + a_p \sum_{i=1}^n x_{pi}^2, \\ s_y^2 = \frac{\sum_{i=1}^n (y_i - a_0 - a_1 x_{1i} - \dots - a_p x_{pi})^2}{n}. \end{array} \right. \quad (2.54)$$

Эта система соответствует системе, полученной по методу наименьших квадратов. То есть ММП и МНК дают одинаковые результаты, если соблюдаются предпосылки МНК. Последнее выражение в системе (2.54) дает оценку дисперсии случайной переменной  $Y_i$  или, что одно и то же, дисперсии случайных остатков. Как было отмечено выше (см. формулу (2.23)), несмещенная оценка дисперсии случайных остатков равна

$$s_e^2 = \frac{\sum_{i=1}^n e_i^2}{n - m - 1}.$$



Аналогичная оценка, полученная с применением ММП (как следует из системы (2.54)), вычисляется по формуле

$$s_{\varepsilon}^2 = \frac{\sum_{i=1}^n e_i^2}{n},$$

т.е. является **смещенной**.

Мы рассмотрели случай применения ММП для нахождения параметров линейной множественной регрессии при условии, что величина  $Y_i$  нормально распределена. Другой подход к нахождению параметров той же регрессии заключается в построении функции максимального правдоподобия для случайных остатков  $\varepsilon_i$ . Для них также предполагается нормальное распределение с параметрами  $(0, \sigma_{\varepsilon})$ . Нетрудно убедиться, что результаты решения в этом случае совпадут с результатами, полученными выше.

## 2.10. Прогнозирование по модели множественной регрессии

Прогнозирование по модели множественной линейной регрессии предполагает оценку ожидаемых значений зависимой переменной при заданных значениях независимых переменных, входящих в уравнение регрессии. Различают точечный и интервальный прогнозы.

**Точечный прогноз** — это расчетное значение зависимой переменной, полученное подстановкой в уравнение множественной линейной регрессии прогнозных (заданных исследователем) значений независимых переменных. Если заданы значения  $x_1^{np}, x_2^{np}, \dots, x_p^{np}$ , то прогнозное значение зависимой переменной (точечный прогноз) будет равно

$$\hat{y}_{np} = a_0 + a_1 x_1^{np} + a_2 x_2^{np} + \dots + a_p x_p^{np}. \quad (2.55)$$

**Интервальный прогноз** — это минимальное и максимальное значения зависимой переменной, в промежутке между

которыми она попадает с заданной долей вероятности и при заданных значениях независимых переменных.

Интервальный прогноз для линейной функции вычисляется по формуле

$$\hat{y}_{\text{пр}} \pm t_T \cdot s_{y_{\text{пр}}}, \quad (2.56)$$

где  $t_T$  — теоретическое значение критерия Стьюдента при  $df = n - m - 1$  степенях свободы;  $s_{y_{\text{пр}}}$  — стандартная ошибка прогноза, вычисляемая по формуле

$$s_{y_{\text{пр}}} = \sqrt{\frac{\sum e^2}{n - m - 1} (1 + \mathbf{X}_{\text{пр}}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_{\text{пр}})}, \quad (2.57)$$

где  $\mathbf{X}$  — матрица исходных значений независимых переменных;  $\mathbf{X}_{\text{пр}}$  — матрица-столбец прогнозных значений независимых переменных вида

$$\mathbf{X}_{\text{пр}} = \begin{pmatrix} 1 \\ x_1^{\text{пр}} \\ x_2^{\text{пр}} \\ \vdots \\ x_p^{\text{пр}} \end{pmatrix}.$$

Найдем прогнозные значения поступления налогов (пример 2.1), при условии, что связь между показателями описывается уравнением

$$y = 3226,44 + 12,45x_1 + 0,06x_2 + 0,31x_3 + e.$$

Зададим прогнозные значения независимых переменных:

- количество занятых  $x_1$ : 500 тыс. человек;
- объем отгрузки в обрабатывающих производствах  $x_2$ : 65 000 млн руб.;
- производство энергии  $x_3$ : 15 000 млн руб.

Найдем точечный и интервальный прогноз поступления налогов.

При заданных значениях независимых переменных поступление налогов в среднем составит

$$\hat{y} = 3226,44 + 12,45 \cdot 500 + 0,06 \cdot 65\,000 + 0,31 \cdot 15\,000 = 18\,013 \text{ млн руб.}$$

Вектор прогнозных значений независимых переменных будет иметь вид

$$X_{\text{пр}} = \begin{pmatrix} 1 \\ 500 \\ 65\,000 \\ 15\,000 \end{pmatrix}.$$

Ошибка прогноза, рассчитанная по формуле (2.57), составила 5556,7. Табличное значение  $t$ -критерия при числе степеней свободы  $df = 44$  и уровне значимости  $\alpha = 0,05$  равно 2,0154. Следовательно, прогнозные значения поступления налогов будут с вероятностью 0,95 находиться в границах:

$$\text{от } 18\,013,69 - 2,0154 \cdot 5556,7 = 6814,1 \text{ млн руб.};$$

$$\text{до } 18\,013,69 + 2,0154 \cdot 5556,7 = 29\,212 \text{ млн руб.}$$

Прогнозирование по нелинейным моделям множественной регрессии также можно осуществлять по формулам (2.55)—(2.57), предварительно линеаризовав указанные модели.

## 2.11. Мультиколлинеарность данных

При построении эконометрической модели предполагается, что независимые переменные воздействуют на зависимую изолированно, т. е. влияние отдельной переменной на результирующий признак не связано с влиянием других переменных. В реальной экономической действительности все явления в той или иной мере связаны, поэтому добиться выполнения этого предположения практически невозможно. Наличие связи между независимыми переменными приводит к необходимости оценки ее влияния на результаты корреляционно-регрессионного анализа.

Различают функциональные и стохастические связи между объясняющими переменными. В первом случае говорят об ошибках спецификации модели, которые должны быть исправлены.

Функциональная связь возникает, если в уравнение регрессии в качестве объясняющих переменных включают, в частности, все переменные, входящие в тождество. Например, можно сказать, что доход  $Y$  складывается из потребления  $C$  и инвестиций  $I$ , т. е. имеет место тождество. Мы предполагаем, что уровень процентных ставок  $r$  зависит от дохода, т.е. модель в общем виде может быть представлена в виде

$$r = \alpha_0 + \alpha_1 Y + \varepsilon;$$

$$Y = C + I.$$

Неопытный исследователь, желая улучшить модель, может включить в уравнение также переменные «потребление» и «инвестиции», что приведет к функциональной связи между объясняющими переменными:

$$r = \alpha_0 + \alpha_1 Y + \alpha_2 C + \alpha_3 I + \varepsilon;$$

$$Y = C + I.$$

Функциональная взаимосвязь столбцов матрицы  $X$  приведет к невозможности найти единственное решение уравнения регрессии, так как  $\mathbf{a} = (X^T X)^{-1} X^T Y$ , а нахождение обратной матрицы  $(X^T X)^{-1}$  предполагает деление алгебраических дополнений матрицы  $(X^T X)$  на ее определитель, который в данном случае будет равен нулю.

Более часто между объясняющими переменными наблюдается стохастическая связь, что приводит к уменьшению величины определителя матрицы  $(X^T X)$ : чем сильнее связь, тем меньше будет определитель. Это приводит к росту не только оценок параметров, полученных с использованием МНК, но и их стандартных ошибок, которые вычисляются по формуле (2.24):

$$S_{a_j} = s_\varepsilon \sqrt{\left[ (X^T X)^{-1} \right]_{jj}},$$

в которой, как мы видим, также используется матрица  $(X^T X)^{-1}$ .

Корреляционная связь может существовать как между двумя объясняющими переменными (*интеркорреляция*), так и между несколькими (*мультиколлинеарность*).

Существует несколько признаков, указывающих на наличие мультиколлинеарности. В частности, такими признаками являются:

- не соответствующие экономической теории знаки коэффициентов регрессии. Например, нам известно, что объясняющая переменная  $x$  оказывает прямое воздействие на объясняемую переменную  $y$ , в то же время коэффициент регрессии при этой переменной меньше нуля;
- значительные изменения параметров модели при небольшом сокращении (увеличении) объема исследуемой совокупности;
- незначимость параметров регрессии, обусловленная высокими значениями стандартных ошибок параметров.

Существование корреляционной связи между независимыми переменными может быть выявлено с помощью показателей корреляции между ними, в частности с помощью парных коэффициентов корреляции  $r_{x_i x_j}$ , которые можно записать в виде матрицы

$$r_{xx} = \begin{pmatrix} r_{x_1 x_1} & r_{x_1 x_2} & \dots & r_{x_1 x_p} \\ r_{x_2 x_1} & r_{x_2 x_2} & \dots & r_{x_2 x_p} \\ \dots & \dots & \dots & \dots \\ r_{x_p x_1} & r_{x_p x_2} & \dots & r_{x_p x_p} \end{pmatrix}. \quad (2.58)$$

Коэффициент корреляции переменной с самой собой равен единице ( $r_{x_i x_i} = 1$ ), а коэффициент корреляции переменной  $x_i$  с переменной  $x_j$  равен коэффициенту корреляции переменной  $x_j$  с переменной  $x_i$  ( $r_{x_i x_j} = r_{x_j x_i}$ ). Следовательно, данная матрица является симметрической, поэтому в ней указывают только главную диагональ и элементы под ней:

$$r_{xx} = \begin{pmatrix} 1 & & & & \\ r_{x_2x_1} & 1 & & & \\ \dots & \dots & \dots & \dots & \\ r_{x_px_1} & r_{x_px_2} & \dots & \dots & 1 \end{pmatrix}.$$

Высокие значения парных линейных коэффициентов корреляции указывают на наличие интеркорреляции, т.е. линейной связи между двумя объясняющими переменными. Чем выше величина  $|r_{x_i x_j}|$ ,  $i \neq j$ , тем выше интеркорреляция. Так как при построении моделей избежать отсутствия связей между объясняющими переменными практически невозможно, существует следующая рекомендация относительно включения двух переменных в модель в качестве объясняющих. Обе переменные можно включить в модель, если выполняются соотношения

$$\begin{cases} r_{yx_i} > r_{x_i x_j}, \\ r_{yx_j} > r_{x_i x_j}, \end{cases} \quad i \neq j, \quad (2.59)$$

т.е. теснота связи результирующей и объясняющей переменных больше, чем теснота связи между объясняющими переменными.

Наличие мультиколлинеарности можно подтвердить, найдя определитель матрицы (2.58). Если связь между независимыми переменными полностью отсутствует, то недиагональные элементы будут равны нулю, а определитель матрицы — единице. Если связь между независимыми переменными близка к функциональной (т.е. является очень тесной), то определитель матрицы  $r_{xx}$  будет близок к нулю.

Еще один метод измерения мультиколлинеарности является следствием анализа формулы стандартной ошибки коэффициента регрессии (2.28):

$$s_{a_j} = \frac{\sigma_y}{\sigma_{x_j}} \sqrt{\frac{1 - R_{yx_1 \dots x_p}^2}{(1 - R_{x_j x_1 \dots x_{j-1} x_{j+1} \dots x_p}^2)(n - m - 1)}}.$$

Как следует из данной формулы, стандартная ошибка будет тем больше, чем меньше будет величина, которую называют

**фактор инфляции дисперсии** (или **фактор вздутия дисперсии**)  $VIF$ :

$$VIF_{x_j} = \frac{1}{(1 - R_{x_j, x_1 \dots x_{j-1} x_{j+1} \dots x_p}^2)},$$

где  $R_{x_j, x_1 \dots x_{j-1} x_{j+1} \dots x_p}^2$  — коэффициент детерминации, найденный для уравнения зависимости переменной  $x_j$  от других переменных  $x_1, \dots, x_p$ , входящих в рассматриваемую модель множественной регрессии.

Так как величина  $R_{x_j, x_1 \dots x_{j-1} x_{j+1} \dots x_p}^2$  отражает тесноту связи между переменной  $x_j$  и прочими объясняющими переменными, то она, по сути, характеризует мультиколлинеарность применительно к данной переменной  $x_j$ . При отсутствии связи показатель  $VIF_{x_j}$  будет равен (или близок) единице, усиление связи ведет к стремлению этого показателя к бесконечности. Считают, что если  $VIF_{x_j} > 3$  для каждой переменной  $x_j$ , то имеет место мультиколлинеарность.

Измерителем мультиколлинеарности является также так называемый **показатель (число) обусловленности** матрицы  $(X^T X)$ . Он равен отношению максимального  $\lambda_{\max}$  и минимального  $\lambda_{\min}$  собственных чисел этой матрицы:

$$\eta = \frac{\lambda_{\max}}{\lambda_{\min}}. \quad (2.60)$$

Считается, что если порядок этого соотношения превышает  $10^5 - 10^6$ , то имеет место сильная мультиколлинеарность<sup>1</sup>.

Проверим наличие мультиколлинеарности в рассматриваемом нами примере 2.1. Матрица парных коэффициентов корреляции имеет вид

$$r_{xx} = \begin{pmatrix} 1 & 0,603 & 0,600 \\ 0,603 & 1 & 0,478 \\ 0,600 & 0,478 & 1 \end{pmatrix}.$$

<sup>1</sup> См., например, Вучков И., Бояджиева Л., Солаков Е. Прикладной регрессионный анализ : пер. с болг. М. : Финансы и статистика, 1987. С. 110.

Можно отметить, что связи между объясняющими переменными достаточно тесные, особенно между переменными  $x_1$  и  $x_2$ ;  $x_1$  и  $x_3$ , что указывает на интеркорреляцию этих переменных. Более слабая связь наблюдается между переменными  $x_2$  и  $x_3$ . Найдем определитель матрицы  $r_{xx}$ .

$$\det(r_{xx}) = 1 + 0,603 \cdot 0,478 \cdot 0,600 + 0,603 \cdot 0,478 \cdot 0,600 - 0,600^2 - 0,603^2 - 0,478^2 = 0,394.$$

Полученное значение ближе к нулю, чем к единице, что указывает на наличие мультиколлинеарности объясняющих переменных.

Проверим обоснованность включения всех трех независимых переменных в модель регрессии, используя правило (2.59). Парные линейные коэффициенты корреляции зависимой и независимых переменных равны

$$r_{yx_1} = 0,736; r_{yx_2} = 0,721; r_{yx_3} = 0,702.$$

Они больше, чем показатели тесноты связи между независимыми переменными, следовательно, правило (2.59) выполняется, все три переменные можно включить в модель регрессии.

Измерим степень мультиколлинеарности переменных с помощью фактора инфляции дисперсии (*VIF*). Для этого необходимо рассчитать коэффициенты детерминации для регрессий:

$$\hat{x}_1 = a_0 + a_2x_2 + a_3x_3 \rightarrow R_{x_1 x_2 x_3}^2;$$

$$\hat{x}_2 = a_0 + a_1x_1 + a_3x_3 \rightarrow R_{x_2 x_1 x_3}^2;$$

$$\hat{x}_3 = a_0 + a_1x_1 + a_2x_2 \rightarrow R_{x_3 x_1 x_2}^2.$$

Для этого к каждой регрессии необходимо применить МНК, оценить ее параметры и рассчитать коэффициент детерминации. Для нашего примера результаты расчетов следующие:



$$\hat{x}_1 = 210,58 + 0,002x_2 + 0,009x_3; R_{x_1 x_2 x_3}^2 = 0,49;$$

$$\hat{x}_2 = -7914,77 + 123,98x_1 + 1,03x_3; R_{x_2 x_1 x_3}^2 = 0,38;$$

$$\hat{x}_3 = 3124,27 + 21,69x_1 + 0,03x_2; R_{x_3 x_1 x_2}^2 = 0,38.$$

Следовательно, фактор инфляции дисперсии для каждой независимой переменной будет равен

$$VIF_{x_1} = \frac{1}{1-0,49} = 1,96;$$

$$VIF_{x_2} = \frac{1}{1-0,38} = 1,62;$$

$$VIF_{x_3} = \frac{1}{1-0,38} = 1,62.$$

Все рассчитанные величины не превысили критического значения, равного трем, следовательно, при построении модели можно пренебречь существованием связей между независимыми переменными.

Для нахождения собственных чисел матрицы  $(X^T X)$  (с целью расчета показателя обусловленности  $\eta$  (2.60)) необходимо найти решение характеристического уравнения

$$|X^T X - \lambda I| = 0.$$

Матрица  $(X^T X)$  для нашего примера имеет вид

$$\begin{pmatrix} n & \sum x_1 & \sum x_2 & \sum x_3 \\ \sum x_1 & \sum x_1^2 & \sum x_1 x_2 & \sum x_1 x_3 \\ \sum x_2 & \sum x_1 x_2 & \sum x_2^2 & \sum x_2 x_3 \\ \sum x_3 & \sum x_1 x_3 & \sum x_2 x_3 & \sum x_3^2 \end{pmatrix},$$

а матрица, модуль определителя которой нужно приравнять нулю, получится следующей:

$$\begin{pmatrix} n-\lambda & \sum x_1 & \sum x_2 & \sum x_3 \\ \sum x_1 & \sum x_1^2-\lambda & \sum x_1x_2 & \sum x_1x_3 \\ \sum x_2 & \sum x_1x_2 & \sum x_2^2-\lambda & \sum x_2x_3 \\ \sum x_3 & \sum x_1x_3 & \sum x_2x_3 & \sum x_3^2-\lambda \end{pmatrix}.$$

Характеристический многочлен в данном случае будет иметь четвертую степень, что затрудняет решение задачи вручную. В данном случае рекомендуется воспользоваться возможностями вычислительной техники. Например, в ППП *EViews* получены следующие собственные числа матрицы  $(X^T X)$ :

$$\begin{aligned} \lambda_1 &= 3,939 \cdot 10^{11}; \\ \lambda_2 &= 1,0107 \cdot 10^9; \\ \lambda_3 &= 2,268 \cdot 10^6; \\ \lambda_4 &= 1,107 \cdot 10^1. \end{aligned}$$

Следовательно, показатель обусловленности  $\eta$  будет равен

$$\eta = \frac{\lambda_{\max}}{\lambda_{\min}} = \frac{3,939 \cdot 10^{11}}{1,107 \cdot 10} = 3,559 \cdot 10^{10},$$

что свидетельствует о наличии в модели сильной мультиколлинеарности.

Методами устранения мультиколлинеарности являются следующие.

1. Анализ связей между переменными, включаемыми в модель регрессии в качестве объясняющих (независимых), с целью отбора только тех переменных, которые слабо связаны друг с другом.
2. Функциональные преобразования тесно связанных между собой переменных. Например, мы предполагаем, что поступление налогов в городах зависит от количества жителей и площади города. Очевидно, что эти переменные будут тесно связаны. Их можно заменить одной относительной переменной «плотность населения».
3. Если по каким-то причинам перечень независимых переменных не подлежит изменению, то можно воспользоваться специальными методами корректировки моделей

с целью исключения мультиколлинеарности: ридж-регрессией (гребневой регрессией), методом главных компонент.

Применение **ридж-регрессии** предполагает корректировку элементов главной диагонали матрицы  $(X^T X)$  на некую произвольно задаваемую положительную величину  $\tau$ . Значение рекомендуется брать от 0,1 до 0,4. Н. Дрейпер, Г. Смит в своей работе приводят один из способов «автоматического» выбора величины  $\tau$ , предложенный Хоэрлом, Кеннардом и Белдвином<sup>1</sup>:

$$\tau = \frac{m \cdot SS_e}{n - m - 1} \cdot \frac{1}{\mathbf{a}^{*T} \mathbf{a}^*}, \quad (2.61)$$

где  $m$  — количество параметров (без учета свободного члена) в исходной модели регрессии;  $SS_e$  — остаточная сумма квадратов, полученная по исходной модели регрессии без корректировки на мультиколлинеарность;  $\mathbf{a}^*$  — вектор-столбец коэффициентов регрессии, преобразованных по формуле

$$a_j^* = a_j \cdot \sqrt{\sum (x_j - \bar{x}_j)^2} \quad (2.62)$$

где  $a_j$  — параметр при переменной  $x_j$  в исходной модели регрессии.

После выбора величины  $\tau$  формула для оценки параметров регрессии будет иметь вид

$$\mathbf{a}_\tau = (X_\tau^T X_\tau + \tau I)^{-1} X_\tau^T Y_\tau \quad (2.63)$$

где  $I$  — единичная матрица;  $X_\tau$  — матрица значений независимых переменных: исходных или преобразованных по формуле (2.64);  $Y_\tau$  — вектор значений зависимой переменной: исходных или преобразованных по формуле (2.65).

При построении ридж-регрессии рекомендуется преобразовывать независимые переменные

$$x_{\tau j} = \frac{x_j - \bar{x}_j}{\sqrt{\sum (x_j - \bar{x}_j)^2}} \quad (2.64)$$

<sup>1</sup> Дрейпер Н., Смит Г. Указ. соч. С. 514.

и результативную переменную

$$y_{\tau} = y - \bar{y}. \quad (2.65)$$

В этом случае после оценки параметров по формуле (2.63) необходимо перейти к регрессии по исходным переменным, используя соотношения

$$a_j = \frac{a_{\tau j}}{\sqrt{\sum (x_j - \bar{x}_j)^2}}, \quad j = 1, 2, \dots, p; \quad a_0 = \bar{y} - \sum_j a_j \bar{x}_j. \quad (2.66)$$

Оценки параметров регрессии, полученные с помощью формулы (2.63), будут смещенными. Однако, так как определитель матрицы  $(X^T X + \tau I)$  больше определителя матрицы  $(X^T X)$ , дисперсия оценок параметров регрессии уменьшится, что положительно повлияет на прогнозные свойства модели.

Рассмотрим применение ридж-регрессии для примера 2.1. Найдем величину  $\tau$  с помощью формулы (2.61). Для этого сначала рассчитаем вектор преобразованных коэффициентов регрессии по формуле (2.62):

$$\begin{aligned} a_1 &= 12,45 \cdot 1742,3 = 21\,688; \\ a_2 &= 0,06 \cdot 437\,747,7 = 26\,319; \\ a_3 &= 0,31 \cdot 77\,228,5 = 23\,963. \end{aligned}$$

Произведение  $\mathbf{a}^{*T} \mathbf{a}^*$  равно  $1,737 \cdot 10^9$ . Следовательно, рекомендуемое  $\tau$  составит

$$\tau = \frac{3 \cdot 1,33 \cdot 10^9}{48 - 3 - 1} \cdot \frac{1}{1,737 \cdot 10^9} = 0,0521.$$

После применения формулы (2.63) и преобразований по формуле (2.66) получим уравнение регрессии

$$y = 3532,46 + 12,35x_1 + 0,06x_2 + 0,30x_3 + e.$$

Применение *метода главных компонент* предполагает переход от взаимозависимых переменных  $x$  к независимым друг от друга переменным  $z$ , которые называют *главными*

**компонентами.** Каждая главная компонента  $z_j$  может быть представлена как линейная комбинация центрированных (или стандартизованных) объясняющих переменных  $t_j$ . Напомним, что центрирование переменной предполагает вычитание из каждого  $i$ -го значения данной  $j$ -й переменной ее среднего значения:

$$t_{ji}^* = x_{ji} - \bar{x}_j, \quad (2.67)$$

а стандартизация (масштабирование) — деление выражения (2.67) на среднее квадратическое отклонение, рассчитанное для исходных значений переменной  $x_j$ :

$$t_{ji} = \frac{x_{ji} - \bar{x}_j}{\sigma_{x_j}}. \quad (2.68)$$

Так как независимые переменные часто имеют разный масштаб измерения, формула (2.68) считается более предпочтительной.

Количество компонент может быть меньше или равно количеству исходных независимых переменных  $p$ . Компоненту с номером  $k$  можно записать следующим образом:

$$z_k = f_{k1}t_1 + f_{k2}t_2 + \dots + f_{kp}t_p. \quad (2.69)$$

Можно показать, что оценки  $f_{kj}$  в формуле (2.69) соответствуют элементам  $k$ -го собственного вектора матрицы  $\frac{1}{n}(\mathbf{T}^T \mathbf{T})$ , где  $\mathbf{T}$  — матрица размером  $(n \times p)$ , содержащая стандартизованные переменные. Нумерация главных компонент не является произвольной. Первая главная компонента имеет максимальную дисперсию, ей соответствует максимальное собственное число матрицы  $\frac{1}{n}(\mathbf{T}^T \mathbf{T})$ ; последняя — минимальную дисперсию и наименьшее собственное число.

Доля дисперсии  $k$ -й компоненты в общей дисперсии независимых переменных рассчитывается по формуле

$$d_k = \frac{\lambda_k}{\sum_j \lambda_j}, \quad (2.70)$$

где  $\lambda_k$  — собственное число, соответствующее данной компоненте; в знаменателе формулы (2.70) приведена сумма всех собственных чисел матрицы  $\frac{1}{n}(\mathbf{T}^T \mathbf{T})$ .

После расчета значений компонент  $z_j$  строят регрессию, используя МНК. Зависимую переменную в регрессии по главным компонентам (2.71) целесообразно центрировать (стандартизовать) по формулам (2.67) или (2.68).

$$t_y = b_1 z_1 + b_2 z_2 + \dots + b_k z_k + \delta, \quad (2.71)$$

где  $t_y$  — стандартизованная (центрированная) зависимая переменная;  $b_1, b_2, \dots, b_k$  — коэффициенты регрессии по главным компонентам;  $z_1, z_2, \dots, z_k$  — главные компоненты, упорядоченные по убыванию собственных чисел  $\lambda_k$ ;  $\delta$  — случайный остаток.

После оценки параметров регрессии (2.71) можно перейти к уравнению регрессии в исходных переменных, используя выражения (2.67)—(2.69).

Рассмотрим применение метода главных компонент на данных примера 2.1. Отметим, что матрица  $\frac{1}{n}(\mathbf{T}^T \mathbf{T})$  для стандартизованных переменных является в то же время матрицей парных линейных коэффициентов корреляции между независимыми переменными. Она уже была рассчитана и равна

$$\frac{1}{n} \mathbf{T}^T \mathbf{T} = \mathbf{r}_{xx} = \begin{pmatrix} 1 & 0,603 & 0,600 \\ 0,603 & 1 & 0,478 \\ 0,600 & 0,478 & 1 \end{pmatrix}.$$

Найдем собственные числа и собственные векторы этой матрицы, используя ППП *Eviews*. Получим следующие результаты.

Собственные числа матрицы  $\frac{1}{n}(\mathbf{T}^T \mathbf{T})$ :

$$\lambda_1 = 2,123; \quad \lambda_2 = 0,522; \quad \lambda_3 = 0,356.$$

Доля дисперсии независимых переменных, отражаемой компонентами, составила

$$d_1 = 0,708; d_2 = 0,174; d_3 = 0,119.$$

Объединим собственные векторы матрицы  $\frac{1}{n}(T^T T)$ , записав их как столбцы приведенной ниже матрицы  $F$ . Они упорядочены по убыванию собственных чисел, т.е. первый столбец является собственным вектором максимального собственного числа ( $\lambda_1 = 2,123$ ) и т.д.:

$$F = \begin{pmatrix} 0,604 & 0,007 & -0,800 \\ 0,564 & 0,702 & 0,434 \\ 0,563 & -0,710 & 0,420 \end{pmatrix}.$$

Следовательно, три компоненты (соответствующие трем собственным векторам) можно записать в виде

$$\begin{aligned} z_1 &= 0,604t_1 + 0,564t_2 + 0,563t_3; \\ z_2 &= 0,007t_1 + 0,702t_2 - 0,710t_3; \\ z_3 &= -0,800t_1 + 0,434t_2 + 0,420t_3. \end{aligned}$$

После стандартизации исходных переменных по формуле (2.68) и расчета значений компонент  $z_1, z_2, z_3$  (по  $n$  значений каждой компоненты) с помощью МНК найдем параметры уравнения (2.71):

$$\begin{aligned} t_y &= 0,587z_1 + 0,022z_2 + 0,059z_3 + \delta. \\ (t) \quad & (11,0) \quad (0,21) \quad (0,46) \end{aligned}$$

В полученном уравнении регрессии значим только параметр при первой компоненте. Это закономерный результат с учетом того, что данная компонента описывает 70,8% вариации независимых переменных. Так как компоненты независимы, при исключении из модели одних компонент параметры уравнения при других компонентах не меняются. Таким образом, имеем уравнение регрессии с одной компонентой:

$$t_y = 0,587z_1 + \delta.$$

Преобразуем полученное выражение в регрессию с исходными переменными

$$t_y = 0,587(0,604t_1 + 0,564t_2 + 0,563t_3) + \delta;$$

$$t_y = 0,354t_1 + 0,331t_2 + 0,330t_3 + \delta.$$

Откуда

$$a_1 = 0,354 \frac{\sigma_y}{\sigma_{x_1}} = 0,354 \frac{10\,188}{251,5} = 14,34;$$

$$a_2 = 0,331 \frac{\sigma_y}{\sigma_{x_2}} = 0,331 \frac{10\,188}{63\,183} = 0,053;$$

$$a_3 = 0,330 \frac{\sigma_y}{\sigma_{x_3}} = 0,330 \frac{10\,188}{11\,147} = 0,302;$$

$$a_0 = \bar{y} - \sum_j a_j \bar{x}_j = 17\,264 - 14,34 \cdot 449,8 - \\ - 0,053 \cdot 63\,291 - 0,302 \cdot 14\,930 = 2950,585.$$

Таким образом, используя метод главных компонент, мы получили уравнение регрессии

$$y = 2950,58 + 14,34x_1 + 0,053x_2 + 0,302x_3 + e.$$

Устранение мультиколлинеарности с помощью ридж-регрессии и метода главных компонент привело к определенному изменению параметров исходной регрессии, которая имела вид

$$y = 3226,44 + 12,45x_1 + 0,06x_2 + 0,31x_3 + e.$$

Отметим, что эти изменения были относительно невелики, что указывает на невысокую степень мультиколлинеарности.

## 2.12. Гетероскедастичность случайных остатков

Постоянство дисперсии случайных остатков называют **гомоскедастичностью** случайных остатков. Напротив, если эта



дисперсия не постоянна, то такое явление называют **гетероскедастичностью** случайных остатков.

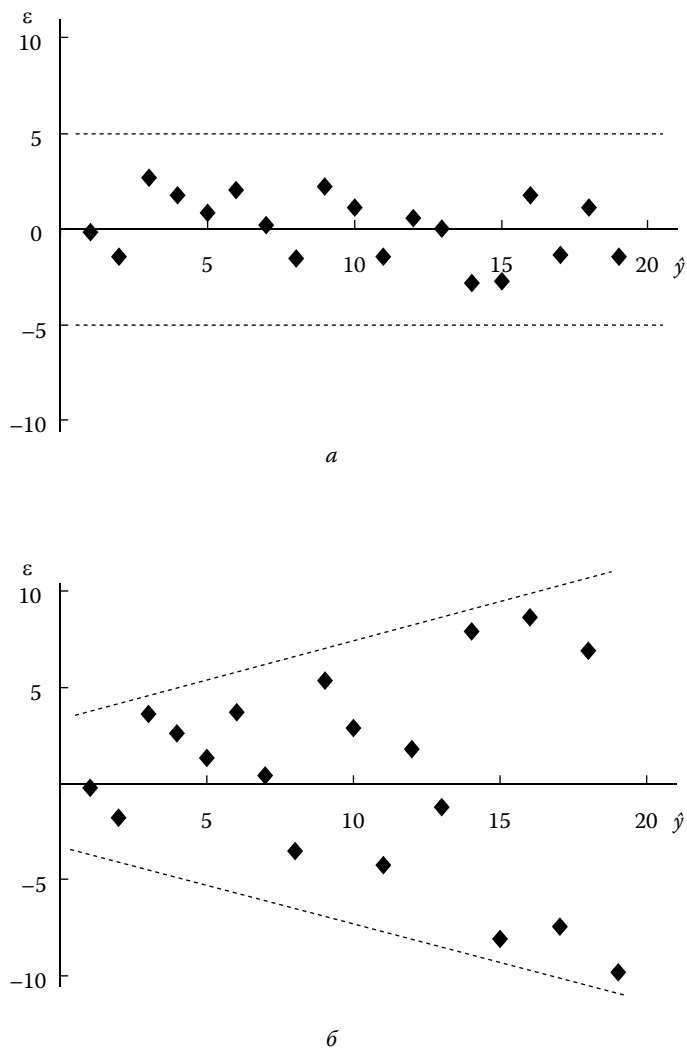
Проверка выполнения требования гомоскедастичности случайных остатков может быть произведена визуально, на основе графика остатков, или с помощью специальных критериев.

Для проведения визуального анализа необходимо построить график зависимости величин случайных остатков от выровненного значения результирующей переменной. В случае гомоскедастичности «облако» случайных остатков находится в области, параллельной оси абсцисс (рис. 2.2, а). Все прочие случаи соответствуют гетероскедастичности случайных остатков (например, как на рис. 2.2, б). Аналогичные графики можно построить также для зависимости случайных остатков от значений конкретных независимых переменных, входящих в регрессию.

К тестам, позволяющим выявить наличие гетероскедастичности случайных остатков, относят тесты Гольдфельда — Квандта, Парка, Глейзера, Уайта, Бреуша — Пагана, ранговой корреляции Спирмена и т.д.

**Тест Гольдфельда — Квандта** применяется, если случайные остатки предполагаются нормально распределенными величинами и объем наблюдений достаточно большой. Процедура проверки следующая.

1. Все наблюдения упорядочивают по мере возрастания какой-либо независимой переменной, которая, как предполагается, оказывает влияние на изменение дисперсии случайных остатков.
2. Упорядоченную совокупность делят на три группы, причем первая и последняя должны быть равного объема, с числом наблюдений, больших, чем число параметров модели регрессии. Пусть в первую и третью группы отобрано по  $k$  наблюдений.
3. По первой и третьей группам находят параметры уравнений регрессии той же структуры, что и исходное уравнение регрессии, и остаточные суммы квадратов по каждой модели.
4. Используя данные об остаточных суммах квадратов моделей первой и третьей групп, рассчитывают фактическое значение  $F$ -критерия Фишера по формуле



**Рис. 2.2.** Зависимость случайных остатков от выровненного значения результирующей переменной:  
*a* — нет зависимости (гомоскедастичность); *б* — дисперсия остатков увеличивается с увеличением выровненного значения результата (один из случаев гетероскедастичности)

$$F = \frac{SS_e^{(\max)}}{SS_e^{(\min)}},$$

где  $SS_e^{(\max)}$  — бóльшая остаточная сумма квадратов;  $SS_e^{(\min)}$  — меньшая остаточная сумма квадратов.

5. Сравнивают фактическое значение  $F$ -критерия с табличным, найденным для  $df_1 = df_2 = k - m - 1$  степеней свободы. Если  $F$ -фактическое больше табличного, то гипотеза об отсутствии гетероскедастичности отклоняется.

Тесты Парка, Глейзера, Уайта и Бреуша — Пагана основываются на предположении, чт.д.сперсия случайных остатков представляет собой определенную функцию от некоторой независимой переменной (или переменных). Перед применением этих тестов по уравнению регрессии необходимо рассчитать случайные остатки  $e_i$ .

Для **теста Парка** строят зависимость вида

$$\ln \varepsilon_i^2 = \alpha_0 + \alpha_j \ln x_{ji} + v_i, \quad (2.72)$$

где  $x_{ji}$  —  $i$ -е значение  $j$ -й независимой переменной, оказывающей влияние на дисперсию остатков;  $v_i$  — случайный остаток.

По **тесту Глейзера** находят параметры целой серии уравнений, задаваемых функцией

$$|\varepsilon_i| = \alpha_0 + \alpha_j x_{ji}^k + v_i, \quad (2.73)$$

где  $k$  — какое-либо число, например  $k = -1; -0,5; 0,5; 1$  и т.п.

**Тест Уайта** заключается в построении квадратичной функции, включающей все независимые переменные, входящие в исходную модель, а также их попарные произведения. Включение попарных произведений независимых переменных является необязательным, их можно опустить. Для случая с двумя переменными эта функция будет иметь вид

$$\varepsilon_i^2 = \alpha_0 + \alpha_{11}x_{1i} + \alpha_{12}x_{1i}^2 + \alpha_{21}x_{2i} + \alpha_{22}x_{2i}^2 + \gamma_{12}x_{1i}x_{2i} + v_i, \quad (2.74)$$

где  $\alpha, \gamma$  — неизвестные параметры.

**Тест Бреуша — Пагана** предполагает исследование влияния на дисперсию остатков нескольких независимых переменных, которые включают в регрессию вида

$$\frac{\varepsilon_i^2}{s_1^2} = \alpha_0 + \alpha_j x_{ji} + \alpha_{j+1} x_{(j+1)i} + \dots + \alpha_{j+k} x_{(j+k)i} + v_i, \quad (2.75)$$

где  $x_{ji}, x_{(j+1)i}, \dots, x_{(j+k)i}$  —  $i$ -е значения  $j$ -й,  $(j+1)$ -й,  $\dots$   $(j+k)$ -й независимых переменных, оказывающих влияние на дисперсию остатков;  $s_1^2$  — оценка дисперсии случайных остатков, рассчитанная по формуле

$$s_1^2 = \frac{\sum_{i=1}^n e_i^2}{n}.$$

Остатки считаются гетероскедастичными, если параметр  $\alpha_j$  в функциях по тесту Парка (2.72) или тесту Глейзера (2.73) значим (для теста Глейзера — хотя бы при одном значении  $k$ ). При проверке по тесту Уайта говорят, что остатки гетероскедастичны, если вся функция (2.74) значима по  $F$ -критерию Фишера.

Проверка гетероскедастичности по тесту Бреуша — Пагана заключается в расчете по функции (2.75) факторной суммы квадратов

$$SS_{\text{факт}} = \sum (\hat{z}_i - \bar{z})^2,$$

где переменная  $z_i = \frac{e_i^2}{s_1^2}$ .

Далее рассчитывается фактическое значение критерия  $\chi^2$  по формуле

$$\chi^2 = \frac{SS_{\text{факт}}}{2},$$

которое сравнивается с табличным  $\chi_{\text{таб}}^2$  (число степеней свободы равно  $df = k + 1$ , т.е. числу независимых переменных в модели (2.75); уровень значимости равен  $\alpha$ ). Нулевая гипотеза о гомоскедастичности случайных остатков отвергается, если

$$\chi^2 > \chi_{\text{таб}}^2.$$

**Тест ранговой корреляции Спирмена**, так же как и ранее рассмотренные тесты, основывается на предположении о зависимости (прямой или обратной) величины дисперсии случайных остатков от значений какой-либо независимой переменной. Для проведения проверки по этому тесту значения случайных остатков, взятые по модулю, и значения этой переменной ранжируют (например, по возрастанию), а затем находят коэффициент корреляции рангов Спирмена

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)},$$

где  $d_i$  — разность между рангами  $i$ -го случайного остатка и  $i$ -го значения независимой переменной.

Полученное значение коэффициента корреляции проверяют на значимость, рассчитывая фактическое значение  $t$ -критерия Стьюдента (2.76) и сравнивая его с табличным значением при числе степеней свободы  $df = n - 2$ .

$$t_\rho = \frac{\rho \sqrt{n-2}}{\sqrt{1-\rho^2}}. \quad (2.76)$$

Если фактическое значение критерия больше табличного, то гипотеза о гомоскедастичности остатков отклоняется.

Проверим на гетероскедастичность модель регрессии из нашего примера:

$$y = 3226,44 + 12,45x_1 + 0,06x_2 + 0,31x_3 + e.$$

Рассчитаем случайные остатки  $e$  для этой модели (табл. 2.5).

**Таблица 2.5.** Расчет случайных остатков для модели регрессии поступления налогов от количества занятых, объема отгрузки в обрабатывающих производствах и производства энергии

$y$	$\hat{y}$	$e$	$y$	$\hat{y}$	$e$
1422,20	4804,33	-3382,13	16 868,50	14 895,12	1973,38
2529,70	5056,17	-2526,47	18 019,40	13 781,67	4237,73
2629,10	5144,80	-2515,70	18 950,30	27 753,87	-8803,57

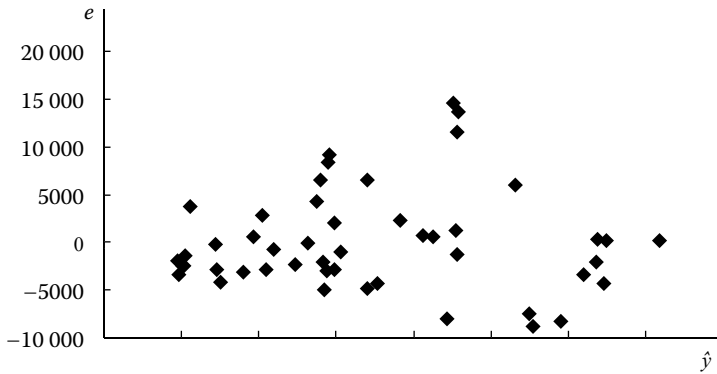
Окончание табл. 2.5

$y$	$\hat{y}$	$e$	$y$	$\hat{y}$	$e$
2764,30	4755,64	-1991,34	19 995,50	27 517,44	-7521,94
3347,50	7553,53	-4206,03	20 445,60	13 948,95	6496,65
3914,20	5263,55	-1349,35	21 220,80	29 518,94	-8298,14
4400,80	7241,83	-2841,03	21 360,00	20 644,76	715,24
5904,00	8992,88	-3088,88	21 418,80	19 152,00	2266,80
6956,70	7161,55	-204,85	21 477,10	22 791,19	-1314,09
7595,10	10 469,73	-2874,63	21 816,30	21 263,08	553,22
9257,80	14 251,91	-4994,11	22 824,90	14 496,10	8328,80
9317,10	5569,87	3747,23	23 579,30	17 021,68	6557,62
9978,80	12 356,21	-2377,41	23 702,60	14 531,28	9171,32
10 144,80	10 929,40	-784,60	24 007,20	22 773,85	1233,35
10 215,40	9619,74	595,66	27 581,20	31 028,32	-3447,12
11 349,50	14 390,38	-3040,88	28 057,50	32 314,79	-4257,29
12 046,90	14 174,13	-2127,23	29 815,50	31 859,41	-2043,91
12 061,40	14 898,60	-2837,20	32 236,50	31 936,40	300,10
12 104,20	17 000,04	-4895,84	32 657,40	32 494,15	163,25
13 042,40	10 214,84	2827,56	32 672,70	26 620,17	6052,53
13 104,30	13 167,07	-62,77	34 351,10	22 852,20	11 498,90
13 396,40	17 660,39	-4263,99	36 050,40	35 892,53	157,87
14 170,30	22 136,25	-7965,95	36 544,30	22 893,37	13 650,93
14 227,00	15 269,09	-1042,09	37 136,90	22 606,54	14 530,36

График зависимости случайных остатков от выровненного значения зависимой переменной имеет вид, представленный на рис. 2.3. Можно отметить определенное увеличение разброса точек в центральной части графика и уменьшение разброса для последних нескольких точек. Такая картина может свидетельствовать о наличии гетероскедастичности остатков.

Применим для анализа дисперсии остатков рассмотренные выше тесты. Так как большинство тестов основано на гипотезе, что известна переменная, вызывающая гетероскедастичность остатков, обратимся сначала к тесту Уайта, в котором рассматриваются все независимые переменные, входящие в модель регрессии.

Используем короткую форму теста Уайта, без включения попарных произведений независимых переменных. Получим следующий результат:



**Рис. 2.3.** Изменение дисперсии случайных остатков с ростом выровненного значения зависимой переменной

$$e^2 = -1,004 \cdot 10^7 - 4755,6x_1 - 36,2x_1^2 + 108,1x_2 - 5,2 \cdot 10^{-4}x_2^2 + 4869x_3 - 0,07x_3^2 + v;$$

$$(t) \quad (-0,04) \quad (-0,42) \quad (0,22) \quad (-0,28) \quad (1,77) \quad (-1,42)$$

$$F = 1,41.$$

Табличное значение  $F$ -критерия равно 2,33 ( $\alpha = 0,05$ ;  $df_1 = 6$ ;  $df_2 = 41$ ). Таким образом, по тесту Уайта нет оснований отвергнуть нулевую гипотезу о гомоскедастичности остатков. Отметим также, что все параметры незначимы, но наибольшее значение  $t$ -критерия (и достаточно близкое к табличному) имеют параметры при переменной  $x_3$  (табличное значение  $t$ -критерия составило 2,02 ( $\alpha = 0,05$ ;  $df = 41$ )). Таким образом, переменная  $x_3$  может быть рассмотрена в других тестах как возможная причина гетероскедастичности.

Тест Бреуша — Пагана позволяет рассматривать различные комбинации переменных в качестве объясняющих гетероскедастичность остатков. Уравнение теста, включающее в себя все три независимые переменные, будет иметь вид

$$z = 0,36 + 2,4 \cdot 10^{-6}x_1 + 5,1 \cdot 10^{-5}x_2 - 10^{-6}x_3;$$

$$\chi^2 = \frac{SS_{\text{факт}}}{2} = \frac{14,82}{2} = 7,41.$$

Табличное значение критерия  $\chi^2$  равно 7,82 ( $\alpha = 0,05$ ;  $df = 3$ ), таким образом, нет оснований отвергнуть нулевую гипотезу о гомоскедастичности случайных остатков. Руководствуясь предположениями, сделанными в ходе анализа теста Уайта, проведем тест Бреуша — Пагана применительно только к переменной  $x_3$ . Получим следующие результаты:

$$\hat{z} = 0,33 + 4,45 \cdot 10^{-5} x_3;$$

$$\chi^2 = \frac{SS_{\text{факт}}}{2} = \frac{11,83}{2} = 5,915.$$

Табличное значение критерия  $\chi^2$  в данном случае равно 3,84 ( $\alpha = 0,05$ ;  $df = 1$ ), таким образом, мы отвергаем нулевую гипотезу о гомоскедастичности случайных остатков. Остатки гетероскедастичны по переменной  $x_3$ . Анализ по тесту Бреуша — Пагана при необходимости можно продолжить, исследуя влияние на дисперсию случайных остатков других независимых переменных. Опираясь на выявленное влияние на дисперсию остатков переменной  $x_3$ , проверим эту связь с помощью других тестов.

Использование критерия Гольдфелда — Квандта предполагает упорядочивание данных, в нашем случае по переменной  $x_3$ .

Общий объем наблюдений составляет 48 регионов, т.е. их можно разделить на три равные группы по 16 наблюдений в каждой или по 18 наблюдений в первой и третьей группах и 12 наблюдений во второй. Так как критерий Гольдфелда — Квандта предполагает построение уравнений регрессии той же структуры, что и исходное уравнение, остановимся на втором варианте деления совокупности как обеспечивающим большую достоверность регрессионного анализа (18 наблюдений на три коэффициента регрессии, т.е. по шесть наблюдений на каждый коэффициент).

Для первой и третьей совокупностей наблюдений найдем параметры уравнений множественной регрессии вида  $\hat{y} = a_0 + a_1x_1 + a_2x_2 + a_3x_3$  и рассчитаем случайные остатки по каждому из них. Получим следующие результаты.

Первая группа (минимальные значения  $x_3$ ):

$$\hat{y} = 711,4 + 3,5x_1 + 0,06x_2 + 0,94x_3,$$



$$SS_e^{(1)} = 35\,799\,250,9.$$

Третья группа (максимальные значения  $x_3$ ):

$$\hat{y} = 6289,4 + 18,4x_1 + 0,027x_2 + 0,17x_3;$$

$$SS_e^{(3)} = 665\,106\,150,5.$$

Разделим большую остаточную сумму квадратов (по третьей группе) на меньшую (по первой группе):  $F = \frac{665\,106\,150,5}{35\,799\,250,9} = 18,58$ . Табличное значение  $F$ -критерия равно 2,48 при  $df_1 = df_2 = 18 - 4 = 14$  степенях свободы и уровне значимости 0,05. Следовательно, дисперсия остатков зависит от величины значений переменной  $x_3$ , гипотеза о гомоскедастичности остатков отвергается.

Расчеты по тестам Парка и Глейзера по переменной  $x_3$  приводят к следующим результатам.

Тест Парка:

$$\ln e^2 = 14,02 + 0,16 \ln x_3 + v, \quad t_{a_3} = 0,43.$$

Тест Глейзера:

$$\text{при } k = 1 \quad |e| = 2640,3 + 0,088x_3 + v, \quad t_{a_3} = 2,0134;$$

$$\text{при } k = -1 \quad |e| = 4399 - 2\,679\,286 \frac{1}{x_3} + v, \quad t_{a_3} = -1,29;$$

$$\text{при } k = 0,5 \quad |e| = 1282,9 + 23,6\sqrt{x_3} + v, \quad t_{a_3} = 2,204;$$

$$\text{при } k = 0,1 \quad |e| = 14,02 + 0,16^{10}\sqrt{x_3} + v, \quad t_{a_3} = 0,43.$$

Табличное значение критерия Стьюдента равно 2,0129 ( $\alpha = 0,05$ ;  $df = 46$ ). Таким образом, по тесту Глейзера при  $k =$

1 и 0,5 гипотеза о гомоскедастичности остатков отвергается, по тесту Парка — нет оснований отвергнуть нулевую гипотезу.

Оценка гетероскедастичности остатков по переменной  $x_3$  с использованием коэффициента ранговой корреляции Спирмена привела к следующему результату:

$$\rho = 0,258, \quad t_\rho = 1,81.$$

То есть коэффициент ранговой корреляции незначим (табличное значение критерия Стьюдента, так же как в тестах Глейзера и Парка, равно 2,0129), нет оснований отвергнуть нулевую гипотезу о гомоскедастичности остатков.

Подводя итоги выявления гетероскедастичности в нашем примере, отметим, что по ряду тестов (Бреуша — Пагана, Гольдфельда — Квандта, Глейзера) гипотеза о гомоскедастичности остатков была отвергнута, т.е. можно утверждать, что на дисперсию случайных остатков оказывает влияние переменная  $x_3$ . То, что гетероскедастичность была выявлена не во всех тестах, связано с тем, что разные тесты опираются на разные предпосылки о форме связи величины случайных остатков и независимой переменной. Исследование по тесту Глейзера показывает, что эта форма может быть описана выражением  $\sigma_\varepsilon = f(x_3)$ , где  $f$  — линейная функция.

Причинами гетероскедастичности случайных остатков могут быть неверная функциональная форма уравнения регрессии (неверная спецификация модели), неоднородность исследуемой совокупности. Соответственно способами устранения гетероскедастичности являются построение модели иной функциональной формы и (или) разбиение совокупности на однородные группы. Если по каким-то причинам это сделать невозможно или нежелательно, т.д.я нахождения параметров уравнения регрессии можно воспользоваться обобщенным методом наименьших квадратов.

## 2.13. Обобщенный метод наименьших квадратов

**Обобщенный метод наименьших квадратов (ОМНК)** применяется в тех случаях, когда нарушены условия Гаусса — Маркова, касающиеся характера случайных остатков, а именно:

- гомоскедастичность (постоянство дисперсии) случайных остатков;
- некоррелированность остатков между собой.

Нарушение этих условий означает, что ковариационная матрица остатков  $\mathbf{\Omega}$  не является скалярной. Она будет иметь вид

$$M(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^T) = \mathbf{\Omega} = \begin{pmatrix} \sigma_{\varepsilon_1}^2 & \text{cov}(\varepsilon_1\varepsilon_2) & \text{cov}(\varepsilon_1\varepsilon_3) & \dots & \text{cov}(\varepsilon_1\varepsilon_n) \\ \text{cov}(\varepsilon_2\varepsilon_1) & \sigma_{\varepsilon_2}^2 & \text{cov}(\varepsilon_2\varepsilon_3) & \dots & \text{cov}(\varepsilon_2\varepsilon_n) \\ \text{cov}(\varepsilon_3\varepsilon_1) & \text{cov}(\varepsilon_3\varepsilon_2) & \sigma_{\varepsilon_3}^2 & \dots & \text{cov}(\varepsilon_3\varepsilon_n) \\ \dots & \dots & \dots & \dots & \dots \\ \text{cov}(\varepsilon_n\varepsilon_1) & \text{cov}(\varepsilon_n\varepsilon_2) & \text{cov}(\varepsilon_n\varepsilon_3) & \dots & \sigma_{\varepsilon_n}^2 \end{pmatrix} \quad (2.77)$$

На главной диагонали этой матрицы расположены дисперсии случайных остатков для различных наблюдений ( $i = 1 - n$ ), не одинаковые по своей величине. Ковариации случайных остатков  $\varepsilon_i$  и  $\varepsilon_j$  в общем случае также ненулевые. Если остальные условия построения классической нормальной линейной модели выполняются (см. параграф 2.2), то модель

$$Y = X\boldsymbol{\alpha} + \boldsymbol{\varepsilon} \quad (2.78)$$

называется **обобщенной линейной моделью**.

Напомним, что формула для расчета вектора-столбца неизвестных параметров с помощью обычного МНК в матричной форме имеет вид (см. параграф 2.2)

$$\boldsymbol{a} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y.$$

Скалярность матрицы  $\mathbf{\Omega}$  использовалась выше при выводе выражения для ковариационной матрицы оценок параметров уравнения регрессии  $\Sigma_{\boldsymbol{a}}$ . Можно показать, что, если матрица  $\mathbf{\Omega}$  не скалярная, то оценки параметров, содержащиеся в векторе  $\boldsymbol{a}$ , являются неэффективными (т.е. их дисперсия

не является наименьшей из дисперсий возможных оценок параметров  $\mathbf{a}$ ). В свою очередь, эти дисперсии являются смещенными и несостоятельными, что приводит к ложным выводам при оценке качества модели (оценке статистической значимости параметров и всей модели в целом) и при проведении прогнозирования по ней.

Сущность обобщенного метода наименьших квадратов состоит в том, чтобы устранить нарушения предпосылок МНК, «скорректировав» расчеты параметров уравнения регрессии с учетом значений ковариационной матрицы остатков. Такая «корректировка» может быть проведена с использованием формулы

$$\mathbf{a}^* = (\mathbf{X}^T \mathbf{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{\Omega}^{-1} \mathbf{Y}, \quad (2.79)$$

где  $\mathbf{\Omega}$  — ковариационная матрица остатков (2.77).

Доказательство эффективности оценок, полученных с помощью обобщенного МНК (ОМНК), содержится в теореме Айткена.

**Теорема Айткена.** В классе линейных несмещенных оценок вектора  $\mathbf{a}$  для обобщенной линейной модели оценка  $\mathbf{a}^*$  (2.79) имеет наименьшую ковариационную матрицу.

Найдем математическое ожидание оценки  $\mathbf{a}^*$ . В выражении для  $\mathbf{a}^*$  заменим вектор  $\mathbf{Y}$  на функцию регрессии (2.78).

$$\begin{aligned} M\mathbf{a}^* &= M[(\mathbf{X}^T \mathbf{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{\Omega}^{-1} \mathbf{Y}] = \\ &= M[(\mathbf{X}^T \mathbf{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{\Omega}^{-1} \mathbf{X} \mathbf{\alpha} + (\mathbf{X}^T \mathbf{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{\Omega}^{-1} \mathbf{\varepsilon}]. \end{aligned}$$

Математическое ожидание суммы равно сумме математических ожиданий:

$$\begin{aligned} M\mathbf{a}^* &= M[(\mathbf{X}^T \mathbf{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{\Omega}^{-1} \mathbf{X} \mathbf{\alpha}] + M[(\mathbf{X}^T \mathbf{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{\Omega}^{-1} \mathbf{\varepsilon}] = \\ &= M(I_n \mathbf{\alpha}) + M[(\mathbf{X}^T \mathbf{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{\Omega}^{-1} \mathbf{\varepsilon}]. \end{aligned}$$

Первое слагаемое после упрощения (перемножение матриц дает единичную матрицу) равно вектору значений параметров  $\mathbf{\alpha}$  в генеральной совокупности, его математическое ожидание (как константы) равно ему самому. Во втором слагаемом произведение, состоящее из матриц  $\mathbf{X}$  и  $\mathbf{\Omega}$ , является неслучайной

величиной, так как  $X$  содержит независимые случайные переменные,  $\Omega$  — ковариационная матрица случайных остатков. Математическое ожидание случайных остатков (согласно одному из условий Гаусса — Маркова) равно нулю, поэтому второе слагаемое равно нулю:

$$Ma^* = \alpha + 0 = \alpha.$$

Таким образом, ОМНК-оценка является несмещенной.

Ковариационная матрица случайных остатков  $\Omega$  является невырожденной симметричной матрицей, поэтому она может быть представлена через некую невырожденную матрицу  $P$  следующим образом:

$$\Omega = PP^T.$$

В свою очередь, обратная матрица  $\Omega^{-1}$  может быть выражена как

$$\Omega^{-1} = (P^{-1})^T P^{-1}.$$

Запишем формулу оценок параметров линейной регрессии по ОМНК, используя вместо матрицы  $\Omega$  матрицу  $P$ .

$$a^* = (X^T (P^{-1})^T (P^{-1}) X)^{-1} X^T (P^{-1})^T (P^{-1}) Y.$$

Если произвести замену переменных

$$X^* = P^{-1} X; \quad Y^* = P^{-1} Y; \quad a^* = P^{-1} \epsilon, \quad (2.80)$$

то оценка  $a^*$  будет иметь вид оценки параметров по обычному МНК:

$$a^* = (X^{*T} X^*)^{-1} X^* Y^*.$$

Иначе говоря, если преобразовать переменные  $X$  и  $Y$  по формулам (2.80), то применение к ним обычного МНК приведет к тому же результату, что и применение ОМНК к исходным переменным.

Преобразование переменных (2.80) приводит к следующему виду обобщенной модели регрессии:

$$Y^* = X^* \alpha + \epsilon^*. \quad (2.81)$$

Покажем, что применение ОМНК позволяет получить гомоскедастичные, некоррелированные случайные остатки. Математическое ожидание случайных остатков модели (2.81) равно нулю, так как в исходной модели эта предпосылка не нарушена ( $M\epsilon = 0$ ):

$$M\epsilon^* = M[(P^{-1})\epsilon] = (P^{-1})M(\epsilon) = 0.$$

Ковариация случайных остатков для модели (2.81) равна (в матричной форме)

$$M(\epsilon^* \epsilon^{*T}) = M\{(P^{-1})\epsilon[(P^{-1})\epsilon]^T\} = M[(P^{-1})\epsilon\epsilon^T(P^{-1})^T].$$

Так как матрица  $(P^{-1})$  формируется на основе ковариационной матрицы остатков  $\Omega$ , т.е. элементы являются постоянными величинами, их можно вынести за знак математического ожидания:

$$M(\epsilon^* \epsilon^{*T}) = (P^{-1})M(\epsilon\epsilon^T)(P^{-1})^T.$$

Математическое ожидание произведения  $\epsilon\epsilon^T$  равно ковариационной матрице случайных остатков  $\Omega$ , которую мы представили как  $PP^T$ :

$$M(\epsilon^* \epsilon^{*T}) = (P^{-1})\Omega(P^{-1})^T = (P^{-1})PP^T(P^{-1})^T.$$

После перегруппировки получаем два одинаковых произведения матрицы  $P$  на обратную ей матрицу  $(P^{-1})$ , которые по определению дают единичные матрицы:

$$M(\epsilon^* \epsilon^{*T}) = [(P^{-1})P][(P^{-1})P]^T = I_n.$$

Напомним, что единичная матрица содержит на главной диагонали единицы, а ее остальные элементы равны нулю. Следовательно, мы доказали, что ковариационная матрица остатков при применении к обобщенной модели линейной

регрессии ОМНК удовлетворяет условиям Гаусса — Маркова о гомоскедастичности и отсутствия автокорреляции случайных остатков. Более того, мы получили, что при применении ОМНК дисперсии случайных остатков всегда будут равны единице (элементы главной диагонали единичной матрицы).

Подводя итоги, отметим, что получены два способа применения ОМНК: с использованием ковариационной матрицы остатков  $\Omega$  и формулы (2.79) и с использованием матрицы  $P$  и модели с преобразованными переменными (2.81), для решения которой достаточно применить обычный МНК.

Если предпосылки МНК о гомоскедастичности и некоррелированности остатков выполняются, то матрица  $\Omega$  превращается в скалярную и обобщенный метод наименьших квадратов дает тот же результат, что и обычный МНК.

Основная проблема, возникающая при использовании ОМНК, заключается в том, что фактические значения элементов матрицы  $\Omega$  неизвестны. Поэтому для применения этого метода используют их *оценки*, полученные на основе исследования имеющихся в распоряжении данных. В этом случае говорят о *доступном ОМНК*.

Для оценки элементов ковариационной матрицы остатков  $\Omega$  выдвигают разные предположения об их характере и структуре.

Как правило, считают, что гетероскедастичность остатков присуща данным пространственных выборок, а автокорреляция остатков — временным рядам. Поэтому в зависимости от характера исходных данных часто предполагают или гетероскедастичность остатков при условии отсутствия автокорреляции, или их автокорреляцию при условии гомоскедастичности.

Если остатки только гетероскедастичны, ковариационная матрица остатков имеет вид

$$\Omega = \begin{pmatrix} \sigma_{\varepsilon_1}^2 & 0 & \dots & 0 \\ 0 & \sigma_{\varepsilon_2}^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_{\varepsilon_n}^2 \end{pmatrix}. \quad (2.82)$$

Таким образом, следует дать оценку не всем элементам матрицы  $\Omega$ , а только ее  $n$  элементам, стоящим на главной диагонали.

Обобщенный метод наименьших квадратов можно с учетом матрицы (2.82) записать в виде системы нормальных уравнений

$$\begin{cases} \sum_i \frac{y_i}{\sigma_{\varepsilon i}^2} = a_0 \sum_i \frac{1}{\sigma_{\varepsilon i}^2} + a_1 \sum_i \frac{x_{1i}}{\sigma_{\varepsilon i}^2} + \dots + a_p \sum_i \frac{x_{pi}}{\sigma_{\varepsilon i}^2}, \\ \sum_i \frac{y_i x_{1i}}{\sigma_{\varepsilon i}^2} = a_0 \sum_i \frac{x_{1i}}{\sigma_{\varepsilon i}^2} + a_1 \sum_i \frac{x_{1i}^2}{\sigma_{\varepsilon i}^2} + \dots + a_p \sum_i \frac{x_{pi} x_{1i}}{\sigma_{\varepsilon i}^2}, \\ \vdots \\ \sum_i \frac{y_i x_{pi}}{\sigma_{\varepsilon i}^2} = a_0 \sum_i \frac{x_{pi}}{\sigma_{\varepsilon i}^2} + a_1 \sum_i \frac{x_{1i} x_{pi}}{\sigma_{\varepsilon i}^2} + \dots + a_p \sum_i \frac{x_{pi}^2}{\sigma_{\varepsilon i}^2}. \end{cases}$$

В данной системе переменные и их попарные произведения как бы взвешиваются путем умножения на  $\frac{1}{\sigma_{\varepsilon i}^2}$ , где  $i$  — номер наблюдения. Поэтому ОМНК для случая гетероскедастичных остатков называют еще **методом взвешенных наименьших квадратов**. Подчеркнем, что весовой коэффициент значений исходных переменных  $x$  и  $y$  для каждого номера наблюдения будет свой. То есть для первого наблюдения все переменные делятся на  $\sigma_{\varepsilon 1}^2$ , для второго — на  $\sigma_{\varepsilon 2}^2$ , для последнего,  $n$ -го наблюдения, — на  $\sigma_{\varepsilon n}^2$ .

Дисперсия случайного остатка  $\varepsilon_i$  (при условии  $M\varepsilon_i = 0$ ) равна

$$M(\varepsilon_i - M\varepsilon_i)^2 = M\varepsilon_i^2.$$

Поэтому в качестве оценок дисперсий случайных остатков могут быть использованы квадраты остатков, полученных при применении обычного МНК.

Процедура применения ОМНК в данном случае предполагает следующие шаги:

- к исходным данным применяется обычный МНК и вычисляются случайные остатки  $e_i$ ;
- делаются предположения относительно функциональной зависимости дисперсии случайных остатков от каких-либо переменных:



$$\hat{\varepsilon}^2 = f(x_1, x_2, \dots). \quad (2.83)$$

В качестве функции, может, например, использоваться функция по тесту Уайта

$$\hat{\varepsilon}^2 = a_0 + a_{11}x_1 + a_{12}x_1^2 + \dots + a_{p1}x_p + a_{p2}x_p^2,$$

или любая другая аналогичная функция;

- с помощью МНК находят параметры модели (2.83), используя в качестве фактических значений зависимой переменной случайные остатки  $\varepsilon_i$ , найденные на первом шаге;
- по модели (2.83) рассчитывают выровненные значения случайных остатков  $\hat{\varepsilon}_i^2$ . Эти значения рассматривают как оценки неизвестных дисперсий случайных остатков  $\sigma_{\varepsilon_i}^2$ , т.е. диагональных элементов матрицы (2.82);
- определяют параметры множественной линейной регрессии с помощью ОМНК в зависимости от алгоритма, выбранного для расчета параметров: по формуле (2.79) или с использованием преобразованных переменных с применением к ним обычного МНК. В последнем случае возникает вопрос о виде матрицы ( $P^{-1}$ ).

Так как в рассматриваемом случае матрица  $\mathbf{\Omega}$  имеет только диагональные элементы, то обратная к ней матрица  $\mathbf{\Omega}^{-1}$  также является диагональной и имеет вид

$$\mathbf{\Omega}^{-1} = \begin{pmatrix} \frac{1}{\sigma_{\varepsilon_1}^2} & 0 & \dots & 0 \\ 0 & \frac{1}{\sigma_{\varepsilon_2}^2} & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \frac{1}{\sigma_{\varepsilon_n}^2} \end{pmatrix}.$$

Тогда матрица ( $P^{-1}$ ) равна

$$P^{-1} = \begin{pmatrix} \frac{1}{\sigma_{\varepsilon_1}} & 0 & \dots & 0 \\ 0 & \frac{1}{\sigma_{\varepsilon_2}} & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \frac{1}{\sigma_{\varepsilon_n}} \end{pmatrix}.$$

С учетом преобразования переменных уравнение множественной линейной регрессии будет иметь вид

$$\frac{y_i}{\sigma_{\varepsilon_i}} = a_0 \frac{1}{\sigma_{\varepsilon_i}} + a_1 \frac{x_{1i}}{\sigma_{\varepsilon_i}} + \dots + a_p \frac{x_{pi}}{\sigma_{\varepsilon_i}} + \frac{e_i}{\sigma_{\varepsilon_i}}.$$

Отметим, что полученное уравнение регрессии не имеет свободного члена.

Применение ОМНК для случая гетероскедастичности остатков еще более упрощается, если предполагается зависимость дисперсии случайных остатков от квадратов значения какой-то одной переменной  $x_j$ :

$$\sigma_{\varepsilon_i}^2 = \sigma^2 x_{ji}^2,$$

где  $i$  — номер наблюдения;  $j$  — номер переменной.

В этом случае ковариационная матрица остатков будет иметь вид

$$\Omega = \sigma^2 \begin{pmatrix} x_{j1}^2 & 0 & \dots & 0 \\ 0 & x_{j2}^2 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & x_{jn}^2 \end{pmatrix}.$$

Диагональные элементы матрицы  $\Omega$  уже известны и равны квадратам фактических значений переменной  $x_j$ . Постоянный

множитель  $\sigma^2$  при расчетах по формуле (2.79) сокращается, поэтому матрицу  $\Omega$  можно заменить матрицей  $\Omega_1$  без этого множителя:

$$\Omega_1 = \begin{pmatrix} x_{j1}^2 & 0 & \dots & 0 \\ 0 & x_{j2}^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & x_{jn}^2 \end{pmatrix}.$$

Матрица ( $P^{-1}$ ) в этом случае равна

$$P^{-1} = \begin{pmatrix} \frac{1}{x_{j1}} & 0 & \dots & 0 \\ 0 & \frac{1}{x_{j2}} & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \frac{1}{x_{jn}} \end{pmatrix}.$$

Использование значений какой-либо независимой переменной в качестве основы для оценки ковариационной матрицы остатков позволяет дать экономическую интерпретацию результатов, полученных по ОМНК. Пусть, например, имеется уравнение множественной линейной регрессии с двумя независимыми переменными

$$y = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \varepsilon. \quad (2.84)$$

Предположим также, что установлена гетероскедастичность случайных остатков  $\varepsilon$  и зависимость их дисперсий от квадрата переменной  $x_2$ . Тогда применение ОМНК предполагает, что каждый элемент уравнения регрессии (2.84) будет разделен на  $x_2$ :

$$\frac{y}{x_2} = \alpha_0 \frac{1}{x_2} + \alpha_1 \frac{x_1}{x_2} + \alpha_2 + \frac{\varepsilon}{x_2}.$$

Преобразованное уравнение регрессии содержит свободный член ( $\alpha_2$ ). Однако его экономическая интерпретация остается

той же, что в исходном уравнении, — это коэффициент регрессии при переменной  $x_2$ . Параметр  $\alpha_0$  является свободным членом и не интерпретируется, а параметр  $\alpha_1$  является коэффициентом регрессии при переменной  $x_1$ .

Если случайные остатки автокоррелированы и предполагается отсутствие гетероскедастичности, то ковариационная матрица остатков примет вид

$$\Omega = \begin{pmatrix} \sigma_\varepsilon^2 & \text{cov}(\varepsilon_1\varepsilon_2) & \dots & \text{cov}(\varepsilon_1\varepsilon_n) \\ \text{cov}(\varepsilon_2\varepsilon_1) & \sigma_\varepsilon^2 & \dots & \text{cov}(\varepsilon_2\varepsilon_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(\varepsilon_n\varepsilon_1) & \text{cov}(\varepsilon_n\varepsilon_2) & \dots & \sigma_\varepsilon^2 \end{pmatrix}. \quad (2.85)$$

В этой матрице неизвестны уже  $n(n-1)/2$  элементов, поэтому, как правило, выдвигают предположения относительно взаимосвязи этих элементов между собой, сокращая тем самым количество неизвестных. Например, предполагают, что случайные остатки связаны автокорреляционной зависимостью первого порядка:

$$\varepsilon_i = \rho_1 \varepsilon_{i-1} + \delta_i, \quad (2.86)$$

где  $\delta_i$  — случайные остатки, удовлетворяющие предпосылкам МНК,  $\rho_1$  — неизвестный параметр.

Уравнение (2.86) не имеет свободного члена, так как предпосылка о равенстве нулю математического ожидания остатков по-прежнему соблюдается. Если применить к выражению (2.86) МНК, то получим следующее выражение для  $\rho_1$ :

$$\rho_1 = \frac{\sum \varepsilon_i \varepsilon_{i-1}}{\sum \varepsilon_{i-1}^2} = \frac{\text{cov}(\varepsilon_i \varepsilon_{i-1})n}{n\sigma_{\varepsilon_{i-1}}^2} = \frac{\text{cov}(\varepsilon_i \varepsilon_{i-1})}{\sigma_\varepsilon^2}.$$

Переход в числителе дроби (2.86) от суммы произведений остатков к их ковариации возможен, так как их средние равны нулю. Этим же объясняется переход от суммы квадратов остатков  $\sum_{i-1} \varepsilon_{i-1}^2$  к дисперсии этой величины. Далее, так как мы предполагаем гомоскедастичность остатков, дисперсия остатка

$\sigma_{\varepsilon_i}^2$  была заменена на постоянную величину  $\sigma_{\varepsilon}^2$ . Отсюда ковариация случайных остатков  $\varepsilon_i$  и  $\varepsilon_{i-1}$  равна

$$\text{cov}(\varepsilon_i \varepsilon_{i-1}) = \rho_1 \sigma_{\varepsilon}^2.$$

Для выявления зависимости между остатком  $\varepsilon_i$  и остатками  $\varepsilon_{i-2}, \varepsilon_{i-3}, \varepsilon_{i-4}, \dots$  будем последовательно применять формулу (2.86):

$$\varepsilon_{i-1} = \rho_1 \varepsilon_{i-2} + \delta_{i-1};$$

$$\varepsilon_i = \rho_1 (\rho_1 \varepsilon_{i-2} + \delta_{i-1}) + \delta_i;$$

$$\varepsilon_i = \rho_1^2 \varepsilon_{i-2} + \rho_1 \delta_{i-1} + \delta_i = \rho_1^2 \varepsilon_{i-2} + \delta_{i-1}^*.$$

где  $\delta_{i-1}^*$  — случайный остаток.

Аналогично

$$\varepsilon_i = \rho_1^3 \varepsilon_{i-3} + \delta_{i-2}^*;$$

$$\varepsilon_i = \rho_1^4 \varepsilon_{i-4} + \delta_{i-3}^*$$

и т.д.

Тогда ковариации между остатком  $\varepsilon_i$  и остатками  $\varepsilon_{i-2}, \varepsilon_{i-3}, \varepsilon_{i-4}, \dots$  будут равны

$$\text{cov}(\varepsilon_i \varepsilon_{i-2}) = \rho_1^2 \sigma_{\varepsilon}^2,$$

$$\text{cov}(\varepsilon_i \varepsilon_{i-3}) = \rho_1^3 \sigma_{\varepsilon}^2;$$

$$\text{cov}(\varepsilon_i \varepsilon_{i-4}) = \rho_1^4 \sigma_{\varepsilon}^2$$

и т.д.

Ковариационная матрица остатков (2.85) примет вид

$$\mathbf{\Omega} = \sigma_{\varepsilon}^2 \begin{pmatrix} 1 & \rho_1 & \rho_1^2 & \dots & \rho_1^{n-1} \\ \rho_1 & 1 & \rho_1 & \dots & \rho_1^{n-2} \\ \rho_1^2 & \rho_1 & 1 & \dots & \rho_1^{n-3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho_1^{n-1} & \rho_1^{n-2} & \rho_1^{n-3} & \dots & 1 \end{pmatrix}.$$

где  $\sigma_\varepsilon^2$  — постоянная дисперсия остатков.

Как и в случае гетероскедастичности остатков, постоянный множитель  $\sigma_\varepsilon^2$  можно опустить, используя для расчетов только матрицу

$$\mathbf{\Omega}_1 = \begin{pmatrix} 1 & \rho_1 & \rho_1^2 & \dots & \rho_1^{n-1} \\ \rho_1 & 1 & \rho_1 & \dots & \rho_1^{n-2} \\ \rho_1^2 & \rho_1 & 1 & \dots & \rho_1^{n-3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho_1^{n-1} & \rho_1^{n-2} & \rho_1^{n-3} & \dots & 1 \end{pmatrix}.$$

Параметр  $\rho$ , как правило, неизвестен. Его оценку производят, исследуя взаимосвязь остатков, полученных после применения к модели регрессии обычного МНК. Более подробно применение ОМНК в случае автокорреляции остатков будет рассмотрено в следующих главах. В данном параграфе приведем пример применения ОМНК к регрессии с гетероскедастичными остатками.

Как было выявлено ранее, случайные остатки в модели регрессии

$$y = 3226,44 + 12,45x_1 + 0,06x_2 + 0,31x_3 + e$$

гетероскедастичны, корень из дисперсии случайных остатков и значения независимой переменной  $x_3$  связаны линейно. То есть мы можем предположить, что  $\sigma_\varepsilon^2 = \sigma^2 x_3^2$ .

В этом случае можно применить метод взвешенных наименьших квадратов, разделив каждую переменную, входящую в уравнение регрессии, на соответствующее значение переменной  $x_3$ :

$$y_i^* = \frac{y_i}{x_{3i}}; \quad x_{1i}^* = \frac{x_{1i}}{x_{3i}}; \quad x_{2i}^* = \frac{x_{2i}}{x_{3i}}.$$

Преобразование переменных затронет также столбец множителей-единиц при свободном члене

$$x_{0i}^* = \frac{1}{x_{3i}}.$$

Модель с преобразованными переменными будет иметь вид

$$\frac{y}{x_3} = \alpha_0 \frac{1}{x_3} + \alpha_1 \frac{x_1}{x_3} + \alpha_2 \frac{x_2}{x_3} + \alpha_3 + \frac{\varepsilon}{x_3},$$

или

$$y^* = \alpha_0 x_0^* + \alpha_1 x_1^* + \alpha_2 x_2^* + \alpha_3 + \varepsilon^*.$$

Ниже приведена табл. 2.6, содержащая преобразованные переменные по данным примера 2.1.

**Таблица 2.6.** Преобразование переменных при применении ОМНК для модели регрессии поступления налогов от количества занятых, объема отгрузки в обрабатывающих производствах и производства энергии

Субъект РФ	$y^*$	$x_0^*$	$x_1^*$	$x_2^*$
Республика Ингушетия	1,940	0,0 013 643	0,146	0,363
Еврейская автономная область	1,240	0,0 004 902	0,040	1,404
Республика Тыва	1,300	0,0 004 943	0,050	0,213
Республика Алтай	2,351	0,0 008 503	0,074	1,044
Карачаево-Черкесская Республика	0,783	0,0 002 339	0,044	2,555
Республика Калмыкия	2,618	0,0 006 689	0,082	0,621
Республика Адыгея	1,467	0,0 003 333	0,062	4,188
Республика Северная Осетия — Алания	1,769	0,0 002 997	0,098	3,323
Магаданская область	0,838	0,0 001 204	0,012	0,299
Кабардино-Балкарская Республика	1,312	0,0 001 727	0,061	3,041
Республика Хакасия	0,525	0,0 000 567	0,014	2,248
Чукотский автономный округ	1,496	0,0 001 606	0,005	0,085
Республика Марий Эл	1,332	0,0 001 335	0,043	6,166
Псковская область	1,799	0,0 001 773	0,057	5,687
Чеченская Республика	1,656	0,0 001 621	0,058	0,094
Республика Карелия	0,773	0,0 000 681	0,023	2,721
Курганская область	0,996	0,0 000 827	0,032	3,168
Республика Мордовия	1,650	0,0 001 368	0,060	8,959

Окончание табл. 2.6

Субъект РФ	$y^*$	$x_0^*$	$x_1^*$	$x_2^*$
Костромская область	0,579	0,0 000 478	0,016	2,415
Камчатский край	1,025	0,0 000 786	0,015	0,884
Орловская область	1,363	0,0 001 040	0,039	3,963
Ивановская область	0,724	0,0 000 540	0,027	2,316
Республика Дагестан	1,127	0,0 000 795	0,088	1,673
Тамбовская область	1,494	0,0 001 050	0,052	5,013
Новгородская область	1,877	0,0 001 112	0,036	9,002
Республика Бурятия	1,438	0,0 000 798	0,031	2,367
Смоленская область	0,435	0,0 000 229	0,012	1,795
Курская область	0,457	0,0 000 229	0,012	1,538
Забайкальский край	1,494	0,0 000 731	0,035	0,578
Липецкая область	1,226	0,0 000 578	0,033	13,218
Ульяновская область	1,297	0,0 000 607	0,038	4,646
Пензенская область	1,776	0,0 000 829	0,053	5,912
Кировская область	1,030	0,0 000 479	0,033	3,651
Чувашская Республика	1,278	0,0 000 586	0,036	5,033
Астраханская область	2,167	0,0 000 949	0,045	3,283
Брянская область	2,242	0,0 000 951	0,054	5,437
Амурская область	1,435	0,0 000 606	0,025	0,994
Калужская область	2,315	0,0 000 964	0,051	15,601
Тульская область	1,131	0,0 000 410	0,031	7,468
Вологодская область	1,210	0,0 000 431	0,027	10,193
Алтайский край	1,202	0,0 000 403	0,045	4,644
Тверская область	0,717	0,0 000 222	0,015	2,294
Белгородская область	1,740	0,0 000 533	0,040	12,444
Владимирская область	1,626	0,0 000 498	0,034	7,110
Мурманская область	0,999	0,0 000 291	0,014	1,427
Воронежская область	0,920	0,0 000 255	0,027	3,200
Рязанская область	1,527	0,0 000 418	0,022	3,991
Калининградская область	2,407	0,0 000 648	0,030	9,565

По преобразованным данным была получена следующая модель:

$$y^* = 524,15x_0^* + 8,018x_1^* + 0,058x_2^* + 0,723 + e^* .$$

(t) (1,29) (2,18) (3,16) (5,41)



Все коэффициенты полученной модели значимы (табличное значение критерия Стьюдента равно 2,015 ( $\alpha = 0,05$ ;  $df = 44$ )). Напомним, что в данном уравнении параметр при переменной  $x_0^*$  является свободным членом исходной модели, а параметр, равный 0,723, — коэффициентом регрессии при переменной  $x_3$ . Таким образом, в исходном виде уравнение регрессии будет иметь вид

$$y = 524,15 + 8,018x_1 + 0,058x_2 + 0,723x_3 + e.$$

Сопоставив полученное уравнение с исходным

$$y = 3226,44 + 12,45x_1 + 0,06x_2 + 0,31x_3 + e,$$

можно отметить существенное изменение практически всех параметров.

Анализ случайных остатков нового уравнения был проведен по тем тестам, которые показали гетероскедастичность остатков, т.е. по тестам Глейзера и Бреуша — Пагана. Получены следующие результаты.

Тест Глейзера:

$$\text{при } k = 1 \quad |e| = 0,34 - 1,06 \cdot 10^{-6} x_3 + v, \quad t_{a_3} = -0,35;$$

$$\text{при } k = 0,5 \quad |e| = 0,37 - 0,0005\sqrt{x_3} + v, \quad t_{a_3} = -0,64;$$

Тест Бреуша — Пагана:

$$\hat{z} = 1,22 - 1,49 \cdot 10^{-5} x_3;$$

$$\chi^2 = \frac{SS_{\text{факт}}}{2} = \frac{1,32}{2} = 0,66.$$

Табличное значение критерия Стьюдента в данном случае равно 2,0129 ( $\alpha = 0,05$ ;  $df = 46$ ), критерия  $\chi^2$  равно 3,84 ( $\alpha = 0,05$ ;  $df = 1$ ), следовательно, по каждому из тестов нет оснований отвергнуть гипотезу о гомоскедастичности остатков.

## Контрольные вопросы и задания

1. Какой вид связей между показателями отражает уравнение регрессии? Запишите уравнение регрессии в матричном виде.
2. Почему конкретные значения параметров уравнения регрессии называют их оценками?
3. Назовите характеристики качества получаемых оценок параметров регрессии.
4. Какая модель называется классической нормальной линейной моделью?
5. Какие требования предъявляются к случайным остаткам модели регрессии?
6. Обоснуйте вывод системы нормальных уравнений.
7. Какие вы знаете показатели силы связи; чем они отличаются друг от друга; как интерпретируются?
8. Докажите утверждение «среднее значение стандартизованной переменной равно нулю».
9. На основе системы нормальных уравнений (2.8) обоснуйте выражение  $a_0 = \bar{y} - a_1x_1 - a_2x_2 - \dots - a_px_p$ .
10. Докажите несмещенность МНК-оценок. Какие предпосылки Гаусса — Маркова используются для доказательства несмещенности?
11. Обоснуйте выражение для расчета дисперсии оценок параметров регрессии.
12. Как рассчитать показатели тесноты связи в множественной линейной регрессии? Как они интерпретируются?
13. Для чего рассчитывается скорректированный коэффициент детерминации?
14. В чем состоят особенности расчета показателей тесноты связи в нелинейных функциях регрессии?
15. Что такое «значимость параметра»? Какой критерий применяется для ее оценки?
16. Кратко опишите алгоритм применения критерия Стьюдента.
17. Опишите процедуру оценки значимости уравнения регрессии. В чем состоит нулевая гипотеза? Какой критерий применяется?
18. Как связана значимость уравнения регрессии со значимостью коэффициента детерминации?

19. Приведите примеры ограничений на коэффициенты множественной регрессии. Как оценить значимость введения линейных ограничений?
20. Каковы особенности анализа нелинейных функций регрессии?
21. Назовите критерии выбора наилучшей функции регрессии. Какая проблема возникает при сопоставлении различных функций регрессии?
22. Опишите алгоритм выбора наилучшего преобразования зависимой переменной по тесту Бокса — Кокса.
23. Используя данные примера 2.1, с помощью теста Бокса — Кокса выберите наилучшее преобразование зависимой переменной при условии  $\lambda = -1; -0,5; 0,5$ . Сравните ответ с результатом, полученным для  $\lambda = 0; 1$  в параграфе 2.8.
24. Чем от теста Бокса — Кокса отличается преобразование Зарембки? Кратко опишите это преобразование.
25. В чем сущность метода максимального правдоподобия?
26. При каких условиях оценки параметров, полученные по ММП и МНК, совпадают?
27. Для какого показателя оценка, полученная при использовании ММП и МНК, различна?
28. Рассчитайте прогнозные значения поступления налогов для  $x_1 = 1200$  тыс. человек;  $x_2 = 237\ 000$  млн руб.;  $x_3 = 45\ 000$  млн руб.
29. Дайте определение мультиколлинеарности.
30. В каких случаях проявляется функциональная связь между независимыми переменными? В чем заключается решение проблемы мультиколлинеарности в этом случае?
31. Чем отличается интеркорреляция от мультиколлинеарности?
32. Каковы признаки наличия мультиколлинеарности в модели регрессии?
33. Назовите количественные критерии, позволяющие выявить мультиколлинеарность переменных.
34. Назовите методы устранения мультиколлинеарности.
35. В чем сущность метода ридж-регрессии?
36. В чем сущность метода главных компонент?
37. Поясните термины «гомоскедастичность» и «гетероскедастичность» случайных остатков.
38. Какие тесты применяют для выявления гетероскедастичности случайных остатков?

39. Как графически проанализировать дисперсии случайных остатков?
40. В чем сходство и в чем различие тестов Парка, Глейзера, Уайта, Бреуша — Пагана?
41. В чем состоит тест ранговой корреляции Спирмена?
42. Кратко опишите тест Гольдфельда — Квандта. В каких случаях его рекомендуется применять?
43. По данным примера 2.1 проверьте гетероскедастичность случайных остатков: а) по переменной  $x_1$ ; б) по переменной  $x_2$ . Какие тесты из рассмотренных в данной главе необходимо применить?
44. В каких случаях применяют обобщенный метод наименьших квадратов?
45. Дайте формулировку теоремы Айткена.
46. В чем заключается проблема применения ОМНК и каковы пути ее решения?
47. Что называют методом взвешенных наименьших квадратов? В чем он состоит?

---

## **ГЛАВА 3 Фиктивные переменные**

### **3.1. Особенности включения в модели регрессии неколичественных показателей**

В предыдущих главах при рассмотрении уравнений парной и множественной регрессии мы предполагали, что все изучаемые переменные являются количественными, т.е. принимают числовые значения. В общем случае это предположение может не выполняться. Значение зависимой переменной  $y$  часто складывается под влиянием факторов, описываемых как количественными, так и неколичественными показателями. Например, на величину надоя с одной фуражной коровы (количественный показатель) оказывает влияние качество кормов — высшее, 1 сорт, 2 сорт (неколичественный показатель). Для отражения влияния неколичественного показателя на результативный признак используют так называемые фиктивные переменные (*dummy variables*).

**Фиктивные переменные** — это переменные бинарного типа, имеющие всего два значения — единица и ноль:

$$z = \begin{cases} 1, \\ 0. \end{cases}$$

Если неколичественный показатель принимает всего два альтернативных значения, то необходимо ввести одну фиктивную переменную. Например, выдвинута гипотеза о разном уровне производительности труда для разных форм собственности — государственной и негосударственной. Тогда каждой форме собственности будет соответствовать одно значение фиктивной переменной. Например:

$$z = \begin{cases} 1, & \text{если собственность государственная,} \\ 0, & \text{если собственность негосударственная.} \end{cases}$$

Возможно и другое закрепление значений фиктивной переменной за категориями неколичественного показателя «форма собственности»:

$$z = \begin{cases} 1, & \text{если собственность негосударственная,} \\ 0, & \text{если собственность государственная.} \end{cases}$$

В общем случае необходимо ввести на единицу меньше фиктивных переменных, чем категорий неколичественного показателя. В частности, для приведенного выше примера с качеством кормов (три возможные категории) необходимо ввести две фиктивные переменные.

Фиктивная переменная равна единице, если неколичественный показатель принял закрепленное за этой переменной значение. Значение неколичественного показателя, не закрепленное ни за одной фиктивной переменной, называют **базой сравнения** значений зависимой переменной при разных состояниях неколичественного фактора. Ей соответствует равенство всех фиктивных переменных нулю:

$$z_1 = z_2 = \dots = z_{k-1} = 0,$$

где  $k$  — количество значений неколичественного показателя.

В качестве базового значения неколичественного показателя целесообразно выбирать такое, которое, как предполагается, обеспечивает минимальные или максимальные значения зависимой переменной. Для примера с качеством кормов

базовой категорией можно выбрать «2 сорт» как фактор, дающий минимальные надои коров:

$$z_1 = \begin{cases} 1, & \text{высшее качество кормов,} \\ 0, & \text{не высшее качество кормов (т.е. 1 или 2 сорт);} \end{cases}$$

$$z_2 = \begin{cases} 1, & \text{корма 1 сорта,} \\ 0, & \text{корма не 1 сорта (т.е. высшего или 2 сорта);} \end{cases}$$

$$z_1 = z_2 = 0, \text{ корма 2 сорта.}$$

Если в модели регрессии необходимо учесть несколько количественных показателей, то для каждого из них вводится свой набор фиктивных переменных. В этом случае фиктивные переменные нумеруются двойным индексом  $z_{ji}$ , где  $j$  — номер неколичественного показателя,  $i$  — номер значения (категории)  $j$ -го неколичественного показателя.

## 3.2. Спецификация моделей регрессии с фиктивными независимыми переменными

В общем случае модель с фиктивными переменными имеет вид

$$y = f(x_1, \dots, x_p, z_{11}, z_{12}, \dots, z_{21}, z_{22}, \dots, z_{j1}, z_{j2}, \dots, \varepsilon),$$

где  $y$  — зависимая переменная;  $x_1, x_2, \dots, x_p$  — количественные независимые переменные;  $z_{11}, z_{12}$  — фиктивные переменные, соответствующие категориям первого неколичественного показателя;  $z_{21}, z_{22}$  — фиктивные переменные, соответствующие категориям второго неколичественного показателя;  $z_{j1}, z_{j2}$  — фиктивные переменные, соответствующие категориям  $j$ -го неколичественного показателя;  $\varepsilon$  — случайный остаток.

Включение фиктивных переменных в модель расширяет круг вопросов, подлежащих решению на этапе спецификации модели.

Учет влияния неколичественного показателя, принимающего три и более значений, означает необходимость ввода в уравнение регрессии двух и более фиктивных переменных, что приводит к значительному увеличению количества параметров и должно быть подкреплено соответствующим объемом наблюдений (не менее семи на один параметр, не считая свободного члена).

При включении в уравнение регрессии фиктивных переменных возникает также вопрос о характере влияния количественных независимых переменных на результирующую переменную при различных значениях неколичественного показателя. Ниже будут рассмотрены различные варианты моделей регрессии с фиктивными переменными.

### 3.3. Модели регрессии с фиктивными переменными сдвига

Рассмотрим линейное уравнение регрессии. Для простоты возьмем в качестве объясняющих одну количественную переменную  $x_1$  и один неколичественный показатель с двумя альтернативными значениями, что требует включения в уравнение одной фиктивной переменной  $z_{11}$ :

$$y = \alpha_0 + \alpha_1 x_1 + \gamma_{11} z_{11} + \varepsilon. \quad (3.1)$$

Из уравнения (3.1) следует, что при  $z_{11} = 1$  значение зависимой переменной  $y$  равно

$$y = (\alpha_0 + \gamma_{11}) + \alpha_1 x_1 + \varepsilon, \quad (3.2)$$

а при  $z_{11} = 0$  равно

$$y = \alpha_0 + \alpha_1 x_1 + \varepsilon. \quad (3.3)$$

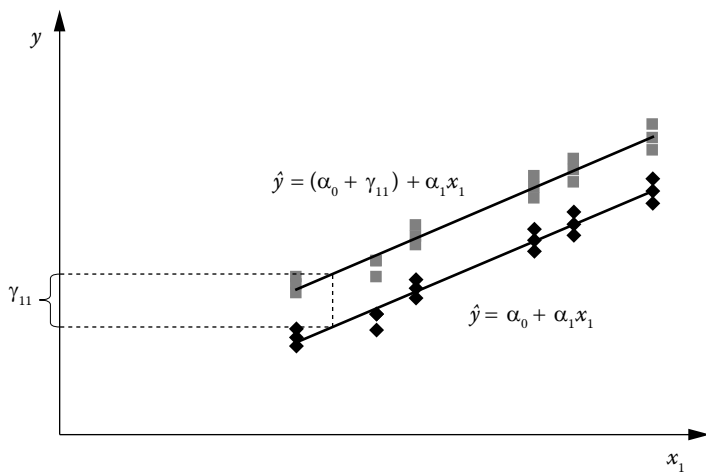
Сравнивая два полученных уравнения (3.2) и (3.3), видим, что они различаются величиной свободного члена. То есть для одного уровня неколичественного показателя значение зависимой переменной всегда в среднем будет на  $\gamma_{11}$  единиц выше или ниже, чем для альтернативного.



Графически эта ситуация соответствует двум параллельным прямым (рис. 3.1). Отметим, что коэффициент  $\alpha_1$  при количественной переменной  $x_1$  остается неизменным. То есть изменение  $x_1$  оказывает одинаковое влияние на изменение переменной  $y$  вне зависимости от значений неколичественного показателя.

Так как изменение значения фиктивной переменной в модели 3.1 приводит к изменению значения зависимой переменной на некую фиксированную среднюю величину, такую переменную еще называют *фиктивной переменной сдвига*. Изменение ее значения приводит к переходу от одной параллельной прямой к другой.

Рассмотрим данные примера 2.1. В качестве зависимой переменной по-прежнему выберем показатель «поступление налогов, сборов и иных обязательных платежей в консолидированный бюджет РФ (без поступлений ЕСН) в 2009 г. (млн руб.)»; в качестве количественной независимой переменной — «объем отгруженных товаров собственного производства, выполненных работ и оказанных услуг собственными силами по виду экономической деятельности «Обрабатывающие производства» в Российской Федерации в 2009 г. (млн руб.)».



**Рис. 3.1.** Сила влияния количественной переменной  $x_1$  на переменную  $y$  не зависит от значения неколичественного показателя

Предположим, что на сбор налогов влияет также географическое положение регионов. Разделим все субъекты РФ, имеющиеся в нашей выборке, на две группы:

- регионы Центрального и Северо-Западного федеральных округов. Для краткости будем называть эти регионы «центральными»;
- регионы прочих федеральных округов («прочие» регионы).

Так как имеется всего две категории (значения) нечисленного показателя «географическое положение», то необходимо ввести только одну фиктивную переменную. За базисную категорию примем «прочие» регионы:

$$z_{21} = \begin{cases} 1, & \text{центральные регионы,} \\ 0, & \text{прочие регионы.} \end{cases}$$

Исходные данные приведены в табл. 3.1.

**Таблица 3.1.** Поступление налогов и объем отгруженных товаров по субъектам РФ

Субъект РФ	Поступление налогов, млн руб.	Отгрузка в обрабатывающих производствах, млн руб.	Региональная принадлежность
	$y$	$x_2$	$z_{21}$
Республика Ингушетия	1422,20	266	0
Еврейская автономная область	2529,70	2865	0
Республика Тыва	2629,10	431	0
Республика Алтай	2764,30	1228	0
Карачаево-Черкесская Республика	3347,50	10 921	0
Республика Калмыкия	3914,20	928	0
Республика Адыгея	4400,80	12 565	0
Республика Северная Осетия — Алания	5904,00	11 088	0
Магаданская область	6956,70	2486	0
Кабардино-Балкарская Республика	7595,10	17 609	0
Республика Хакасия	9257,80	39 640	0
Чукотский автономный округ	9317,10	531	0
Республика Марий Эл	9978,80	46 180	0
Псковская область	10 144,80	32 074	1
Чеченская Республика	10 215,40	579	0

Окончание табл. 3.1

Субъект РФ	Поступление налогов, млн руб.	Отгрузка в обработки-вающих производствах, млн руб.	Региональная принадлежность
	$y$	$x_2$	$z_{21}$
Республика Карелия	11 349,50	39 962	1
Курганская область	12 046,90	38 308	0
Республика Мордовия	12 061,40	65 507	0
Костромская область	12 104,20	50 532	1
Камчатский край	13 042,40	11 245	0
Орловская область	13 104,30	38 089	1
Ивановская область	13 396,40	42 865	1
Республика Дагестан	14 170,30	21 031	0
Тамбовская область	14 227,00	47 738	1
Новгородская область	16 868,50	80 915	1
Республика Бурятия	18 019,40	29 660	0
Смоленская область	18 950,30	78 278	1
Курская область	19 995,50	67 241	1
Забайкальский край	20 445,60	7910	0
Липецкая область	21 220,80	228 812	1
Ульяновская область	21 360,00	76 523	0
Пензенская область	21 418,80	71 307	0
Кировская область	21 477,10	76 151	0
Чувашская Республика	21 816,30	85 926	0
Астраханская область	22 824,90	34 576	0
Брянская область	23 579,30	57 187	1
Амурская область	23 702,60	16 412	0
Калужская область	24 007,20	161 769	1
Тульская область	27 581,20	182 031	1
Вологодская область	28 057,50	236 267	1
Алтайский край	29 815,50	115 197	0
Тверская область	32 236,50	103 158	1
Белгородская область	32 657,40	233 608	1
Владимирская область	32 672,70	142 867	1
Мурманская область	34 351,10	49 081	1
Воронежская область	36 050,40	125 343	1
Рязанская область	36 544,30	95 522	1
Калининградская область	37 136,90	147 573	1

После применения МНК к уравнению (3.1) получим:

$$\hat{y} = 9404,1 + 0,098x_2 + 3713,6z_{21} + e. \quad (3.4)$$

(t) (6,26) (4,81) (1,43)

Табличное значение  $t$ -критерия Стьюдента составляет 2,014 при числе степеней свободы  $df = n - m - 1 = 48 - 2 - 1 = 45$  и уровне значимости  $\alpha = 0,05$ . Следовательно, коэффициент при фиктивной переменной  $z_{21}$  незначим. Как было показано выше, он равен величине, на которую значение зависимой переменной в среднем отличается для разных значений неколичественного показателя. То есть он показывает, насколько (в среднем) поступление налогов в центральных регионах больше, чем в прочих регионах. Для нашего примера эта величина составляет 3713,6 млн руб., однако мы не можем утверждать это с высокой степенью вероятности. Так как этот параметр незначим, возможно, региональные различия в поступлении налогов близки к нулю.

Если не принимать в расчет незначимость параметра при фиктивной переменной, уравнения регрессии для двух групп регионов, построенные на основе уравнения (3.4), будут иметь вид

$$z_{21} = 1: \quad \hat{y} = 13117,7 + 0,098x_2;$$

$$z_{21} = 0: \quad \hat{y} = 9404,1 + 0,098x_2.$$

Отметим, что если находить параметры этих уравнений регрессии отдельно, применяя к исходным данным МНК, получим не совсем тот же результат:

для центральных регионов:

$$\hat{y} = 15794,6 + 0,073x_2; \quad (3.5)$$

(t) (4,76) (2,79)

для прочих регионов:

$$\hat{y} = 6822,7 + 0,186x_2. \quad (3.6)$$

(t) (4,62) (5,45)

Свободные члены этих двух моделей различаются на величину 8971,9 что существенно больше, чем «прибавка» к свободному члену в модели с фиктивной переменной. Более того, все параметры в последних двух моделях статистически значимы. Можно заметить также, что в этих моделях существенно (более чем в два раза) отличаются коэффициенты

при количественной переменной  $x_2$ . То есть можно предположить, что влияние объема отгрузки в обрабатывающих производствах в центральных и прочих регионах на собираемость налогов различно. Это соображение приводит к необходимости построения моделей с **фиктивными переменными наклона**, рассмотренными в следующем параграфе.

### 3.4. Модели регрессии с фиктивными переменными наклона

Если влияние количественной независимой переменной на результативную меняется для разных значений неколичественного показателя, то можно записать

$$\hat{y} = \alpha_0 + \alpha_{11}x_1, \text{ если } z_{11} = 0; \quad (3.7a)$$

$$\hat{y} = \alpha_0 + \alpha_{12}x_1, \text{ если } z_{11} = 1; \quad (3.7б)$$

$$\alpha_{11} \neq \alpha_{12}.$$

В таком случае говорят, что имеют место структурные изменения в исследуемой зависимости количественных переменных. Для их учета в уравнении регрессии фиктивную переменную вводят как множитель при количественной переменной:

$$\hat{y} = \alpha_0 + \alpha_1x_1 + \varphi_{111}x_1z_{11}^1. \quad (3.7в)$$

Действительно, если рассмотреть это уравнение для  $z_{11} = 1$  и для  $z_{11} = 0$ , получим соответственно:

$$\text{при } z_{11} = 0 \quad \hat{y} = \alpha_0 + \alpha_1x_1;$$

$$\text{при } z_{11} = 1 \quad \hat{y} = \alpha_0 + (\alpha_1 + \varphi_{111})x_1.$$

Следовательно, коэффициент  $\alpha_{12}$  из модели (3.7б) будет равен  $(\alpha_{11} + \varphi_{111})$ .

Графически модель (3.7в) можно представить в виде двух прямых с разным углом наклона, отражающих зависимость

<sup>1</sup> Так как параметр  $\varphi$  объединяет две переменные —  $x_1$  и  $z_{11}$ , то он имеет тройной индекс —  $\varphi_{111}$ .

переменной  $y$  от количественной переменной  $x_1$  при разных значениях фиктивной переменной (рис. 3.2). Так как речь идет о фиктивной переменной, включение которой позволяет изменить угол наклона прямой, такую переменную называют **фиктивной переменной наклона**.

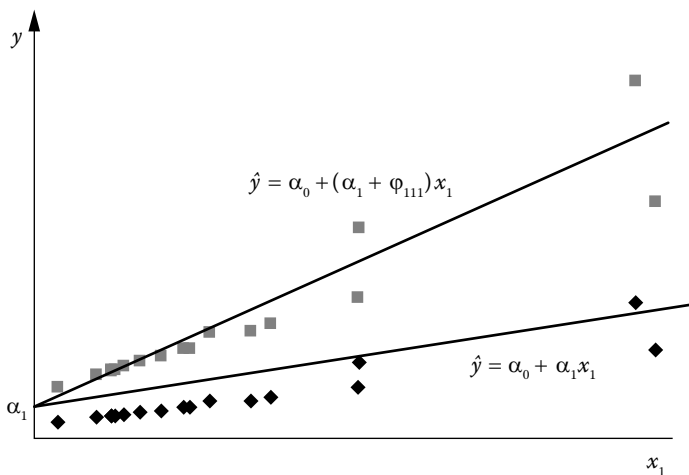
Параметр  $\alpha_1$  интерпретируется как сила влияния количественной переменной при значении неколичественного показателя, соответствующего  $z_{11} = 0$ , а параметр  $\varphi_{111}$  — как среднее изменение силы влияния количественной переменной при переходе от одной категории неколичественного показателя к другой (при переходе от  $z_{11} = 0$  к  $z_{11} = 1$ ).

Построим уравнение регрессии с фиктивной переменной наклона для нашего примера. После применения МНК к выражению (3.7в) получим выражение

$$\hat{y} = 9436,96 + 0,145x_2 - 0,028x_2z_{21}.$$

$$(t) \quad (5,87) \quad (3,53) \quad (-0,76)$$

Несмотря на большое различие коэффициентов регрессии по моделям, построенным по отдельным группам регионов (в предыдущем параграфе), мы получили, что коэффициент при фиктивной переменной  $z_{21}$  статистически незначим. То есть



**Рис. 3.2.** Сила влияния количественной переменной  $x_1$  на переменную  $y$  зависит от значения неколичественного показателя

различия между коэффициентами регрессии для данных по центральному и прочим регионам несущественны. В противном случае мы могли бы интерпретировать значение  $-0,028$  как более низкую «налоговую отдачу» с рубля отгруженной продукции по обрабатывающим производствам (на  $0,028$  руб. с 1 руб. отгрузки) в центральных регионах, чем в других регионах РФ.

Модели, полученные для разных значений фиктивной переменной  $z_{21}$ , будут иметь следующий вид:

для центральных регионов

$$\hat{y} = 9436,96 + 0,117x_2; \quad (3.8)$$

для прочих регионов

$$\hat{y} = 9436,96 + 0,145x_2. \quad (3.9)$$

Выражения (3.8) и (3.9) не соответствуют моделям, полученным по каждой группе регионов по отдельности, что приводит к выводу о необходимости дальнейшего совершенствования модели регрессии с фиктивными переменными. Это можно сделать с помощью общей модели.

### 3.5. Общий вид модели регрессии с фиктивными переменными

Модели, рассмотренные выше, можно считать частными случаями общей модели. Очевидно, что это будет модель, в которой для разных категорий неколичественного показателя будут разные параметры уравнения регрессии. Например, для простейшего случая с одной количественной  $x_1$  и одной фиктивной переменной  $z_{11}$  необходимо объединить в одну модель два уравнения регрессии:

$$\text{при } z_{11} = 0 \quad \hat{y} = \alpha_{01} + \alpha_{11}x_1; \quad (3.10)$$

$$\text{при } z_{11} = 1 \quad \hat{y} = \alpha_{02} + \alpha_{12}x_1.$$

Используем для этого модели (3.1) и (3.7в):

$$\hat{y} = \alpha_0 + \alpha_1 x_1 + \varphi_{111} x_1 z_{11} + \gamma_{11} z_{11}. \quad (3.10a)$$

Действительно,

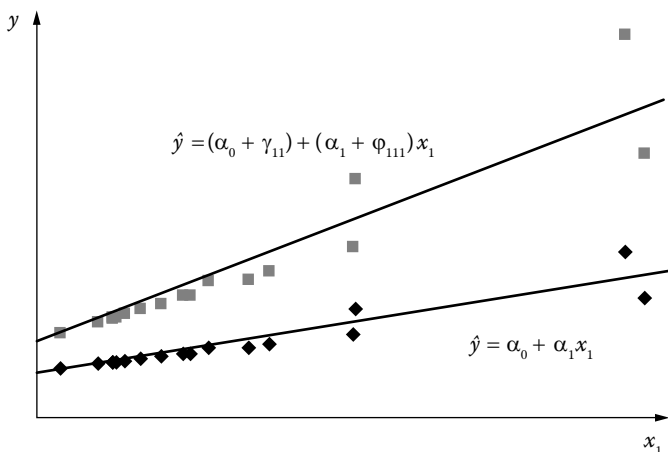
$$\text{при } z_{11} = 0 \quad \hat{y} = \alpha_0 + \alpha_1 x_1;$$

$$\text{при } z_{11} = 1 \quad \hat{y} = (\alpha_0 + \gamma_{11}) + (\alpha_1 + \varphi_{111}) x_1.$$

Параметры уравнения (3.10) можно выразить через параметры уравнения (3.10a) следующим образом:

$$\begin{aligned} \alpha_{01} &= \alpha_0; \quad \alpha_{11} = \alpha_1; \\ \alpha_{02} &= (\alpha_0 + \gamma_{11}); \quad \alpha_{12} = (\alpha_1 + \varphi_{111}). \end{aligned}$$

Графически модель (3.10a) может быть представлена в виде двух прямых, отражающих зависимость переменной  $y$  от количественной переменной при разных значениях фиктивной переменной. Параметры этих прямых различны, поэтому они имеют разный угол наклона (разный коэффициент при количественной переменной) и пересекают ось ординат в разных точках (разный свободный член, рис. 3.3).



**Рис. 3.3.** Влияние нечисленного показателя полностью изменяет параметры регрессии



Продолжим рассмотрение нашего примера. Найдем параметры модели (3.10а), используя МНК. Получим следующее уравнение регрессии:

$$\hat{y} = 6822,7 + 0,186x_2 - 0,112x_2z_{21} + 8971,9z_{21}.$$

(t) (3,82) (4,50) (-2,40) (2,71)

Все параметры статистически значимы (табличное значение  $t$ -критерия равно 2,01 при  $df = 48 - 3 - 1 = 44$  степенях свободы и уровне значимости  $\alpha = 0,05$ ), следовательно, общая форма модели с фиктивными переменными хорошо описывает исследуемые нами связи. Этот вывод подтверждает достаточно высокое значение коэффициента детерминации ( $R^2 = 0,59$ ), а также то, что уравнение регрессии статистически значимо ( $F = 7,79$  при  $F_{\text{табл}} = 2,82$ ).

Каковы причины статистической незначимости параметров при фиктивных переменных в моделях с фиктивными переменными сдвига и фиктивными переменными наклона, в то время как в общей модели все параметры стали значимыми? Ответ, очевидно, следует искать, сопоставляя значения параметров, полученных в общей модели.

Мы получили, что центральные регионы в среднем обеспечивают более высокое поступление налогов, чем прочие регионы (на статистически значимую величину, равную 8971,9 млн руб.). В то же время налоговая отдача от объема отгруженной продукции обрабатывающих производств в центральных регионах, наоборот, меньше (на статистически значимую величину, равную -0,112 млн руб. с каждого миллиона рублей отгруженной продукции). То есть имеет место разнонаправленное воздействие неколичественного показателя «географическое положение» на объем сбора налогов. Оно опосредовано объемом отгрузки (связь обратная) и прочими факторами, которые в совокупности дают прямую связь. В моделях с фиктивными переменными сдвига и фиктивными переменными наклона мы фактически накладываем ограничения на параметры общей модели в виде их равенства нулю:

- в модели с фиктивными переменными сдвига:  $\varphi_{221} = 0$ ;
- в модели с фиктивными переменными наклона:  $\gamma_{21} = 0$ .

Из-за неучета значимой переменной разнонаправленное влияние неколичественного показателя отражается в одном коэффициенте регрессии (при фиктивной переменной), что приводит

к его незначимости. Отметим, что если бы влияние географического положения на сбор налогов, опосредованное количественными переменными, имело бы одинаковую направленность, то все три модели с фиктивными переменными имели бы значимые параметры.

Проинтерпретируем параметры полученной общей модели с фиктивными переменными.

Коэффициент 0,186 при переменной  $x_2$  означает, что в регионах, не являющихся центральными ( $z_{21} = 0$ ), сбор налогов составляет в среднем 0,186 млн руб. с 1 млн руб. отгруженной продукции обрабатывающих производств (или 18,6 коп. с 1 руб. отгруженной продукции). В центральных регионах ( $z_{11} = 1$ ) отдача с рубля отгруженной продукции обрабатывающих производств меньше на 0,112 руб. (11,2 коп.), т.е. составляет  $0,186 - 0,112 = 0,074$  руб. Действие других факторов приводит к тому, что в целом по центральным регионам средний сбор налогов выше, чем в прочих, на величину 8971,9 млн руб.

Рассмотрим, как соотносятся параметры общей модели, полученной по данным примера, с параметрами уравнений регрессии, полученных с помощью МНК, по каждой совокупности регионов в отдельности (уравнения (3.5) и (3.6)). Оформим результаты сравнения в табл. 3.2.

**Таблица 3.2.** Модель с фиктивными переменными и модели, построенные по отдельным частям совокупности

Регионы	Общая модель с фиктивными переменными	Модели, построенные по каждой совокупности отдельно
Центральные регионы ( $z_{21} = 1$ )	$\hat{y} = 15794,7 + 0,074x_2$	$\hat{y} = 15794,6 + 0,073x_2$
Прочие регионы ( $z_{21} = 0$ )	$\hat{y} = 6822,7 + 0,186x_2$	$\hat{y} = 6822,7 + 0,186x_2$

Отметим, что параметры моделей для каждой группы регионов, полученные различным способом, совпадают. Небольшое расхождение для центральных регионов объясняется ошибками округления.

Модели (3.1), (3.7в), (3.10а) были построены для простейшего случая с двумя объясняющими переменными: одной

фиктивной и одной количественной. Запишем каждую из этих моделей в общей форме.

### 1. Модель с фиктивными переменными сдвига

$$\hat{y} = \alpha_0 + \alpha_1 x_1 + \dots + \alpha_p x_p + \gamma_{11} z_{11} + \gamma_{12} z_{12} + \dots + \gamma_{21} z_{21} + \gamma_{22} z_{22} + \dots + \gamma_{j1} z_{j1} + \gamma_{j2} z_{j2} + \dots, \quad (3.11)$$

где  $\alpha_0$  — свободный член;  $\alpha_1, \dots, \alpha_p$  — коэффициенты при количественных переменных  $x_1, \dots, x_p$ ;  $\gamma_{11}, \gamma_{12}, \dots$  — коэффициенты при фиктивных переменных  $z_{11}, z_{12}, \dots$ , введенных для первого нечисленного показателя;  $\gamma_{21}, \gamma_{22}, \dots$  — коэффициенты при фиктивных переменных  $z_{21}, z_{22}, \dots$ , введенных для второго нечисленного показателя;  $\gamma_{j1}, \gamma_{j2}, \dots$  — коэффициенты при фиктивных переменных  $z_{j1}, z_{j2}, \dots$ , введенных для  $j$ -го нечисленного показателя.

### 2. Модель с фиктивными переменными наклона

$$\begin{aligned} \hat{y} = & \alpha_0 + \alpha_1 x_1 + \dots + \alpha_p x_p + \\ & + \varphi_{111} x_1 z_{11} + \varphi_{112} x_1 z_{12} + \dots + \varphi_{1j1} x_1 z_{j1} + \varphi_{1j2} x_1 z_{j2} + \dots + \\ & + \varphi_{211} x_2 z_{11} + \varphi_{212} x_2 z_{12} + \dots + \varphi_{2j1} x_2 z_{j1} + \varphi_{2j2} x_2 z_{j2} + \dots + \\ & + \dots + \\ & + \varphi_{p11} x_p z_{11} + \varphi_{p12} x_p z_{12} + \dots + \varphi_{pj1} x_p z_{j1} + \varphi_{pj2} x_p z_{j2} + \dots \end{aligned} \quad (3.12)$$

Для удобства чтения уравнение (3.12) записано в несколько строк. В первой приведено обычное уравнение множественной регрессии с количественными независимыми переменными  $x_1 \dots x_p$ . В другие строки, начиная со второй, включены (с параметрами  $\varphi$ ) попарные произведения количественных и фиктивных переменных.

### 3. Общая форма модели с фиктивными переменными объединяет модели (3.11) и (3.12):

$$\begin{aligned} \hat{y} = & \alpha_0 + \alpha_1 x_1 + \dots + \alpha_p x_p + \\ & + \gamma_{11} z_{11} + \gamma_{12} z_{12} + \dots + \gamma_{21} z_{21} + \gamma_{22} z_{22} + \dots + \gamma_{j1} z_{j1} + \gamma_{j2} z_{j2} + \dots + \\ & + \varphi_{111} x_1 z_{11} + \varphi_{112} x_1 z_{12} + \dots + \varphi_{1j1} x_1 z_{j1} + \varphi_{1j2} x_1 z_{j2} + \dots + \\ & + \varphi_{211} x_2 z_{11} + \varphi_{212} x_2 z_{12} + \dots + \varphi_{2j1} x_2 z_{j1} + \varphi_{2j2} x_2 z_{j2} + \dots + \\ & + \dots + \\ & + \varphi_{p11} x_p z_{11} + \varphi_{p12} x_p z_{12} + \dots + \varphi_{pj1} x_p z_{j1} + \varphi_{pj2} x_p z_{j2} + \dots \end{aligned}$$

Частным случаем модели с фиктивными переменными является модель, не содержащая количественных независимых переменных. Рассмотрим этот вид моделей на примере модели с одной фиктивной переменной ( $z_{11}$ ). Очевидно, что здесь не может быть фиктивной переменной наклона, так как отсутствует количественная независимая переменная. Поэтому модель может быть единственного вида — с фиктивной переменной сдвига

$$\hat{y} = \alpha_0 + \gamma_{11}z_{11}.$$

Так как фиктивная переменная принимает всего два значения, может быть рассчитано всего два выровненных значения зависимой переменной: при  $z_{11} = 0$  и при  $z_{11} = 1$ . Это будут, соответственно, величины  $\alpha_0$  и  $(\alpha_0 + \gamma_{11})$ . Их можно рассматривать как средние значения переменной  $y$  по группам, образованным двумя значениями фиктивной переменной. Модели с фиктивными переменными, не содержащие количественных независимых переменных, называются **ANOVA-моделями (моделями дисперсионного анализа)**. Модели с фиктивными переменными, содержащие независимые количественные переменные, называются **ANCOVA-моделями (моделями ковариационного анализа)**.

Существенным недостатком моделей с фиктивными переменными является значительное число параметров, требующее соответствующего увеличения числа наблюдений. Например, если неколичественный показатель имеет три возможных значения, то для линейной модели с фиктивными переменными сдвига требуется два дополнительных слагаемых. Если же рассматриваются модели с несколькими количественными факторами и фиктивными переменными наклона, то число параметров еще больше возрастает. Как правило, объем наблюдений, находящихся в распоряжении исследователя, весьма незначителен и не может обеспечить статистическую значимость столь «длинной» модели.

Решением может быть объединение нескольких значений неколичественного показателя в одно, если эти значения не различаются по силе влияния на результат. Например, местоположение квартиры определяется ее адресом: городом, улицей, номером дома и номером квартиры. Это индивидуальная

характеристика рассматриваемого объекта. Она отражена в неколичественном показателе «местоположение», имеющем столько значений, сколько адресов. Такие индивидуальные характеристики можно укрупнить, объединив квартиры по принадлежности к номеру дома, еще крупнее — к номеру квартала, к району города или, наконец, разделить все квартиры на те, которые расположены в центре города, и находящиеся вне него.

Модели с фиктивными переменными могут применяться для анализа временных рядов, в частности для изучения периодических колебаний. Если рассматриваются сезонные колебания продаж (например, меховой одежды), можно построить модель зависимости объема продаж от квартала, месяца и т.п. Чем шире будет интервал времени, рассматриваемый в качестве фактора, тем короче модель, что повышает ее статистическую надежность. Однако укрупнение временных интервалов ведет к значительному усреднению характеристик влияния неколичественного показателя, что понижает аналитическую ценность модели.

### 3.6. Исследование структурных изменений с помощью теста Чоу

Как было показано выше, при разных значениях фиктивной переменной получаются разные уравнения регрессии. Целесообразность применения двух уравнений регрессии вместо одного можно оценить, не прибегая к вводу фиктивных переменных. Для этого используется *тест Чоу*.

Пусть имеется  $n$  наблюдений, позволяющих охарактеризовать зависимость результативной переменной от одной или нескольких количественных независимых переменных. У исследователя есть основания предполагать, что совокупность неоднородна с точки зрения числовых характеристик этой зависимости. Предполагается также, что однородность может быть достигнута в случае, если разбить эту совокупность по определенному критерию на две части.

Для проверки выдвинутого предположения находят параметры трех уравнений регрессии. Первое уравнение строится для всей совокупности наблюдений, второе и третье — для соответствующих выделенных подмножеств совокупности

наблюдений. Для каждого из этих уравнений находят остаточную сумму квадратов  $SS_e$ :

$$SS_e = \sum (y - \hat{y})^2 = \sum e^2.$$

Обозначим остаточную сумму квадратов, рассчитанную по общему уравнению регрессии, через  $SS_{e0}$ , по уравнениям регрессий для подмножеств наблюдений — через  $SS_{e1}$  и  $SS_{e2}$ .

Тогда равенство  $SS_{e0} = SS_{e1} + SS_{e2}$  выполняется, если параметры всех трех уравнений регрессии равны.

В противном случае  $SS_{e0} > SS_{e1} + SS_{e2}$ . Чем больше разница между двумя частями этого неравенства, тем больше различия между двумя подмножествами с точки зрения параметров уравнений регрессии. Существенность различий проверяют с помощью  $F$ -критерия.

Фактическое значение  $F$ -критерия находят по формуле

$$F = \frac{[SS_{e0} - (SS_{e1} + SS_{e2})]}{(SS_{e1} + SS_{e2})} \cdot \frac{n - m_1 - m_2 - 2}{m_1 + m_2 + 1 - m}, \quad (3.13)$$

где  $m_1$  и  $m_2$  — количество параметров (без свободного члена) в уравнениях, построенных по подмножествам;  $m$  — количество параметров (без свободного члена) для уравнения, построенного по всей совокупности;  $n$  — число наблюдений по всей совокупности.

Табличное значение  $F$ -критерия находят для степеней свободы  $df_1 = m_1 + m_2 + 1 - m$  и  $df_2 = n - m_1 - m_2 - 2$ . Если фактическое значение окажется больше табличного, то имеют место структурные сдвиги и целесообразно строить уравнение регрессии с соответствующей фиктивной переменной.

Часто при применении теста Чоу рассматривают уравнения одинаковой структуры, т.е.  $m_1 = m_2 = m$ . Тогда формула (3.13) имеет вид

$$F = \frac{[SS_{e0} - (SS_{e1} + SS_{e2})]}{(SS_{e1} + SS_{e2})} \cdot \frac{n - 2m - 2}{m + 1}.$$

Табличное значение  $F$ -критерия в этом случае находят для степеней свободы  $df_1 = m + 1$  и  $df_2 = n - 2m - 2$ .

Рассмотрим применение теста Чоу на нашем примере.

Нам необходимо найти параметры уравнения

$$\hat{y} = \alpha_0 + \alpha_2 x_2$$

для следующих массивов данных:

- для всех данных ( $n = 48$ );
- для данных по центральным регионам ( $n = 21$ );
- для данных по прочим регионам ( $n = 27$ ).

Применяя МНК, получим уравнения регрессии и значения сумм квадратов остатков.

По всем регионам:

$$\hat{y} = 9910,39 + 0,116x_2; \quad SS_{e0} = 2\,394\,978\,286.$$

По центральным регионам:

$$\hat{y} = 15\,794,6 + 0,073x_2; \quad SS_{e1} = 1\,239\,199\,040.$$

По прочим регионам:

$$\hat{y} = 6822,7 + 0,186x_2; \quad SS_{e2} = 786\,754\,752,4.$$

Фактическое значение  $F$ -критерия равно

$$F = \frac{2\,394\,978\,286 - (1\,239\,199\,040 + 786\,754\,752,4)}{(1\,239\,199\,040 + 786\,754\,752,4)} \cdot \frac{48 - 2 - 2}{1 + 1} \approx 4,01.$$

Табличное значение  $F$ -критерия равно 3,21 (при  $\alpha = 0,05$  и  $df_1 = 1 + 1 = 2$  и  $df_2 = 48 - 2 - 2 = 44$  степенях свободы).

Так как фактическое значение  $F$ -критерия больше табличного, следует признать существенность различия характеристик налоговой отдачи с 1 руб. отгруженной продукции обрабатывающих производств от местоположения региона. Тот же вывод мы получили, исследуя общую модель с фиктивными переменными.

Тест Чоу можно использовать также, если исследуется зависимость уровня ряда от времени, т.е. если строится уравнение тренда вида  $\hat{y} = f(t)$ . В этом случае деление совокупности наблюдений на две части производится относительно определенного момента времени, в который, по мнению исследователя, произошли какие-либо структурные изменения.

## Контрольные вопросы и задания

1. Можно ли учесть в уравнении регрессии неколичественные показатели? Каким образом?
2. Дайте определение фиктивной переменной.
3. Сколько фиктивных переменных нужно ввести, если имеются два неколичественных фактора, причем один из них имеет три возможных значения, а другой — два?
4. Как интерпретируется коэффициент регрессии при фиктивной переменной сдвига?
5. Как интерпретируется коэффициент регрессии при фиктивной переменной наклона?
6. Каков общий вид модели регрессии с одной количественной и одной фиктивной переменной?
7. В какой модели (моделях) с фиктивными переменными учитывается различие влияния количественных переменных на результативную переменную при разных значениях неколичественного фактора?
8. Назовите достоинства и недостатки моделей с фиктивными переменными.
9. Пусть имеется уравнение регрессии с одной количественной переменной и тремя фиктивными переменными, описывающими один неколичественный показатель. Сколько возможных значений у неколичественного показателя? Как на основе заданного уравнения регрессии найти уравнения парной регрессии, содержащие только количественную переменную? Сколько будет таких уравнений и почему?
10. Какова область применения теста Чоу?
11. Какие показатели сравниваются между собой по тесту Чоу? Какой статистический критерий в этом случае используется?
12. Опишите методику применения теста Чоу.



# **Системы эконометрических уравнений**

Изученные в предыдущих главах методы работы с регрессиями опирались на предположение о том, что исследуемый экономический показатель может быть описан некоторой регрессионной зависимостью от других показателей. В настоящей главе мы рассмотрим ситуации, возникающие, как правило, при изучении сложных экономических систем, когда изучаемое экономическое явление или процесс описываются несколькими переменными, некоторые из которых могут быть взаимозависимыми.

Целью настоящей главы является знакомство с различными видами систем эконометрических уравнений, изучение методов идентификации и оценивания этих систем.

## **4.1. Виды систем эконометрических уравнений и методы их оценивания**

При изучении систем эконометрических уравнений принято различать эндогенные, экзогенные, лаговые и предопределенные переменные.

*Эндогенными* переменными называют зависимые переменные, стоящие, как правило, в левой части системы; их число должно быть равно числу уравнений в системе (требование полноты).

**Экзогенными** переменными называют переменные, влияющие на эндогенные, но не зависящие от них и не коррелирующие с ошибками.

**Лаговыми** называют переменные, принимающие значения не текущего, а предыдущих периодов. Например, если в системе есть переменные  $y_t, y_{t-1}, y_{t-2}$ , то они обозначают переменную  $y$  в моменты или интервалы времени  $t, t-1$  (т.е. в момент или интервал, непосредственно предшествующий моменту или интервалу  $t$ ) и  $t-2$  (т.е. в момент или интервал, непосредственно предшествующий моменту или интервалу  $t-1$ ) соответственно.

**Предопределенными** называют экзогенные (текущие и лаговые) и лаговые эндогенные переменные, т.е. переменные, значения которых к рассматриваемому моменту или интервалу времени определены (либо потому, что они экзогенные, т.е. используются как заданные извне, либо потому, что они лаговые, т.е. используются уже известные значения за предыдущие периоды). Предполагается, что предопределенные переменные, в отличие от эндогенных, не коррелируют с ошибкой.

Простейшая система уравнений с  $n$  эндогенными переменными  $y$  и  $m$  экзогенными переменными  $x$  может быть записана в следующем виде (некоторые из коэффициентов могут быть равны нулю)<sup>1</sup>:

$$\begin{cases} y_1 = a_{10} + a_{11}x_1 + a_{12}x_2 + \dots + a_{1m}x_m + \varepsilon_1, \\ y_2 = a_{20} + a_{21}x_1 + a_{22}x_2 + \dots + a_{2m}x_m + \varepsilon_2, \\ \dots \\ y_n = a_{n0} + a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nm}x_m + \varepsilon_n. \end{cases}$$

В том случае, если ошибки  $\varepsilon_1, \dots, \varepsilon_n$  не коррелируют друг с другом, эта система называется **системой независимых эконометрических уравнений**: эндогенные переменные не зависят друг от друга. Такие уравнения могут оцениваться по отдельности методами, изученными в предыдущих главах (при выполнении традиционных предпосылок, обычным МНК).

Если же ошибки  $\varepsilon_1, \dots, \varepsilon_n$  коррелируют друг с другом, то это случай **уравнений, кажущихся несвязанными**; он будет рассмотрен ниже.

<sup>1</sup> Мы ограничиваемся рассмотрением только линейных систем.

Следующий пример системы эконометрических уравнений — случай, когда зависимая переменная у одного уравнения выступает в качестве фактора в другом уравнении (находится в его правой части), но зависимая переменная второго уравнения не является фактором первого. Это — **система рекурсивных уравнений**. Например:

$$\begin{cases} y_1 = a_{10} + a_{11}x_1 + a_{12}x_2 + \dots + a_{1m}x_m + \varepsilon_1, \\ y_2 = b_{21}y_1 + a_{20} + a_{21}x_1 + a_{22}x_2 + \dots + a_{2m}x_m + \varepsilon_2, \\ y_3 = b_{31}y_1 + b_{32}y_2 + a_{30} + a_{31}x_1 + a_{32}x_2 + \dots + a_{3m}x_m + \varepsilon_3, \\ \dots \\ y_n = b_{n1}y_1 + b_{n2}y_2 + \dots + b_{n, n-1}y_{n-1} + a_{n0} + a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nm}x_m + \varepsilon_n. \end{cases}$$

Как и в системе независимых уравнений, каждое уравнение системы рекурсивных уравнений может оцениваться по отдельности обычным МНК (при выполнении предпосылок теоремы Гаусса — Маркова).

Наконец, **системы одновременных уравнений** содержат взаимозависимые эндогенные переменные. В общем случае это системы вида

$$\begin{cases} y_1 = b_{12}y_2 + b_{13}y_3 + \dots + b_{1n}y_n + a_{10} + a_{11}x_1 + a_{12}x_2 + \dots + a_{1m}x_m + \varepsilon_1, \\ y_2 = b_{21}y_1 + b_{23}y_3 + \dots + b_{2n}y_n + a_{20} + a_{21}x_1 + a_{22}x_2 + \dots + a_{2m}x_m + \varepsilon_2, \\ \dots \\ y_n = b_{n1}y_1 + b_{n2}y_2 + \dots + b_{n, n-1}y_{n-1} + a_{n0} + a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nm}x_m + \varepsilon_n. \end{cases} \quad (4.1)$$

Данная система называется также **структурной формой** модели. В отличие от независимых и рекурсивных систем, система одновременных уравнений в структурной форме содержит взаимозависимые уравнения, которые не могут быть оценены по отдельности обычным МНК, так как в этом случае корреляция эндогенных переменных с ошибками приводит к смещенным и несостоятельным оценкам. При отсутствии корреляции (и автокорреляции) ошибок в разных уравнениях (и при условии идентифицируемости уравнений, о чем будет идти речь ниже) системы одновременных уравнений оцениваются двухшаговым методом наименьших квадратов (ДМНК), и в некоторых случаях возможно применение косвенного МНК (КМНК). Если же ошибки в разных уравнениях коррелируют между

собой, то для повышения эффективности оценок применяется трехшаговый МНК (ТМНК). В случае же наличия автокорреляции ошибок и присутствия в модели лагов эндогенных переменных для обеспечения состоятельности оценок требуется использовать метод Фэа<sup>1</sup>.

## 4.2. Системы одновременных уравнений

### 4.2.1. Идентификация и оценивание

Оценивание системы одновременных уравнений обычным МНК дает смещенные и несостоятельные оценки по той причине, что эндогенные регрессоры коррелируют с ошибками. Однако если записать систему в *приведенной форме*, в которой все эндогенные переменные объясняются только predetermined переменными, то эта система может быть оценена состоятельным образом обычным МНК, поскольку predetermined переменные не коррелируют с ошибкой.

В системе (4.1) predetermined переменные представлены только экзогенными (лагов эндогенных переменных нет), поэтому ее приведенная форма имеет вид

$$\begin{cases} y_1 = A_1 + B_{11}x_1 + B_{12}x_2 + \dots + B_{1m}x_m + u_1, \\ y_2 = A_2 + B_{21}x_1 + B_{22}x_2 + \dots + B_{2m}x_m + u_2, \\ \dots \\ y_n = A_n + B_{n1}x_1 + B_{n2}x_2 + \dots + B_{nm}x_m + u_n. \end{cases}$$

где  $u_1, u_2, \dots, u_n$  — соответствующие ошибки.

Состоятельные оценки коэффициентов приведенной формы в некоторых случаях позволяют получить состоятельные оценки коэффициентов уравнений структурной формы. Такие уравнения называются *идентифицируемыми*.

Если же какое-то уравнение структурной формы содержит слишком много коэффициентов по сравнению с уравнением приведенной формы, то оно может оказаться *неидентифицируемым*

<sup>1</sup> См.: Fair R. C. «The Estimation of Simultaneous Equation Models with Lagged Endogenous Variables and First Order Serially Correlated Errors» *Econometrica*. Vol. 38. May 1970. Pp. 507—516.

(оценки его коэффициентов не могут быть получены из оценок коэффициентов приведенной формы).

Модель называется точно идентифицируемой, если все ее структурные коэффициенты определяются однозначно по коэффициентам приведенной формы. В этом случае число оцениваемых коэффициентов в структурной и приведенной формах одинаково. Здесь для оценки коэффициентов структурной модели применяется косвенный МНК (КМНК), состоящий в том, что МНК оценивают приведенную систему, а затем выражают оценки коэффициентов структурной формы через оценки коэффициентов приведенной формы.

Соответствие между коэффициентами приведенной и структурной форм модели может быть также неединственным. Уравнение оказывается сверхидентифицируемым, если число его приведенных коэффициентов больше числа структурных. В этом случае на основе коэффициентов приведенной формы можно с помощью КМНК получить более одного значения структурного коэффициента, и поскольку предпочтение одного из вариантов этих значений приводит к потере информации, оценки КМНК оказываются неэффективными, и здесь необходимо использовать двухшаговый МНК (ДМНК). Оценивание уравнения ДМНК состоит в следующем.

1. Обычным МНК оценивается приведенная форма модели, и на ее основе рассчитывают теоретические значения эндогенных переменных, стоящих в правой части оцениваемого уравнения структурной формы.
2. Уравнение в структурной форме оценивается обычным МНК, но вместо значений эндогенных переменных, стоящих в правой части, используются их теоретические значения.

Если все уравнения системы сверхидентифицируемы, то для оценки структурных коэффициентов каждого уравнения используется ДМНК. Если же в системе есть и точно идентифицируемые уравнения, то структурные коэффициенты по ним можно найти и с помощью КМНК. Для точно идентифицируемых уравнений результаты, полученные КМНК и ДМНК, совпадают.

**(Точно) идентифицируемой** называется модель, в которой каждое уравнение (точно) идентифицируемо. **Сверхидентифицируемой** называется модель, в которой

каждое уравнение идентифицируемо и хотя бы одно уравнение сверхидентифицируемо. Модель называется *неидентифицируемой*, если неидентифицируемо хотя бы одно ее уравнение. В этом случае число приведенных коэффициентов оказывается меньше числа структурных, из-за чего структурные коэффициенты не могут быть получены из приведенных.

Существуют необходимые и достаточные условия идентифицируемости уравнений; они проверяются для каждого уравнения структурной модели.

Необходимое условие идентификации (сверхидентификации) уравнения, или порядковое условие, заключается в том, что число предопределенных переменных, отсутствующих в уравнении, плюс один должно быть равно (больше) числа эндогенных переменных, присутствующих в уравнении. При невыполнении этого условия уравнение считается неидентифицируемым (хотя в некоторых случаях отдельные его коэффициенты могут быть идентифицированы). Следует заметить, что данное условие эквивалентно тому, что для идентифицируемости (сверхидентифицируемости) уравнения число его коэффициентов в приведенной форме должно быть равно (больше) числу его коэффициентов в структурной форме.

Необходимое и достаточное условие идентифицируемости требует вычисления ранга подматрицы матрицы коэффициентов приведенной формы. Однако если полагать эту подматрицу матрицей полного ранга (что, как правило, и имеет место на практике), то это условие будет эквивалентно порядковому.

В качестве примера рассмотрим модель денежного и товарного рынков<sup>1</sup>:

$$R_t = a_1 + b_{12}Y_t + b_{14}M_t + e_1 \text{ (функция денежного рынка);}$$

$$Y_t = a_2 + b_{21}R_t + b_{23}I_t + b_{25}G_t + e_2 \text{ (функция товарного рынка);}$$

$$I_t = a_3 + b_{31}R_t + e_3 \text{ (функция инвестиций),}$$

где  $R$  — процентные ставки;  $Y$  — реальный ВВП;  $M$  — денежная масса;  $I$  — внутренние инвестиции;  $G$  — реальные государственные расходы.

<sup>1</sup> Задача взята из кн.: Практикум по эконометрике (+CD) : учеб. пособие / И. И. Елисеева, С. В. Курьшева, Н. М. Гордеенко [и др.] ; под ред. И. И. Елисеевой. 2-е изд., перераб. и доп. М. : Финансы и статистика, 2007.

Здесь все три анализируемые переменные  $R$ ,  $Y$ ,  $I$  зависят друг от друга. Поэтому это система одновременных уравнений, и каждое из этих уравнений не может оцениваться по отдельности обычным МНК, так как оценки при этом получают-ся смещенные и несостоятельные.

В системе три эндогенные переменные  $R$ ,  $Y$ ,  $I$  и две predetermined переменные  $M$ ,  $G$ . Предetermined переменные представлены только экзогенными переменными, лаговых эндогенных переменных нет.

Проверим выполнение необходимого условия идентифицируемости.

Номер уравнения	Число predetermined переменных, отсутствующих в уравнении, $D$	Число эндогенных переменных, присутствующих в уравнении, $H$	$D + 1 \neq H$	Идентифицируемость
1	1	2	$D + 1 = H$	Точная
2	1	3	$D + 1 < H$	Нет
3	2	2	$D + 1 > H$	Сверх

Таким образом, выявляем точную идентифицируемость первого уравнения, неидентифицируемость второго и сверхидентифицируемость третьего.

Тот же результат получим, взглянув на приведенную форму системы

$$\begin{cases} R_t = A_1 + B_{11}M_t + B_{12}G_t + u_1, \\ Y_t = A_2 + B_{21}M_t + B_{22}G_t + u_2, \\ I_t = A_3 + B_{31}M_t + B_{32}G_t + u_3. \end{cases}$$

По трем коэффициентам уравнения приведенной формы (например, для первого уравнения это  $A_1$ ,  $B_{11}$ ,  $B_{12}$ ) требуется оценить в первом уравнении три параметра  $a_1$ ,  $b_{12}$ ,  $b_{14}$  (уравнение точно идентифицируемо); во втором — четыре параметра  $a_2$ ,  $b_{21}$ ,  $b_{23}$ ,  $b_{25}$  ( $4 > 3$ , следовательно, уравнение неидентифицируемо); в третьем — два параметра  $a_3$ ,  $b_{31}$  ( $2 < 3$ , следовательно, уравнение сверхидентифицируемо).

Таким образом, система в целом неидентифицируема, так как в ней имеется неидентифицируемое уравнение. Однако параметры первого и третьего уравнений могут быть оценены КМНК/ДМНК и ДМНК соответственно.

Применим КМНК к первому уравнению. Оценим обычным МНК параметры каждого уравнения приведенной выше системы.

Для этого нам нужны значения  $R_t$ ,  $Y_t$ ,  $I_t$ ,  $M_t$ ,  $G_t$ , где  $t = 1, 2, \dots, n$ . Воспользуемся условными данными из следующей таблицы<sup>1</sup>.

$t$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
$R_t$	6	5	5	6	5	5	4	5	6	5	5	4	4	5	5
$Y_t$	175	179	185	179	183	188	197	181	158	173	176	194	194	178	185
$I_t$	14	15	15	15	14	15	17	14	12	14	15	16	15	15	15
$M_t$	10	9	10	11	11	12	8	10	11	9	9	8	8	9	9
$G_t$	99	93	102	95	106	104	96	103	94	94	93	101	107	92	101

$t$	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30
$R_t$	5	5	7	4	5	5	6	5	4	6	5	5	5	4	5
$Y_t$	186	186	160	180	202	196	196	187	198	180	177	179	171	191	179
$I_t$	15	15	12	15	16	17	17	16	17	14	15	16	13	16	14
$M_t$	8	9	10	10	10	10	10	11	10	10	9	12	11	12	9
$G_t$	100	101	96	92	111	99	99	98	98	104	93	88	98	99	98

Оценим коэффициенты уравнений приведенной формы обычным МНК:

$$\begin{cases} R_t = 4,50 + 0,15M_t - 0,01G_t + \hat{u}_1, \\ Y_t = 96,77 - 1,44M_t + 1,02G_t + \hat{u}_2, \\ I_t = 15,10 - 0,13M_t + 0,01G_t + \hat{u}_3. \end{cases}$$

Теперь мы должны получить оценки коэффициентов первого уравнения. Оно точно идентифицируемо, поэтому решение должно быть единственным. Здесь есть несколько эквивалентных способов решения.

**Первый способ.** Заметим, что в первом уравнении приведенной формы присутствуют переменные  $M$  и  $G$ , тогда как в структурной форме —  $Y$  и  $M$ . Следовательно, нужно избавиться от переменной  $G$  в приведенном уравнении с помощью переменной  $Y$ . Выразим переменную  $G$  из второго уравнения приведенной формы и подставим в первое уравнение. Получим выражения

<sup>1</sup> Данные сгенерированы с помощью надстройки «Пакет анализа» MS Excel (см. в главном меню Сервис => Анализ данных => Генерация случайных чисел).



$$R_t = \left( A_1 - A_2 \frac{B_{12}}{B_{22}} \right) + \frac{B_{12}}{B_{22}} Y_t + \left( B_{11} - \frac{B_{12}B_{21}}{B_{22}} \right) M_t + u_1^* \quad \Rightarrow$$

$$a_1 = A_1 - A_2 \frac{B_{12}}{B_{22}}; \quad b_{12} = \frac{B_{12}}{B_{22}}; \quad b_{14} = B_{11} - \frac{B_{12}B_{21}}{B_{22}}.$$

**Второй способ.** Вместо каждой эндогенной переменной из правой части первого структурного уравнения подставим ее выражение в приведенной форме; затем приравняем правые части приведенного и структурного уравнений:

$$A_1 + B_{11}M_t + B_{12}G_t + u_1 = a_1 + b_{12}(A_2 + B_{21}M_t + B_{22}G_t + u_2) + b_{14}M_t + e_1.$$

Далее нужно в правой части раскрыть скобки, сгруппировать коэффициенты при переменных  $M$  и  $G$  и свободные члены, а затем приравнять соответствующие коэффициенты правой и левой частей. Получится система из трех линейных уравнений с тремя неизвестными. Ответ будет тот же:

$$a_1 = A_1 - A_2 \frac{B_{12}}{B_{22}}; \quad b_{12} = \frac{B_{12}}{B_{22}}; \quad b_{14} = B_{11} - \frac{B_{12}B_{21}}{B_{22}}.$$

**Третий способ.** Воспользоваться ДМНК. Результат:

$$a_1 = 5,4; \quad b_{12} = -0,0^{12}; \quad b_{14} = 0,1.$$

Таким образом, мы идентифицировали первое уравнение

$$R_t = 5,4 - 0,0Y_t + 0,1M_t + \text{остаток}.$$

Попробуем теперь идентифицировать второе уравнение (мы заранее знаем, что это не получится). Воспользуемся вторым из перечисленных выше способов:

$$A_2 + B_{21}M_t + B_{22}G_t + u_2 = a_2 + b_{21}(A_1 + B_{11}M_t + B_{12}G_t + u_1) + b_{23}(A_3 + B_{31}M_t + B_{32}G_t + u_3) + b_{25}G_t + e_2.$$

Далее нужно в правой части раскрыть скобки, сгруппировать коэффициенты при переменных  $M$  и  $G$  и свободные члены,

<sup>1</sup> Заметим, что неправомерно считать этот коэффициент равным нулю, не протестировав его значимость.

а затем приравнять соответствующие коэффициенты правой и левой частей. Получится система из трех линейных уравнений с четырьмя неизвестными  $a_2, b_{21}, b_{23}, b_{25}$ :

$$\begin{cases} a_2 + A_1 b_{21} + A_3 b_{23} = A_2, \\ B_{11} b_{21} + B_{31} b_{23} = B_{21}, \\ B_{12} b_{21} + B_{32} b_{23} + b_{25} = B_{22}. \end{cases}$$

По основной теореме линейной алгебры эта система совместна (ранг основной матрицы равен рангу расширенной матрицы, т.е. трем), а так как в ней четыре неизвестных и только три уравнения, то она имеет бесчисленное множество решений. Мы не можем из этого бесчисленного множества решений предпочесть какое-то одно решение всем остальным. Поэтому второе уравнение считается неидентифицируемым.

Если мы попробуем идентифицировать третье уравнение (сверхидентифицируемое) КМНК, т.е. выразить его коэффициенты через коэффициенты приведенной формы, то получим выражение

$$A_3 + B_{31}M_t + B_{32}G_t + u_3 = a_3 + b_{31}(A_1 + B_{11}M_t + B_{12}G_t + u_1) + e_3.$$

После группировки получим систему из трех линейных уравнений с двумя неизвестными  $a_3, b_{31}$ :

$$\begin{cases} a_3 + A_1 b_{31} = A_3, \\ B_{11} b_{31} = B_{31}, \\ B_{12} b_{31} = B_{32}. \end{cases}$$

По основной теореме линейной алгебры эта система несовместна (ранг основной матрицы меньше ранга расширенной матрицы):  $b_{31} = B_{31}/B_{11} = -0,9$  и, в то же время,  $b_{31} = B_{32}/B_{12} = -1,0$ . Так что для параметра  $b_{31}$  мы имеем две оценки, для каждой из которых будет также свое значение параметра  $a_3$  (соответственно 19,0 и 19,6). Это и есть сверхидентифицируемость. Поэтому здесь применяется не КМНК, а ДМНК, заключающийся в том, что все эндогенные переменные, стоящие в правой части сверхидентифицируемого структурного уравнения, заменяются их теоретическими значениями, полученными из соответствующих уравнений приведенной формы; затем структурное уравнение оценивается обычным МНК.

В нашем случае необходимы теоретические значения переменной  $R$ . Они находятся для каждого  $t = 1, 2, \dots, n$  по формуле

$$\hat{R}_t = A_1 + B_{11}M_t + B_{12}G_t.$$

Получим следующие значения:

$t$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
$\hat{R}_t$	5,01	4,92	4,98	5,2	5,09	5,26	4,74	4,97	5,21	4,91	4,92	4,69	4,63	4,93	4,84
$t$	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30
$\hat{R}_t$	4,7	4,84	5,04	5,08	4,89	5,01	5,01	5,17	5,02	4,96	4,92	5,42	5,17	5,31	4,87

Затем обычным МНК оцениваются параметры уравнения

$$I_t = a_3 + b_{31}\hat{R}_t + e_3^*.$$

Получим выражение

$$I_t = 19,3 - 0,9\hat{R}_t + \text{остаток}.$$

Предлагаем читателю в качестве упражнения оценить первое уравнение ДМНК и убедиться, что в этом случае он дает те же оценки, что и КМНК.

В том случае, когда ошибки в различных (идентифицируемых) уравнениях коррелируют между собой, существует возможность повысить эффективность оценивания, приняв во внимание дополнительную информацию об этой корреляции. Трехшаговый МНК (ТМНК) повторяет первые два шага ДМНК, затем на основе остатков оценивается ковариационная матрица и система переоценивается заново обобщенным МНК.

Продолжим рассмотрение примера. Мы оценили первое и третье уравнения системы в предположении о том, что ошибки различных уравнений не коррелируют друг с другом. Если же теперь предположить, что корреляция имеет место, то имеет смысл использовать ТМНК. Читателю предлагается оценить систему из первого и третьего уравнений с использованием в качестве инструментов переменных  $M$  и  $G$  в компьютерном пакете

(например, в *EViews* или *Gauss*) и убедиться в том, что в данном примере оценки, полученные ДМНК и ТМНК, почти одинаковы, а стандартные ошибки в последнем случае меньше (т.е. оценки эффективнее).

#### 4.2.2. Тестирование на экзогенность

Как отмечалось выше, для оценивания систем одновременных уравнений применяется ДМНК, потому что корреляция регрессоров с ошибками делает оценки обычного МНК несостоятельными. Однако если корреляция некоторых переменных с ошибкой вызывает сомнение, то имеет смысл протестировать их на экзогенность, поскольку трактовка экзогенных переменных как эндогенных хотя и не приводит к потере состоятельности оценок, но снижает их эффективность (увеличивает дисперсию, делая оценки более колеблемыми относительно их математического ожидания).

**Тест Хаусмана — *Vu* на экзогенность** устроен следующим образом. Если в рассматриваемом уравнении имеется несколько переменных, экзогенность которых подлежит проверке, то сначала для каждой такой переменной строится инструмент. В качестве инструмента берутся теоретические значения данной переменной из регрессии ее на предопределенные переменные (т.е. из приведенной формы). Далее строится регрессия объясняемой переменной исходного уравнения на все переменные этого уравнения и построенные инструменты. Нулевая гипотеза состоит в том, что все коэффициенты при инструментах равны нулю, альтернативная — что хотя бы один из них значимо отличен от нуля. Для проверки нулевой гипотезы применяется *F*-критерий Фишера. Если же тестируется значимость только одного коэффициента, то применение *F*-критерия эквивалентно использованию *t*-критерия Стьюдента.

Проверим, являются ли экзогенными переменные *Y* и *R* соответственно в первом и третьем уравнениях примера из 4.2.1.

Рассмотрим первое уравнение. В правой части присутствуют одна эндогенная переменная *Y* и одна экзогенная переменная *M*. Добавим в это уравнение инструмент  $\hat{Y}$  (теоретические значения, рассчитанные по второму уравнению приведенной формы) и оценим полученное уравнение (в скобках приведены значения *p-value*):

$$R = 5,4 - 0,0Y + 0,1M + 0,0\hat{Y}. \\ (0,2039)(0,0003)(0,1617)(0,1046)$$

Для проверки экзогенности переменной  $Y$  тестируется гипотеза о незначимости инструмента  $\hat{Y}$ . Поскольку  $p$ -value превосходит 10%-ный уровень значимости<sup>1</sup>, то нулевая гипотеза об экзогенности  $Y$  не отвергается. Следует отметить, что это, однако, еще не означает, что  $Y$  — экзогенная переменная для данного уравнения. Это означает лишь, что используемые данные не противоречат гипотезе об экзогенности переменной  $Y$  в рамках данной модели.

Читателю предлагается протестировать на экзогенность переменную  $R$  в третьем уравнении, а также попробовать оценить регрессию  $Y$  на  $R$ ,  $I$ ,  $G$  и соответствующие инструменты и убедиться в том, что полная коллинеарность регрессоров не позволит этого сделать.

### 4.3. Уравнения, кажущиеся несвязанными

Система может состоять из уравнений, кажущихся несвязанными, например из уравнений, описывающих спрос на различные виды продукции в зависимости от их цен. Если продукты конкурирующие (взаимозаменяемые) или дополняющие, то спрос на них будет зависеть от цен на другие продукты. Если же продукты не конкурируют и не дополняют друг друга, но, тем не менее, покупаются в пределах одной экономики, то спрос на них может испытывать общие экономические шоки, отраженные в коррелирующих ошибках. Так что уравнения будут казаться совершенно несвязанными (например, содержать непересекающиеся наборы переменных), но связь будет незримо присутствовать в коррелирующих ошибках.

Оценивать такие уравнения можно и по отдельности, но при этом будет потеря в эффективности, вызванная тем, что информация о корреляции ошибок не учитывается. Чтобы

<sup>1</sup> В данном случае  $p$ -value = 0,1046 означает, что если нулевая гипотеза (об экзогенности) верна, то абсолютные значения  $t$ -статистики, большие или равные абсолютному значению фактической  $t$ -статистики (здесь оно равно 1,68), будут получаться (для различных выборок) с вероятностью 0,1046 = 10,46%.

получить эффективные оценки, нужно объединить уравнения системы в одно и оценивать его обобщенным МНК. Оценивание представляет собой двухшаговую процедуру. На первом шаге уравнения оцениваются по отдельности обычным МНК, и остатки из этих регрессий используются для оценивания ковариационной матрицы ошибок. Далее на втором шаге применяется обобщенный МНК для оценивания объединенного уравнения. В левой части этого уравнения стоит объясняемая переменная, последовательно принимающая значения всех объясняемых переменных. В правой части стоят все регрессоры, которые тоже последовательно принимают все свои значения, если они присутствуют в уравнении, соответствующем объясняемой переменной данного наблюдения, и нуль в противном случае. Данная процедура позволяет получить состоятельные и асимптотически эффективные оценки.

Описанный метод дает идентичные с обычным МНК оценки в двух случаях: когда корреляция между ошибками равна нулю и когда во всех уравнениях используются одни и те же объясняющие переменные.

Предположим, что исследуется модель следующего вида:

$$R_t = a_1 + b_{14}M_t + e_1 \text{ (функция денежного рынка);}$$

$$Y_t = a_2 + b_{23}I_t + b_{25}G_t + e_2 \text{ (функция товарного рынка).}$$

Данные два уравнения, представляющие собой ограниченную версию примера из 4.2.1, кажутся несвязанными. Однако их ошибки могут коррелировать между собой. Уравнения содержат только predetermined (экзогенные) переменные в правой части<sup>1</sup>, поэтому обычный МНК даст состоятельные оценки. Рассчитаем их, используя данные из 4.2.1 (в скобках приводятся стандартные ошибки коэффициентов):

$$R_t = 3,52 + 0,15M_t + \text{остаток;}$$

$$(s.e.) (1,1055) (0,1117)$$

$$Y_t = -12,45 + 6,82I_t + 0,95G_t + \text{остаток.}$$

$$(s.e.) (6,6825) (0,2333) (0,0596)$$

<sup>1</sup> Здесь мы считаем, что переменная  $I$  экзогенна и не зависит от  $R$  (пример условный).

Если корреляция ошибок также имеет место, то эти оценки становятся неэффективными.

Оценим эту же систему в компьютерном пакете *EViews 6* с использованием *SUR*-метода<sup>1</sup> (в скобках приводятся стандартные ошибки коэффициентов):

$$R_t = 4,82 + 0,02M_t + \text{остаток};$$

(s.e.) (0,9370) (0,0944)

$$Y_t = -6,70 + 6,49I_t + 0,94G_t + \text{остаток}.$$

(s.e.) (5,5483) (0,1938) (0,0494)

Как видим, оценки коэффициентов изменились, причем стандартные ошибки понизились. Таким образом, учет дополнительной информации о корреляции ошибок помогает повысить точность оценивания.

## Контрольные вопросы и задания

1. Какие типы переменных принято выделять в системах эконометрических уравнений?
2. Назовите основные виды систем эконометрических уравнений.
3. Поясните, почему нельзя использовать обычный МНК для нахождения параметров системы одновременных уравнений.
4. Что называют структурной формой модели?
5. Для чего необходима приведенная форма модели? Какой вид она имеет?
6. Что такое идентификация модели?
7. Какие классы моделей можно выделить с точки зрения их идентификации?
8. В чем состоят необходимое и достаточное условия идентификации?
9. Какие методы могут быть использованы для нахождения параметров системы эконометрических уравнений? Какова область их применения?
10. В чем состоит тест Хаусмана — Vu на экзогенность?
11. Опишите особенности оценки параметров системы, состоящей из уравнений, кажущихся несвязанными.

<sup>1</sup> *SUR* — *Seemingly Unrelated Regression*.

# **Моделирование изолированного динамического ряда**

В целях прогнозирования широко используются эконометрические модели по временным (динамическим) рядам. Это связано с тем, что прогнозирование по пространственной (статической) информации осложнено необходимостью построения дополнительных прогнозов для объясняющих переменных. Например, по совокупности регионов за один год можно построить модель сбережений населения от их доходов. Однако воспользоваться этой моделью для прогноза можно, лишь зная прогноз доходов, который в свою очередь будет зависеть от прогноза инфляции и развития экономической деятельности в регионе.

## **5.1. Компоненты динамического ряда**

Модели по рядам динамики могут строиться на основе:

- изолированного динамического ряда, т.е. изучается один динамический ряд, например по данным о численности занятых за несколько лет строится модель динамики численности занятых;
- системы взаимосвязанных рядов динамики, т.е. когда один из рядов рассматривается как моделируемый объект, а другие — как его факторы, например строится модель прибыли в зависимости от объема реализации, численности работающих, фондовооруженности труда и т.п.



При построении моделей по временным рядам необходимо учитывать компоненты (составные части) динамического ряда.

Уровни динамического ряда в конкретный период времени  $t$  принимают те или иные значения в результате действия разных факторов. Одни из них являются основными, формирующими величину уровня  $y_t$  на данном этапе исторического развития, а другие — случайными, несущественными с точки зрения содержания его материальной природы. Фактическую величину уровня динамического ряда  $y_t$  можно представить как функцию трех компонент:

- **тенденции** ряда, обусловленной влиянием общих факторов, определяющих основное направление развития явления за длительный период времени — **тренд** ряда;
- **периодических колебаний**, вызванных особенностями существования явления в одни периоды по сравнению с другими (циклические — период колебаний несколько лет, сезонные — внутригодовые колебания);
- **случайных колебаний**, связанных с действием разного рода второстепенных факторов, — случайная компонента.

Символически функцию можно представить в виде

$$y_t = f(T, P, \xi),$$

где  $y_t$  — фактический уровень динамического ряда в период времени  $t$ ;  $T$  — тренд ряда;  $P$  — периодические колебания (циклические, сезонные);  $\xi$  — случайная составляющая.

Рассматриваемые компоненты динамического ряда обязательно присущи каждому временному ряду. Могут быть ряды динамики, в которых отсутствуют как тенденция, так и периодические колебания. В этом случае уровни ряда являются функцией случайной компоненты:  $y_t = f(\xi)$ . Они колеблются вокруг среднего уровня, что характерно для так называемого **стационарного ряда**. На графике такой ряд представляет собой ломаную линию, параллельную оси времени (рис. 5.1).

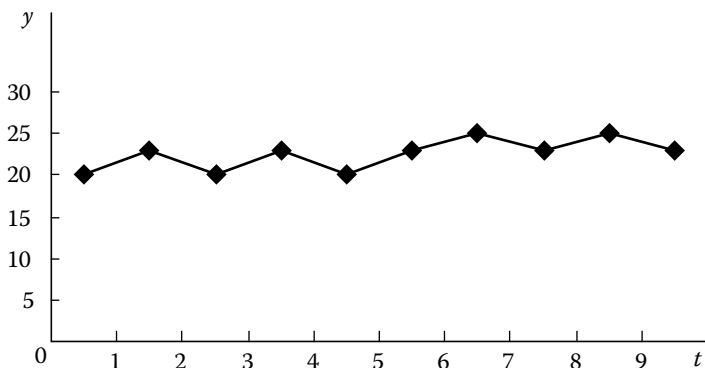


Рис. 5.1. Ряд без тенденции и периодических колебаний

Модель уровня такого динамического ряда имеет вид

$$y_t = \bar{y} + \xi,$$

где  $y_t$  — уровни динамического ряда;  $\bar{y}$  — средний за период времени уровень ряда;  $\xi$  — случайная составляющая, определяемая как  $\xi = y_t - \bar{y}$ .

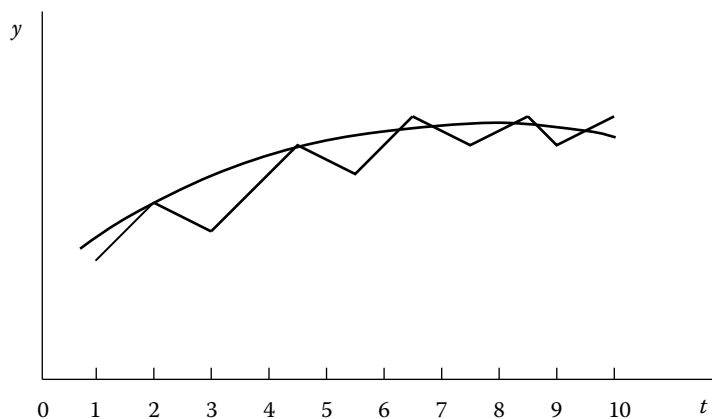
Такие ряды в экономике сравнительно редки. Чаще имеют место ряды с тенденцией. В основном ряды без тенденции наблюдаются при изучении динамики показателей из относительных и средних величин. Например, доля социальных платежей во внебюджетные фонды России в процентах от фонда оплаты труда на предприятиях представляет собой во многих случаях подобный ряд, так как при наличии дифференциации выплат с подавляющей их части страховые платежи взимаются с предприятий по единой ставке.

Большинство динамических рядов в экономике характеризуются тенденцией и случайными колебаниями (рис. 5.2).

Модель уровня такого ряда имеет вид

$$y_t = f(T) + \xi,$$

где  $f(T)$  — математическая функция, характеризующая закономерность развития явления во времени, т.е. описывающая тенденцию развития явления — тренд ряда;  $\xi$  — случайные колебания.



**Рис. 5.2.** Ряд с тенденцией (плавная линия) и случайными колебаниями (ломаная линия)

Если обозначить теоретическое значение уровня ряда, соответствующее определенной математической функции тренда, как  $\hat{y}_t$ , то случайные колебания  $\xi$  составят величину

$$\xi_t = y_t - \hat{y}_t.$$

Например, за 2000—2008 гг. динамический ряд среднемесячной номинальной начисленной заработной платы работника по России составил (в тыс. руб.):

2000	2001	2002	2003	2004	2005	2006	2007	2008
2,2	3,2	4,4	5,5	6,7	8,6	10,6	13,6	17,3

Наблюдающаяся тенденция к росту может быть описана уравнением тренда вида

$$\hat{y}_t = 1,9299e^{0,2466t},$$

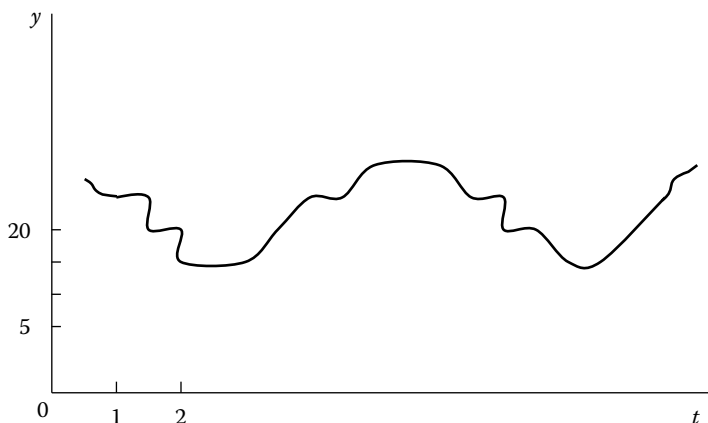
где  $t$  принимает значения 1, 2, ..., 9.

Согласно этой тенденции теоретическое значение уровня ряда в 2008 г.  $\hat{y}_t$  составило 17,8 тыс. руб. (в уравнение подставлено значение  $t = 9$ ). Так как фактическое значение в этот год

$y_t$  составило 17,3 тыс. руб., то величина случайной составляющей  $\xi$  окажется равной 0,5 тыс. руб. Конкретные экономические выводы о росте оплаты труда должны быть скорректированы на уровень инфляции.

При изучении динамики явления за продолжительный период времени уровни ряда могут обнаруживать регулярные колебания, повторяющиеся через равные промежутки времени: спады или подъемы. Такие колебания принято называть **периодическими** (рис. 5.3 и 5.4).

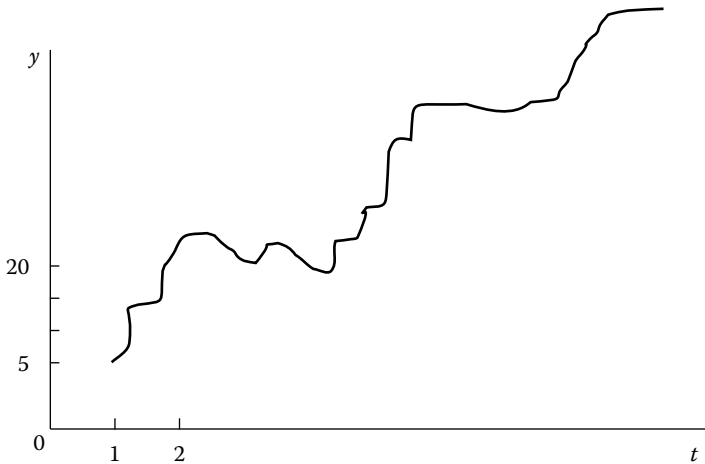
Если период колебаний насчитывает несколько лет, то такие периодические колебания считают циклическими. Например, солнечная активность проявляется с периодом 10—11 лет. В сфере предпринимательства могут иметь место экономические циклы, включающие в себя рост, спад, свертывание и затем оживление экономической деятельности. Длина цикла зависит от вида деятельности и охватывает нередко 3—12 лет. Так, в ряде стран производство свинины и ее цена подвергаются регулярным колебаниям, цикл которых длится около 3 лет<sup>1</sup>.



**Рис. 5.3.** Ряд с периодическими и случайными колебаниями

$$y_t = f(P, \xi)$$

<sup>1</sup> Ланге О. Введение в эконометрику : пер. с польского. М. : Прогресс, 1964. С. 145.



**Рис. 5.4.** Ряд с тенденцией, периодическими и случайными колебаниями  $y_t = f(T, P, \xi)$

Регулярные колебания в течение года называются **сезонными** (обозначаются  $S$ ). Например, к сезонным относятся колебания спроса на одежду с изменением сезона года (весна, лето, осень, зима), колебания цен на сельскохозяйственную продукцию и т.п. Цикл колебаний равен году. Наличие сезонных колебаний означает, что на протяжении ряда лет в одни и те же кварталы (месяцы) года наблюдается рост или снижение уровня ряда. Так, рождественские праздники обуславливают рост товарооборота в декабре и январе.

В отличие от периодических случайные колебания не носят регулярный характер и связаны с действием разного рода случайных причин.

Рассматриваемые компоненты динамического ряда позволяют представить уровень динамического ряда в виде **аддитивной** или **мультипликативной** моделей:

$$y_t = T + S + \xi \text{ — аддитивная модель;}$$

$$y_t = T \cdot S \cdot \xi \text{ — мультипликативная модель.}$$

Выбор вида модели зависит от характера периодических колебаний. Если амплитуда, например, сезонных колебаний

остаётся во времени постоянной, то применяется аддитивная модель. Если же амплитуда колебаний изменяется во времени, то рассматривается мультипликативная модель (рис. 5.5 и 5.6).

На рис. 5.5 показаны тенденция к увеличению (прямая линия) и периодические колебания — отклонения от тенденции, равные по всей длине динамического ряда, т.е. с одинаковой амплитудой волны.

На рис. 5.6 также присутствуют тенденция уровней ряда и периодические колебания, амплитуда которых возрастает во времени.

В аддитивной модели компоненты ряда выражены в тех же единицах измерения, что и рассматриваемый в динамике признак. Так, если  $y_t$  выражается в тысячах тонн, то и составные части ряда тоже выражены в тысячах тонн. Пусть  $y_t = 45$  тыс. т, а согласно тенденции  $\hat{y}_t = 40$  тыс. т. Если имеют место периодические колебания, то разница  $y_t - \hat{y}_t = 5$  тыс. т характеризует периодическую и случайную составляющие. Предположим, что периодическая компонента составила 12 тыс. т, т.е. тренд вместе с периодической составляющей равен 52 тыс. т. Однако с учетом случайных колебаний фактическое значение  $y_t = 45$  тыс. т. Следовательно, случайная компонента составила величину  $\xi = 45 - 52 = -7$ , или иначе  $\xi = 5 - 12 = -7$  тыс. т.

При мультипликативной модели периодическая и случайная составляющие выражены в относительных величинах. Так,

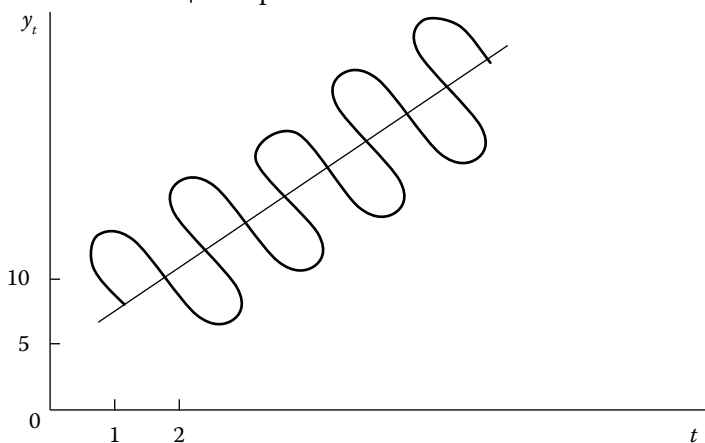
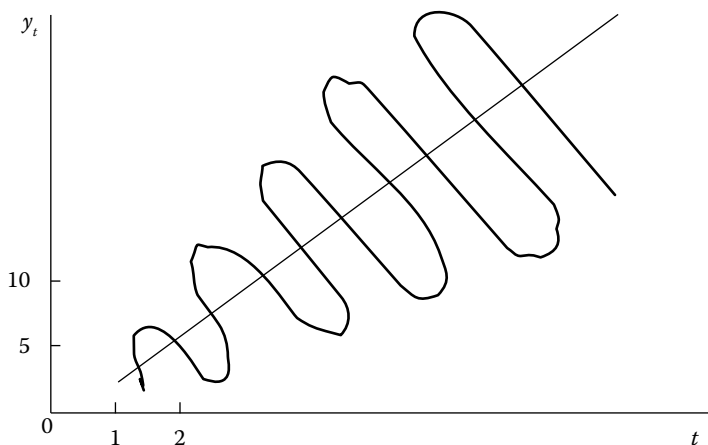


Рис. 5.5. Аддитивная модель



**Рис. 5.6.** Мультипликативная модель

при сезонных колебаниях  $S$  — это индекс сезонности. Пусть  $y_t = 38$  тыс. т, а согласно тенденции  $\hat{y}_t = 31$  тыс. т. Индекс сезонности для периода  $t$  составил, например 119,4%. Тогда тренд с учетом сезонности окажется равным 37 тыс. т ( $31 \cdot 1,194$ ), а случайная компонента по абсолютной величине составит 1 тыс. т ( $38 - 37$ ). Ее можно представить в виде относительной величины:  $\xi_t = \frac{38}{37} = 1,027$ . В этом случае мультипликативное разложение уровня динамического ряда окажется следующим:

$$y_t = \hat{y}_t \cdot S \cdot \xi, \text{ т.е. } 38 = 31 \cdot 1,194 \cdot 1,027$$

Если случайную компоненту представить абсолютной величиной, то получим модель смешанного типа:

$$y_t = \hat{y}_t \cdot S \cdot \xi, \text{ т.е. } 38 = 31 \cdot 1,194 + 1.$$

Рассмотренные компоненты динамического ряда учитываются как при построении модели изолированного временного

ряда, так и при построении регрессионных моделей на основе системы взаимосвязанных рядов динамики, что будет изложено далее.

## 5.2. Автокорреляция уровней динамического ряда и характеристика его структуры

При наличии тенденции в ряде динамики уровни ряда характеризуются автокорреляцией, т.е. каждый последующий уровень ряда зависит от предыдущего. Например, цена на товар сегодня, как правило, зависит от цены вчерашнего дня. Корреляционная связь между последовательными значениями уровней динамического ряда называется **автокорреляцией уровней динамического ряда**.

Для измерения автокорреляции уровней динамического ряда используется коэффициент автокорреляции уровней

$$r_{y_t, y_{t-\tau}} = \frac{(\overline{y_t y_{t-\tau}} - \bar{y}_t \bar{y}_{t-\tau})}{\sigma_{y_t} \sigma_{y_{t-\tau}}}, \quad (5.1)$$

где  $y_t$  — фактические уровни динамического ряда;  $y_{t-\tau}$  — уровни того же динамического ряда, но сдвинутые на  $\tau$  шагов во времени;  $\tau$  — величина лага (сдвига во времени), принимающая значения 1, 2, 3, ... и определяющая порядок коэффициента автокорреляции.

При  $\tau = 1$  рассчитывается коэффициент автокорреляции **первого порядка**, т.е. измеряется корреляция текущих значений уровней динамического ряда  $y_t$  с предшествующими уровнями  $y_{t-1}$ .

При  $\tau = 2$  изучается зависимость текущих уровней ряда  $y_t$  с уровнями этого же ряда, сдвинутыми на 2 временных шага  $y_{t-2}$ , т.е. рассчитывается коэффициент автокорреляции **второго порядка**, а при  $\tau = 3$  — соответственно **третьего порядка**, при  $\tau = k$  — коэффициент автокорреляции  **$k$ -го порядка**. Чем длиннее динамический ряд, тем выше может быть порядок коэффициента автокорреляции уровней.

Коэффициент автокорреляции уровней ряда практически рассчитывается по формуле линейного коэффициента



корреляции. Поэтому его величина изменяется в пределах от  $-1$  до  $+1$ . Чем ближе его величина к  $\pm 1$ , тем сильнее зависимость текущих уровней динамического ряда от предыдущих.

Если ряд характеризуется четко выраженной тенденцией, то для него коэффициент автокорреляции первого порядка приближается к  $+1$ . Так, для рассмотренного ранее ряда динамики заработной платы работника коэффициент автокорреляции уровней первого порядка составил  $0,9987$ , демонстрируя тесную связь последующих уровней ряда от предыдущих.

Поскольку в примере рассчитывается коэффициент автокорреляции первого порядка, т.е. когда  $\tau = 1$ , формула его расчета приобретает вид

$$r_{y_t, y_{t-1}} = \frac{(\overline{y_t y_{t-1}} - \bar{y}_t \bar{y}_{t-1})}{\sigma_{y_t} \sigma_{y_{t-1}}}, \quad (5.2)$$

где  $y_t$  — уровни ряда в момент времени  $t$ ;  $y_{t-1}$  — те же уровни ряда, но сдвинутые на год, т.е. уровни ряда в момент времени  $(t - 1)$  (предыдущий год).

Так как оба ряда ( $y_t$  и  $y_{t-1}$ ) для расчета коэффициента автокорреляции должны быть одинаковой длины, то первое значение по ряду  $y_t$  в расчетах не участвует. По нашему примеру необходимые суммы для подсчета отдельных элементов формулы коэффициента автокорреляции уровней составили

$$\bar{y}_t = \sum_{t=2}^n y_t / (n-1) = 69,9 / 8 = 8,7375;$$

$$\bar{y}_{t-1} = \sum_{t=2}^n y_{t-1} / (n-1) = 54,8 / 8 = 6,85;$$

$$\bar{y}_t \bar{y}_{t-1} = \sum_{t=2}^n y_t y_{t-1} / (n-1) = 610,39 / 8 = 76,29875;$$

$$\sigma_{y_t} = \sqrt{\sum_{t=2}^n y_t^2 / (n-1) - (\bar{y}_t)^2} = \sqrt{\frac{775,31}{8} - 8,7375^2} = 4,535399;$$

$$\sigma_{y_{t-1}} = \sqrt{\sum_{t=2}^n y_{t-1}^2 / (n-1) - (\bar{y}_{t-1})^2} = \sqrt{\frac{480,86}{8} - 6,85^2} = 3,631116.$$

Соответственно коэффициент автокорреляции уровней составит

$$r_{y_t, y_{t-1}} = \frac{76,29875 - 8,7375 \cdot 6,85}{4,535399 \cdot 3,631116} = 0,9987.$$

Методика расчета коэффициентов автокорреляции более высоких порядков та же, но при этом число коррелируемых пар уменьшается. В нашем примере их восемь ( $ct = 2$  по  $t = 9$ ). Если же увеличим лаг до 2 лет, т.е.  $\tau = 2$ , то останется семь коррелируемых пар ( $ct = 3$  по  $t = 9$ ), при  $\tau = 3$  будет шесть коррелируемых пар ( $ct = 4$  по  $t = 9$ ). Ввиду уменьшения числа наблюдений при расчете коэффициента автокорреляции уровней, увеличение величины лага не беспредельно: принято считать, что максимальная величина лага должна быть не более чем  $n/4$  ( $n$  — длина динамического ряда). Для нашего примера при  $n = 9$  максимальная величина лага составит 2 года ( $\tau = 2$ ).

Для расчета коэффициента автокорреляции второго порядка составим таблицу.

**Таблица 5.1.** Расчет коэффициента автокорреляции уровней второго порядка (для ряда динамики заработной платы работника)

$t$	$y_t$	$y_{t-2}$	$y_t y_{t-2}$	$y_t^2$	$y_{t-2}^2$
1	2,2	—	—	—	—
2	3,2	—	—	—	—
3	4,4	2,2	9,68	19,36	4,84
4	5,5	3,2	17,6	30,25	10,24
5	6,7	4,4	29,48	44,89	19,36
6	8,6	5,5	47,3	73,96	30,25
7	10,6	6,7	71,02	112,36	44,89
8	13,6	8,6	116,96	184,96	73,96
9	17,3	10,6	183,38	299,29	112,36
$\Sigma$	66,70*	41,2	475,42	765,07	295,9

\* Подсчитано без первых двух строк

Так как теперь в расчете участвует семь коррелируемых пар  $y_t$  и  $y_{t-2}$ , то первые две строки табл. 5.1 не принимаются во внимание. Коэффициенты автокорреляции разных порядков принято обозначать  $r_1, r_2, r_3, \dots, r_k$ , где  $1, 2, \dots, k$  указывает на номер порядка коэффициента автокорреляции. Формула расчета коэффициента автокорреляции второго порядка следующая:

$$r_2 = \frac{\overline{y_t y_{t-2}} - \bar{y}_t \bar{y}_{t-2}}{\sigma_{y_t} \sigma_{y_{t-2}}}$$

где  $\bar{y}_t = \sum_{t=3}^n y_t / (n-2) = 66,70 / 7 = 9,528571$ ;

$$\bar{y}_{t-2} = \sum_{t=3}^n y_{t-2} / (n-2) = 41,2 / 7 = 5,88571;$$

$$\bar{y}_t \bar{y}_{t-2} = \sum_{t=3}^n y_t y_{t-2} / (n-2) = 475,42 / 7 = 67,91714;$$

$$\sigma_{y_t} = \sqrt{\sum_{t=3}^n y_t^2 / (n-2) - (\bar{y}_t)^2} = \sqrt{\frac{765,07}{7} - 9,528571^2} = 4,3014;$$

$$\sigma_{y_{t-2}} = \sqrt{\sum_{t=3}^n y_{t-2}^2 / (n-2) - (\bar{y}_{t-2})^2} = \sqrt{\frac{295,9}{7} - 5,885714^2} = 2,762209.$$

Соответственно коэффициент автокорреляции равен

$$r_2 = \frac{67,91714 - 9,528571 \cdot 5,8857147}{4,3014 \cdot 2,762209} = 0,996072.$$

В рассмотренном примере уровни динамического ряда имеют тенденцию к возрастанию, и коэффициенты автокорреляции приближаются к +1. Аналогичная картина будет наблюдаться и при тенденции к уменьшению уровней динамического ряда. Например, лесовосстановление в России за 1995—2002 гг. характеризуется тенденцией к снижению. Уровни ряда (в тыс. га) составили:

1995	1996	1997	1998	1999	2000	2001	2002
1454	1110	1092	1019	964	973	960	887

Коэффициенты автокорреляции первого и второго порядков оказались равными  $r_1 = 0,812$  и  $r_2 = 0,885$ , что подтверждает наличие тенденции в ряду динамики. При этом  $r_1 > 0$  и  $r_2 > 0$ , хотя ряд и имеет тенденцию к снижению. Чем тенденция по ряду динамики более четкая, тем ближе  $r_1$  и  $r_2$  к +1.

Для стационарного динамического ряда с небольшими колебаниями уровней,  $r_1$  достаточно близок к нулю и может принимать небольшое отрицательное значение. Так, предположим, что уровни ряда приняли следующие значения (последовательно во времени):

3; 1; 2; 1; 2; 1; 3; 3; 2; 3; 1; 2; 1; 1; 3; 3; 2; 2; 1; 3; 3; 2; 2; 3; 1; 2; 2; 1; 3; 1.

Коэффициент автокорреляции первого порядка составил  $-0,209$ , а коэффициент автокорреляции второго порядка составил  $0,056$ .

Серию коэффициентов автокорреляции уровней ряда с последовательным увеличением величины лага принято называть **автокорреляционной функцией** (АКФ).

Для стационарного временного ряда с увеличением величины лага взаимосвязь  $y_t$  и  $y_{t-\tau}$  ослабевает и АКФ характеризуется монотонным убыванием, что графически должно представлять затухающую кривую (рис. 5.7).

По стационарному ряду АКФ оценивается исходя из формулы коэффициента автокорреляции

$$r(\tau) = \frac{\frac{1}{n-\tau} \sum_{t=1}^{n-\tau} (y_t - \bar{y})(y_{t+\tau} - \bar{y})}{\frac{1}{n} \sum_{t=1}^n (y_t - \bar{y})^2}, \quad (5.3)$$

где  $n$  — длина временного ряда  $y_1, y_2, \dots, y_n$ ;  $\tau$  — временной сдвиг;  $\bar{y}$  — средняя арифметическая по исходному ряду  $(\sum_{t=1}^n \frac{y_t}{n})$ .

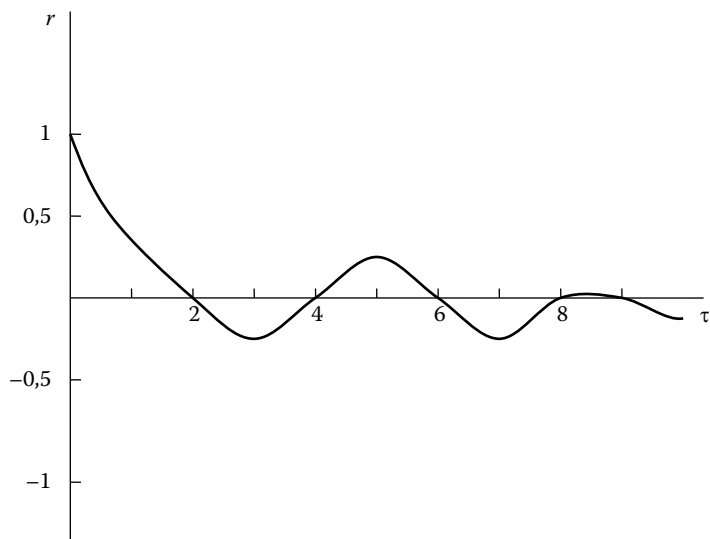
В нашем примере АКФ для стационарного ряда составила:  $r_1 = -0,209$ ;  $r_2 = 0,056$ ;  $r_3 = -0,114$ ;  $r_4 = -0,356$ ;  $r_5 = 0,057$ ;  $r_6 = -0,074$ ;  $r_7 = -0,003$ . Однако при ограниченной длине динамического ряда поведение АКФ в виде рис. 5.7 не всегда соблюдается.

АКФ дает представление о внутренней структуре динамического ряда. С помощью АКФ можно определить наличие или отсутствие в ряду динамики периодических колебаний и соответственно величину периода колебаний: она равна той величине лага  $\tau$ , при которой коэффициент автокорреляции уровней наибольший.

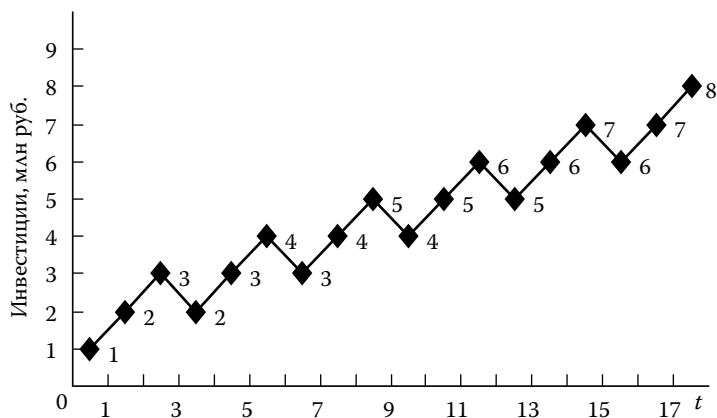
Предположим, что объем продаж товара за 18 мес. характеризуют следующим образом (рис. 5.8).

График показывает наличие тенденции, а также периодических колебаний. Это подтверждает и АКФ:

$$\begin{array}{lll} r_1 = 0,863; & r_2 = 0,829; & r_3 = 1; \\ r_4 = 0,812; & r_5 = 0,755; & r_6 = 1. \end{array}$$



**Рис. 5.7.** Коррелограмма АКФ



**Рис. 5.8.** Динамика объема продаж

Достаточно высокое значение коэффициента автокорреляции первого порядка ( $r_1 = 0,863$ ) означает наличие тенденции в ряде динамики. Вместе с тем максимальное значение

коэффициента автокорреляции наблюдается при лаге 3 и кратном ему лаге 6, т.е. для ряда характерна регулярная колеблемость уровней через 3 мес.: подъем в течение 3 мес. сменяется спадом в следующий месяц. Иными словами, волнообразное изменение объема продаж повторяется через 3 мес., что и демонстрирует АКФ. Для динамического ряда с монотонной тенденцией к возрастанию (или уменьшению) уровней АКФ имеет значения, близкие к +1, которые медленно снижаются с возрастанием величины лага. Например, за 60 кварталов динамика объема продаж характеризовалась уравниванием тренда

$$\hat{y}_t = 101,659 - 0,9254t,$$

где  $y$  — объем продаж в тыс. руб.;  $t = 1, 2, \dots, 60$ .

Коэффициент детерминации для него составил 0,973, характеризуя хорошее качество описания тенденции ряда: отклонения фактических уровней ряда от теоретических, обусловленных тенденцией, составляют всего 2,7%. АКФ для данного ряда оказалась следующей:  $r_1 = 0,991$ ;  $r_2 = 0,984$ ;  $r_3 = 0,980$ ;  $r_4 = 0,979$ ;  $r_5 = 0,973$ ;  $r_6 = 0,968$ ;  $r_7 = 0,963$ ;  $r_8 = 0,965$ ;  $r_9 = 0,963$ ;  $r_{10} = 0,962$ ;  $r_{11} = 0,959$ ;  $r_{12} = 0,957$ ;  $r_{13} = 0,952$ ;  $r_{14} = 0,955$ ;  $r_{15} = 0,943$ .

Если ряд характеризуется сменой тенденций, то АКФ примет значения, стремительно уменьшающиеся с возрастанием величины лага, сопровождаемые иногда сменой знака коэффициента автокорреляции. Так, например динамический ряд описывается параболой второго порядка (рис. 5.9).

АКФ оказывается следующей:

Лаг	$r$	Коррелограмма		
		-1	0	+1
1	0,887			xxxxxxxx
2	0,535			xxxxx
3	0,192			xx
4	-0,104	x		
5	-0,350	xxxx		

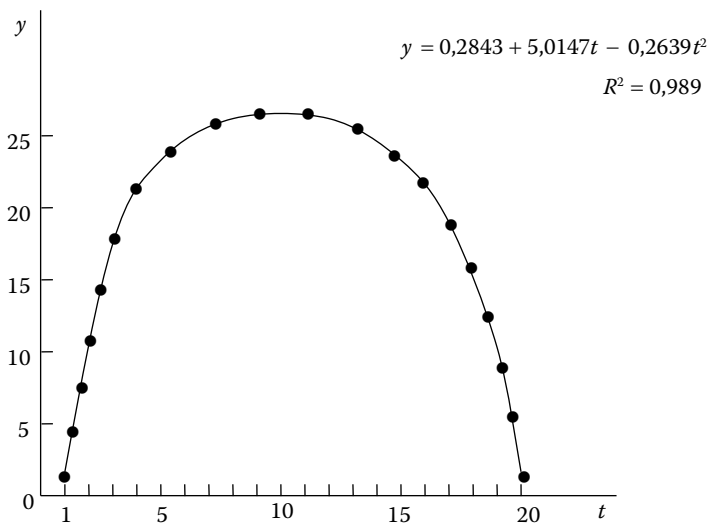


Рис. 5.9. Тренд в виде параболы второй степени

Похожая ситуация имеет место, например, при анализе динамики числа раненых в ДТП (на 100 тыс. человек населения) за 1999—2008 гг. по Тюменской области. Тенденция описывается параболой вида  $y = 80,537 + 45,756t - 3,5053t^2$ . Коэффициенты автокорреляции уровней с увеличением величины лага составили: 0,831; 0,588; 0,179; -0,544.

Иными словами, знание АКФ может помочь при подборе модели рассматриваемого динамического ряда.

### 5.3. Модели тенденции развития

#### 5.3.1. Общая характеристика моделей тенденции

Закономерность изменения уровней динамического ряда во времени может быть представлена в виде модели тенденции. При ее построении уровни динамического ряда рассматриваются как функция времени  $t$  и случайной компоненты

$\xi$ . Тогда модель уровня динамического ряда можно представить в виде

$$y_t = \bar{y} + (\hat{y}_t - \bar{y}) + (y_t - \hat{y}_t), \quad (5.4)$$

где  $y_t$  — фактический уровень динамического ряда в период времени  $t$ ;  $\bar{y}$  — средний уровень динамического ряда за весь период времени;  $\hat{y}_t$  — теоретический уровень динамического ряда, обусловленный тенденцией развития, т.е. тренд ряда.

В этой модели  $(\hat{y}_t - \bar{y})$  характеризует эффект тенденции, а  $(y_t - \hat{y}_t)$  — случайную составляющую  $\xi$ . Ввиду того, что  $\bar{y} + (\hat{y}_t - \bar{y}) = \hat{y}_t$ , данную модель уровня временного ряда можно представить в виде

$$y_t = \hat{y}_t + \xi_t, \quad (5.5)$$

где  $\hat{y}_t$  — модель тенденции, когда уровни ряда рассматриваются как функция времени  $t$ :  $\hat{y}_t = f(t)$ .

Совершенно ясно, что практическая значимость модели тенденции будет тем выше, чем меньше будут остаточные колебания (случайная составляющая  $\xi_t = y_t - \hat{y}_t$ ).

Построение модели тенденции (уравнения тренда) включает в себя следующие этапы работы:

- выбор математической функции, описывающей тенденцию;
- оценку параметров модели;
- проверку адекватности выбранной функции и оценку точности модели;
- расчет точечного и интервального прогнозов.

В настоящее время компьютерные программы анализа временных рядов содержат достаточно широкий набор математических функций для построения уравнения тренда. Все многообразие их можно свести в три группы:

- функции с монотонным характером возрастания (убывания) и отсутствием пределов роста (снижения);
- кривые с насыщением, т.е. устанавливается нижняя или верхняя граница изменения уровней ряда;
- S-образные кривые, т.е. кривые с насыщением, имеющие точку перегиба.



В первую группу функций входят полиномы  $k$ -й степени

$$\hat{y}_t = a + b_1t + b_2t^2 + \dots + b_kt^k. \quad (5.6)$$

При  $k = 1$  получаем линейный тренд:  $\hat{y}_t = a + b_1t$ , который часто записывают как  $\hat{y}_t = a + bt$ .

По содержанию линейный тренд означает, что уровни динамического ряда изменяются с одинаковой скоростью, т.е. с равным абсолютным приростом (параметр « $b$ »). В этом можно убедиться, подставляя в уравнение линейного тренда порядковые значения  $t$  (1, 2, 3, ...,  $k$ ): теоретические уровни ряда  $\hat{y}_t$  будут изменяться на величину параметра  $b$ , т.е. в арифметической прогрессии.

Например, уравнение тренда для индексов потребительских цен за 12 мес. года составило:  $\hat{y}_t = 99,9 + 1,9t$ , где  $t = 1, 2, \dots, 12$ . Из уравнения очевидно, что ежемесячно цены возрастали в среднем на 1,9 процентных пункта.

При  $k = 2$  получаем параболу второй степени

$$\hat{y}_t = a + b_1t + b_2t^2. \quad (5.7)$$

Данная функция рекомендуется для моделирования тенденции, если временной ряд характеризуется постоянным абсолютным ускорением, т.е. постоянными являются вторые разности (приросты абсолютных приростов). В этом случае скорость ряда изменяется линейно:

$t$	$\hat{y}_t = a + b_1t + b_2t^2$	Скорость $\Delta = y_t - y_{t-1}$	Ускорение $\Delta'' = \Delta y_t - \Delta y_{t-1}$
0	$a$	—	—
1	$a + b_1 + b_2$	$b_1 + b_2$	—
2	$a + 2b_1 + 4b_2$	$b_1 + 3b_2$	$2b_2$
3	$a + 3b_1 + 9b_2$	$b_1 + 5b_2$	$2b_2$
4	$a + 4b_1 + 16b_2$	$b_1 + 7b_2$	$2b_2$
5	$a + 5b_1 + 25b_2$	$b_1 + 9b_2$	$2b_2$

Как видим, параметр  $a$  означает начальный уровень ряда динамики при  $t = 0$ . Параметр  $b_1$  представляет собой средний

абсолютный прирост за весь период времени, если  $t$  обозначено так, что  $\Sigma t = 0$  (при обозначении  $t$  в виде ряда натуральных чисел, что наиболее распространено при компьютерной обработке, параметр  $b_1$  такой интерпретации не имеет). Параметр  $b_2$  характеризует половину абсолютного ускорения динамического ряда.

Например, динамика численности детей в возрасте 7 лет характеризуется по району за последние 15 лет уравнением тренда

$$\hat{y}_t = 323,7 + 10,8t - 1,6t^2,$$

где  $y$  — тыс. человек;  $t = 1, 2, \dots, 15$ .

Следовательно, ежегодно численность детей сокращается в среднем с ускорением в 3,2 тыс. человек.

При  $k = 3$  имеем параболу третьей степени

$$\hat{y}_t = a + b_1t + b_2t^2 + b_3t^3. \quad (5.8)$$

Этот вид тренда предполагает, что по временному ряду стабильны третьи разности ( $\Delta'''$ ), т.е. приросты вторых приростов ( $\Delta'' = \Delta''_t - \Delta''_{t-1}$ ), а абсолютные ускорения имеют линейную тенденцию:

$t$	$\hat{y}_t = a + b_1t + b_2t^2 + b_3t^3$	$\Delta' = y_t - \Delta y_{t-1}$	$\Delta'' = \Delta_t - \Delta_{t-1}$	$\Delta''' = \Delta''_t - \Delta''_{t-1}$
0	$a$	—	—	—
1	$a + b_1 + b_2 + b_3$	$b_1 + b_2 + b_3$	—	—
2	$a + 2b_1 + 4b_2 + 8b_3$	$b_1 + 3b_2 + 7b_3$	$2b_2 + 6b_3$	—
3	$a + 3b_1 + 9b_2 + 27b_3$	$b_1 + 5b_2 + 19b_3$	$2b_2 + 12b_3$	$6b_3$
4	$a + 4b_1 + 16b_2 + 64b_3$	$b_1 + 7b_2 + 37b_3$	$2b_2 + 18b_3$	$6b_3$
5	$a + 5b_1 + 25b_2 + 125b_3$	$b_1 + 9b_2 + 61b_3$	$2b_2 + 24b_3$	$6b_3$

Полиномы высоких степеней требуют достаточно длинных динамических рядов, чтобы параметры тренда были статистически надежными: на каждый параметр при  $t$  должно

приходиться не менее 6—7 временных единиц. Следовательно, парабола уже третьей степени должна содержать не менее 20 лет (если уровни ряда представлены по годам), что предполагает достаточно стабильную экономику.

Чаще отдают предпочтение функциям с меньшим числом параметров. Среди них широкое применение находит показательная функция

$$\hat{y}_t = ab^t, \quad (5.9)$$

или равносильная ей экспонента

$$\hat{y}_t = e^{a+bt} \text{ (либо } y = ae^{bt}), \quad (5.10)$$

которые характеризуются стабильным коэффициентом (темпом) роста:

$t$	1	2	3	4	5
$y = ab^t$	$ab$	$ab^2$	$ab^3$	$ab^4$	$ab^5$
Коэффициент роста	—	$b$	$b$	$b$	$b$

Например, за ряд лет динамика прибыли характеризуется уравнением вида  $\hat{y}_t = 13,5 \cdot 1,5^t$ , где  $t = 1, 2, \dots, n$ . Следовательно, ежегодно прибыль возрастает в среднем на 50% (коэффициент роста 1,5). Данный тренд в виде экспоненты примет выражение  $y = e^{2,603+0,405t}$ , где  $e^{2,603} = 13,5$  и  $e^{0,405} = 1,5$ . Рост по экспоненте означает геометрическую прогрессию уровней динамического ряда, что в экономике возможно сравнительно небольшой период времени (ограничены ресурсы, меняются условия рынка).

Если стабильными оказываются коэффициенты опережения темпов роста, то динамический ряд может быть описан **логарифмической параболой**

$$\hat{y}_t = ab^t c^{t^2}. \quad (5.11)$$

Свое название данная функция получила ввиду того, что прологарифмировав ее, получим параболу второй степени  $\lg y = \lg a + t \lg b + t^2 \lg c$ .

Для этой функции темпы роста изменяются в одно и то же число раз ( $c^2$ ):

$t$	$\hat{y}_t = ab^t c^{t^2}$	Коэффициент роста $k_t$	Коэффициент опережения ( $k_t/k_{t-1}$ )
1	$abc$	—	—
2	$ab^2c^4$	$bc^3$	—
3	$ab^3c^9$	$bc^5$	$c^2$
4	$ab^4c^{16}$	$bc^7$	$c^2$
5	$ab^5c^{25}$	$bc^9$	$c^2$

Например, дебиторская задолженность за ряд лет характеризуется уравнением  $\hat{y}_t = 1,47 \cdot 1,30^t \cdot 1,05^{t^2}$ . Следовательно, имеет место ускоренное увеличение дебиторской задолженности с коэффициентом опережения темпов роста  $1,05^2$ , т.е. 1,1025. Другими словами, темпы роста ежегодно возрастали в среднем в 1,1025 раза.

При моделировании тенденции используются и другие функции, приводимые к линейному виду. Так, при замедленном росте уровней ряда может использоваться полулогарифмическая кривая

$$\hat{y}_t = a + b \ln t. \quad (5.12)$$

В 1990-е гг. по этой функции развивалось в стране потребление картофеля.

Предполагая разную меру пропорциональности изменений уровней во времени, может быть использована степенная функция

$$\hat{y}_t = at^b. \quad (5.13)$$

При  $b > 0$  она характеризует непрерывный рост уровней с падающими темпами роста, а при  $b < 0$  — их ускоренное снижение. Величина  $t^b$  означает базисный коэффициент роста:

$t$	$y = at^b$	Базисный коэффициент роста
1	$a$	1
2	$a2^b$	$2^b$

Окончание табл.

$t$	$y = at^b$	Базисный коэффициент роста
3	$a3^b$	$3^b$
4	$a4^b$	$4^b$
5	$a5^b$	$5^b$

Поэтому степенная функция практически сообщает о величине среднего коэффициента роста

$$\bar{K} = \sqrt[t-1]{t^b}. \quad (5.14)$$

Например, обеспеченность городского населения республики Коми жильем ( $m^2$  общей площади на человека) за 1990—1999 гг. характеризовалась уравнением вида  $\hat{y}_t = 15,876t^{0,08}$ , где  $t = 1, 2, \dots, 10$ . Следовательно, за весь период обеспеченность населения жильем выросла в 1,202 раза ( $10^{0,08}$ ), т.е. ежегодно она возрастала в среднем на 2,07% ( $\bar{K} = \sqrt[9]{1,202} = 1,027$ ).

К кривым с насыщением можно отнести гиперболы вида

$$y = a + \frac{b}{t}; \quad (5.15)$$

$$y = a + \frac{b}{c+t}. \quad (5.16)$$

Равносторонняя гипербола ( $y = a + \frac{b}{t}$ ) при  $b > 0$  означает, что уровни ряда снижаются во времени и асимптотически приближаются к параметру  $a$ .

Например, индексы потребительских цен (декабрь к декабрю предыдущего года) за 1998—2003 гг. по России изменялись по гиперболе вида

$$\hat{y}_t = 95,557 + \frac{87,126}{t}; \quad r^2 = 0,9897.$$

Уравнение характеризует падающую тенденцию индекса потребительских цен (ИПЦ), при которой ИПЦ не может быть меньше 95,6%. Тренд описывает 99% вариации ИПЦ и лишь 1% ее связан с действием случайных факторов.

Если  $b < 0$ , то уравнение тренда  $y = a + \frac{b}{t}$  характеризует тенденцию к росту уровней ряда с асимптотической границей равной параметру « $a$ ». Так, численность мужчин старше трудоспособного возраста в Санкт-Петербурге за 1979—1995 гг. характеризовалась повышающейся тенденцией:  $\hat{y}_t = 296,92 - \frac{89,90}{t}$ , из которой следует, что численность мужчин этой возрастной категории за данный период не превышала 296,9 тыс. человек. Этот максимум выдерживался и для 1996 и 1997 гг., а в 1998 г. он превысил эту величину, составив 303,1 тыс. человек.

Гипербола  $y = a + \frac{b}{c+t}$  при  $b > 0$  и  $c > 0$  всегда характеризует падающую тенденцию с нижней асимптотой, равной параметру  $a$ . При  $b < 0$  данная кривая означает рост уровней ряда, который происходит до определенного предела, описываемого параметром  $a$ . Рассматриваемая гипербола предполагает более плавное замедление изменения уровней, чем равносторонняя гипербола  $y = a + \frac{b}{t}$ .

Среди гипербол нередко используется так называемая **обратная функция**  $y = \frac{1}{a+bt}$ . Свое название она получила в связи с тем, что при сведении ее к линейному виду используются обратные значения  $y$ , т.е.  $\frac{1}{y} = a + bt$ . Следует отметить, что если  $\frac{1}{y}$  имеет экономический смысл, то параметры данной функции интерпретируются аналогично линейному тренду. Например, предположим, что динамика трудоемкости продукции  $y$  характеризуется уравнением  $y = \frac{1}{1+2t}$ . Оно означает снижение трудоемкости и рост производительности труда на 2 ед. В иных случаях параметры обратной функции экономически не интерпретируются.

При  $b > 0$  ряд характеризуется понижающейся тенденцией, а при  $b < 0$  — повышающейся.

Среди кривых с насыщением может использоваться модифицированная экспонента

$$y = c \pm ab^t, \quad (5.17)$$

где  $c$  — асимптота (верхняя для функции  $y = c - ab^t$  и нижняя для функции  $y = c + ab^t$ ). Так, при изучении тенденции роста уровня механизации труда целесообразно учитывать ограничение роста (показатель уровня механизации труда не может быть больше 100%). Если изучается динамика детской смертности, то можно установить нижнюю асимптоту, т.е. минимальное значение детской смертности исходя из достигнутых условий жизни.

Модифицированная экспонента характеризуется постоянным отношением последовательных во времени приростов. Величина этого отношения равна параметру  $b$ :

$t$	$y = c + ab^t$	$\Delta_i$	$\frac{\Delta_i}{\Delta_{i-1}}$
0	$c + a$	—	—
1	$c + ab$	$a(b - 1)$	—
2	$c + ab^2$	$ab(b - 1)$	$b$
3	$c + ab^3$	$ab^2(b - 1)$	$b$
4	$c + ab^4$	$ab^3(b - 1)$	$b$

Так, модифицированная экспонента роста уровня механизации труда  $\hat{y}_t = 100 - 12,7 \cdot 0,895^t$  означает, что ежегодно скорость ряда снижается в 0,895 раз или на 10,5%. Верхняя граница уровня механизации труда 100%.

Величина  $(100 - y)$  характеризует уровень использования ручного труда. Поэтому в уравнении интерпретируется и параметр  $a$ :  $a = 12,7\%$  означает начальный уровень ручного труда. Соответственно 87,3% составит начальный уровень механизированного труда.

Модифицированная экспонента служит базовой кривой для других кривых с насыщением, а именно S-образных кривых: логистической кривой и кривой Гомперца. Тенденция развития явления в S-образных кривых охватывает три этапа: вначале довольно медленный рост, который затем убыстрится, далее сменяется уменьшением роста и приближением уровня ряда к предельному значению, т.е. к уровню насыщения.

Если в модифицированной экспоненте вместо  $y$  ввести обратную величину  $\frac{1}{y}$ , то получим логистическую кривую вида

$$y = \frac{1}{c + ab^t}, \quad (5.18)$$

которую называют кривой Перла — Рида. В ней верхняя асимптота составит величину  $\frac{1}{c}$  (рис. 5.10).

Точка перегиба у этой кривой равна  $t_p = \frac{1}{\ln b} \ln \frac{c}{a}$ . Значение  $y$  в точке перегиба равно  $\frac{1}{2c}$ . При практических расчетах исследователь может не иметь в полном виде S-образную кривую. Тогда точка перегиба находится за пределами наблюдаемых величин уровней ряда. В этом случае верхняя асимптота является теоретическим максимумом и ориентироваться на него в дальнейшем прогнозе достаточно проблематично.

Однако чаще сегодня применяется логистическая кривая вида

$$y = \frac{c}{1 + be^{-at}}, \quad (5.19)$$

где  $c$  — верхняя асимптота;  $b$  и  $a$  — параметры функции;  $e$  — основание натурального логарифма.

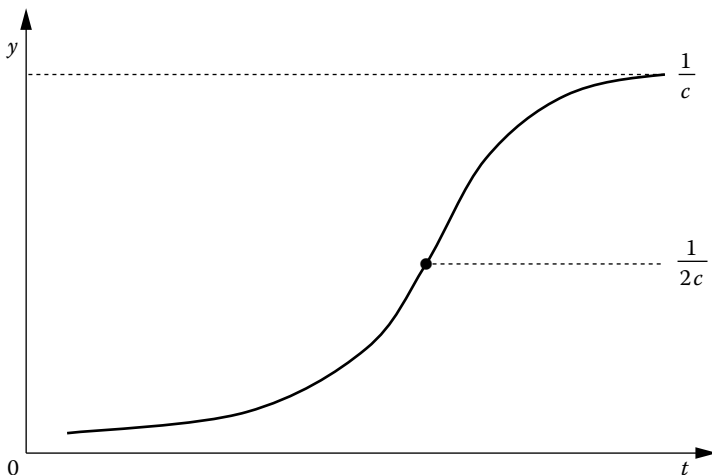


Рис. 5.10. Логистическая кривая Перла — Рида



Механизм развития производства новых товаров описывается иногда этой кривой.

Г. Тинтнер<sup>1</sup> применил данную функцию для описания тенденции роста численности населения Швеции за 100 лет по десятилетним интервалам с 1850 по 1950 г.:

$$y = \frac{10\,328\,806}{1 + 2,1176e^{-0,14t}}.$$

Согласно этой кривой верхняя асимптота роста численности населения Швеции составила 10 328 806 человек (по данным статистики в 2005 г. население Швеции составляло 9,0 млн человек).

Максимальное значение показателя  $s$  соответствует на графике отрезку кривой, параллельному оси абсцисс. Минимальное значение функции, равное нулю при  $t$ , стремящемся к  $-\infty$ , обычно отсутствует при использовании модели тенденции в прогнозных расчетах.

К классу  $S$ -образных кривых относится также кривая Гомперца

$$y = ca^{b^t}. \quad (5.20)$$

Она нашла применение в страховых расчетах и экстраполяции численности населения.

Верхняя асимптота соответствует значению параметра  $c$ , а нижняя равна нулю, если  $\ln a < 0$  (рис. 5.11, 5.12).

Если  $\ln a > 0$ , то кривая имеет нижнюю асимптоту, равную величине параметра  $c$  (рис. 5.13, 5.14).

Кривая Гомперца основана на модифицированной экспоненте. Прологарифмировав уравнение кривой Гомперца, получим после замены переменных уравнение модифицированной экспоненты

$$\begin{aligned} \lg y &= \lg c + b^t \lg a; \lg y = Y; \\ \lg c &= c' \text{ и } \lg a = a' \rightarrow Y = c' + a'b^t. \end{aligned}$$

<sup>1</sup> Тинтнер Г. Введение в эконометрию : пер. с нем. М. : Статистика. 1965. С. 291.

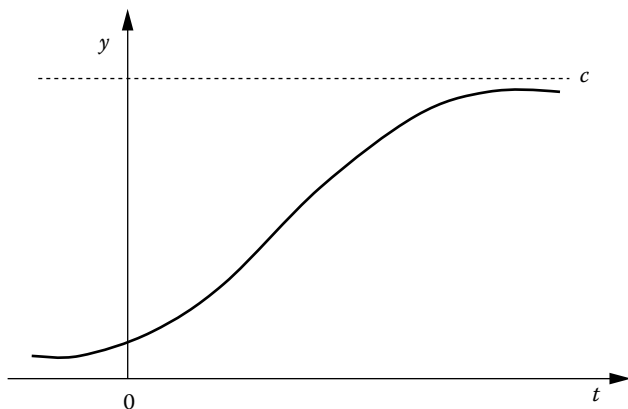


Рис. 5.11. Кривая Гомперца при  $\ln a < 0; b < 1$

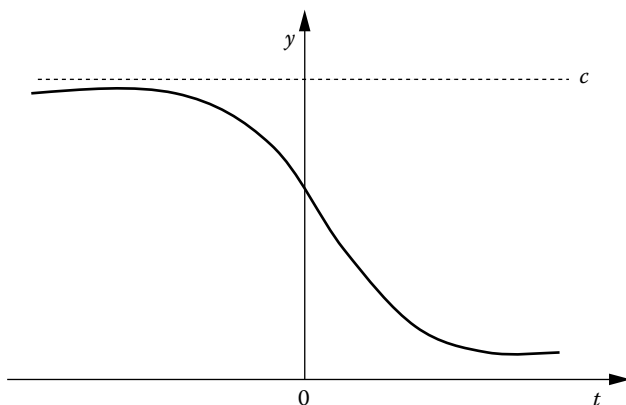


Рис. 5.12. Кривая Гомперца при  $\ln a < 0; b > 1$

Параметр  $c'$  будет характеризовать уровень насыщения. Точкой перегиба данной кривой будет точка

$$t_p = \frac{1}{\ln b} \ln\left(-\frac{1}{\ln a}\right)$$

со значением функции  $y_{t_p} = \frac{c}{e}$ , где  $e$  — основание натурального логарифма.

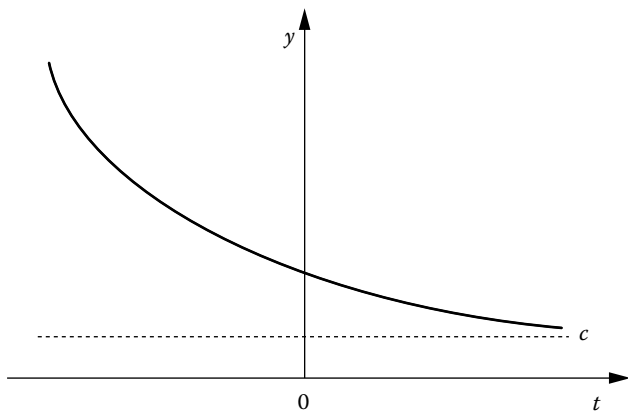


Рис. 5.13. Кривая Гомперца при  $\ln a > 0; b < 1$

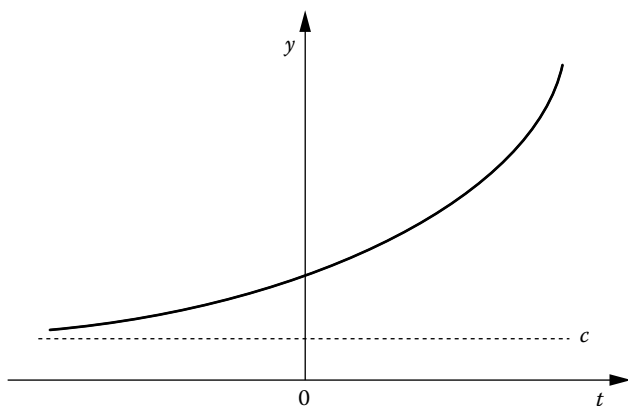


Рис. 5.14. Кривая Гомперца при  $\ln a > 0; b > 1$

Например, затраты на строительство автомобильных дорог описаны в работе К. Д. Льюис<sup>1</sup> в виде кривой Гомперца  $\hat{y}_t = 4644,5 \cdot 0,0961435^{0,93176^t}$ . Уравнение тренда показывает предельное значение затрат 4644,5 ден. ед. Точка перегиба составляет 12 лет; ей соответствуют затраты в 1708,6 ден. ед. Далее прирост затрат постепенно падает.

<sup>1</sup> Льюис К. Д. Методы прогнозирования экономических показателей : пер. с англ. М. : Финансы и статистика. 1986. С. 111—112.

### 5.3.2. Оценка параметров уравнения тренда

При использовании полиномов разных степеней оценка параметров уравнения тренда производится методом наименьших квадратов (МНК) точно так же, как оценки параметров уравнения регрессии на основе пространственных данных. В качестве зависимой переменной рассматриваются уровни динамического ряда, а в качестве независимой переменной — фактор времени  $t$ , который обычно выражается рядом натуральных чисел  $1, 2, \dots, n$ .

Оценка параметров нелинейных функций проводится МНК после линеаризации, т.е. приведения их к линейному виду. Рассмотрим применение МНК для некоторых нелинейных функций, которые не излагались подробно в главе, посвященной регрессии.

Для оценки параметров показательной кривой  $y = ab^t$  или экспоненты  $y = e^{a+bt}$  (либо  $y = ae^{bt}$ ) путем логарифмирования функции приводятся к линейному виду  $\ln y = \ln a + t \ln b$  или экспоненты:  $\ln y = a + bt$ . Далее строится система нормальных уравнений

$$\begin{cases} \sum \ln y = n \ln a + \ln b \sum t, \\ \sum t \ln y = \ln a \sum t + \ln b \sum t^2. \end{cases} \quad (5.21)$$

#### Пример 5.1

Число зарегистрированных ДТП (на 100 000 человек населения) по Новгородской области за 2000—2008 гг. характеризуется данными:

2000	2001	2002	2003	2004	2005	2006	2007	2008
105,7	105,3	156	158,1	160,1	178	191,5	274,6	287,3

Исходя из графика была выбрана показательная кривая  $y = ab^t$ . Для построения системы нормальных уравнений были рассчитаны вспомогательные величины

$$\begin{aligned} \sum \ln y &= 46,217; \sum_{1}^6 t = 45; \\ \sum t \ln y &= 238,496; \sum_{1}^6 t^2 = 285. \end{aligned}$$

Система нормальных уравнений составила

$$\begin{cases} 46,217 = 9\ln a + 45\ln b, \\ 238,496 = 45\ln a + 285\ln b. \end{cases}$$

Решая ее, получим значения

$$\ln a = 4,517598; a = e^{4,5176} = 91,61524;$$

$$\ln b = 0,123523; b = e^{0,12352} = 1,131476.$$

Соответственно имеем экспоненту  $y = 91,615e^{0,1235t}$  или показательную кривую  $Y = 91,615 \cdot 1,1315^t$ .

За период с 2000 по 2008 г. число дорожно-транспортных происшествий возрастало в среднем ежегодно на 13,5%. Экспонента достаточно хорошо описывает тенденцию исходного временного ряда: коэффициент детерминации составил 0,9202. Следовательно, данный тренд объясняет 92% колеблемости уровней ряда и лишь 8% ее связаны со случайными факторами.

---

Некоторую специфику имеет оценка параметров кривых с насыщением: модификационной экспоненты, логистической кривой, кривой Гомперца, гиперболы вида  $y = a + \frac{b}{c+t}$ .

По этим функциям должна быть сначала определена асимптота. Если она может быть задана исследователем на основе анализа временного ряда, то другие параметры могут быть оценены по МНК. В этих случаях данные функции приводятся к линейному виду. Рассмотрим оценку параметров этих кривых на отдельных примерах, начиная с модифицированной экспоненты.

### Пример 5.2

Уровень механизации труда (в %) характеризуется динамическим рядом (табл. 5.2)

**Таблица 5.2.** Расчет параметров модифицированной экспоненты  
 $y = c - ab^t$

Годы	$y$	$Y = c - y$	$\ln Y$	$t$	$t \ln Y$	$t^2$	$\hat{y}_t$
2001	82	18	2,890	1	2,890	1	82,4
2002	85	15	2,708	2	5,416	4	85,6
2003	89	11	2,399	3	7,197	9	88,2
2004	91	9	2,197	4	8,788	16	90,3
2005	92	8	2,079	5	10,395	25	92,0
2006	93	7	1,946	6	11,676	36	93,5
2007	94	6	1,792	7	12,544	49	94,6
2008	96	4	1,386	8	11,088	64	95,6
$\Sigma$	722	—	17,397	36	69,994	204	722,2

Так как уровень механизации труда не может превышать 100%, то имеется объективно заданная верхняя асимптота  $c = 100$ . Для оценки параметров  $a$  и  $b$  приведем рассматриваемую функцию к линейному виду  $c - y = ab^t$ ; обозначим  $(c - y)$  через  $Y$  и прологарифмируем:  $\ln Y = \ln a + t \ln b$ . Далее применим МНК и получим систему нормальных уравнений

$$\begin{cases} \Sigma \ln Y = n \ln a + \ln b \Sigma t, \\ \Sigma t \ln Y = \ln a \Sigma t + \ln b \Sigma t^2. \end{cases}$$

Для нашего примера, исходя из данных итоговой строки табл. 3, имеем систему уравнений

$$\begin{cases} 17,397 = 8 \ln a + 36 \ln b, \\ 69,994 = 36 \ln a + 204 \ln b. \end{cases}$$

Решив ее, получим  $\ln a = 3,06311$ ;  $\ln b = -0,19744$ . Соответственно потенцируя, получим:  $a = 21,394$ ,  $b = 0,8208$ , т.е. уравнение  $Y = 21,394 \cdot 0,8208^t$ .

Если перейти от  $Y$  к исходным уровням ряда, уравнение модифицированной экспоненты составит  $\hat{y}_t = 100 - 21,394 \times \times 0,8208^t$ , где параметр  $b = 0,8208$  показывает средний коэффициент снижения уровня использования ручного труда за 1998—2005 гг. Расчетные значения  $y$ , т.е.  $\hat{y}_t$  могут быть найдены путем подстановки в уравнение  $\hat{y}_t = 100 - 21,394 \times$

$\times 0,8208^t$  соответствующих значений  $t$ . Либо на основе уравнения  $\ln Y = 3,06311 - 0,19744 t$  при компьютерной обработке определяется  $\ln Y$  и далее  $100 - e^{\ln y}$ . Так, при  $t = 8$   $\ln Y = 1,48363$  и  $100 - e^{1,48363} = 100 - 4,40892 = 95,59108 = 95,6$  (см. последнюю графу таблицы). Ввиду некоторой смещенности оценок (так как МНК применяется к логарифмам)  $\sum y_t \neq \sum \hat{y}_t$ , хотя в примере эти величины достаточно близки друг другу.

Если асимптота  $c$  не задана, то оценка параметров модифицированной экспоненты усложняется. В этих случаях могут использоваться разные методы оценивания: метод трех сумм, метод трех точек<sup>1</sup>, с помощью регрессии<sup>2</sup>, метод Брианта<sup>3</sup>. Рассмотрим применение метода регрессии для оценки параметров модифицированной экспоненты вида  $y = c - ab^t$ .

### Пример 5.3

В таблице представлены данные о расходах предприятия на рекламу за 10 мес. года.

**Таблица 5.3.** Данные о расходах предприятия на рекламу за 10 мес. года (в тыс. руб.)

$t$	$y$	$z$	$\ln z$	$t - 1$	$c$
1	121				516,3118
2	196	75	4,317488	1	516,8105
3	256	60	4,094345	2	516,3499
4	305	49	3,89182	3	516,2838
5	345	40	3,688879	4	516,4648
6	377	32	3,465736	5	516,1502
7	403	26	3,258097	6	515,9256
8	425	22	3,091042	7	516,6434
9	442	17	2,833213	8	516,3721
10	456	14	2,639057	9	516,3558

<sup>1</sup> Четыркин Е. М. Статистические методы прогнозирования. М. : Статистика. 1975. С. 114—122.

<sup>2</sup> Там же. С. 125—130.

<sup>3</sup> Льюис К. Д. Указ. соч. С. 107—109.

Найдем по нашему ряду цепные абсолютные приросты  $z$  и представим их через параметры нашей функции, т.е.  $z = c - ab^t - c + ab^{t-1} = ab^{t-1}(1 - b)$ . Известно, что для модифицированной экспоненты логарифм абсолютных приростов линейно зависит от фактора времени  $t$ . Следовательно, можно записать, что  $\ln z = \ln a + (t - 1) \ln b + \ln(1 - b)$ . Обозначим  $\ln a + \ln(1 - b)$  через  $d$ . Тогда  $\ln z = d + (t - 1) \ln b$ , т.е. линейное в логарифмах уравнение. Применяя МНК, получим оценки параметров  $d$ ,  $\ln b$ , а соответственно и параметра  $b$ . В рассматриваемом примере на основании граф табл. 5.3  $\ln z$  и  $(t - 1)$  было найдено уравнение регрессии:  $\ln z = 4,519641 - 0,20882(t - 1)$ . Исходя из него получаем  $\ln b = -0,20882$ ;  $b = 0,811538$ .  $4,519641 = \ln a + \ln(1 - b) = \ln[a(1 - b)]$ . Тогда  $a(1 - b) = e^{4,519641}$ , откуда параметр  $a = 91,80264 / (1 - 0,811538) = 487,1145$ .

Далее можно найти оценку параметра  $c$  как среднее значение из величин  $c = y + ab^t$ , найденных для каждого месяца (см. последнюю графу табл. 5.3). Предельная величина расходов на рекламу составит 516,4 тыс. руб. Искомое уравнение тренда примет вид

$$y = 516,4 - 487,1 \cdot 0,81^t.$$

Рассмотренный метод применим, если абсолютные приросты — величины положительные. Если же некоторые приросты окажутся меньше нуля, то нужно проводить сглаживание уровней временного ряда методом скользящей средней.

Для логистической кривой Перла — Рида ( $y = \frac{1}{c + ab^t}$ ) аналогично параметры  $a$  и  $b$  могут быть найдены МНК, если асимптота  $c$  задана. Тогда данная функция преобразовывается в линейную из логарифмов  $\frac{1}{y} = c + ab^t$ ;  $\frac{1}{y} - c = ab^t$ ; обозначим  $\frac{1}{y} - c$  через  $Y$  и прологарифмируем, т.е.  $\ln Y = \ln a + t \ln b$ . Далее параметры  $a$  и  $b$  определяются МНК, как и в примере по табл. 5.3.

Для логистической кривой вида  $y = \frac{c}{1 + be^{-at}}$  параметры  $a$  и  $b$  могут быть оценены МНК, если асимптота  $c$  задана, так как в этом случае функция линеаризуема:  $1 + be^{-at} = \frac{c}{y}$ ;



$\frac{c}{y} - 1 = be^{-at}$  обозначим через  $Y$  величину  $(\frac{c}{y} - 1)$  и прологарифмируем:  $Y = be^{-at}$  и  $\ln Y = \ln b - at$ . Далее, применяя МНК, оцениваем параметры  $a$  и  $b$ .

При практических расчетах значение верхней асимптоты логистической кривой может быть определено исходя из существа развития явления, различного рода ограничений для его роста (нормативы потребления, законодательные акты), а также графически.

Если верхняя асимптота не задана, то для оценки параметров могут использоваться разные методы: Фишера, Юла, Родса, Нейра и др. Сравнительная оценка и обзор этих методов изложены в работе Е. М. Четыркина<sup>1</sup>.

Покажем на примере расчет параметров логистической кривой по методу Фишера.

#### Пример 5.4

Производство продукции характеризуется данными, представленными в табл. 5.4.

**Таблица 5.4.** Расчет параметров логистической кривой

$t$	$y_t$	$z_t = \frac{1}{2} \ln \frac{y_{t+1}}{y_{t-1}}$	$Y = \frac{403}{y_t} - 1$	$\ln Y$	$\hat{y}_t$
1	12	—	32,583	3,484	12,2
2	28	0,788	13,393	2,595	26,2
3	58	0,661	5,948	1,783	54,3
4	105	0,572	2,838	1,043	104,2
5	182	0,453	1,214	0,194	176,8
6	260	0,282	0,55	-0,598	256,4
7	320	0,163	0,259	-1,349	321,0
8	360	0,086	0,119	-2,125	361,7
9	380	—	0,061	-2,805	383,5
$\Sigma$	1705	—	—	2,222	1696,3

<sup>1</sup> Четыркин Е. М. Указ. соч. С. 126—133.

Метод Фишера основан на определении производной для логистической кривой  $y_t = \frac{c}{1 + be^{-at}}$ . Дифференцируя данную функцию по  $t$ , получим уравнение

$$\frac{dy_t}{dt} = ay_t \left(1 - \frac{y_t}{c}\right) \rightarrow \frac{dy_t}{dt} \cdot \frac{1}{y_t} = a \left(1 - \frac{y_t}{c}\right) = a - \frac{a}{c} y_t.$$

Обозначим темп прироста логистической кривой  $\frac{dy_t}{dt} \cdot \frac{1}{y_t}$  через  $z_t$ . Тогда  $z_t = a - \frac{a}{c} y_t$ , т.е. для  $z_t$  имеем линейную функцию с параметрами  $a$  и  $\frac{a}{c}$ . Чтобы найти решение, необходимо оценить  $z_t$ . Предполагая, что интервалы между уровнями в ряду динамики равны, Фишер предложил приближенно оценивать  $z_t$  в виде уравнения  $z_t = \frac{1}{2} \ln \frac{y_{t+1}}{y_{t-1}}$ , где  $t = 2, 3, \dots, n-1$ . Для нашего примера значения  $z_t$  представлены в графе 3 табл. 5.4. Далее применяем МНК к уравнению  $z_t = a - \frac{a}{c} y_t$ , т.е. строим регрессию  $z_t$  от  $y_t$ , беря данные от  $t = 2$  до  $t = 8$ . Уравнение регрессии запишется в виде  $z_t = 0,806 - 0,002y_t$ . Исходя из него находим параметры  $a$  и  $c$  для логистической кривой. Параметр  $a = 0,806$ . Данное уравнение статистически значимо:  $F$ -критерий равен 689,6;  $R^2 = 0,996$ . Соответственно для него значимы и параметры:  $t$ -критерий для параметра  $a$  равен 47,2 и для параметра  $-\frac{a}{c}$  равен -26,2. Так как  $-\frac{a}{c} = -0,002$ , то  $\frac{a}{c} = 0,002$  и  $c = \frac{a}{0,002} = \frac{0,806}{0,002} = 403$ , т.е. верхняя асимптота производства продукции составляет 403 ед.

После того, как найдены параметры  $a$  и  $c$ , находим параметр  $b$ . Для этого функцию  $y_t = \frac{c}{1 + be^{-at}}$  представим как  $\frac{c}{y_t} - 1 = be^{-at}$ . Обозначим через  $Y$  выражение в левой части равенства, т.е.  $Y = \frac{c}{y_t} - 1$ . Тогда имеем уравнение  $Y = be^{-at}$ . Прологарифмируем его:  $\ln Y = \ln b - at$ . В этом уравнении свободным членом является  $\ln b$ . Его можно определить из первого уравнения системы нормальных уравнений, а именно  $\ln b = \overline{\ln Y} + a\bar{t}$ . Для нашего примера имеем уравнение  $\ln b = \frac{1}{9} 2,222 + 0,806 \cdot 5 = 4,2769$ . Соответственно  $b = e^{4,2769} = 72,016$ . Таким образом, логистическая кривая запишется в виде

$$y_t = \frac{403}{1 + 72,016e^{-0,806t}}.$$

Теоретические значения данной функции представлены в графе 6 табл. 5.4 (найжены путем подстановки соответствующих значений  $t$ ). Они достаточно близко подходят к исходным данным: коэффициент корреляции между ними равен 0,999;  $\sum y_t \neq \sum \hat{y}_t$  ввиду того, что в расчетах использовались логарифмы. Если предположить, что предельное значение объема производства продукции равно 400 ед., т.е. применить МНК к уравнению  $\ln Y = \ln b - at$ , то получим  $\ln b = 4,212$  и  $b = =67,5$ ; параметр  $a$  при компьютерной обработке определяется как  $-a = -0,8$ . Соответственно уравнение тренда запишется в виде  $y = 400 / (1 + 67,5e^{-0,8t})$ . Результаты двух уравнений достаточно близки.

Параметры кривой Гомперца  $y_t = ca^{bt}$  также могут быть оценены МНК, если асимптота  $c$  задана, так как в этом случае данная функция сводима к линейному виду  $\frac{y_t}{c} = a^{bt}$ . Прологарифмировав ее, получим уравнение  $\lg(\frac{y_t}{c}) \stackrel{c}{=} b^t \lg a$ . Вторично прологарифмировав, получим уравнение  $\lg(\lg \frac{y_t}{c}) = t \lg b + \lg(\lg a)$ . Обозначив  $\lg(\lg \frac{y_t}{c})$  через  $y^*$ ,  $\lg b$  через  $B$  и  $\lg(\lg a)$  через  $A$ , запишем кривую Гомперца в линейном виде  $y^* = A + Bt$ , для оценки параметров которой применим МНК.

При практическом применении кривой Гомперца могут возникнуть некоторые сложности по динамическому ряду с повышающейся тенденцией. В этом случае задается верхняя асимптота  $c$  и логарифмы  $(y/c) < 0$ . При повторном логарифмировании в расчетах используются лишь положительные значения  $\lg(y/c)$ . Продемонстрируем возможность оценки параметров кривой Гомперца с верхней асимптотой на примере динамики по предприятию товарных запасов на начало каждого месяца (тыс. долл.).

**Таблица 5.5.** Расчет параметров кривой Гомперца

$t$	$y$	$y/250 = Y$	$\lg Y$	$(-1) \lg Y$	$Y^* = \lg[(-1) \lg Y]$	$\hat{y}_t$
1	50	0,2	-0,69897	0,69897	-0,15554	44,75
2	80	0,32	-0,49485	0,49485	-0,30553	76,87
3	114	0,456	-0,34104	0,341035	-0,4672	111,38

Окончание табл. 5.5

$t$	$y$	$y/250 = Y$	$\lg Y$	$(-1) \lg Y$	$Y^* = \lg[(-1) \lg Y]$	$\hat{y}_t$
4	144	0,576	-0,23958	0,239578	-0,62055	143,63
5	170	0,68	-0,16749	0,167491	-0,77601	170,98
6	190	0,76	-0,11919	0,119186	-0,92377	192,68
7	207	0,828	-0,08197	0,08197	-1,08635	209,12
8	219	0,876	-0,0575	0,057496	-1,24036	221,20
9	228	0,912	-0,04001	0,040005	-1,39788	229,88
10	234	0,936	-0,02872	0,028724	-1,54175	236,03
11	242	0,968	-0,01412	0,014125	-1,85002	240,33
12	244	0,976	-0,01055	0,01055	-1,97674	243,33

Исходя из экономических соображений о нецелесообразности чрезмерного увеличения запасов выдвинута гипотеза, что верхняя асимптота не превысит 250 тыс. долл., т.е. в примере  $c = 250$ .

Используя данные предпоследней графы табл. 5.5 и применяя МНК, получим уравнение  $Y^* = 0,037383 - 0,16398t$ , в котором параметры выражены в логарифмах  $A = \lg(\lg a) = 0,037383$ ;  $B = \lg b = -0,16398$ . Далее переходим к искомым параметрам  $a$ ;  $b$ . Потенцируя, определим параметр  $b = 10^{-0,16398} = 0,685522$ . Аналогично определим  $\lg a = 10^{0,037383} = 1,089891$  и для оценки параметра  $a$  потенцируем величину  $(-\lg a)$ . Соответственно получим, что  $a = 10^{-1,089891} = 0,081303$ , и кривая Гомперца примет вид

$$Y = 250 \cdot 0,0813^{0,6855^t}$$

Расчетные (теоретические) значения для этой кривой приведены в последней графе табл. 5.5. В рассматриваемом примере верхняя асимптота может быть определена методом регрессии аналогично, как и для модифицированной экспоненты. В этих целях необходимо преобразовать кривую Гомперца. Прологарифмируем уравнение

$$\ln y = \ln c + \ln a \cdot b^t$$

и получим модифицированную экспоненту  $Y = C + A \cdot b^t$ , где  $Y = \ln y$ ;  $C = \ln c$ ;  $A = \ln a$ . Далее найдем абсолютные приросты  $\ln y$  и выразим их через параметры модифицированной экспоненты

$$\Delta Y = C + Ab^t - C - Ab^{t-1} = Ab^{t-1} \cdot (b - 1).$$

Вновь прологарифмировав это выражение, придем к выражению  $\ln \Delta Y = \ln A + (t - 1) \ln b + \ln (b - 1)$ . Выразив в нем  $\ln A + \ln (b - 1)$  через  $d$ , получим выражение  $\ln \Delta Y = d + (t - 1) \ln b$ , т.е. линейное в логарифмах уравнение, в котором параметры  $d$  и  $b$  могут быть оценены МНК. В нашем примере уравнение принимает вид

$$\ln \Delta Y = -0,34366 - 0,36132 (t - 1).$$

Отсюда  $b = e^{-0,36132} = 0,696755$ ;  $d = 0,34366$ .

Параметр  $d = \ln A + \ln (b - 1) = \ln (\ln a) + \ln (b - 1) = \ln [\ln a \times (b - 1)]$ . Соответственно  $\ln a \cdot (b - 1) = e^{0,34366} = 0,70917$ ;  $\ln a = 0,70917 / (0,696755 - 1) = -2,3386$  и параметр  $a = 0,096463$ . Зная параметры  $a$  и  $b$ , найдем для каждой строки таблицы параметр  $C$  как  $C = Y - Ab^t$ . Далее оценим среднее значение  $C$  и на его основе определим верхнюю асимптоту  $C = \ln c$ ;  $c = e^C$ . Для нашего примера параметр  $c = e^{5,523} = 250,4$ , что практически совпало с выдвинутой ранее гипотезой. Кривая Гомперца в этом варианте расчетов примет вид  $y = 250,4 \cdot 0,0965^{0,69^t}$ .

Некоторую специфику имеет также оценка параметров гиперболы вида  $y = a + \frac{b}{c+t}$ . Если асимптота  $a$  задана (например, исходя из графика временного ряда), то рассматриваемая функция достаточно легко преобразовывается в линейный вид, что позволяет по МНК оценить параметры  $b$  и  $c$ . Преобразование гиперболы  $y = a + \frac{b}{c+t}$  в линейную функцию сводится к следующему. Обозначим  $(y - a)$  через  $Y$ ; т.е.  $Y = b / (c + t)$ . Тогда имеем уравнение  $(1/Y) = (c + t) / b$ . Заменим  $(1/Y)$  на  $z$ , получим  $z = (c + t) / b$ ;  $zb = c + t$ ; откуда  $t = bz - c$ . Далее традиционно найдем оценку параметров  $b$  и  $c$ . Если асимптота  $a$  не задана заранее, то рассматриваемая гипербола может быть сведена к линейной множественной регрессии, что позволяет оценить ее параметры по МНК.

Из уравнения  $y = a + b / (c + t)$  следует, что  $y(c + t) = a(c + t) + b$  или  $cy + ty = ac + at + b$ . В этом равенстве  $ty$  представляет собой при расчетах ряд значений  $t^*y$ , которые обозначим через  $Y$ . Тогда имеем  $Y = (ac + b) + at - cy$ . Далее обозначим  $(ac + b)$  через  $A$  и получим линейную множественную регрессию  $Y = A +$

$+ at - cy$ , в которой параметры  $A$ ,  $a$ ,  $c$  могут быть найдены по МНК. Далее можно оценить параметр  $b$  как  $b = A - ac$ . Если в первом варианте расчетов асимптота  $a$  была задана правильно, то оба метода дают одинаковые оценки параметров. Однако второй подход предпочтителен, ибо его результаты не зависят от субъективизма исследователя. Вместе с тем следует отметить, что использование формулы гиперболы не всегда приводит к лучшим результатам, чем описание тенденции с помощью равностоярной гиперболы  $y = a + b/t$ , параметры которой можно оценить по МНК.

### 5.3.3. Оценка адекватности модели тенденции

Модель тенденции считается адекватной реальному процессу, если теоретические (найденные по уравнению тренда) уровни ряда достаточно близко подходят к фактическим их значениям, т.е.  $y_t$  и  $\hat{y}_t$  мало отличаются друг от друга. Для оценки адекватности модели проводится анализ остатков ( $e_t = y_t - \hat{y}_t$ ).

Модели тенденции можно сравнивать по величине остаточной суммы квадратов

$$S^2 = \sum (y_t - \hat{y}_t)^2. \quad (5.22)$$

Чем меньше эта величина, тем в большей мере уравнение тренда подходит для описания тенденции временного ряда. Обратимся к примеру 5.1 динамики числа зарегистрированных ДТП (на 100 000 человек населения) по Новгородской области за 2000—2008 гг.

Предположим, что было рассчитано не только уравнение экспоненты, но и уравнение линейного тренда  $Y = 69,1 + 22,1t$ .

Для линейного тренда остаточная сумма квадратов составила 3874,62, а для экспоненты 2617,701. Следовательно, экспонента лучше описывает тенденцию ряда.

Другим показателем при выборе функции тренда является коэффициент детерминации  $R^2$ . Чем выше  $R^2$ , тем соответственно выше вероятность того, что данная модель тенденции описывает исходные данные. В нашем примере  $R^2$  для экспоненты составил 0,9202, а для линейного тренда 0,8832, подтверждая

еще раз, что экспонента в большей мере подходит для описания тенденции.

Величина  $1 - R^2$  отражает влияние случайной составляющей, т.е. показывает, какая доля вариации уровней динамического ряда не связана с тенденцией. Так, в нашем примере для экспоненты  $1 - R^2 = 0,0798$ . Это означает, что лишь 8% вариации уровней динамического ряда не связаны с рассматриваемой тенденцией.

Однако рассмотренные критерии адекватности модели тенденции ( $S^2$  и  $R^2$ ) могут привести к неправильным выводам, если не учитывать статистическую значимость параметров уравнения тренда. Общеизвестно, что можно по ряду из  $n$  точек построить полином степени  $(n - 1)$ , и он пройдет через все точки ряда. В этом случае  $S^2 = 0$ , а  $R^2 = 1$ . Однако в таком случае уравнение  $y_t = a + b_1t + b_2t^2 + \dots + b_{n-1}t^{n-1}$  описывает исходные уровни ряда, но не является моделью тенденции, ибо отражает не только тенденцию, но и влияние случайной компоненты.

Если для нашего примера построить параболу пятой степени, то получим уравнение  $\hat{y}_t = 288,6 - 356,111t + 223,3635t^2 - 56,9569t^3 + 6,4075t^4 - 0,2609t^5$ .

Однако по  $t$ -критерию Стьюдента при 5%-ном уровне существенности параметры этой параболы оказываются статистически незначимыми. Поэтому, хотя для этой функции коэффициент детерминации и выше (0,98585), ее нельзя считать лучшей формой уравнения тренда.

Уравнение тренда хорошо описывает тенденцию, если отсутствует автокорреляция в остатках ( $y_t - \hat{y}_t$ ), т.е. остатки текущего периода не коррелируют с остатками предыдущего периода.

Измерить автокорреляцию в остатках можно с помощью **коэффициента автокорреляции остатков**

$$r_{a_t} = \frac{\bar{e}_t \bar{e}_{t-1} - \bar{e}_t \bar{e}_{t-1}}{\sigma_{e_t} \sigma_{e_{t-1}}}, \quad (5.23)$$

где  $e_t = y_t - \hat{y}_t$ , т.е. остатки текущего периода;  $e_{t-1}$  — остатки предыдущего периода.

Иными словами, автокорреляция в остатках оценивается так же, как и автокорреляция уровней ряда, с тем лишь

отличием, что в расчетах используются остаточные величины  $e_t$ , а не уровни динамического ряда  $y_t$ .

### Пример 5.5

Динамика численности детей в возрасте от 8 до 13 лет в  $N$ -ской области за последние 15 лет характеризуется параболой второй степени

$$\hat{y}_t = 110,5 + 10,8t - 0,8t^2,$$

где  $y_t$  — численность детей, тыс. человек;  $t = 1, 2, 3, \dots, 15$ ;  $R^2 = 0,967$ ; все параметры уравнения тренда статистически значимы.

Для рассматриваемой модели остаточные величины за 15 лет составили:

$t$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
$e$	6	0,2	-2,6	-4	-5	-3,1	0,7	2,8	2,5	2,9	3,3	1,6	-1,0	-4	-0,3

Исходя из формулы (5.23) коэффициент автокорреляции остатков окажется равным 0,627824.

Его величина не столь мала, чтобы утверждать об отсутствии автокорреляции остатков. Очевидно, данное уравнение тренда не является наилучшим, ибо нарушена предпосылка МНК об отсутствии автокорреляции остатков.

Для комплексной оценки адекватности модели тенденции можно пользоваться и другими характеристиками, которые обычно рассматриваются при оценке качества регрессионных моделей, например средней ошибкой аппроксимации, показателями асимметрии и эксцесса для остаточных величин.

### 5.3.4. Тесты на наличие автокорреляции остатков первого порядка

Проверка модели на автокоррелированность остатков обычно проводится с помощью критерия Дарбина — Уотсона. Предполагается, что при наличии автокорреляции остатков верно соотношение



$$e_t = \rho e_{t-1} + v_t, \quad (5.24)$$

где  $\rho$  — коэффициент автокорреляции остатков;  $v_t$  — независимые случайные величины. Если  $\rho = 0$ , то отсутствует автокорреляция остатков и оценка параметров МНК является наилучшей линейной несмещенной оценкой (НЛНО). Если  $\rho \neq 0$ , то оценка параметров МНК является смещенной и выводы по модели не корректны. Таким образом, при тестировании на автокорреляцию остатков выдвигается нулевая гипотеза  $H_0: \rho = 0$  против альтернативной гипотезы  $H_1: \rho \neq 0$ .

Статистика Дарбина — Уотсона определяется соотношением

$$D-W = \frac{\sum_{t=2}^n (e_t - e_{t-1})^2}{\sum_{t=1}^n e_t^2}. \quad (5.25)$$

Критерий Дарбина — Уотсона и коэффициент автокорреляции остатков связаны между собой соотношением

$$D-W \cong 2(1-\rho). \quad (5.26)$$

Убедиться в этом можно, если преобразуем формулу коэффициента автокорреляции остатков и раскроем числитель формулы критерия Дарбина — Уотсона.

Параметры уравнения тренда обычно оцениваются МНК и в соответствии с предпосылками МНК  $\sum_{t=1}^n e_t = 0$  и  $\bar{e}_t = \frac{\sum_{t=1}^n e_t}{n} = 0$ .

Поскольку ряд  $e_{t-1}$  сдвинут по отношению к ряду  $e_t$  на один временной интервал, то вариации в этих рядах мало отличаются друг от друга и можно предположить, что  $\sigma_{e_t} \cong \sigma_{e_{t-1}}$ . Тогда коэффициент автокорреляции остатков можно представить в следующем виде:

$$\rho = \frac{\sum_{t=2}^n e_t e_{t-1}}{\sum_{t=1}^n e_t^2}. \quad (5.27)$$

Данная формула является приближенной, ибо если  $e_1$  и  $e_n$  сильно различаются между собой, то и  $\sigma_{e_t} \neq \sigma_{e_{t-1}}$ .

Для нашего примера имеем соотношения

$$\sum_2^{15} e_t e_{t-1} = 79,07; \sum_1^{15} e_t^2 = 110,94 + 6^2 = 146,94.$$

Соответственно коэффициент автокорреляции остатков окажется равным

$$\rho = \frac{79,07}{146,94} = 0,538111,$$

что несколько ниже найденного ранее значения (0,627824), но также демонстрирует заметную связь соседних величин остатков, т.е. их автокорреляцию.

Далее преобразуем формулу критерия Дарбина — Уотсона

$$\begin{aligned} D-W &= \frac{\sum_{t=2}^n (e_t - e_{t-1})^2}{\sum_{t=1}^n e_t^2} = \frac{\sum_{t=2}^n e_t^2 - 2 \sum_{t=2}^n e_t e_{t-1} + \sum_{t=2}^n e_{t-1}^2}{\sum_{t=1}^n e_t^2} = \\ &= \frac{\sum_{t=2}^n e_t^2 - e_1^2 + \sum_{t=2}^n e_{t-1}^2 - e_n^2 - 2 \sum_{t=2}^n e_t e_{t-1}}{\sum_{t=1}^n e_t^2}. \end{aligned}$$

Разделим каждое слагаемое числителя на  $\sum_{t=1}^n e_t^2$  и предположим, что  $\sum e_t^2 = \sum e_{t-1}^2$ :

$$1 - \frac{e_1^2}{\sum_{t=1}^n e_t^2} + 1 - \frac{e_n^2}{\sum_{t=1}^n e_t^2} - 2 \frac{\sum_{t=1}^n e_t e_{t-1}}{\sum_{t=1}^n e_t^2} = 2 \left( 1 - \frac{\sum_{t=1}^n e_t e_{t-1}}{\sum_{t=1}^n e_t^2} \right) - \frac{e_1^2 + e_n^2}{\sum_{t=1}^n e_t^2}.$$

При большем числе наблюдений ( $e_1^2 + e_n^2$ ) значительно меньше  $\sum_{t=1}^n e_t^2$ . Поэтому приближенно можно считать, что  $D-W \cong 2(1-\rho)$ . Из этого соотношения очевидно, что при полной положительной автокорреляции остатков ( $\rho=1$ ) критерий  $D-W=0$ , а при полной отрицательной автокорреляции ( $\rho=-1$ ) критерий  $D-W=4$ . Если же автокорреляция в остатках отсутствует, т.е.  $\rho=0$ , то  $D-W=2$ . Иными словами критерий Дарбина — Уотсона изменяется в пределах:  $0 \leq D-W \leq 4$ .

Дарбин и Уотсон разработали пороговые значения показателя  $D - W$ , позволяющие принять или отвергнуть гипотезу об отсутствии автокорреляции в остатках.

При заданном числе наблюдений  $n$  (длина динамического ряда) и  $m$  параметров при  $t$  в уравнении тренда (или  $m$  объясняющих переменных в уравнении регрессии) установлены при 5%-ом уровне значимости верхняя  $D - W_u$  ( $u$  — upper) и нижняя  $D - W_l$  ( $l$  — low) границы критерия.

Фактическое значение критерия  $D - W$  сравнивается с табличными значениями. Если  $D - W < 2$ , то возможны следующие варианты:

- 1) при  $D - W < D - W_l$  (нижней границы) нулевая гипотеза об отсутствии автокорреляции отвергается и делается вывод о наличии положительной автокорреляции в остатках;
- 2) при  $D - W > D - W_u$  (верхней границы) нулевая гипотеза об отсутствии автокорреляции не отвергается, т.е. делается вывод об отсутствии корреляционной связи последующих остатков с предыдущими;
- 3) при  $D - W_l \leq D - W \leq D - W_u$  нельзя ни отвергнуть, ни принять нулевую гипотезу об отсутствии автокорреляции в остатках, т.е. значение  $D - W$  попало в область неопределенности и необходимы дальнейшие исследования, например по большему числу наблюдений.

Если фактическое значение  $D - W > 2$ , что означает отрицательную автокорреляцию, то с пороговыми табличными значениями сравнивается величина  $4 - D - W$ .

При этом возможны следующие варианты:

- 1) при  $4 - DW < D - W_l$  отвергается нулевая гипотеза об отсутствии автокорреляции и делается вывод о наличии отрицательной автокорреляции в остатках;
- 2) при  $4 - DW > D - W_u$  нулевая гипотеза об отсутствии автокорреляции в остатках принимается;
- 3) при  $D - W_l \leq 4 - DW \leq D - W_u$  нельзя сделать определенного вывода о наличии или отсутствии автокорреляции в остатках по имеющимся данным.

По величине критерия Дарбина — Уотсона можно определить величину коэффициента автокорреляции остатков исходя из соотношения:  $D - W \cong 2(1 - \rho)$ . Отсюда  $\frac{D - W}{2} \cong 1 - \rho$

и соответственно  $\rho \cong 1 - \frac{D-W}{2}$ . Поэтому если  $D-W > 2$ , то  $\rho < 0$ , а при  $D-W < 2$  имеем  $\rho > 0$ .

Таким образом, если фактическое значение критерия Дарбина — Уотсона не слишком отличается от 2, то можно сделать вывод об отсутствии автокорреляции в остатках. Использование критерия Дарбина — Уотсона можно показать графически (рис. 5.15).

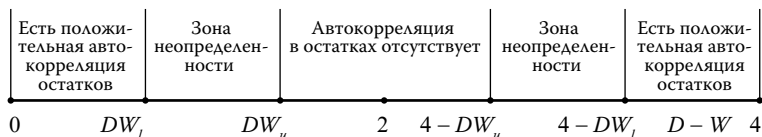


Рис. 5.15. Схема применения критерия Дарбина — Уотсона

Обратимся к рассмотренному ранее примеру и найдем критерий Дарбина — Уотсона, используя формулу (5.25):

$$D - W = \frac{\sum_{t=2}^{15} (e_t - e_{t-1})^2}{\sum_{t=1}^{15} e_t^2}.$$

В примере  $\sum (e_t - e_{t-1})^2 = 99,65$ ;  $D - W = \frac{99,65}{146,64} = 0,678$ .

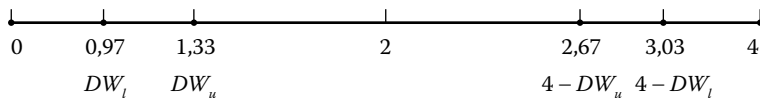
При  $n = 15$  и  $m = 2D - W_l = 0,95$ . Фактическое значение критерия меньше его нижней границы, что подтверждает вывод о наличии в остатках автокорреляции и означает плохой выбор модели тенденции.

По временным рядам экономических данных автокорреляция остатков чаще всего бывает положительной, т.е. сохраняется знак остатка при переходе к следующему периоду времени (за положительными остатками следуют в основном положительные, а за отрицательными — отрицательные). В этом случае альтернативная гипотеза  $H_1$  предполагает, что  $\rho > 0$ . Вместе с тем может иметь место и отрицательная автокорреляция, когда положительные и отрицательные остатки чередуются.

Предположим, что динамика просроченной задолженности за 12 мес. года характеризовалась на предприятии уравнением тренда  $\hat{y}_t = 410,1 - 6,5t$  при  $t = 1, 2, \dots, 12$ ;  $R^2 = 0,921$ . Критерий Дарбина — Уотсона составил 2,226. Требуется оценить автокорреляцию в остатках.

Так как фактическое значение  $D - W > 2$ , то речь идет об оценке отрицательной автокорреляции в остатках. Приближенно

$\rho = 1 - \frac{2,226}{2} = -0,113$ , что означает слабую автокорреляцию в остатках. Убедимся в этом через сравнение фактического значения  $4 - DW$  с табличными значениями при 5%-ном уровне значимости  $4 - DW = 4 - 2,226 = 1,774$ . При  $n = 12$  и  $m = 1$   $D - W_l = 0,97$  и  $D - W_u = 1,33$ . Так как  $4 - DW > D - W_u$  (верхней границы), то автокорреляция в остатках отсутствует и линейный тренд достаточно хорошо описывает исходный динамический ряд. Ответ на данный вопрос можно дать и графически. С этой целью необходимо нанести на числовую ось в интервале  $[0; 4]$  табличные значения критерия  $D - W$ , а именно  $D - W_l = 0,97$ ;  $D - W_u = 1,33$ ;  $4 - DW_u = 2,67$ ;  $4 - DW_l = 3,03$ , и посмотреть, в какой интервал попадает фактическое значение критерия  $D - W = 2,226$  (рис. 5.16).



**Рис. 5.16.** Интервалы для критерия Дарбина — Уотсона

Фактическое значение  $D - W = 2,226$  попадает в промежуток от  $D - W_u$  до  $4 - DW_u$ . Следовательно, можно принять нулевую гипотезу об отсутствии автокорреляции в остатках.

Критерий Дарбина — Уотсона используется и для оценки автокорреляции по регрессионным моделям, построенным по временным рядам. Однако его нельзя использовать, если в модели среди объясняющих переменных содержатся лаговые значения результативного признака. Например, инвестиции в основной капитал в текущем периоде рассматриваются не только как функция разных экономических переменных

(прибыли, процентной ставки), но и достигнутого в предыдущий период размера инвестиций.

Критерий Дарбина — Уотсона является наиболее распространенным в современной эконометрике. Вместе с тем для проверки некоррелированности ошибок могут использоваться асимптотические тесты. Одним из них является тест Бреуша — Годфри (*Breusch — Godfrey*). Выдвигается гипотеза о независимости случайных величин  $E_t$ , разделенных  $K$  периодами времени и очищенных от влияния промежуточных значений остатков. Иными словами, модель зависимости остатков имеет вид

$$E_t = a_1 E_{t-1} + a_2 E_{t-2} + \dots + a_k E_{t-k} + v_t,$$

где  $E_t, E_{t-1}, E_{t-2}, \dots, E_{t-k}$  — остатки по линейной регрессионной модели, построенной по временным рядам ( $y_t = A_0 + A_1 x_{1t} + \dots + A_p x_{pt} + E_t$ );  $v_t$  — независимые случайные величины. В модель остатков ( $E_t$ ) могут быть включены также  $X_1, X_2, \dots, X_p$ , и она примет вид  $E_t = b_1 X_1 + \dots + b_p X_p + a_1 E_{t-1} + a_2 E_{t-2} + \dots + a_k E_{t-k} + v_t$ .

Как видим, в этом критерии рассматривается автокорреляция остатков  $k$ -го порядка в отличие от автокорреляции первого порядка, оцениваемой критерием Дарбина — Уотсона. В тесте Бреуша — Годфри проверяется нулевая гипотеза  $H_0: a_1 = a_2 = \dots = a_k = 0$ . Если эта гипотеза верна, то при большом количестве наблюдений ( $n$ ) статистика критерия  $nR^2$  (где  $R^2$  — коэффициент детерминации по модели регрессии для остатков) имеет распределение, близкое к распределению хи-квадрат ( $\chi^2$ ) с  $K$  степенями свободы. Нулевая гипотеза отвергается, если вычисленное значение  $nR^2$  превышает критическое (табличное) значение  $\chi^2$  при заданном уровне значимости  $\alpha$ , т.е. если имеем соотношение  $nR^2 > \chi^2_{1-\alpha(K)}$ .

При  $K = 1$  тест Бреуша — Годфри оценивает автокорреляцию остатков первого порядка и поэтому сопоставим с критерием Дарбина — Уотсона. Выводы по этим критериям практически не расходятся. Так, в работе М. Вербик<sup>1</sup> приводится пример изучения спроса на мороженое в зависимости от семейного дохода, цены и температуры воздуха по данным за 34 недели. Фактическое значение статистики Дарбина — Уотсона по моде-

<sup>1</sup> Вербик М. Путеводитель по современной эконометрике. М : Научная книга. 2008. С. 174—179.

ли регрессии составило 1,0212, что меньше нижнего критического значения и означает наличие в остатках автокорреляции первого порядка. Этот же вывод подтверждает и тест Бреуша — Годфри. При  $R^2 = 0,149$  для модели по остаткам статистика Бреуша — Годфри составит 5,066, значение которой превышает 5%-ное критическое значение 3,84 из распределения  $\chi^2$  с одной степенью свободы. Учитывая, что рассматриваемый временной ряд не слишком велик, Вербик в расчетах использует величину  $(n - 1) R^2$ , но выводы от этого не меняются.

Несмотря на то, что критерий Бреуша — Годфри является асимптотическим, а критерий Дарбина — Уотсона точным и поэтому более широко используемым, в ряде случаев целесообразно применение критерия Бреуша — Годфри. Так, его можно использовать и в моделях регрессии с лаговыми значениями результативного признака среди объясняющих переменных. Этот критерий учитывает наличие в остатках автокорреляции не только первого порядка. Поэтому его называют также *тестом серий*.

Другим асимптотическим тестом для оценки наличия в остатках автокорреляции может выступать  $Z$ -статистика. Полагаем, что истинное значение коэффициента автокорреляции  $\rho$  неизвестно и исследователь имеет дело с выборочными оценками  $\rho_B$ , которые при большом объеме выборки  $n$  имеют асимптотически нормальное распределение  $\sqrt{n} \cdot (\rho_B - \rho) \rightarrow N(0, 1 - \rho^2)$ , где  $\rightarrow$  обозначает «асимптотически распределено как» с математическим ожиданием для  $\rho$ , равным нулю, и дисперсией  $(1 - \rho^2)/n$ . Тогда  $Z$ -статистика составит

$$Z = (\rho_B - \rho) / [(1 - \rho^2)/n]^{0,5}.$$

$Z$  приближенно имеет стандартное нормальное распределение. Если нулевая гипотеза  $H_0$  верна и  $\rho = 0$ , то  $Z$ -статистика окажется равной  $Z = \rho_B \cdot \sqrt{n}$ , что и используется в качестве критической величины. Нулевая гипотеза отклоняется на 5%-ном уровне значимости, если фактическое значение  $Z > 1,96$  и делается вывод о наличии в остатках автокорреляции первого порядка. По рассмотренному ранее примеру из работы Вербик имеем при  $\rho_B = 0,401$  и  $n = 34$   $Z = 2,34$ , что больше 1,96 и подтверждает сделанный ранее вывод о наличии в остатках автокорреляции первого порядка.

Рассмотренные асимптотические тесты предполагают длинные ряды динамики, а критерий Дарбина — Уотсона основан на теории малых выборок. Поэтому он и получил наиболее широкое применение. Кроме изложенных тестов для оценки в остатках автокорреляции могут применяться и другие тесты. Так, в компьютерном пакете «*Econometric Views*» приводится не только тест Бреуша — Годфри, но и  $Q$ -тест Льюинга — Бокса, который основан на изучении частной автокорреляционной функции. Ее исследование важно в авторегрессионных моделях, которые рассматриваются в других главах учебника.

### 5.3.5. Доверительные интервалы прогноза по трендовым моделям

Модели тенденций широко используются в экстраполяции. Прогнозирование по ним означает нахождение интервалов прогноза. Поскольку модели тенденций по сути есть модели регрессии, в которых в качестве объясняющей переменной используется фактор времени, то интервалы прогноза по трендовым моделям могут быть даны так же, как и по модели регрессии. Это значит, что по уравнению тренда строится точечный прогноз  $y_p$  и далее определяется ошибка прогноза  $S_p$ . Ошибка прогноза зависит от колеблемости уровней ряда  $y_t$  вокруг тренда  $\hat{y}_t$ , а также от длины временного ряда и периода упреждения. Колеблемость уровней ряда вокруг тренда измеряется в виде стандартного отклонения  $S_y$ , которое представляет собой корень квадратный из остаточной дисперсии на одну степень свободы

$$S_y = \sqrt{\frac{\sum (y_t - \hat{y}_t)^2}{n - m - 1}},$$

где  $y_t$  — фактические значения уровней временного ряда;  $\hat{y}_t$  — расчетные значения уровней временного ряда, найденные по уравнению тренда;  $n$  — длина временного ряда;  $m$  — число параметров при факторе времени  $t$ .

Чем больше стандартное отклонение  $S_y$ , тем шире должен быть интервал прогноза для тренда. Доверительный интервал прогноза для тренда составит  $\hat{y}_t \pm t_\alpha S_y$ . При прогнозе индивидуального значения уровня временного ряда ошибка прогноза



$S_p$  определяется как  $S_p = S_y Q$ , где  $Q$  — поправочный коэффициент, величина которого зависит от длины временного ряда  $n$  и периода упреждения  $L$ .

Чем длиннее временной ряд, тем меньше ошибка прогноза  $S_p$ . Чем больше период упреждения (период, на который дается прогноз), тем больше ошибка прогноза  $S_p$ . Величина  $Q$  зависит также от вида уравнения тренда. Для линейного тренда  $y = a + bt$  поправочный коэффициент  $Q$  определяется по формуле

$$Q = 1 + \frac{1}{n} + \frac{(t_p - \bar{t})^2}{\sum (t - \bar{t})^2}. \quad (5.28)$$

Так как  $t$  представляет собой ряд натуральных чисел, то  $\bar{t} = (n + 1)/2$ ;  $\sum (t - \bar{t})^2 = n(n^2 - 1)/12$ ;  $t_p = n + L$ ;  $t_p - \bar{t} = 0,5(n + 2L - 1)$ . Отсюда получаем формулу  $Q = \sqrt{\frac{n+1}{n} + \frac{3(n+2L-1)^2}{n(n^2-1)}}$ .

Из этой формулы очевидно, что величина  $Q$  зависит от длины временного ряда  $n$  и периода упреждения  $L$ .

Аналогично будет определяться величина  $Q$  и для нелинейных функций, сводимых при преобразовании к линейному виду. Для полиномов второй и более высоких степеней также возможно найти алгебраическое выражение для поправочного коэффициента  $Q$ .

В результате получим громоздкие формулы, которые при компьютерной обработке излишни. Ошибку прогноза  $S_p$  независимо от вида уравнения тренда определим матричным методом

$$S_p = S_y (1 + T_L^* (T^* T)^{-1} T_L)^{0,5},$$

где  $T_L$  — вектор прогнозных значений фактора времени, т.е.  $T_L = 1$ ;  $t_L$ ;  $T_L^*$  — вектор-строка прогнозных значений фактора времени;  $(T^* T)$  — матрица значений коэффициентов при факторе времени в системе нормальных уравнений;  $(T^* T)^{-1}$  — обратная матрица значений коэффициентов при факторе времени.

Доверительный интервал прогноза уровня ряда для конкретного периода времени составит  $y_p \pm t_{\alpha} S_p$ . Рассмотрим построение прогноза для параболы второй степени.

**Пример 5.6**

За 7 мес. ставка межбанковского кредита изменялась следующим образом:

$t$	1	2	3	4	5	6	7
$y$	21,5	22	21	20	18	14	9,5

Уравнение параболического тренда составило

$$Y = 19,85714 + 2,035714t - 0,5t^2, R^2 = 0,996415.$$

$$t \quad 37,40 \quad 6,69 \quad -13,45$$

Как видим, уравнение параболического тренда хорошо описывает исходные данные и все параметры статистически значимы. Чтобы дать интервальный прогноз, найдем стандартное отклонение  $S_y$ . При компьютерной обработке его величина выдается в результатах регрессии  $S_y = 0,340693$ . Прогноз на 8-й мес. означает, что транспонированный вектор прогнозных значений примет вид (1; 8; 64).

Матрица значений коэффициентов при  $t$  составила ( $T^T T$ )

7	28	140
28	140	784
140	784	4676

Ее определитель равен 16 464. Обратная матрица  $(T^T T)^{-1}$  имеет вид

2,428571	-1,28571	0,142857
-1,28571	0,797619	-0,09524
0,142857	-0,09524	0,011905

В результате ошибка прогноза окажется равной

$$S_p = 0,340693 \cdot (1 + 2,428597)^{0,5} = 0,630843.$$

Точечный прогноз на 8-й мес. составит 4,142857%, а интервальный прогноз с вероятностью 0,95 при  $t_\alpha = 2,7764$  окажется в диапазоне  $2,391\% < y_p < 5,894\%$ .

---

## 5.4. Моделирование периодических колебаний

При исследовании длинных экономических временных рядов может возникать потребность выделять жизненные циклы в развитии явления. Это могут быть циклы в сфере производства, строительства, рынка товаров и услуг с длительностью от несколько месяцев до нескольких лет. Иными словами, в ряду динамики может иметь место циклическая или сезонная составляющая, которая должна быть учтена при построении модели динамического ряда. Пути изучения периодических колебаний могут быть разными. Одним из них является разложение временного ряда в ряд Фурье, которое в дальнейшем привело к развитию аппарата спектрального анализа.

### 5.4.1. Ряд Фурье

Ряд Фурье представляет декомпозицию динамического ряда на составляющие, которые связаны с частотой колебаний уровней. Его построение зависит от наличия или отсутствия тенденции в ряду динамики. При отсутствии тенденции, т.е. при стационарном динамическом ряде, методика построения ряда Фурье применяется непосредственно к уровням динамического ряда. Если же в ряде динамики наблюдается тенденция, то ряд Фурье применяется к отклонениям от тенденции. Соответственно эти различия учитываются и при прогнозировании:

- по стационарному временному ряду прогноз дается по ряду Фурье (см. 5.4.1.1);
- по ряду с тенденцией производится суммарный прогноз, т.е. сначала строится прогноз исходя из тенденции развития уровней ряда и далее к нему прибавляется прогноз по ряду Фурье отклонений от тренда (см. 5.4.1.2).

### 5.4.1.1. Ряд Фурье по стационарному ряду

Стационарный ряд с периодическими колебаниями представлен на рис. 5.17.

Уровни ряда варьируют вокруг среднего значения  $\bar{y}$ , а их колебания (волны) повторяются. Интервал времени, необходимый, чтобы динамический ряд начал повторяться, называется **периодом** и обозначен на графике  $P$ . Его величина (расстояние между пиками или впадинами) составляет на графике 10 мес. (12 – 2). Если ряд имеет период  $P$ , то он, как правило, имеет также период  $2P$ ,  $3P$  и т. п. В общем случае для стационарного периодического временного ряда справедливо равенство:  $y_t = y_{t+P}$ , где  $c = 1, 2, \dots$

Величина, обратная периоду, называется **частотой динамического ряда** ( $f$ ):  $f = 1/P$ . Частота указывает на число повторений цикла в единицу времени:  $f = 1/10$  в месяц (по графику).

Отклонение от среднего уровня до пика (или впадины) называется **амплитудой временного ряда** (на графике —  $A$ ).

Расстояние между началом отсчета времени ( $t = 0$ ) и ближайшим пиковым значением называется **фазой** ( $F$ ).

Стационарный периодический временной ряд можно задать четырьмя параметрами: периодом  $P$  или частотой  $f$ , амплитудой  $A$ , фазой  $F$  и средним значением  $\bar{y}$ , что может быть представлено в виде

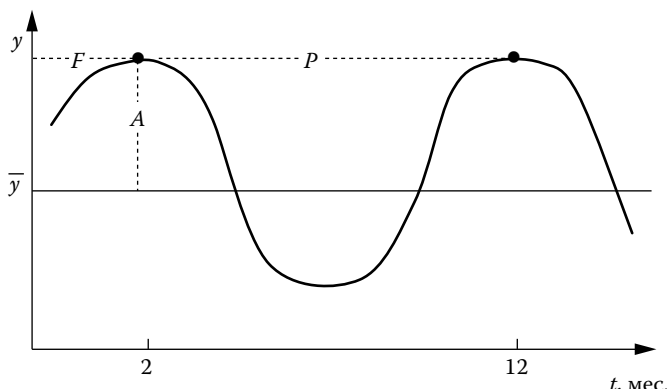


Рис. 5.17. Стационарный ряд с периодическими колебаниями

$$y_t = \bar{y} + A \cos W(t - F), \quad (5.29)$$

где  $W$  — угловая частота, измеряемая в радианах в единицу времени и равная  $W = 2\pi f$ ;  $0 \leq W \leq 2\pi$ ;  $F$  — фаза.

Пусть, например имеем ряд, представленный на рис. 5.18.

Для этого ряда  $\bar{y} = 20$ ;  $P = 12$  мес.;  $f = 1/12$ ;  $A = y_{\max} - \bar{y} = 30 - 20 = 10$ ;  $F = 6$  мес. Учитывая, что  $y_t = \bar{y} + A \cos W(t - F)$ , где  $W = 2\pi f = 2\pi \cdot 1/12 = \pi/6$ , тогда для  $t = 6$  мес. имеем  $y_t = 20 + 10 \cos \frac{\pi}{6}(6 - 6) = 20 + 10 \cos 0 = 20 + 10 \cdot 1 = 30$ ; для  $t = 12$  мес. имеем  $y_t = 20 + 10 \cos \frac{\pi}{6}(12 - 6) = 20 + 10 \cos \pi = 20 + 10(-1) = 10$ ; для  $t = 18$  мес. имеем  $y_t = 20 + 10 \cos \frac{\pi}{6}(18 - 6) = 20 + 10 \cos 2\pi = 20 + 10 \cdot 1 = 30$ .

Рассмотренное выражение  $y_t = \bar{y} + A \cos W(t - F)$  называется **гармоническим представлением ряда** и часто записывается через синусы и косинусы без упоминания о фазе:

$$y_t = \bar{y} + a \cos Wt + b \sin Wt, \quad (5.30)$$

где  $a = A \cos F$  и  $b = A \sin F$ .

Ввиду того, что  $\cos^2 x + \sin^2 x = 1$ ,  $a^2 + b^2 = A^2$ , т.е. существует взаимосвязь между амплитудой колебаний и параметрами гармоника, параметры гармоника также связаны с фазой ряда

$$\operatorname{tg} F = \frac{\sin F}{\cos F} = \frac{b}{a}, \text{ или } \operatorname{arctg} \frac{b}{a} = F.$$

Теоретически стационарный временной ряд с периодическими колебаниями может быть представлен как сумма среднего

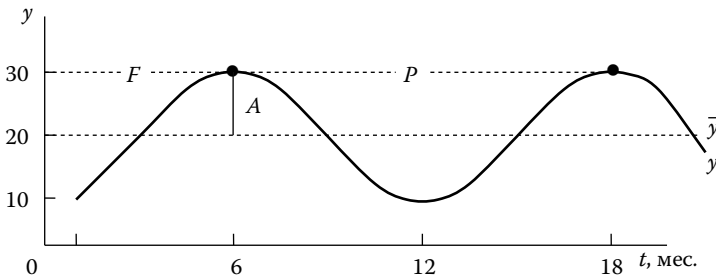


Рис. 5.18. Периодический стационарный ряд

значения и ряда синусоид и косинусоид, что и называется **рядом Фурье**:

$$y_t = \bar{y} + \sum_{i=1}^{\infty} a_i \cos W_i t + \sum_{i=1}^{\infty} b_i \sin W_i t. \quad (5.31)$$

Анализируемые ряды динамики обычно имеют конечную длину  $N$ . Поэтому ряд Фурье приобретает вид

$$y_t = \bar{y} + \sum_{i=1}^n a_i \cos W_i t + \sum_{i=1}^n b_i \sin W_i t, \quad (5.32)$$

где  $n = N/2$  ( $N$  — длина временного ряда).

При замене  $\bar{y}$  параметром  $a_0$  ряд Фурье принимает вид

$$y_t = a_0 + \sum_{i=1}^n a_i \cos W_i t + \sum_{i=1}^n b_i \sin W_i t. \quad (5.33)$$

Оценка параметров данного уравнения обычно дается МНК. Покажем его применение для случая одной гармоник:

$$y_t = a_0 + a_1 \cos t + b_1 \sin t, \quad (5.34)$$

где  $t$  принимает значения от нуля с последующим увеличением на  $\frac{2\pi}{N}$ .

Система нормальных уравнений примет вид

$$\begin{cases} Na_0 + a_1 \sum \cos t + b_1 \sum \sin t = \sum y_t, \\ a_0 \sum \cos t + a_1 \sum \cos^2 t + b_1 \sum \sin t \cos t = \sum y_t \cos t, \\ a_0 \sum \sin t + a_1 \sum \cos t \sin t + b_1 \sum \sin^2 t = \sum y_t \sin t. \end{cases} \quad (5.35)$$

В этой системе  $\sum \cos t = \sum \sin t = 0$ . Поэтому из первого уравнения системы получаем, что  $a_0 = \frac{\sum y_t}{N} = \bar{y}$ . Так как  $\sum \sin t \cos t = 0$  (табл. 5.6), то из второго уравнения системы получим оценку параметра  $a_1$ , а из третьего — параметра  $b_1$ .

$$a_1 = \frac{\sum y_t \cos t}{\sum \cos^2 t}; \quad (5.36)$$

$$b_1 = \frac{\sum y_t \sin t}{\sum \sin^2 t}. \quad (5.37)$$

Так как  $\cos^2 t = \frac{1 + \cos 2t}{2}$ , то

$$\sum \cos^2 t = \frac{1}{2} \sum (1 + \cos 2t) = \frac{1}{2} (N + \sum \cos 2t) = \frac{1}{2} (N + 0) = \frac{N}{2}.$$

Аналогично  $\sum \sin^2 t = \frac{N}{2}$ , так как  $\sum \sin^2 t = \sum (1 - \cos^2 t) = N - \frac{N}{2} = \frac{N}{2}$ . Следовательно, параметры гармонии определяются как

$$a_1 = \frac{2 \sum y_t \cos t}{N} \text{ и } b_1 = \frac{2 \sum y_t \sin t}{N}.$$

Ряд Фурье с двумя гармониками имеет вид

$$y_t = a_0 + a_1 \cos t + b_1 \sin t + a_2 \cos 2t + b_2 \sin 2t. \quad (5.38)$$

При этом параметры  $a_0$ ,  $a_1$  и  $b_1$  соответствуют тем значениям, которые были найдены при рассмотрении одной гармоники. Параметры  $a_2$  и  $b_2$  найдем аналогично:

$$a_2 = \frac{2}{N} \sum y_t \cos 2t \text{ и } b_2 = \frac{2}{N} \sum y_t \sin 2t.$$

В общем виде для  $i$  гармоник параметры ряда Фурье определяются по формулам

$$a_i = \frac{2}{N} \sum y_t \cos W_i t \text{ и } b_i = \frac{2}{N} \sum y_t \sin W_i t.$$

Чаще всего описание временного ряда не превышает четырех гармоник.

### Пример 5.7

Производство товара «К» по месяцам характеризуется данными (в ед.)

№ месяца	1	2	3	4	5	6	7	8	9	10
$y_t$	22	24	23	14	6	5	6	8	15	17
№ месяца	11	12	13	14	15	16	17	18	19	20
$y_t$	24	25	24	18	8	5	9	14	19	23

Графическое представление этого временного ряда дано на рис. 5.19.

Перед нами стационарный динамический ряд, для которого  $\bar{y} = 15,45$ ;  $\sigma_y^2 = 52,1475$ .

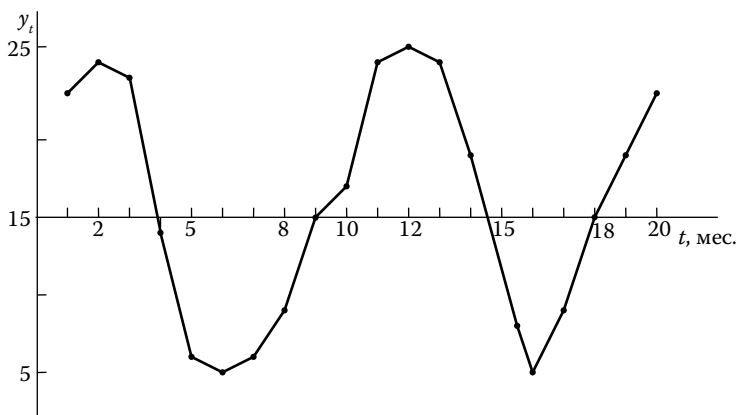


Рис. 5.19. Периодический ряд динамики производства товара  $\mathcal{K}$

Расчеты для определения параметров ряда Фурье представлены в табл. 5.6.

Таблица 5.6. Расчет параметров по ряду Фурье

№ п/п	$y_t$	$t$	$\cos t$	$\sin t$	$\cos 2t$	$\sin 2t$	$\cos 3t$	$\sin 3t$	$\cos 4t$	$\sin 4t$
1	22	0	1	0	1	0	1	0	1	0
2	24	$0,1\pi$	0,951	0,309	0,809	0,588	0,588	0,809	0,309	0,951
3	23	$0,2\pi$	0,809	0,588	0,309	0,951	-0,309	0,951	-0,809	0,588
4	14	$0,3\pi$	0,588	0,809	-0,309	0,951	-0,951	0,309	-0,809	-0,588
5	6	$0,4\pi$	0,309	0,951	-0,809	0,588	-0,809	-0,588	0,309	-0,951
6	5	$0,5\pi$	0	1	-1	0	0	-1	1	0
7	6	$0,6\pi$	-0,309	0,951	-0,809	-0,588	0,809	-0,588	0,309	0,951
8	8	$0,7\pi$	-0,588	0,809	-0,309	-0,951	0,951	0,309	-0,809	0,588
9	15	$0,8\pi$	-0,809	0,588	0,309	-0,951	0,309	0,951	-0,809	-0,588
10	17	$0,9\pi$	-0,951	0,309	0,809	-0,588	-0,588	0,809	0,309	-0,951
11	24	$\pi$	-1	0	1	0	-1	0	1	0
12	25	$1,1\pi$	-0,951	-0,309	0,809	0,588	-0,588	-0,809	0,309	0,951
13	24	$1,2\pi$	-0,809	-0,588	0,309	0,951	0,309	-0,951	-0,809	0,588



Окончание табл. 5.6

№ п/п	$y_t$	$t$	$\cos t$	$\sin t$	$\cos 2t$	$\sin 2t$	$\cos 3t$	$\sin 3t$	$\cos 4t$	$\sin 4t$
14	18	$1,3\pi$	-0,588	-0,809	-0,309	0,951	0,951	-0,309	-0,809	-0,588
15	8	$1,4\pi$	-0,309	-0,951	-0,809	0,588	0,809	0,588	0,309	-0,951
16	5	$1,5\pi$	0	-1	-1	0	0	1	1	0
17	9	$1,6\pi$	0,309	-0,951	-0,809	-0,588	-0,809	0,588	0,309	0,951
18	14	$1,7\pi$	0,588	-0,809	-0,309	-0,951	-0,951	-0,309	-0,809	0,588
19	19	$1,8\pi$	0,809	-0,588	0,309	-0,951	-0,309	-0,951	-0,809	-0,588
20	23	$1,9\pi$	0,951	-0,309	0,809	-0,588	0,588	-0,809	0,309	-0,951
$\Sigma$	309	—	0	0	0	0	0	0	0	0

Отсчет  $t$  ведется с нуля, прибавляя каждый раз величину  $\frac{2\pi}{N}$ , т.е. в нашем случае  $0,1\pi$  (графа  $t$ ). Таблица содержит значения  $\cos t$ ,  $\sin t$ ,  $\cos 2t$ ,  $\sin 2t$ ,  $\cos 3t$ ,  $\sin 3t$ ,  $\cos 4t$ ,  $\sin 4t$  для расчета параметров уравнения с четырьмя гармониками

$$y_t = a_0 + a_1 \cos t + b_1 \sin t + a_2 \cos 2t + b_2 \sin 2t + a_3 \cos 3t + b_3 \sin 3t + \dots + a_4 \cos 4t + b_4 \sin 4t.$$

Чтобы воспользоваться ранее приведенными формулами  $a_i$  и  $b_i$ , были найдены по данным табл. 5.6 следующие значения:

$$\begin{aligned} \Sigma y \cos t &= 6,667; & \Sigma y \sin t &= -17,948 \\ \Sigma y \cos 2t &= 92,883; & \Sigma y \sin 2t &= 26,577; \\ \Sigma y \cos 3t &= -2,698; & \Sigma y \sin 3t &= -10,568 \\ \Sigma y \cos 4t &= -16,753; & \Sigma y \sin 4t &= 11,274. \end{aligned}$$

Ввиду того, что в нашем примере  $\frac{2}{N} = 0,1$ , то  $a_i$  и  $b_i$  составят

$$\begin{aligned} a_1 &= 0,6667; & b_1 &= -1,7948 \\ a_2 &= 9,2883; & b_2 &= -2,6577; \\ a_3 &= -0,2698 & b_3 &= -1,0568; \\ a_4 &= -1,6753; & b_4 &= 1,1274. \end{aligned}$$

Соответственно ряд Фурье представит собой следующее выражение:

$$y_t = \bar{y} + \sum_{i=1}^4 C_{it},$$

где  $C_{it}$  — гармоники вида

$$C_{it} = a_i \cos W_i t + b_i \sin W_i t.$$

Для нашего примера соответствующие гармоники даны в табл. 5.7.

**Таблица 5.7.** Четыре периодические составляющие динамического ряда производства продукции  $K$

Номер гармоники	Гармоническая функция	$R^2$
1	$0,6667 \cos t - 1,7948 \sin t$	0,0351
2	$9,2883 \cos 2t - 2,6577 \sin 2t$	0,930
3	$-0,2698 \cos 3t - 1,0568 \sin 3t$	0,942
4	$-1,6753 \cos 4t + 1,1274 \sin 4t$	0,976

Ряд Фурье с одной гармоникой имеет вид

$$\hat{y}_t = 15,45 + 0,6667 \cos t - 1,7948 \sin t,$$

а с двумя:

$$\hat{y}_t = 15,45 + 0,6667 \cos t - 1,7948 \sin t + 9,2883 \cos 2t - 2,6577 \sin 2t.$$

Аналогично записывается модель с тремя и четырьмя гармониками.

Как очевидно из последней графы табл. 5.7, модель с двумя гармониками хорошо описывает исходный динамический ряд. Хотя  $R^2$  при увеличении числа гармоник и возрастает, но параметры модели для третьей гармоники по  $t$ -критерию Стьюдента оказываются статистически незначимыми. Поэтому при выборе модели лучше предпочесть модель с двумя гармониками.

Как очевидно из графика (см. рис. 5.19), для рассматриваемого временного ряда амплитуда колебаний  $A$  приближается к 10, что и имеет место для уравнения с двумя гармониками

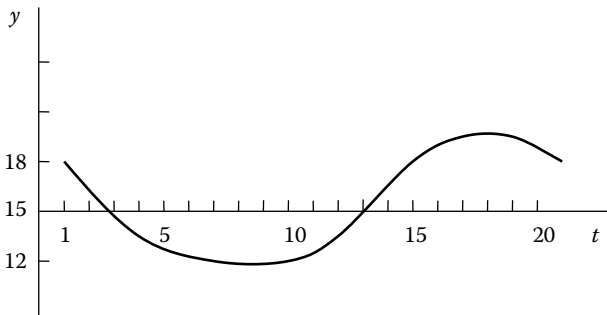
$$A_2 = \sqrt{a_2^2 + b_2^2} = \sqrt{9,2883^2 + 2,6577^2} = 9,66,$$

т.е. интервал, через который ряд начинает повторяться, равен 10 мес.

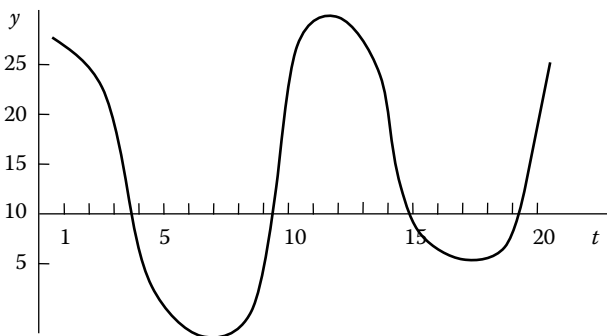
Для уравнения с одной гармоникой период повторения составит 20 мес. и, естественно, выровненный динамический ряд (ряд  $\hat{y}_t$ ) плохо аппроксимирует исходные данные (рис. 5.20).

Для прогноза, используя ряд Фурье с двумя гармониками, в уравнение подставляем следующее по порядку значение  $t$ : в примере на 21-й мес.  $t = 2\pi$ . С учетом того, что  $\cos 2\pi = \cos 4\pi = 1$ , а  $\sin 2\pi = \sin 4\pi = 0$ , прогноз окажется следующим:

$$y_p = 15,45 + 0,6667\cos 2\pi - 1,7948 \sin 2\pi + 9,2883\cos 4\pi - 2,6577 \times \sin 4\pi = 25,4 \text{ ед.}$$



а



б

**Рис. 5.20.** Ряды с одной (а) и ряд с двумя гармониками (б)

### 5.4.1.2. Ряд Фурье по ряду с тенденцией

В экономике чаще встречаются динамические ряды с тенденцией. В этом случае при наличии периодических колебаний ряд Фурье может быть использован, если привести ряд к стационарному виду. Для этой цели можно найти линейный тренд  $\hat{y}_t = a + bt$  и применить ряд Фурье к остаточным величинам  $e_t = y_t - \hat{y}_t$ . Возможен и иной путь: ряд Фурье строится по первым разностям, что равносильно учету линейного тренда. Иными словами, по ряду динамики определяются цепные абсолютные приросты:  $\Delta_t = y_t - y_{t-1}$ , которые далее используются как информационная база для построения ряда Фурье.

#### Пример 5.8

Динамика средних цен на товар по месяцам года характеризуется графиком, представленным на рис. 5.21.

Как очевидно из графика, ряд имеет тенденцию. Уравнение линейного тренда составило  $\hat{y}_t = 6,683 + 0,388t$ , где  $t = 1, 2, \dots, 12$ . Оно описывает 73,7% вариации средних цен  $R^2 = 0,737$  и статистически значимо, ибо  $F = 28$  при табличном значении 4,84 (для 5%-ного уровня значимости). Подставив в данное уравнение соответствующие значения  $t$ , получим расчетные величины средних цен  $\hat{y}_t$  и остатки  $e_t = y_t - \hat{y}_t$ . График остатков представлен на рис. 5.22.

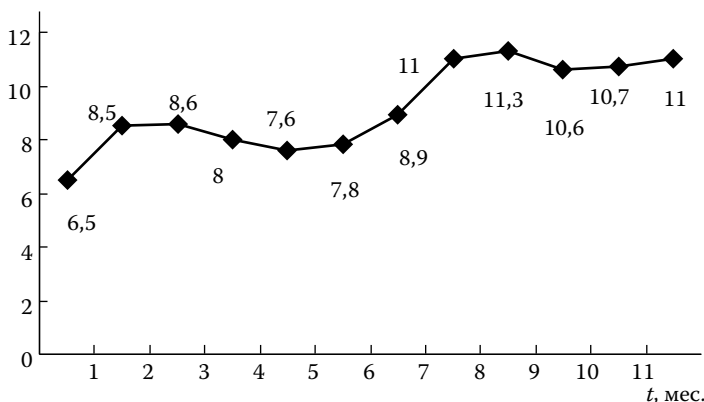
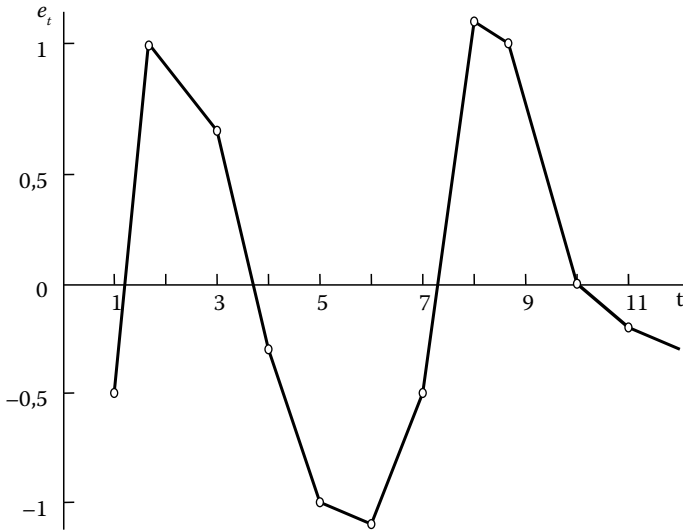


Рис. 5.21. Ряд с тенденцией и периодической составляющей



**Рис. 5.22.** График остатков от линейного тренда

Остатки  $e_t$  представляют собой стационарный ряд и хорошо описываются рядом Фурье с двумя гармониками

$$e_t = 0,123 \cos t - 0,296 \sin t - 0,137 \cos 2t + 1,005 \sin 2t; R^2 = 0,880.$$

В данном уравнении свободный член  $a_0$  отсутствует, так как  $\sum e_t = 0$ .

Модель рассматриваемого динамического ряда представит собой систему уравнений

$$\begin{cases} y_t = a + bt + e_t, \\ e_t = a_1 \cos t + b_1 \sin t + a_2 \cos 2t + b_2 \sin 2t. \end{cases}$$

В нашем примере модель ряда имеет вид

$$\begin{cases} y_t = 6,683 + 0,388t + e_t, \\ e_t = 0,123 \cos t - 0,296 \sin t - 0,137 \cos 2t + 1,005 \sin 2t. \end{cases}$$

Подставив в данную систему значения  $t$  (для тренда: 1, 2, ..., 12; для остатков:  $0, \frac{\pi}{6}, \frac{2\pi}{6}, \dots, \frac{11\pi}{6}$ ), найдем теоретические

значения уровней ряда динамики, которые тесно коррелируют с исходными данными ( $R = 0,9855$ ).

Аналогично поступим и для прогноза на следующий месяц года, т.е. строится суммарный прогноз: прогноз по тренду и плюс прогноз по ряду Фурье для остаточных величин.

В рассматриваемом примере прогноз на 14-й мес. составит:

а) по тренду:  $6,683 + 0,388 \cdot 14 = 12,122$ ;

б) по остаткам:  $0,123 \cos \frac{\pi}{6} - 0,296 \sin \frac{13\pi}{6} - 0,137 \cos \frac{13\pi}{3} +$   
 $+ 1,005 \sin \frac{13\pi}{3} = 0,123 \cdot 0,866 - 0,296 \cdot 0,5 - 0,137 \cdot 0,5 +$   
 $+ 1,005 \cdot 0,866 = 0,760$

в) итого: 12,88.

Ряд Фурье может использоваться также для отображения и прогнозирования динамики с сезонными колебаниями. При этом амплитуда колебаний не должна превышать четыре квартала или 12 мес., ибо сезонные — это внутригодовые колебания. Вместе с тем, сезонные колебания могут изучаться и с помощью иных моделей, позволяющих не только учесть сезонность, но и измерить ее количественно, что имеет, несомненно, практическое значение.

#### 5.4.2. Аддитивная модель сезонности

Аддитивная модель предполагает агрегирование отдельных компонент уровней динамического ряда на основе информации за несколько лет. В зависимости от того, есть или нет тенденция в ряду динамики, она может иметь следующий вид:

$$y_t = \bar{y} + S + \xi \quad \text{— при отсутствии тенденции;}$$

$$y_t = \hat{y}_t + S + \xi \quad \text{— при наличии тенденции,}$$

где  $y_t$  — уровень динамического ряда в период времени  $t$ ;  $\bar{y}$  — средний уровень динамического ряда;  $\hat{y}_t$  — теоретический уровень ряда согласно тенденции;  $S$  — сезонная составляющая, измеренная в тех же единицах, что и уровень ряда;  $\xi$  — случайная компонента, измеренная в тех же единицах, что и уровень ряда.

### 5.4.2.1. Аддитивная модель при отсутствии тенденции

При отсутствии тенденции в ряду динамики общая колеблемость уровней ряда раскладывается на две составляющие: влияние сезонности  $S$  и влияние случайности  $\xi$ . Тогда имеем равенство

$$(y_t - \bar{y}) = (\bar{y}_{S_j} - \bar{y}) + (y_t - \bar{y}_{S_j}), \quad (5.39)$$

где  $\bar{y}_{S_j}$  — средний уровень ряда соответствующего периода внутри года (месяца, квартала) за ряд лет.

В данном равенстве величина  $(\bar{y}_{S_j} - \bar{y})$  отражает влияние сезонности, а величина  $(y_t - \bar{y}_{S_j})$  характеризует влияние случайной компоненты. Влияние сезонной компоненты можно оценить, если уровень ряда представить с помощью линейной модели с фиктивными переменными

$$y_t = a + b_1 z_1 + b_2 z_2 + b_3 z_3 + \xi_t, \quad (5.40)$$

где  $z_1, z_2, z_3$  — фиктивные переменные для кварталов I, II и III, принимающие значение 1 для рассматриваемого квартала и 0 — для остальных.

Так,  $z_1 = 1$  только для I квартала,  $z_2 = 1$  — для II квартала и  $z_3 = 1$  — для III квартала. Применяя к матрице исходных данных  $(y_t, z_1, z_2, z_3)$  МНК, получим оценку параметров  $(a, b_1, b_2, b_3)$ . В данной модели сравнение ведется с IV кварталом, для которого  $z = 0$ .

Параметры модели интерпретируются следующим образом. Параметр  $a = y_S$  для IV квартала, т.е.  $\bar{y}_{S_4}$ ; параметры  $b_j = \bar{y}_{S_j} - \bar{y}_{S_4}$ , т.е. показывают, насколько средний уровень  $j$ -го квартала ниже или выше среднего уровня за IV квартал. Такая интерпретация параметров обусловлена спецификой фиктивных переменных и применением к модели МНК.

Использование МНК в нашем случае приводит к системе нормальных уравнений

$$\begin{cases} \sum y_t = na + b_1 \sum z_1 + b_2 \sum z_2 + b_3 \sum z_3, \\ \sum y_t z_1 = a \sum z_1 + b_1 \sum z_1^2 + b_2 \sum z_1 z_2 + b_3 \sum z_1 z_3, \\ \sum y_t z_2 = a \sum z_2 + b_1 \sum z_1 z_2 + b_2 \sum z_2^2 + b_3 \sum z_2 z_3, \\ \sum y_t z_3 = a \sum z_3 + b_1 \sum z_1 z_3 + b_2 \sum z_2 z_3 + b_3 \sum z_3^2. \end{cases} \quad (5.41)$$

В этой системе  $n$  — число кварталов-лет;  $\sum z_1 = \sum z_2 = \sum z_3 = k$  (число исследуемых лет);  $\sum z_1 z_2 = \sum z_1 z_3 = \sum z_2 z_3 = 0$ ;  $\sum y_t z_1 = \sum y_{j=1}$  (итог по I кварталу за ряд лет);  $\sum y_t z_2 = \sum y_{j=2}$  (итог по II кварталу за ряд лет);  $\sum y_t z_3 = \sum y_{j=3}$  (итог по III кварталу за ряд лет);  $\sum z^2$  по переменным  $z = \sum z$  по ним.

Вычтем из первого уравнения три последующих  $\sum y_{j=4} = ka$  и получим

$$a = \frac{\sum y_{j=4}}{k} = \bar{y}_{S_{j=4}},$$

т.е. параметр  $a$  отражает средний уровень за IV квартал.

Разделив уравнения 2, 3 и 4 на  $k$  и подставив значение параметра  $a$ , найдем оценки параметров  $b_1$ ,  $b_2$  и  $b_3$ , а именно

$$b_1 = \bar{y}_{S_{j=1}} - \bar{y}_{S_{j=4}};$$

$$b_2 = \bar{y}_{S_{j=2}} - \bar{y}_{S_{j=4}};$$

$$b_3 = \bar{y}_{S_{j=3}} - \bar{y}_{S_{j=4}}.$$

Предположим, что по данным в поквартальном разрезе за три года о реализации товара была построена аддитивная модель

$$y_t = 30 - 3z_1 + 96z_2 + 144z_3 + \xi_t.$$

Уравнение показывает, что в IV квартале средний размер реализации за год составил 30 ед. В I квартале средний размер реализации за год был ниже, чем в IV квартале, на 3 ед., а во II и III кварталах — наоборот, выше, чем в IV квартале, на 96 и 144 ед. соответственно.

На основании модели  $\hat{y}_t = 30 - 3z_1 + 96z_2 + 144z_3$  могут быть найдены средние значения для каждого квартала  $\bar{y}_{S_j}$ , а именно  $\bar{y}_{S_j} = b_j + a$ , т.е.

$$\bar{y}_{S_I} = b_1 + a = -3 + 30 = 27;$$

$$\bar{y}_{S_{II}} = b_2 + a = 96 + 30 = 126;$$

$$\bar{y}_{S_{III}} = b_3 + a = 144 + 30 = 174.$$



Иными словами, имея модель уровней динамического ряда, одновременно имеем значения средних уровней для каждого квартала. Это позволяет исходя из модели оценить сезонные колебания по их абсолютной величине  $S_j = \bar{y}_{S_j} - \bar{y}$ . Учитывая, что  $\bar{y} = 89,25$  сезонная компонента составит  $S_1 = -62,25$ ;  $S_2 = 36,75$ ;  $S_3 = 84,75$ ;  $S_4 = -59,25$ . Различия в величине сезонной составляющей I и IV кварталов не велики, что подтверждается  $t$ -статистикой для параметра  $b_1$  ( $t = -0,6$ ). Влияние сезонного фактора во II и III кварталах весьма ощутимо: увеличение объемов продаж статистически значимо по  $t$ -критерию Стьюдента (18,7 и 28,1 для  $b_2$  и  $b_3$  соответственно). В целом рассматриваемая модель достаточно хорошо представляет исходный ряд  $R^2 = 0,993$ , т.е. на долю случайной колеблемости приходится 0,7%.

#### 5.4.2.2. Аддитивная модель при наличии тенденции

При наличии тенденции в ряду динамики общая колеблемость уровней ряда раскладывается на три составляющие:

$$(y_t - \bar{y}) = (\hat{y}_t - \bar{y}) + (y_S - \hat{y}_t) + (y_t - y_S), \quad (5.42)$$

общая
влияние
влияние
влияние  
вариация
тенденции
сезонности
случайности

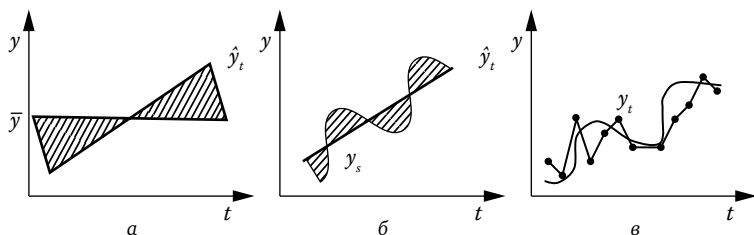
где  $y_S$  — тренд с учетом сезонности, т.е. уровень ряда, обусловленный одновременно влиянием тенденции и сезонности.

Графически влияние этих составляющих представлено на рис. 5.23.

Чем больше угол наклона линии тренда  $\hat{y}_t$  к среднему значению ряда  $\bar{y}$ , тем больше влияние тенденции (рис. 5.23, а). Чем больше плавная кривая  $y_S$  отклоняется от линии тренда  $\hat{y}_t$ , тем значительнее влияние сезонности (рис. 5.23, б). Чем ближе фактические уровни ряда  $y_t$  подходят к плавной линии точек  $y_S$ , тем меньше влияние случайности (рис. 5.23, в).

Существуют разные подходы расчета отдельных составляющих рассматриваемой аддитивной модели, которые зависят от того, как найдены выровненные данные  $\hat{y}_t$ , отражающие тенденцию, а именно:

- а) путем исключения сезонности из данных;



**Рис. 5.23.** Разложение колеблемости уровней динамического ряда на составляющие:  
 а — влияние тенденции; б — влияние сезонности; в — влияние случайности

б) включая сезонность, т.е. выравнивая непосредственно исходные уровни динамического ряда.

Чаще предпочтение отдается первому подходу, при котором вначале производится выравнивание динамического ряда методом скользящих средних для выделения сезонных колебаний  $S$ , а далее, исключив их, определяется тренд без сезонных колебаний  $T = \hat{y}_t$ .

Построение модели включает в себя следующие расчеты.

1. Нахождение сглаженных уровней динамического ряда методом скользящих средних  $\tilde{y}_t$ .
2. Оценка сезонной компоненты  $S_j$  и ее корректировка  $\hat{S}_j$ .
3. Элиминирование сезонной компоненты из исходных данных временного ряда  $(y_t - \hat{S}_j)$ , т.е. проводится десезонализация уровней динамического ряда.
4. Построение уравнения линейного тренда по уровням ряда с элиминированием сезонности.
5. Расчет выровненных значений трендовой составляющей  $\hat{y}_t$ .
6. Расчет теоретических уровней ряда с учетом сезонности  $(\hat{y}_t + \hat{S}_j)$ .
7. Расчет случайной компоненты  $\xi$ , позволяющий оценить далее качество построенной модели.

Данная методика обычно рассматривается в учебниках по эконометрике при разложении уровней динамического ряда по аддитивной модели<sup>1</sup>.

<sup>1</sup> Эконометрика : учебник / под ред. И. И. Елисеевой. М. : Проспект. 2009. С. 172—175.

Аддитивная модель уровней динамического ряда при наличии тенденции и сезонности может быть построена как модель регрессии с включением в нее фактора времени  $t$  и фиктивных переменных  $z$ . Так, при квартальном разрезе информации модель примет вид

$$y_t = a + bt + c_1 z_1 + c_2 z_2 + c_3 z_3 + \xi_t. \quad (5.43)$$

Фактор времени  $t$  в этой модели позволит учесть влияние тенденции. Сезонный фактор представлен фиктивными переменными. Предположим, что  $z_1$ ,  $z_2$  и  $z_3$  соответствуют учету сезонного фактора в I, II и III кварталах соответственно, т.е.

$$z_1 = \begin{cases} 1 & \text{— для I квартала,} \\ 0 & \text{— для остальных;} \end{cases}$$
$$z_2 = \begin{cases} 1 & \text{— для II квартала,} \\ 0 & \text{— для остальных;} \end{cases}$$
$$z_3 = \begin{cases} 1 & \text{— для III квартала,} \\ 0 & \text{— для остальных.} \end{cases}$$

Так как фиктивные переменные  $z$  принимают только значения 1 и 0, то практически мы имеем модель тенденции для каждого квартала:

$$y_t = (a + c_1) + bt + \xi_t \quad \text{— для I квартала;}$$

$$y_t = (a + c_2) + bt + \xi_t \quad \text{— для II квартала;}$$

$$y_t = (a + c_3) + bt + \xi_t \quad \text{— для III квартала;}$$

$$y_t = a + bt + \xi_t \quad \text{— для IV квартала.}$$

Иными словами, параметры при фиктивных переменных  $c_1$ ,  $c_2$ ,  $c_3$  отражают изменение уровня ряда соответствующего квартала под воздействием сезонности по сравнению с IV кварталом.

Предположим, что по данным за три года о численности безработных района (в тыс. человек) была построена аддитивная модель

$$\hat{y}_t = 12,417 - 0,344t - 2,031z_1 - 3,688z_2 - 5,010z_3;$$

$t$	38,5	-11	-6,7	-12,5	-17,3
-----	------	-----	------	-------	-------

$$R^2 = 0,984; F = 108,25.$$

Параметр  $b = -0,344$  указывает на тенденцию снижения уровней ряда при элиминировании сезонности, т.е. ежеквартально независимо от влияния сезонности численность безработных снижалась в среднем на 344 человека.

Параметры  $c_1, c_2, c_3$  показывают, что в I, II и III кварталах уровни ряда независимо от влияния тенденции были в среднем ниже, чем в IV квартале, на соответствующие величины.

Параметр  $a = 12,417$  характеризует уровень IV квартала нулевого года. В рассматриваемой аддитивной модели параметры при фиктивных переменных  $c_1, c_2, c_3$  фиксируют не уровни сезонности соответствующего квартала, а их отличие от воздействия сезонности в IV квартале.

Вместе с тем, зная параметры  $c_1, c_2$  и  $c_3$ , можно определить показатели сезонности. Исходя из содержания параметров при фиктивных переменных, имеем

$$c_1 = S_1 - S_4;$$

$$c_2 = S_2 - S_4;$$

$$c_3 = S_3 - S_4,$$

где  $S_1, S_2, S_3$  и  $S_4$  — показатели сезонности соответствующих кварталов.

$$\text{Тогда } c_1 + c_2 + c_3 = S_1 + S_2 + S_3 - 3S_4.$$

Но так как  $\sum_{i=1}^4 S_i = 0$ , то можно записать, что  $c_1 + c_2 + c_3 = 0 - S_4 - 3S_4$ , или  $c_1 + c_2 + c_3 = 0 - 4S_4$ . Отсюда получим значение сезонной компоненты для IV квартала

$$S_4 = -\frac{1}{4}(c_1 + c_2 + c_3).$$

Для нашего примера, используя значения параметров модели, получим  $S_4 = -\frac{1}{4}(-2,031 - 3,688 - 5,010) = 2,682$ . Далее определяем сезонные компоненты I, II и III кварталов

$$S_1 = c_1 + S_4;$$

$$S_2 = c_2 + S_4;$$

$$S_3 = c_3 + S_4.$$

По данным примера имеем значения компонент

$$S_1 = -2,031 + 2,682 = 0,651;$$

$$S_2 = -1,006 + 2,682 = -1,006;$$

$$S_3 = -5,010 + 2,682 = -2,328;$$

$$S_4 = 2,682;$$

$$\Sigma S = -0,001 \text{ (при более точном подсчете получим нуль).}$$

Часто в практических исследованиях сезонность изучают по месячным данным, ибо сезонность может проявлять себя и внутри квартала. В этом случае применение модели с фиктивными переменными (а их будет 11) потребует информации не менее, чем за 7—8 лет, чтобы на каждый параметр модели приходилось достаточное число степеней свободы и можно было получить надежные оценки параметров. При ограниченной по числу лет информации изучение сезонности по месячным данным целесообразно вести по аддитивной модели, основанной на разложении уровней динамического ряда по компонентам.

#### 5.4.3. Мультипликативная модель сезонности

В мультипликативной модели уровень динамического ряда рассматривается как произведение его компонент:

$$y_t = \hat{y}_t \cdot K_S \cdot E_t, \quad (5.44)$$

где  $y_t$  — фактические уровни динамического ряда;  $\hat{y}_t$  — теоретические значения уровней динамического ряда согласно тенденции;  $K_S$  — коэффициент сезонности;  $E_t$  — коэффициент случайной компоненты.

В данной модели  $\hat{y}_t \cdot K_S$  представляет собой тренд с учетом сезонной волны  $y_S$ , т.е. уровень ряда, обусловленный влиянием как тенденции, так и сезонности  $y_S = \hat{y}_t \cdot K_S$ . Используя величину  $y_S$ , мультипликативную модель можно представить как

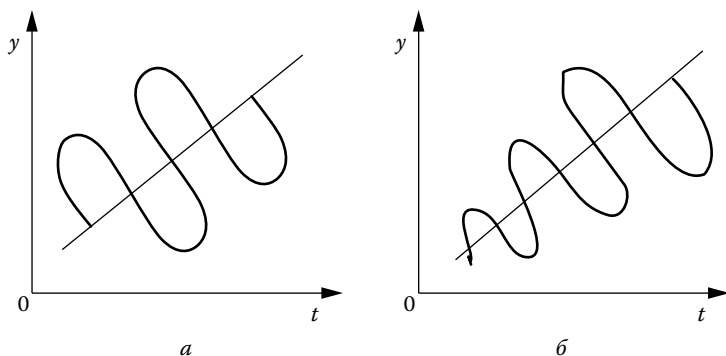
$$y_t = \hat{y}_t \cdot \frac{y_s}{\hat{y}_t} \cdot \frac{y_t}{y_s}, \quad (5.45)$$

где  $\frac{y_s}{\hat{y}_t} = K_s$  и  $\frac{y_t}{y_s} = E_t$ .

Как видим, отличие мультипликативной модели от аддитивной состоит в том, что в мультипликативной модели сезонная и случайная составляющие определены в виде относительных величин (коэффициентов), а в аддитивной модели — в виде абсолютных величин (тысяч рублей, тонн, человек и т.п.). И хотя сезонная компонента в анализе может быть определена и как абсолютная величина  $S$ , и как относительная  $K_s$ , это не означает, что по одному и тому же динамическому ряду обе модели одинаково возможны. Эти модели в практических расчетах дадут близкие результаты, если амплитуда колебаний уровней ряда слабо изменяется во времени.

Ввиду того, что в мультипликативной модели сезонность выражена в процентах, при наличии тенденции в ряду динамики амплитуда сезонных колебаний меняющаяся. Так, если коэффициент сезонности примет значение для I квартала 1,2 или 120%, то при повышающейся тенденции в ряду динамики прирост в 20% будет для I квартала каждого года представлять увеличивающуюся сезонную волну (рис. 5.24).

Мультипликативная модель сезонности традиционно строится на разложении уровней динамического ряда по компонентам.



**Рис. 5.24.** Сезонность: сравнение аддитивной (а) и мультипликативной моделей с линейной тенденцией (б)

Как правило, алгоритмы построения мультипликативных моделей сезонности различаются в зависимости от того, как найдены коэффициенты сезонности, а именно:

- на основе линейного тренда для исходных уровней динамического ряда;
- на основе центрированных скользящих средних;
- на основе простых средних.

В курсах по эконометрике чаще используется второй подход. В этом случае алгоритм построения модели практически тот же, что и по аддитивной модели.

1. Нахождение сглаженных уровней динамического ряда методом скользящих средних  $\hat{y}_t$ .
2. Оценка сезонной составляющей в виде коэффициентов сезонности  $K_{S_j}$  и их корректировка  $\hat{K}_{S_j}$ .
3. Элиминирование сезонной компоненты из исходных данных временного ряда  $\frac{y_t}{\hat{K}_{S_j}}$ , т.е. проведение десеозонлизации уровней динамического ряда.
4. Построение уравнения линейного тренда по уровням ряда без сезонности.
5. Расчет выровненных значений трендовой составляющей  $\hat{y}_t$ .
6. Расчет теоретических уровней ряда с учетом сезонности  $y_S(\hat{y}_t \cdot \hat{K}_{S_j})$ .
7. Расчет случайной компоненты  $E_t(\frac{y_t}{y_S})$ . Для аналитических целей случайная компонента может быть найдена и по абсолютной величине:

$$\xi_t = y_t - y_S.$$

Пример построения подобной модели см. в учебнике по эконометрике<sup>1</sup>.

Прогноз по рассматриваемой мультипликативной модели дается как прогноз по тренду, скорректированному на коэффициент сезонности и коэффициент случайной компоненты:

$$y_p = (\hat{Y} = a + bt_p) \cdot \hat{K}_{S_1} \cdot E_t.$$

<sup>1</sup> Эконометрика : учебник под ред. И. И. Елисеевой. С. 180—182.

Если, например, прогноз по тренду составил 400 тыс. руб.; коэффициент сезонности по расчетам равен 94%, а случайная компонента равна 98%, то прогноз по мультипликативной модели представит величину

$$y_p = 400 \cdot 0,94 \cdot 0,98 = 368,48 \text{ тыс. руб.}$$

Основной недостаток рассмотренной модели состоит в том, что в ней не дается оценка качества модели и статистической значимости отдельных ее компонент. Чтобы оценить существенность влияния отдельных составляющих мультипликативной модели, надо включить в модель каждую из них. С этой целью может быть построена мультипликативная модель сезонности как регрессионная модель с фиктивными переменными. Для отражения влияния тенденции используется фактор времени  $t$ , а для учета сезонности — фиктивные переменные  $z$ , при квартальных данных  $z_1, z_2, z_3$ , принимающие соответственно для I, II, III кварталов значения 1, и 0 — для IV квартала.

Мультипликативная модель сезонности имеет следующий вид:

$$y_t = a \cdot b^t \cdot c_1^{z_1} \cdot c_2^{z_2} \cdot c_3^{z_3} \cdot E_t,$$

где  $y_t$  — исходные уровни динамического ряда;  $t$  — фактор времени, учитывающий влияние тенденции (принимает значения как ряд натуральных чисел);  $z_1, z_2, z_3$  — фиктивные переменные, учитывающий влияние сезонности;  $E_t$  — случайная компонента.

Прологарифмировав данное уравнение, приходим к линейному виду модели

$$\ln y_t = \ln a + t \ln b + z_1 \ln c_1 + z_2 \ln c_2 + z_3 \ln c_3 + \ln E_t,$$

оценка параметров которой может быть дана МНК. Далее, потенцировав, оценим исходные параметры модели  $a, b, c_1, c_2$  и  $c_3$ .

Предположим, что динамика прибыли за последние 12 кварталов характеризуется уравнением

$$y_t = 8,558 \cdot 1,135^t \cdot 0,624^{z_1} \cdot 0,7057^{z_2} \cdot 0,0767^{z_3}; R^2 = 0,982.$$



Уравнение показывает, что в ряду динамики имеется четкая тенденция: ежеквартально независимо от влияния сезонности прибыль возрастает в среднем на 13,5%. Параметры при фиктивных переменных  $z_1, z_2, z_3$  показывают соотношение прибыли соответствующего квартала к прибыли IV квартала, взятого за базу сравнения. В примере прибыль в IV квартале была выше, чем в остальных кварталах, и поэтому величина параметров  $c_1, c_2, c_3$  оказалась меньше единицы.

На основании параметров при фиктивных переменных можно оценить коэффициенты сезонности. Известно, что при квартальных данных сумма коэффициентов сезонности должна быть равна четырем. В нашем примере  $c_1 + c_2 + c_3 + 1 = 3,09687$ . Соответственно поправочный коэффициент составит 1,291628. Умножая его на параметр  $c_j$ , найдем коэффициент сезонности: для I квартала — 0,806; для II квартала — 0,911; для III квартала — 0,991; для IV квартала — 1,292. Они носят аналитический характер, позволяя проводить десеозонлизацию данных. Для прогнозирования достаточно пользоваться параметрами модели.

## Контрольные вопросы и задания

1. Как используются показатели динамики для обоснования типа модели тренда?
2. Дайте определение автокорреляции уровней и поясните, как она используется при моделировании динамического ряда.
3. Что такое автокорреляционная функция и как рассчитывается выборочная оценка коэффициента автокорреляции?
4. Как интерпретируются параметры модели в виде показательной кривой?
5. Как оцениваются параметры S-образных кривых?
6. Как оценить соответствие модели характеру тренда?
7. Что такое автокорреляция в остатках и как она измеряется?
8. В чем причины появления автокорреляции в остатках?
9. С какой целью используется критерий Дарбина — Уотсона?
10. В каких пределах принимает значения критерий Дарбина — Уотсона?
11. Опишите алгоритм построения аддитивной модели временного ряда.

12. Дайте определение сезонной декомпозиции и корректировки временных рядов.
13. В чем состоит различие аддитивных моделей при наличии и отсутствии тренда в ряду динамики?
14. В чем отличие подходов к оцениванию сезонности в аддитивной и мультипликативной моделях?
15. Как используются фиктивные переменные при моделировании сезонных колебаний?
16. Как изменяется ряд Фурье для оценки периодических колебаний?
17. Как оцениваются параметры модели в виде ряда Фурье?
18. Как используется ряд Фурье при наличии тенденции в ряду динамики?

# **Модели регрессии по временным рядам**

## **6.1. Специфика изучения взаимосвязей по рядам динамики**

Временные ряды как источник информации накладывают отпечаток на методологию построения регрессионных моделей. Это связано с возможной ложной корреляцией и ложной регрессией. Высокая корреляция между уровнями временных рядов может иметь место и при отсутствии реальной связи между явлениями. Если, например прокоррелировать динамические ряды заработной платы и дорожно-транспортных происшествий, то можно обнаружить коэффициент детерминации, близкий к единице. Установление связи там, где ее на самом деле нет, означает, что имеет место ложная корреляция. Соответственно и уравнение связи в этом случае будет представлять собой ложную регрессию.

Наличие ложной корреляции и ложной регрессии связано с повышающейся тенденцией каждого из рядов динамики, с автокорреляцией их уровней. Поэтому при построении регрессионных моделей по рядам динамики требуется их предварительная специальная обработка.

Если ряды динамики характеризуются наличием тренда, то при построении модели регрессии надо учесть тренд, например исключить его. В противном случае корреляция уровней рядов динамики будет преувеличена (коэффициент корреляции будет близок к  $+1$  при одинаковой тенденции в рядах и будет

близок  $\kappa - 1$  при противоположной тенденции). Предположим, что строится регрессия личных сбережений граждан от доходов населения по данным за ряд лет. Коэффициент детерминации при этом составил 0,95. Может показаться, что получен хороший результат и уравнение регрессии пригодно для прогноза. Однако анализируя остатки ( $y - \hat{y}_x$ ), мы обнаружим наличие в них автокорреляции. Следовательно, наше уравнение регрессии содержит систематическую погрешность, так как не учитывает влияние тенденции. Высокое значение коэффициента детерминации указывает лишь на то, что обоим рядам свойственна тенденция к повышению уровней. Наличие в двух рядах динамики детерминированного тренда приведет, естественно, к ложной регрессии, которой нельзя воспользоваться для прогнозирования, так как остатки будут автокоррелированы и оценка параметра регрессии окажется неэффективной и не являющейся оценкой теоретического параметра связи между переменными  $x_t$  и  $y_t$ . Модель регрессии по временным рядам пригодна для прогноза, только если остатки ( $y - \hat{y}_x$ ) представляют собой стационарный ряд.

Если ряды динамики характеризуются не только тенденцией, но и периодическими колебаниями, то при построении модели регрессии следует учесть обе компоненты динамических рядов. В этом случае можно из первоначальных данных исключить как тенденцию, так и периодическую составляющую. Модель регрессии может быть построена либо по остаточным величинам, либо с включением в нее обоих компонент динамического ряда наряду с экономическими переменными.

Кроме того, изучая параллельные временные ряды, можно столкнуться с таким явлением, как *временной лаг*, т.е. запаздывание уровней одного ряда относительно другого. Например, спрос на товары длительного пользования может зависеть от доходов предыдущих лет. Инвестиции в основной капитал обычно зависят не только от прибыли текущего года, но и от прибыли прошлых периодов. Поэтому при изучении связи по рядам динамики сначала рассчитывается взаимная корреляционная функция, представляющая собой множество коэффициентов корреляции между уровнями рядов  $y_t$  и  $x_t$ , сдвинутыми относительно друг друга на  $\tau$  моментов времени. Величина лага определяется по наибольшему коэффициенту корреляции. Если временной лаг существует, то он должен быть учтен в модели регрессии (см. гл. 7).

Определенные трудности при построении модели регрессии по временным рядам возникают в связи с проблемой мультиколлинеарности факторов, когда за счет тенденции объясняющие переменные оказываются тесно связанными между собой. Выходом из создавшегося положения может явиться построение модели регрессии по отклонениям от тренда.

Однако можно строить регрессию и по уровням рядов динамики, если удастся при этом устранить автокорреляцию в остатках, применяя, например, обобщенный метод наименьших квадратов. Устранение автокорреляции в остатках возможно также путем изменения спецификации модели, включая, например, в правую часть модели регрессии *лагированные* (запаздывающие) переменные.

## 6.2. Учет тенденции при построении модели регрессии

Методы учета тенденции при построении модели регрессии по временным рядам делятся на две группы:

- методы исключения тенденции из уровней динамического ряда и построение модели по остаточным величинам;
- включение в модель регрессии фактора времени.

### 6.2.1. Методы исключения тенденции

Теоретически возможны два подхода для исключения тенденции из уровней временного ряда:

- метод последовательных разностей;
- метод отклонений от тренда.

Наиболее точным из них является *метод отклонений от тренда*, ибо тенденция учитывается в виде уравнения тренда, описывающего закономерность изменения уровней ряда во времени. *Метод последовательных разностей* учитывает тенденцию, представленную полиномом соответствующей степени. Так, если тенденция линейная, то регрессия строится по первым разностям, т.е. абсолютным приростам; если же тенденция характеризуется параболой второй степени, то для модели регрессии используются вторые разности, т.е. абсолютные ускорения.

Поскольку тренд может быть описан любой математической функцией, а не только полиномом  $k$ -го порядка, то теоретически более оправданным является учет тенденции в модели регрессии методом отклонений от тренда. Вместе с тем построение модели регрессии по последовательным разностям как наиболее простой способ учета тенденции находит практическое применение. Последовательные разности используются также при построении модели *ARIMA* (см. гл. 7).

### 6.2.1.1. Метод последовательных разностей

Если в ряде динамики имеется четко выраженная линейная тенденция, то ее можно устранить, перейдя от исходных уровней ряда  $y_t$  к цепным абсолютным приростам  $\Delta_t$ , т.е. первым разностям. Объясняется это тем, что линейный тренд характеризуется постоянным абсолютным приростом. Его величина в уравнении  $\hat{y}_t = a + bt$  соответствует параметру  $b$ . Первые разности в линейном тренде будут варьировать за счет случайной составляющей  $\xi_t$  вокруг своей константы — параметра  $b$ . Тенденция в уровнях временного ряда будет устранена. Так, цепной абсолютный прирост можно представить как

$$\Delta y_t = y_t - y_{t-1}, \text{ и если } y_t = a + bt + \xi_t,$$

$$\text{то } \Delta y_t = (a + bt + \xi_t) - (a + b(t-1) + \xi_{t-1}) = b + (\xi_t - \xi_{t-1}) \quad (6.1)$$

Если ряд динамики характеризуется тенденцией в виде параболы второй степени, то для ее устранения можно заменить исходные уровни ряда на вторые разности  $\Delta''$ , т.е. на величину абсолютных ускорений. Как известно, парабола второй степени характеризуется постоянным абсолютным ускорением (вторыми разностями), а первые разности имеют линейную тенденцию. Поэтому для динамического ряда с тенденцией в виде параболы  $\hat{y}_t = a + bt + ct^2$  вторые разности  $\Delta''$  будут колебаться вокруг величины  $2c$  за счет случайной ошибки  $\xi_t$ , ибо всегда  $y_t = \hat{y}_t + \xi_t$ . Соответственно тенденция из исходных данных временного ряда будет устранена.

Представим абсолютное ускорение как

$$\Delta''_t = \Delta_t - \Delta_{t-1}, \text{ и если } y_t = a + bt + ct^2 + \xi_t,$$

$$\text{то } \Delta_t = y_t - y_{t-1} = (a + bt + ct^2 + \xi_t) - [a + b(t-1) + c(t-1)^2 + \xi_{t-1}] = (b-c) + 2ct + (\xi_t - 2\xi_{t-1} + \xi_{t-2}), \quad (6.2)$$

т.е. первые разности  $\Delta_t$  являются линейной функцией от времени  $t$ . Вторые разности  $\Delta_t''$  окажутся равными

$$\Delta_t'' = (b-c) + 2ct + (\xi_t - \xi_{t-1}) - [(b-c) + 2c(t-1) + (\xi_{t-1} - \xi_{t-2})] = 2c + (\xi_t - 2\xi_{t-1} + \xi_{t-2}). \quad (6.3)$$

Они не зависят от фактора времени  $t$  и могут быть использованы для построения регрессии по временным рядам.

Аналогично можно показать, что если тенденция характеризуется полиномом третьей степени, то для модели регрессии следует использовать третьи разности, чтобы исключить тенденцию из уровней временного ряда. Однако модели регрессии по вторым, третьим разностям мало информативны с точки зрения их интерпретации и последующего использования в прогнозировании. Поэтому ограничимся рассмотрением регрессии по первым разностям.

При исследовании двух динамических рядов с линейными тенденциями модель линейной регрессии примет вид

$$\Delta y_t = a + b\Delta x_t + u_t, \quad (6.4)$$

где  $\Delta y_t, \Delta x_t$  — первые разности;  $u_t$  — случайная ошибка.

Модель (6.4) по существу является моделью скорости роста. Она строится как обычная модель регрессии, но не по уровням динамических рядов, а по их приростам, т.е. по продифференцированным рядам.

Параметр  $b$  в модели характеризует среднее изменение скорости ряда  $y_t$  с изменением абсолютного прироста ряда  $x_t$  на единицу.

Следует заметить, что если модель будет характеризовать высоким показателем  $R^2$  и отсутствием автокорреляции в остатках, то для прогнозирования конкретных значений  $y_t$  можно перейти к уравнению вида

$$y_p = y_n + a + b(x_p - x_n), \quad (6.5)$$

где  $y_p$  — прогнозное значение динамического уровня ряда  $y_t$ ;  $y_n$  — конечный уровень динамического ряда  $y_t$ ;  $x_p$  — прогнозное значение уровня ряда  $x$ ,  $x_n$  — конечный уровень ряда  $x$ .

В данном уравнении величина  $x_p - x_n = \Delta x_p$  оценивает прогнозное значение скорости ряда  $x$ , а  $y_p - y_n = \Delta y_p$  — прогнозное значение скорости ряда  $y$ .

Прогнозное значение фактора  $x_p$  может быть дано либо по модели  $x_t = f(z_t)$ , где  $z_t$  — объясняющая переменная  $x_t$ , либо по тренду  $\hat{x}_p = a + bt$ . От того, насколько хорошо спрогнозировано значение фактора  $x_p$ , зависит качество прогноза  $y_p$ .

### Пример 6.1

Затраты электроэнергии ( $y_t$  — тыс. кВт · час) и объем выпущенной продукции  $K$  ( $x_t$  — тыс. ед.) характеризуются по предприятию за девять кварталов следующими данными:

Кварталы	1	2	3	4	5	6	7	8	9
$x_t$	3,4	3,9	4,2	4,6	5	5,5	5,8	6,1	6,5
$y_t$	1,7	2	2,2	2,4	2,7	3,1	3,3	3,5	3,9

Если к этим данным применить МНК, то получим уравнение регрессии

$$y_t = -0,7478 + 0,7007x_t + \xi_t; R^2 = 0,9948.$$

Значение  $R^2$ , близкое к единице, обусловлено наличием линейной тенденции в рассматриваемых временных рядах. Модель регрессии по критерию Дарбина — Уотсона ( $D - W = 1,15$ ) не позволяет отклонить гипотезу об отсутствии автокорреляции в остатках: табличные значения его при 5%-ном уровне значимости составило 0,82 — ниже и 1,32 — выше.

Чтобы устранить из данных тенденцию, найдем первые разности (табл. 6.1).



**Таблица 6.1.** Первые разности временных рядов затрат электроэнергии  $\Delta y_t$  и объема выпущенной продукции  $\Delta x_t$ 

Годы	1	2	3	4	5	6	7	8	9
$\Delta x$	—	0,5	0,3	0,4	0,4	0,5	0,3	0,3	0,4
$\Delta y$	—	0,3	0,2	0,2	0,3	0,4	0,2	0,2	0,4

Используя МНК, получим уравнение регрессии:

$$\Delta y_t = -0,0231 + 0,7692\Delta x_t + u_t; R^2 = 0,5245; D - W = 1,4005.$$

Коэффициент регрессии  $b = 0,7692$  показывает, что увеличение скорости роста объема продукции на 1 тыс. ед. приводит в среднем к увеличению абсолютного прироста затрат электроэнергии на 769,2 кВт. · ч.

В данном уравнении регрессии отсутствует автокорреляция в остатках: критерий Дарбина — Уотсона превышает верхнее табличное значение 1,33.

Прогноз на 10-й год выполним по уравнению  $y_p = y_n + a + b(x_p - x_n)$ , где  $y_n = 3,9$ ;  $x_n = 6,5$ ;  $x_p$  найдем по уравнению линейного тренда

$$\hat{x}_t = 3,075 + 0,385t; R^2 = 0,997; \text{ при } t = 10 \quad x_p = 6,925.$$

$$\text{Соответственно } y_p = 3,9 - 0,0231 + 0,7692(6,925 - 6,5) = 4,2 \text{ тыс. кВт. ч.}$$

### 6.2.1.2. Метод отклонений от тренда

Как уже указывалось, метод отклонений от тренда является более точным методом исключения тенденции из данных временного ряда. Это связано не только с тем, что тенденция выражается в виде уравнения тренда любой математической функции. Рассматриваемые для модели регрессии ряды динамики могут иметь разные тенденции. Например, ряд  $x_t$  описывается гиперболой, а ряд  $y_t$  — параболой. В этом случае метод

отклонений от тренда позволяет исключить из каждого временного ряда соответствующую ему тенденцию.

Алгоритм построения регрессии при применении метода отклонений следующий.

1. Для каждого временного ряда определяются уравнение тренда и теоретические значения  $\hat{y}_t$ ;  $\hat{x}_t$ .
2. По каждому из рядов находятся остаточные величины  $dy = y_t - \hat{y}_t$ ;  $dx = x_t - \hat{x}_t$ .
3. Строится модель регрессии

$$dy = f(dx). \quad (6.6)$$

В линейной регрессии  $dy = a + bdx + u_t$  параметр  $b$  показывает как в среднем изменяется величина случайных отклонений по ряду  $y_t$  с изменением случайных колебаний ряда  $x_t$  на единицу. Если при этом оба ряда характеризуются линейной тенденцией, то параметр  $a = 0$ , так как  $\sum dy = \sum dx = 0$ . Тогда модель линейной регрессии примет вид  $dy = bdx + u_t$  и параметр  $b$  будет выступать коэффициентом пропорциональности. Его величина будет показывать, во сколько раз случайные отклонения по ряду  $y_t$  в среднем выше (ниже) случайных отклонений по ряду  $x_t$ .

Для прогноза конкретных значений  $y$  можно перейти к уравнению, связывающему между собой уровни временных рядов. С этой целью в модель регрессии  $dy = a + bdx$  подставим значения  $dy$  и  $dx$ , раскрыв их содержание, т.е.  $dy = y_t - \hat{y}_t$  и  $dx = x_t - \hat{x}_t$ .

Тогда имеем, например, для линейной регрессии  $\hat{dy} = a + bdx$ , т.е.  $y_t - \hat{y}_t = a + b(x_t - \hat{x}_t)$ , или  $y_t = \hat{y}_t + a + b(x_t - \hat{x}_t)$ .

Данную модель можно использовать для прогноза

$$y_p = \hat{y}_{t=p} + a + b(x_p - \hat{x}_{t=p}), \quad (6.7)$$

где  $y_p$  — прогнозное значение  $y$ ;  $\hat{y}_{t=p}$  — прогноз  $y$  по тренду при  $t = p$ ;  $x_p$  — прогнозное значение  $x$ , найденное либо по модели регрессии, либо как  $x_n + \Delta x_n$ ;  $\hat{x}_{t=p}$  — прогноз  $x$  исходя из уравнения тренда при  $t = p$ .

Результат прогноза зависит от качества прогноза фактора  $x$  и от качества трендовых моделей, используемых в прогнозировании.

### Пример 6.2

По данным за 10 мес. рассматривается зависимость прибыли предприятия ( $y_t$  — тыс. руб.) от затрат на рекламу ( $x_t$  — тыс. руб.), табл. 6.2.

**Таблица 6.2.** Расчет остаточных величин для построения модели регрессии  $\hat{y}_t = a + bdx + u_t$

Период $t$	$y_t$	$\hat{y}_t$	$\hat{d}y$	$x_t$	$\hat{x}_t$	$\hat{d}x$
1	250	259,7	-9,7	14	15,1	-1,1
2	305	294,1	10,9	20	18	2
3	314	316,3	-2,3	18	20	-2
4	338	333,1	4,9	23	21,5	1,5
5	354	346,7	7,3	24	22,8	1,2
6	363	358,2	4,8	25	23,9	1,1
7	375	368,3	6,7	27	24,8	2,2
8	376	377,2	-1,2	24	25,7	-1,7
9	377	385,3	-8,3	25	26,5	-1,5
10	380	392,6	-12,6	26	27,2	-1,2

Каждый из рядов имеет повышающуюся тенденцию, которая достаточно хорошо описывается степенной функцией. Для ряда прибыли уравнение тренда составило  $\hat{y}_t = 259,7t^{0,1795}$ ;  $R^2 = 0,9657$ , а для ряда затрат на рекламу —  $\hat{x}_t = 15,077t^{0,2566}$ ;  $R^2 = 0,8555$ . Автокорреляция в остатках отсутствует:  $r_{a_1} = 0,116$  для ряда  $y_t$  и  $r_{a_1} = -0,129$  для ряда  $x_t$ . Если коррелировать исходные уровни рядов динамики  $x_t$  и  $y_t$ , то  $R^2$  составит 0,925, а уравнение регрессии запишется в виде  $\hat{y}_t = 116,621 + 10,026x$ . Однако ввиду наличия в каждом из рядов четкой тенденции можно предположить, что результаты регрессионно-корреляционного анализа завышены. Поэтому используем метод устранения тенденции, найдя отклонения от тренда  $\hat{d}y = y_t - \hat{y}_t$  и  $\hat{d}x = x_t - \hat{x}_t$  (см. табл. 6.2).

Применяя к рядам отклонений  $\hat{d}y$  и  $\hat{d}x$  МНК, получим уравнение линейной регрессии  $\hat{d}y = -0,1398 + 4,0589\hat{d}x$ ;  $R^2 = 0,7022$  (при компьютерной обработке в расчетах использованы  $\hat{d}y$  и  $\hat{d}x$  с точностью до 0,00001). Полученное уравнение регрессии

показывает, что при устранении из исходных уровней временных рядов тенденции имеет место связь между остаточными величинами. Поэтому данную модель можно использовать в прогнозировании. На это указывает и отсутствие в модели автокорреляции остатков. Коэффициент автокорреляции равен  $-0,1372$ , а критерий Дарбина — Уотсона равен  $1,679$ . Сравнивая с табличными значениями величину  $(4-DW)$  при  $\alpha = 0,05$  и числе степеней свободы, равном  $10$ , увидим, что фактическое значение критерия превышает верхнюю границу  $1,32$ .

Для прогноза на 11-й мес. воспользуемся уравнением вида

$$y_p = \hat{y}_{t=p} + a + b(x_p - \hat{x}_{t=p}),$$

$$\text{где } \hat{y}_{t=p} = 259,7t^{0,1795} = 259,7 \cdot 11^{0,1795} = 399,39;$$

$$x_p = 27 \text{ (взято как } x_n = 26 \text{ и плюс } \Delta_n = 1);$$

$$\hat{x}_{t=p} = 15,077t^{0,2566} = 15,077 \cdot 11^{0,2566} = 27,9.$$

Тогда  $y_p$  составит

$$y_p = 399,39 + (-0,1398) + 4,0589(27 - 27,9) = 395,6 \text{ тыс. руб.}$$

Рассмотренный метод построения модели регрессии по временным рядам был основным методом в первой половине XX в. Дальнейшее его развитие привело к моделям, в которых устранение тенденции производится путем включения в модель регрессии фактора времени.

### 6.2.2. Включение в модель регрессии фактора времени

Модель регрессии по двум временным рядам *с включением* в нее как отдельной независимой переменной *фактора времени*  $t$  имеет вид

$$y = a + bx + ct + \xi, \quad (6.8)$$

где  $t = 1, 2, 3, \dots, n$ .

Включая в регрессию фактор времени  $t$ , устраняем линейную тенденцию из уровней временных рядов. Это объясняется спецификой множественной регрессии: коэффициенты

регрессии показывают изолированное влияние на результат соответствующего фактора при неизменном уровне других факторов. В рассматриваемом двухфакторном уравнении регрессии (6.8) коэффициент регрессии  $b$  характеризует «чистое» воздействие переменной  $x$  на результат  $y$  в условиях неизменной тенденции, т.е. при ее устранении.

Математически доказано, что если временные ряды характеризуются линейной тенденцией, то включение в модель фактора времени  $t$  равносильно построению модели регрессии по отклонениям от трендов с последующим переходом от нее к исходным уровням временного ряда зависимой переменной  $y$ . Это связано с тем, что уравнение регрессии  $\hat{y} = a + bx + ct$  может быть построено двумя путями:

- применяя метод наименьших квадратов, получаем оценки параметров  $a$ ,  $b$  и  $c$  (именно так строится данная модель при компьютерной обработке);
- последовательно включаем в модель линейную тенденцию ряда  $y$  и линейную регрессию остаточных величин  $dy = bdx + \xi$ , где  $dy$ ,  $dx$  — остаточные величины от линейных тенденций.

В целях уяснения проблемы рассмотрим второй подход построения линейной модели регрессии с включением фактора времени  $t$ . Алгоритм построения модели следующий.

1. Строится линейное уравнение тренда для ряда  $y_t$

$$\hat{y}_t = a^* + c^* t.$$

2. Строится линейное уравнение тренда для ряда  $x_t$

$$\hat{x}_t = A + Bt.$$

3. Находятся остаточные величины  $dy$  и  $dx$

$$dy = y_t - (a^* + c^* t);$$

$$dx = x_t - (A + Bt).$$

4. Строится регрессия по отклонениям от трендов

$$dy = bdx + \xi.$$

5. Определяется модель для ряда  $y_t$

$$y_t = \hat{y}_t + dy,$$

или  $y_t = a^* + c^* t + b(x_t - A - Bt) + \xi.$

Откуда имеем уравнение

$$y_t = (a^* - bA) + (c^* - bB)t + bx + \xi. \quad (6.9)$$

Данное уравнение соответствует уравнению регрессии (6.8):

$$y_t = a + bx + ct + \xi,$$

где

$$a = a^* - bA; \quad (6.10)$$

$$c = c^* - bB. \quad (6.11)$$

Рассмотренный подход к построению модели регрессии позволяет понять, что уравнение регрессии  $y_t = a + bx + ct + \xi$  с включением фактора времени  $t$  учитывает линейные тенденции для временных рядов  $y_t$  и  $x_t$ . Кроме того, строя регрессию по отклонениям от линейных трендов, мы получаем остатки  $\xi_t$  те же, что и в регрессии с включением линейного фактора времени  $t$ . Поэтому при наличии в рядах линейных тенденций целесообразно строить модель регрессии по исходным уровням рядов динамики, включая в нее фактор времени  $t$ . В этом случае модель регрессии по отклонениям от трендов не информативна. Тем более надо учесть, что регрессия по отклонениям от линейных трендов является составной частью регрессии с включением фактора времени  $t$  (шаг 4 при втором подходе).

В регрессии  $\hat{y}_t = a + bx + ct$  параметр  $b$  показывает, на сколько единиц изменяется в среднем  $y$  при изменении  $x$  на одну единицу в условиях неизменной тенденции; параметр  $c$  показывает средний абсолютный прирост  $y$  в условиях неизменно уровня объясняющей переменной  $x$ .

### Пример 6.3

По промышленному предприятию имеются данные за 3 года в поквартальном разрезе об уровне производительности труда ( $y$  — в тыс. руб. на одного работника) и доле активной части основных фондов ( $x$  — в %):

$t$	1	2	3	4	5	6	7	8	9	10	11	12
$y$	5	6	6	7	8	10	11	11	13	12	13	15
$x$	9,9	18,9	19,8	27,9	22,2	29,7	38,7	36	46	37,8	45	54

Модель регрессии с включением в нее фактора времени  $t$  оказалась следующей:

$$y_t = 2,943 + 0,104x + 0,533t; \quad R^2 = 0,9767; \quad F = 188,3.$$

$t$ -критерий 5,47    2,43    3,44.

В целом уравнение регрессии значимо, как и его параметры (табличное значение  $F$ -критерия при  $\alpha = 0,05$  и числе степеней свободы 2 и 9 равно 4,26, а табличное значение  $t$ -критерия Стьюдента при  $\alpha = 0,05$  и  $df = 9$  равно 2,26). Эти же результаты получим, применяя метод последовательного включения в модель регрессии линейной тенденции ряда  $y_t$  и линейной регрессии отклонений от трендов.

Так, в рассматриваемом примере уравнение линейного тренда для ряда  $y_t$  составило

$$\hat{y}_t = 3,954 + 0,892t,$$

где  $t = 1, 2, \dots, 12$ .

Уравнение хорошо описывает тенденцию:  $R^2 = 0,961$ , коэффициент автокорреляции в остатках равен 0,0083. Линейное уравнение тренда для ряда  $x_t$  составило

$$\hat{x}_t = 9,729 + 3,451t; \quad R^2 = 0,9104.$$

Для каждого ряда были рассчитаны отклонения от трендов

$$dy = y_t - \hat{y}_t \quad \text{и} \quad dx = x_t - \hat{x}_t.$$

Получено уравнение регрессии по отклонениям трендов

$$dy = 0,104dx + \xi; \quad R^2 = 0,396.$$

Далее, с использованием формул (6.10), (6.11), были получены оценки параметров регрессии  $y_t = a + bx + ct + \xi$

$$a = 3,954 - 0,104 \cdot 9,729 = 2,943;$$

$b = 0,104$  (совпадает с величиной коэффициента регрессии по уравнению  $dy = bdx + \xi$ );

$$c = 0,892 - 0,104 \cdot 3,451 = 0,533.$$

Данные оценки совпадают с тем, что было получено ранее по модели с включением в нее фактора времени как объясняющей переменной.

Остаточные величины  $\xi_t$  по модели для отклонений от трендов совпадают с остатками для регрессии с включением фактора времени. Автокорреляция в остатках небольшая:  $r_{a_1} = 0,296$ .

В рассматриваемой модели параметр  $b$  показывает, что рост доли активной части основных фондов на 1%-ный пункт в условиях неизменной тенденции способствует росту уровня производительности труда на 0,104 тыс. руб. Параметр  $c$  характеризует среднеквартальный прирост производительности труда независимо от изменения доли активной части основных фондов, т.е. обусловленный влиянием других факторов, не учитываемых в регрессии.

В примере рассмотрены два динамических ряда. Принцип введения в модель фактора времени сохраняется и при изучении трех и более связанных рядов динамики. Так, если строится регрессия  $y = f(x_1, x_2, x_3)$ , то включение в нее фактора времени  $t$  приводит чаще всего к модели вида

$$\hat{y}_t = a + b_1x_1 + b_2x_2 + b_3x_3 + ct. \quad (6.12)$$

В ней параметры  $b_1$ ,  $b_2$  и  $b_3$  показывают изолированное воздействие каждой объясняющей переменной на результат  $y$ , а параметр  $c$  — средний абсолютный прирост  $y$  в условиях неизменности значений переменных  $x_1$ ,  $x_2$  и  $x_3$ .

Время в качестве независимой переменной часто вводится в виде линейного члена даже если другие переменные подвергаются логарифмированию или иному преобразованию. Например, производственная функция с включением фактора времени часто записывается как

$$P = aK^{b_1}L^{b_2}e^{ct}, \quad (6.13)$$



где  $P$  — объем продукции;  $K$  — основной капитал;  $L$  — занятость;  $e$  — основание натурального логарифма;  $t$  — фактор времени, взятый как ряд натуральных чисел  $1, 2, \dots, n$ .

Линеаризуем данную зависимость, прологарифмировав производственную функцию:

$$\ln P = \ln a + b_1 \ln K + b_2 \ln L + ct. \quad (6.14)$$

Здесь фактор времени  $t$  введен в модель линейно. Рассматриваемая производственная функция не линейна относительно оцениваемых параметров. В ней параметры  $b_1$  и  $b_2$  являются коэффициентами эластичности, показывая, на сколько процентов повышается объем продукции при увеличении соответствующего фактора ( $K$  и  $L$ ) на 1% в условиях неизменной тенденции.

Параметр  $c$  обычно интерпретируется как автономный рост объема продукции в условиях неизменности факторов производства  $K$  и  $L$ . Так, если  $c = 0,0175$ , то  $e^{0,0175} = 1,01765$  и, следовательно, ежегодно (если  $t$  — годы) объем продукции возрастает в среднем при неизменных уровнях затрат капитала и труда в 1,01765 раза, или на 1,765%.

Если тенденция в рядах динамики характеризуется полиномом второй и более высоких степеней, то в модель регрессии вводятся  $t$  и  $t^2$ , а иногда  $t$  в более высокой степени.

В этом случае рассматривается регрессия вида

$$\hat{y} = a + bx + ct + dt^2 \quad (6.15)$$

при двух временных рядах

$$\text{или} \quad \hat{y} = a + \sum_{j=1}^p b_j x_j + ct + dt^2 \quad (6.16)$$

при  $p$  временных рядах.

Вводя в модель регрессии фактор времени в виде  $t, t^2, \dots, t^k$ , предполагаем, что коэффициенты при переменных остаются во времени неизменными и характеризуют силу связи результата  $y$  с соответствующей объясняющей переменной  $x$ .

Если предполагается, что в регрессии коэффициенты при независимой переменной подвержены изменению во времени,

то в модель можно ввести преобразованные переменные  $tx$  (где  $t$  — время). Оценка параметров модели дается МНК.

Модель регрессии с включением в нее фактора времени как независимой переменной не всегда эффективна ввиду возможной мультиколлинеарности факторов. Если временные ряды, используемые в регрессии, характеризуются четкой тенденцией ( $R^2 > 0,9$ ), то корреляция  $t$  и  $x_j$  может превышать корреляцию  $x_j$  с  $y$ , и параметры регрессии при объясняющих переменных  $x$  оказываются ненадежными и экономически не интерпретируемыми.

Время может быть учтено в регрессии и через использование лаговых переменных, т.е. запаздывающих переменных, сдвинутых на определенный интервал времени. Например, спрос на недвижимость в значительной мере определяется доходом не текущего, а предыдущих периодов. Вопросы, связанные с построением моделей регрессии с лаговыми переменными, рассматриваются в гл. 7.

Рассмотренные пути учета тенденции при построении модели регрессии по временным рядам не всегда дают желаемые результаты. Регрессия по отклонениям от тренда зачастую имеет низкий показатель детерминации. Регрессия с включением фактора времени нередко сводится лишь к модели тенденции ввиду статистической незначимости коэффициентов регрессии при объясняющих переменных. Но даже при статистической значимости модели регрессии и ее параметров может остаться автокоррелированность ошибок. Одним из методов ее устранения является обобщенный метод наименьших квадратов.

### **6.3. Обобщенный метод наименьших квадратов при построении модели регрессии по временным рядам**

Методы устранения автокорреляции в остатках могут быть разные. Они зависят от причин автокорреляции. Автокорреляция в остатках может быть следствием неправильной спецификации модели: не учтена важная объясняющая переменная, неправильно выбрана форма связи. В этом случае можно попытаться изменить математическую функцию регрессии (например, линейную на степенную), уточнить набор

объясняющих переменных. Однако если эти попытки не увенчались успехом и автокорреляция в остатках имеет место, то для ее устранения можно применить обобщенный метод наименьших квадратов (ОМНК).

ОМНК можно использовать как для парной, так и для множественной регрессии. Для простоты и уяснения сути проблемы рассмотрим регрессию двух временных рядов

$$y_t = a + bx_t + \xi_t. \quad (6.17)$$

Для периода времени  $(t - 1)$  справедливо равенство

$$y_{t-1} = a + bx_{t-1} + \xi_{t-1}. \quad (6.18)$$

Если имеет место автокорреляция в остатках, т.е. последующие по времени остатки зависят от предыдущих, то регрессия остатков может быть представлена как

$$\xi_t = c + d\xi_{t-1} + V_t, \quad (6.19)$$

где  $V_t$  — случайная ошибка для линейной регрессии остатков.

Но так как  $\xi_t = y_t - \hat{y}_t$ , где  $\hat{y}_t$  — теоретические уровни линейной регрессии  $y$  от  $x$ , то  $\sum \xi_t = 0$  и  $\bar{\xi}_t = 0$ . Полагая, что  $\bar{\xi}_{t-1} \cong \bar{\xi}_t$ , имеем  $c = \bar{\xi}_t - d\bar{\xi}_{t-1} = 0$ . Тогда регрессия остатков примет вид

$$\xi_t = d\xi_{t-1} + V_t. \quad (6.20)$$

Параметр  $d$  определим по формуле

$$d = \frac{\text{cov}(\xi_t, \xi_{t-1})}{\sigma_{\xi_{t-1}}^2}, \quad (6.21)$$

где  $\text{cov}(\xi_t, \xi_{t-1}) = \overline{\xi_t \xi_{t-1}} - \bar{\xi}_t \bar{\xi}_{t-1} = \overline{\xi_t \xi_{t-1}}$ .

В результате получим, что  $d = \frac{\overline{\xi_t \xi_{t-1}}}{\sigma_{\xi_{t-1}}^2}$ . Предполагая, что  $\sigma_{\xi_t}^2 \cong \sigma_{\xi_{t-1}}^2$ , можно записать, что

$$d = \frac{\overline{\xi_t \xi_{t-1}}}{\sigma_{\xi_{t-1}}^2} = \frac{\sum \xi_t \xi_{t-1}}{\sum \xi_t^2}, \quad (6.22)$$

т.е.  $d$  — коэффициент автокорреляции остатков первого порядка. Обозначим его через  $\rho$ . Тогда регрессия остатков примет вид

$$\xi_t = \rho \xi_{t-1} + V_t, \quad (6.23)$$

где  $\rho$  — коэффициент автокорреляции остатков первого порядка;  $V_t$  — случайная ошибка, удовлетворяющая всем предпосылкам МНК.

Уравнение (6.23), как было показано ранее, используется при тестировании моделей на автокорреляцию остатков.

Предполагая, что  $\rho$  известен, вычтем из уравнения (6.17) уравнение (6.18), умноженное на  $\rho$ :

$$y_t - \rho y_{t-1} = a(1 - \rho) + b(x_t - \rho x_{t-1}) + (\xi_t - \rho \xi_{t-1}). \quad (6.24)$$

В итоге имеем линейное уравнение с новыми переменными  $(y_t - \rho y_{t-1})$ ,  $(x_t - \rho x_{t-1})$ . Обозначим новую зависимую переменную  $(y_t - \rho y_{t-1})$  через  $y_t^*$ , а объясняющую переменную  $(x_t - \rho x_{t-1})$  через  $x_t^*$ . Примем также, что  $a(1 - \rho) = a^*$ . Учитывая, что  $\xi_t - \rho \xi_{t-1} = V_t$ , получим уравнение

$$y_t^* = a^* + b x_t^* + V_t, \quad (6.25)$$

где  $V_t$  — независимые случайные величины, имеющие нормальное распределение.

Так как ошибки  $V_t$  удовлетворяют предпосылкам МНК (они не содержат автокорреляцию), то оценки  $a^*$  и  $b$  будут обладать свойствами несмещенных оценок и могут быть получены обычным МНК.

Уравнение (6.25) возможно только при  $t \geq 1$ , так как при  $t = 1$  отсутствует лаговая переменная. Чтобы не уменьшать число степеней свободы рекомендуется для первого периода времени ( $t = 1$ ) использовать *поправку Прайса — Уинстена*

$$x_1^* = \sqrt{1 - \rho^2} \cdot x_1; \quad y_1^* = \sqrt{1 - \rho^2} \cdot y_1. \quad (6.26)$$

Таким образом, ОМНК предполагает, что вместо исходных переменных  $y_t$  и  $x_t$  используются взвешенные переменные

$PY = y^*$  и  $PX = x^*$ , где  $P$  — веса. В матричном виде модель регрессии принимает вид  $PY = PXB + P\xi$ . В ней матрица весов  $P$  составит

$$P = \begin{bmatrix} \sqrt{1-\rho^2} & 0 & 0 & \dots & 0 & 0 & 0 \\ -\rho & 1 & 0 & \dots & 0 & 0 & 0 \\ 0 & -\rho & 1 & \dots & 0 & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & -\rho & 1 & 0 \\ 0 & 0 & 0 & \dots & 0 & -\rho & 1 \end{bmatrix}.$$

Иными словами, матрица исходных данных трансформируется

$$y^* = \begin{bmatrix} y_1\sqrt{1-\rho^2} \\ y_2 - \rho y_1 \\ y_3 - \rho y_2 \\ \dots \\ y_n - \rho y_{n-1} \end{bmatrix}; \quad x^* = \begin{bmatrix} x_1\sqrt{1-\rho^2} \\ x_2 - \rho x_1 \\ x_3 - \rho x_2 \\ \dots \\ x_n - \rho x_{n-1} \end{bmatrix}.$$

Для длинных динамических рядов поправка Прайса — Уинстена может не применяться. Тогда матрица весов не содержит первую строку рассмотренной матрицы  $P$ , и в расчетах используется  $(n-1)$  преобразованных наблюдений  $y_t^*$  и  $x_t^*$ .

К преобразованным переменным  $y^*$  и  $x^*$  применяется традиционный МНК и оцениваются параметры  $a^*$  и  $b$ . Далее из соотношения  $a^* = a(1-\rho)$  можно найти параметр  $a$  как

$$a = a^* / (1-\rho). \quad (6.27)$$

ОМНК распространяется и на случай множественной регрессии

$$y_t = a + b_1 x_{1t} + b_2 x_{2t} + \dots + b_p x_{pt} + \xi_t.$$

Если имеет место автокорреляция остатков и  $\xi_t = \rho \xi_{t-1} + V_t$ , то

$$y_t - \rho y_{t-1} = a(1-\rho) + b_1(x_{1t} - \rho x_{1t-1}) + b_2(x_{2t} - \rho x_{2t-1}) + \dots + b_p(x_{pt} - \rho x_{pt-1}) + (\xi_t - \rho \xi_{t-1}), \quad (6.28)$$

где  $\xi_t = \rho \xi_{t-1} + V_t$ .

Или, исходя из прежней символики, строим модель вида

$$y_t^* = a^* + b_1 x_{1t}^* + b_2 x_{2t}^* + \dots + b_p x_{pt}^* + V_t. \quad (6.29)$$

Применяя к переменным  $y_t^*, x_{1t}^*, x_{2t}^*, \dots, x_{pt}^*$  традиционный МНК, найдем оценки параметров  $b_1, b_2, \dots, b_p$ . Свободный член модели определим как  $a = a^* / (1 - \rho)$ . Далее можно написать искомую модель регрессии  $\hat{y}_t = a + b_1 x_{1t} + \dots + b_p x_{pt}$ , в которой устранена автокорреляция остатков.

Иными словами, применение ОМНК к регрессии с автокоррелированными остатками сводится к двухшаговой процедуре:

- преобразование исходных уровней динамических рядов с помощью известного значения коэффициента автокорреляции остатков первого порядка  $\rho$ ;
- применение к преобразованным данным обычного МНК.

#### Пример 6.4

По данным за 12 кварталов по предприятию рассматривается зависимость уровня рентабельности ( $y$  — в %) от оборачиваемости активов ( $x$  — в разах).

**Таблица 6.3.** Исходные данные и результаты анализа

$t$	$y_t$	$x_t$	$\xi_t$	$\xi_{t-1}$	$y_t^*$	$x_t^*$	$\hat{y}$
1	9	2	-2,109	—	7,725	1,717	13,5
2	11,5	2,6	-1,379	-2,109	6,883	1,574	14,8
3	14	3	-0,06	-1,379	8,101	1,666	15,7
4	16,8	3,5	1,265	-0,06	9,618	1,961	16,8
5	18	4	0,99	1,265	9,382	2,205	17,8
6	18,3	4,1	0,995	0,99	9,066	2,048	18,1
7	18,5	4,2	0,9	0,995	9,112	2,097	18,3
8	19	4,4	0,809	0,9	9,51	2,245	18,7
9	19,2	4,6	0,419	0,809	9,453	2,343	19,2
10	20,8	5	0,839	0,419	10,95	2,64	20
11	22	6	-0,912	0,839	11,33	3,435	22,2
12	24,4	7,1	-1,757	-0,912	13,114	4,022	24,6

Непосредственное применение регрессионного анализа к исходным данным дало следующие результаты:

$$\hat{y}_t = 5,2078 + 2,9506x; \quad R^2 = 0,9216;$$

$$F = 117,6;$$

$$t \quad 4,3 \quad 10,8.$$

Уравнение регрессии статистически значимо, как и его параметры. Исследуем автокорреляцию остатков. Значения остатков  $\xi_t$  и  $\xi_{t-1}$  приведены в табл. 6.3. Коэффициент автокорреляции остатков найдем по формуле

$$r_{a_1} = \frac{\sum e_t e_{t-1}}{\sum e_t^2} = \frac{8,3037}{16,1867} = 0,513.$$

Его величина не столь мала, а по критерию Дарбина — Уотсона не может быть отвергнута нулевая гипотеза об отсутствии автокорреляции в остатках (величина  $DW = 0,508$  ниже нижней границы  $0,97$  при уровне значимости  $0,05$ ). Для устранения автокорреляции в остатках применим ОМНК.

Определим преобразованные значения зависимой и объясняющей переменных, взвесив их на коэффициент автокорреляции остатков. В качестве весов  $\rho$  будем использовать полученное значение  $r_{a_1} = 0,513$ . Для  $t = 1$  преобразованные значения составят

$$y_1^* = y_1 \sqrt{1 - \rho^2} = 9 \sqrt{1 - (0,513)^2} = 7,725;$$

$$x_1^* = x_1 \sqrt{1 - \rho^2} = 2 \sqrt{1 - (0,513)^2} = 1,717.$$

Остальные переменные со значениями  $t = 2, 3, \dots, 12$  преобразуются по формулам  $y_t^* = y_t - \rho y_{t-1}$ ;  $x_t^* = x_t - \rho x_{t-1}$ . Результаты расчета новых переменных приведены в табл. 6.3.

К преобразованным переменным  $y^*$  и  $x^*$  применяем обычный МНК. Уравнение регрессии оказалось следующим:

$$y^* = 4,4474 + 2,1778x^*; \quad R^2 = 0,9016; \quad F = 91,66.$$

Хотя в модели с преобразованными переменными коэффициент детерминации ниже, чем в первоначальной модели, но автокорреляция в остатках отсутствует (коэффициент

автокорреляции остатков равен 0,047). В данном уравнении величина 4,4474 — параметр  $a^*$ . Перейдем от него к параметру  $a$ :

$$a = a^* / (1 - \rho) = 9,1322.$$

В окончательном виде уравнение регрессии составит

$$\hat{y} = 9,1322 + 2,1778x$$

(рассчитанные по нему теоретические значения  $y$  представлены в последней графе табл. 6.3). В данном уравнении коэффициент регрессии свободен от автокоррелированности ошибок соответственно не является результатом ложной регрессии и позволяет делать более корректные выводы относительно связи результативного признака с объясняющей переменной.

Рассмотренный ОМНК базируется на предположении, что коэффициент автокорреляции остатков  $\rho$  известен. Однако на практике точное значение  $\rho$  неизвестно и используются его оценки  $\hat{\rho}$ .

В качестве оценки  $\hat{\rho}$  может использоваться значение  $r_{e_t, e_{t-1}}$ , полученное на основе критерия Дарбина — Уотсона для модели по исходным данным. Как было показано ранее,  $DW \cong 2(1 - r_{e_t, e_{t-1}})$ .

Соответственно имеем уравнение

$$r_{e_t, e_{t-1}} \cong 1 - \frac{DW}{2}. \quad (6.30)$$

В нашем примере до коррекции выводов величина  $DW = 0,508$ . Соответственно  $\hat{\rho}$  составит 0,746.

Данный метод оценивания дает неплохие результаты при достаточно большом числе наблюдений. Используя  $\hat{\rho} = 0,746$  и процедуру ОМНК к табл. 6.3, получим уравнение с некоррелированными остатками

$$\hat{y} = 3,314 + 1,670x; R^2 = 0,735,$$

но коэффициент автокорреляции остатков несколько выше (-0,212), хотя нулевая гипотеза об отсутствии автокорреляции остатков принимается. Практически эта модель



по коэффициенту регрессии отличается от предыдущей и исходной моделей.

Считается, что более точную оценку коэффициентов регрессии при ОМНК дает *двухшаговая процедура Дарбина*. Суть ее заключается в следующем. Уравнение (6.24) можно записать в виде

$$y_t = a(1 - \rho) + bx_t + \rho y_{t-1} - \rho x_{t-1} + V_t, \quad (6.31)$$

где  $V_t = \xi_t - \rho \xi_{t-1}$ .

Тогда имеем модель регрессии, в которой  $\rho$  входит в число оцениваемых параметров.

К уравнению (6.28) можно применить обычный МНК, так как остатки  $V_t$  не содержат автокорреляции. Первый шаг процедуры Дарбина состоит в применении к модели (6.31) традиционного МНК для определения оценки коэффициента автокорреляции остатков  $\rho$  при переменной  $y_{t-1}$ .

Далее на втором шаге оценка  $\rho$  используется для вычисления преобразованных переменных  $y_t^* = y_t - \rho y_{t-1}$  и  $x_t^* = x_t - \rho x_{t-1}$ . К этим преобразованным переменным применяется обычный МНК, т.е. строится уравнение (6.22), в котором коэффициент при  $x^*$  служит оценкой коэффициента регрессии  $b$ , а величина  $a^* / (1 - \rho)$  оценивает параметр  $a$ .

Для нашего примера двухшаговая процедура Дарбина приводит к следующим результатам: на первом шаге получено уравнение регрессии обычным МНК

$$y_t = 3,981 + 1,641x_t - 0,771x_{t-1} + 0,602y_{t-1} + V_t;$$

соответственно считаем, что коэффициент автокорреляции остатков  $\rho = 0,602$ ; на втором шаге находим преобразованные переменные  $x_t^*$  и  $y_t^*$  с учетом поправки Прайса — Уинстена и обычным МНК получаем регрессию

$$y_t^* = 4,1298 + 1,9814x_t^*; R^2 = 0,8797;$$

t            8,5            8,6.

Уравнение и его параметры статистически значимы. Автокорреляция остатков отсутствует ( $\rho = -0,226$ ;  $4 - DW =$

= 1,558, что больше верхней границы его критического значения 1,33 при уровне значимости 0,05 и  $n = 12$ ). Параметр  $a$  составит  $-10,3765$ , т.е. в окончательном виде уравнение регрессии имеет вид

$$y_t = 10,3765 + 1,984x_t + V_t.$$

Другим методом, позволяющим оценить  $\hat{\rho}$ , является *итеративная процедура Кохрейна — Оркатта*. Применительно к модели парной регрессии  $y_t = a + bx_t + \xi_t$  она состоит в следующем.

На первом шаге оцениваются параметры регрессии обычным МНК, а также остаточные величины  $\xi_t$  и по ним на основе модели авторегрессии остатков (6.20)  $\xi_t = \rho\xi_{t-1} + V_t$  оценивают МНК параметр  $\rho$ . Если по критерию Дарбина — Уотсона гипотеза о нулевой автокорреляции ошибок отвергается, то переходят далее ко второму шагу итерации и строится модель по преобразованным переменным  $y_t^* = y_t - \rho y_{t-1}$  и  $x_t^* = x_t - \rho x_{t-1}$ , т.е. рассматривается модель

$$\hat{y}_t^* = a' + b' x_t^*, \quad (6.32)$$

оценка параметров которой дается МНК (при этом в расчетах не используются данные при  $t = 1$ , т.е. поправка Прайса — Уинстена). Остатки этой регрессии вновь проверяют на наличие автокорреляции. Используя модель (6.23), получают новую оценку коэффициента автокорреляции остатков  $\rho$ . Процесс продолжается до тех пор, пока разность между последующей и предыдущей оценками  $\rho$  не будет меньше заданного числа, фиксирующего отсутствие автокорреляции ошибок.

Далее полагая, что  $a = a' / (1 - \rho)$  и  $b = b'$ , получим окончательный вид линейного уравнения регрессии.

Итерационная процедура Кохрейна — Оркатта может быть распространена и на множественную регрессию. Она представляет собой метод коррекции статистических выводов относительно коэффициентов регрессии для модели по динамическим рядам при наличии автокорреляции ошибок.

Применим итерационную процедуру Кохрейна — Оркатта к нашему примеру. Ранее уже было приведено уравнение регрессии по исходным данным и соответственно остатки

по нему  $\xi_t$  (см. табл. 6.3). Применение к модели (6.23) МНК приводит к формуле расчета коэффициента автокорреляции остатков  $\rho$

$$\rho = \frac{\sum_{t=2}^n \xi_t \xi_{t-1}}{\sum_{t=2}^n \xi_{t-1}^2}, \quad (6.33)$$

что не совпадает с расчетом коэффициента автокорреляции остатков первого порядка по формуле, используемой в критерии Дарбина — Уотсона. Для нашего примера  $\rho$ , найденное по формуле (6.33), равно 0,634, а исходя из формулы (6.22) равно 0,513. Результаты по длинным рядам динамики, как правило, близки, так как  $\sum \xi_{t-1}^2$  мало отличаются от  $\sum \xi_t^2$ . Полагая  $\rho = 0,634$ , найдем значения преобразованных переменных  $y^*$  и  $x^*$ , которые будут достаточно близки к данным, полученным по двухшаговой процедуре Дарбина. Уравнение регрессии для них составит  $\hat{y}^* = 3,986 + 1,9088x^*$ . Для него  $R^2 = 0,8636$ ;  $F = 63,3$ ;  $\rho = -0,297$ ;  $4 - DW = 1,42$ , т.е. полученная оценка коэффициента регрессии  $b$ , не подверженная автокорреляции ошибок.

Таким образом, по данным табл. 6.3 получены следующие оценки коэффициента регрессии  $b$ :

2,951 — по исходным данным без устранения автокорреляции остатков;

2,178 — по ОМНК;

1,670 — по ОМНК с коэффициентом автокорреляции остатков, исходя из критерия Дарбина — Уотсона;

1,984 — по двухшаговой процедуре Дарбина;

1,909 — по итерационной процедуре Кохрейна — Оркатта.

Как видим, устранение автокорреляции остатков по двухшаговой процедуре Дарбина и по итерационной процедуре Кохрейна — Оркатта приводит к примерно одинаковым оценкам коэффициента регрессии.

Практически процедура Кохрейна — Оркатта равносильна применению ОМНК. Что касается оценки автокорреляции остатков, то поскольку ее истинное значение исследователю неизвестно, то можно использовать более простую процедуру оценивания в виде двухшаговой процедуры Дарбина.

Ввиду того, что в расчетах используется не истинное значение  $\rho$ , а его оценка, описанный метод получил название **до-ступного** (или **практически реализуемого**) **ОМНК**.

При прогнозировании по модели, полученной с помощью ОМНК, рекомендуется учитывать возможную автокорреляцию остатков, т.е. прогноз осуществляется по модели

$$y_p = a + bx_p + \rho (y_n - a - bx_n), \quad (6.34)$$

где  $y_p$  — прогнозируемое значение  $y$  на  $(n + 1)$  период при длине динамического ряда  $n$ , т.е.  $p = n + 1$ ;  $x_p$  — прогнозируемое значение  $x$ ;  $y_n$  и  $x_n$  — конечные уровни динамических рядов.

Для нашего примера это означает, что прогноз строится по модели  $y_p = 9,132 + 2,178x_p + 0,513 (y_n - 9,132 - 2,178x_n)$ .

Если мы предположим, что  $\rho = 1$ , т.е. имеет место полная положительная автокорреляция в остатках, то ОМНК будет сведен к методу последовательных разностей.

Обратимся к уравнению (6.24). При  $\rho = 1$  оно может быть записано в виде

$$y_t - y_{t-1} = b(x_t - x_{t-1}) + (\xi_t - \xi_{t-1}), \quad (6.35)$$

где  $y_t - y_{t-1} = \Delta y_t$ ;  $x_t - x_{t-1} = \Delta x_t$ ;  $\xi_t - \xi_{t-1} = \Delta \xi_t$ , т.е. мы имеем уравнение регрессии по первым разностям

$$\Delta y_t = b \Delta x_t + V_t, \quad (6.36)$$

где  $V_t = \Delta \xi_t$ .

Из уравнения (6.36) по МНК оценивается коэффициент регрессии  $b$ .

Таким образом, если  $\rho$  стремится к единице, а величина критерия Дарбина — Уотсона стремится к нулю, то использование регрессии по первым разностям устраняет автокорреляцию в остатках.

Если предположить, что  $\rho = -1$ , т.е. имеет место полная отрицательная автокорреляция в остатках, то уравнение (6.21) примет вид

$$y_t + y_{t-1} = 2a + b(x_t + x_{t-1}) + (\xi_t + \xi_{t-1}), \quad (6.37)$$

$$\text{или } (y_t + y_{t-1})/2 = a + b(x_t + x_{t-1})/2 + (\xi_t + \xi_{t-1})/2. \quad (6.38)$$

В уравнении (6.38)  $(y_t + y_{t-1})/2$  и  $(x_t + x_{t-1})/2$  — средние величины за два смежных периода. Используя их в качестве

новых преобразованных переменных  $y_t^*$  и  $x_t^*$ , параметры  $a$  и  $b$  можно оценить МНК. В этом случае ОМНК приводит к модели регрессии по скользящим средним.

Однако предположение, что  $\rho = \pm 1$ , представляет собой довольно редкое на практике явление. Поэтому в большинстве случаев ОМНК применяют либо используя величину  $\rho$ , исходя из критерия Дарбина — Уотсона, либо по двухшаговой процедуре Дарбина, либо по процедуре Кохрейна — Оркатта, либо по процедуре Хилдрета — Лу<sup>1</sup>.

ОМНК позволяет строить модель регрессии по исходным уровням временных рядов, не используя методы исключения тенденции или включения в модель фактора времени.

## 6.4. Учет сезонности при построении модели регрессии

При отсутствии тенденции в рядах динамики и при наличии сезонной компоненты в модели регрессии зависимая переменная  $y_t$  может рассматриваться как функция объясняющих переменных  $x_{1t}, x_{2t}, \dots, x_{pt}$  и сезонного фактора  $z$ , который можно ввести в модель в виде фиктивных переменных 0 и 1: например, для квартальных данных в модель вводятся фиктивные переменные  $z_1, z_2$  и  $z_3$ , учитывающие информацию I, II и III кварталов:

$$z_1 = 1 \text{ для I квартала;}$$

$$z_2 = 1 \text{ для II квартала;}$$

$$z_3 = 1 \text{ для III квартала;}$$

$$\text{в остальных случаях } z_1 = z_2 = z_3 = 0.$$

В общем виде линейная модель регрессии имеет вид

$$y_t = a + \sum_{j=1}^p b_j x_{jt} + \sum_{i=1}^k c_i z_{it}, \quad (6.39)$$

<sup>1</sup> Плохотников К. Э. Основы эконометрики в пакете STATISTICA. М. : Вузовский учебник, 2010. С. 118

где  $k = 1, 2, 3$  (для квартальных данных);  $k = 1, 2, 3, \dots, 11$  (для месячных данных). Построение подобного рода моделей требует достаточно длинных рядов динамики, чтобы параметры модели были статистически надежны.

Коэффициенты регрессии  $b_j$  ( $j = 1, 2, \dots, p$ ), показывают воздействие соответствующего фактора  $x_j$  на  $y$  при устранении сезонности. Влияние сезонности оценивается через коэффициенты  $c_i$ , которые представляют собой разность среднего уровня  $y$  для  $i$ -го квартала по сравнению с базовым, т. е. если за базу сравнения взят IV квартал, то  $c_1 = \bar{y}_I - \bar{y}_{IV}$  при условии, что факторы  $x_1, x_2, \dots, x_p$  неизменны. Соответственно  $c_2 = \bar{y}_{II} - \bar{y}_{IV}$ ,  $c_3 = \bar{y}_{III} - \bar{y}_{IV}$ . Чтобы оценить влияние сезонности I квартала по сравнению со II, можно использовать разность  $(c_2 - c_1)$ , а влияние III квартала по сравнению со II —  $(c_3 - c_2)$  и т.п.

Практическое использование рассматриваемой модели затруднено следующими обстоятельствами:

- ограничено число объясняющих переменных  $x$  ввиду сравнительно коротких динамических рядов в экономических исследованиях. Так, по данным за 5—7 лет в поквартальном разрезе в модели может быть учтена одна объясняющая переменная;
- сезонность по  $y$  и по  $x$  может не совпадать по кварталам (например, пик доходов может приходиться на зиму, а пик расходов — на лето);
- коэффициент регрессии при факторе  $x$  может подвергаться сезонности, а в модели он предполагается независимым от сезонной компоненты;
- ряды динамики обнаруживают тенденцию.

Рассмотрим содержание линейной модели  $y_t = a + bx_t + c_1z_1 + c_2z_2 + c_3z_3 + \xi_t$  на следующем примере.

### Пример 6.5

Объем продаж товара фирмой ( $y$  — в тыс. ед.) исследуется в зависимости от объема продаж его дочерним предприятием ( $x$  — в тыс. ед.) по данным за 5 лет.

**Таблица 6.4.** Объем продаж фирмы за 5 лет

Годы	Кварталы								Итого	
	I		II		III		IV			
	y	x	y	x	y	x	y	x	y	x
2005	11	9	15	10	6	8	12	9	44	36
2006	11	10	16	9	4	3	13	11	44	33
2007	10	7	14	10	7	8	12	11	43	36
2008	10	12	16	9	8	11	13	12	47	44
2009	11	8	18	16	7	6	12	12	48	42
Итого	53	46	79	54	32	36	62	55	226	191
В среднем	10,6	9,2	15,8	10,8	6,4	7,2	12,4	11,0	45,2	38,2

Таблица показывает, что объем продаж практически стабилен по годам, но сильно варьирует по кварталам. Если строить модель, не обращая внимания на сезонность, то результаты моделирования окажутся плохими:

$$Y_t = 2,674 + 0,903x_t + \xi_t; R^2 = 0,4466.$$

$$t \quad 1,1 \quad 3,8$$

Они существенно улучшаются, если в модели будет учтена сезонность. Введем в модель сезонную компоненту с помощью фиктивных переменных. Взяв за базу сравнения IV квартал, построим модель зависимости объема продаж фирмой от объема продаж дочерним предприятием с учетом сезонности:

$$\hat{y}_t = a + bx_t + c_1z_1 + c_2z_2 + c_3z_3,$$

где  $z_1 = 1$  для I квартала и 0 для остальных;  $z_2 = 1$  для II квартала и 0 для остальных;  $z_3 = 1$  для III квартала и 0 для остальных.

Результаты регрессии оказались следующими:

$$\hat{y}_t = 8,896 + 0,319x_t - 1,227z_1 + 3,464z_2 - 4,789z_3;$$

$$t \quad 8,3 \quad 3,5 \quad -2,15 \quad 6,3 \quad -7,4$$

$$R^2 = 0,955; D - W = 1,998.$$

Как видим, существенно повысился коэффициент детерминации. Все параметры модели статистически значимы и отсутствует автокорреляция в остатках. Параметр при переменной  $x$  означает, что независимо от сезонного фактора с ростом объема продаж дочерним предприятием на 1 тыс. ед. объем продаж головной организации возрастает в среднем на 0,319 тыс. ед. Для отдельных кварталов модель регрессии составит:

$$\text{для I квартала } \hat{y}_t = 7,669 + 0,319x_t;$$

$$\text{для II квартала } \hat{y}_t = 12,360 + 0,319x_t;$$

$$\text{для III квартала } \hat{y}_t = 4,107 + 0,319x_t;$$

$$\text{для IV квартала } \hat{y}_t = 8,896 + 0,319x_t.$$

Коэффициенты регрессии при  $z$  показывают влияние сезонности: изменение объема продаж в головной организации в соответствующем квартале по сравнению с IV кварталом, когда  $z_1 = z_2 = z_3 = 0$ , в условиях неизменности объема продаж в дочернем предприятии. Все коэффициенты регрессии при сезонном факторе  $z$  согласовываются в примере с изменением средних уровней результата по отдельным кварталам в сравнении с IV кварталом:

По модели	По итоговой строке таблицы
$c_1 = -1,227$	$\bar{y}_I - \bar{y}_{IV} = 10,6 - 12,4 = -1,8;$
$c_2 = +3,464$	$\bar{y}_{II} - \bar{y}_{IV} = 15,8 - 12,4 = 3,4;$
$c_3 = -4,789$	$\bar{y}_{III} - \bar{y}_{IV} = 6,4 - 12,4 = -6.$

Наибольшие изменения в средних уровнях наблюдаются в III и II кварталах (по сравнению с IV кварталом). Автокорреляция в остатках по модели практически отсутствует:  $r_{a_1} = -0,165$ . Если в рядах динамики наблюдается тенденция и сезонная компонента одновременно, то модель регрессии может иметь вид (при двух временных рядах)

$$\hat{y}_t = a + bx_t + ct + d_1z_1 + d_2z_2 + d_3z_3. \quad (6.40)$$



Следует заметить, что при наличии четкой тенденции как в ряду зависимой, так и в рядах объясняющих переменных модель с сезонной составляющей в виде линейно введенных фиктивных переменных не всегда может давать хорошие результаты ввиду мультиколлинеарности факторов, особенно, если корреляция объясняющих переменных с фактором времени выше их корреляции с результивным признаком. В этом случае коэффициенты регрессии часто статистически не значимы или имеют неинтерпретируемые знаки.

Предположим, что строится модель регрессии объема продаж фирмой туристических путевок в Австрию от цены на них за последние 3 года (данные характеризуются тенденцией и сезонностью; пользуются спросом горнолыжные маршруты зимой и термальные озера летом). Располагая информацией по месяцам, вводим в модель фактор времени  $t$  (1, 2, 3, ..., 36) и сезонную компоненту в виде фиктивных переменных  $z_1$  и  $z_2$ , принимающих значение 1 соответственно для зимних и летних периодов года и 0 — для остальных. Модель имеет вид  $y_t = a + bx_t + ct + d_1z_1 + d_2z_2 + \xi_t$ . Параметр  $b$  показывает, что независимо от тенденции и сезонности рост цен на 1 ден. ед. соответствует изменению числа проданных путевок на  $b$  ед. Ежемесячно под воздействием тенденции и независимо от периода времени объем продаж возрастает в среднем на  $c$  ед. При абстрагировании от влияния тенденции и фактора цен сезонность проявляет себя в значениях параметров  $c_1$  и  $c_2$ : в зимний период объем продаж возрастает на  $c_1$ , а летом — на  $c_2$  ед. (все коэффициенты регрессии при сезонной компоненте больше нуля и статистически значимы по  $t$ -критерию Стьюдента). Сезонность можно рассматривать и в более крупном плане: например, осенне-зимний период и весенне-летний период. В этом случае фиктивная переменная может принимать значение 1, например, для I и IV кварталов, а для II и III кварталов — значение 0 (или наоборот). Такой подход уменьшает число вводимых в модель фиктивных переменных, что может быть полезным при сравнительно коротких рядах динамики. Этот подход возможен также, если в каком-то из периодов влияние сезонности оказалось по  $t$ -критерию Стьюдента статистически незначимым.

---

## Контрольные вопросы и задания

1. Что такое ложная корреляция?
2. Почему регрессия по временным рядам требует специальных методов построения модели?
3. В чем состоит смысл метода первых разностей?
4. Опишите алгоритм построения модели регрессии по отклонениям от тренда.
5. Как учитывается тенденция при встроенности в нее фактора времени?
6. В каком случае модель по отклонениям от тренда совпадает с моделью с включением в нее фактора времени?
7. В чем смысл использования обобщенного метода наименьших квадратов (ОМНК) при построении модели регрессии по временным рядам?
8. Какова цель поправки Прайса — Уинстена?
9. Что собой представляет метод Кохрейна — Оркатта?
10. Опишите алгоритм обобщенного метода наименьших квадратов для построения модели регрессии по временным рядам.

### 7.1. Общая характеристика

До сих пор мы рассматривали модели по временным рядам, в которых  $y_t = f(x_t)$ . Между тем в моделях временных рядов зависимая переменная  $y_t$  может быть связана не только со значениями объясняющих переменных  $x$  в момент времени  $t$ , но и с их значениями в предыдущие моменты времени. Так, например, потребление товаров длительного пользования зачастую зависит от доходов не только текущего, но и предыдущих периодов. Аналогично величина основных производственных фондов зависит от размера инвестиций не только текущего года, но и предыдущих лет. В этом случае строятся модели с лаговыми объясняющими переменными. Например,

$$c_t = a + b_1 y_t + b_2 y_{t-1} + \xi_t,$$

где  $c_t$  — потребление в период времени  $t$ ;  $y_t$  — доход в период времени  $t$ ;  $y_{t-1}$  — доход в предыдущий период  $t - 1$ .

В данной модели лаговой является переменная  $y_{t-1}$ , т.е. доход за предыдущий период времени. Возможна ситуация, когда объясняющая переменная  $x$  влияет на результат  $y$  не сразу же, а с определенным запаздыванием во времени, превышающем

один временной интервал. Так, выпуск специалистов высшей квалификации зависит от приема в вузы четырех- или пяти-летней давности.

Объясняющие переменные, взятые в модели регрессии с запаздыванием во времени, называются **лаговыми переменными**. Величина интервала запаздывания называется **лагом**. Так, в модели  $y_t = a + bx_{t-4} + \xi_t$  лаговая переменная  $x$  взята с лагом, равным четыре.

Вместе с тем в правой части модели лаговой может быть и зависимая переменная. Например, спрос на товар может зависеть не только от дохода, но и от достигнутого спроса на него в предыдущий период времени. Или ставка банковского кредита может зависеть не только от объема денежной массы в наличии, но и от достигнутого ранее процента банковского кредита. В этом случае строятся модели с лаговой зависимой переменной. Например,

$$c_t = a + b_1 y_t + b_2 c_{t-1} + \xi_t,$$

где  $c_t$  — потребление в период времени  $t$ ;  $y_t$  — доход в период времени  $t$ ;  $c_{t-1}$  — потребление в предыдущий период времени  $t - 1$ .

Модели регрессии по временным рядам с лаговыми переменными принято называть **динамическими моделями**. Их можно подразделить на три класса

1. Модели с лаговыми объясняющими переменными — **модели с распределенными лагами**

$$y_t = a + b_0 x_t + b_1 x_{t-1} + \dots + b_k x_{t-k} + \xi_t.$$

2. Модели с лаговыми зависимыми переменными — **модели авторегрессии**

$$y_t = a + bx_t + c_1 y_{t-1} + \dots + c_k y_{t-k} + \xi_t.$$

3. Модели с лаговыми зависимыми и независимыми переменными — **авторегрессионные модели с распределенными лагами**

$$y_t = a + b_1 y_{t-1} + \dots + b_k y_{t-k} + c_0 x_t + c_1 x_{t-1} + \dots + c_k x_{t-k} + \xi_t.$$

Центральным вопросом при построении моделей с лаговыми переменными является выбор величины лага и числа лаговых переменных. Теоретически трудно определить величину лага. Определенную помощь может оказать взаимная корреляционная функция: рассчитывается множество коэффициентов корреляции между уровнями временных рядов  $y_t$  и  $x_t$ , сдвинутыми относительно друг друга на последовательно увеличивающиеся интервалы времени. Величина лага определяется по максимальному значению коэффициента корреляции. Например, продажа товара за две декады двумя филиалами фирмы характеризуется данными, представленными в таблице (тыс. ден. ед.).

Числа месяца	Филиал № 1	Филиал № 2	Числа месяца	Филиал № 1	Филиал № 2
1	5	9	11	9,5	13
2	4,5	10,8	12	8	14
3	4	13,5	13	7,6	15
4	4,1	14,5	14	7,5	20
5	5	16	15	7,6	24
6	7	14,7	16	10	25
7	8	14	17	12,2	26
8	9,7	12	18	15	26,3
9	10	11,9	19	15,6	26,4
10	11	12	20	16	27,1

Примем объем продаж филиалом № 1 за  $y_t$ , а филиалом № 2 — за  $x_t$ . Если прокоррелировать  $y_t$  и  $x_t$ , то коэффициент корреляции между ними составит 0,6912. При последовательном сдвиге уровней ряда  $x_t$  на один временной интервал получатся коэффициенты корреляции, представленные в следующей таблице.

Величина лага						
1	2	3	4	5	6	7
0,7738	0,867	0,9445	0,9553	0,8562	0,5977	0,0724

Следовательно, объем продаж филиалом № 1 в наибольшей мере коррелирует с объемом продаж по филиалу № 2

с интервалом в 4 дня. Уравнение регрессии принимает вид  $\hat{y}_t = -0,4575 + 0,6977x_{t-4}$ ;  $R^2 = 0,913$ ;  $F = 146,2$ , что статистически значимо. Оно позволяет по данным филиала № 2, взятым на четыре дня раньше, предсказывать объем продаж по филиалу № 1. Так, например, при объеме продаж за 2-е число в 10,8 тыс. ден. ед. по филиалу № 2 объем продаж по филиалу № 1 составит 6-го числа 7,1 тыс. ден. ед. Соответственно подставляя в уравнение регрессии информацию об объеме продаж филиалом № 2 за 3—16-е числа, получим объем продаж по филиалу № 1 на 7—20-е числа.

Выбор величины лага и количества лагов проводится обычно экспериментально: строятся модели с разным числом лагов и их величиной и изучается значимость коэффициентов регрессии при лаговых переменных; останавливаются на модели, для которой все коэффициенты регрессии при лаговых переменных будут статистически значимыми по  $t$ -критерию Стьюдента.

Построение моделей с лаговыми переменными имеет свою специфику. Дело не только в выборе величины лага и их числа. Во многих случаях оценка параметров моделей с лаговыми переменными не может быть проведена с помощью традиционного МНК ввиду нарушения ряда его предпосылок и требует специальных методов оценивания. При наличии двух и более лаговых переменных возникает проблема мультиколлинеарности факторов, ибо, как правило,  $x_t, x_{t-1}, x_{t-2}, \dots, x_{t-k}$  или  $y_{t-1}, y_{t-2}, \dots, y_{t-k}$  связаны между собой, особенно при наличии тенденции в рядах динамики. Это снижает точность оценок коэффициентов при лаговых переменных и требует изменять приемы оценивания.

## 7.2. Модели с распределенными лагами

### 7.2.1. Интерпретация параметров модели с распределенными лагами

Модели с распределенными лагами бывают двух типов:

— с конечным числом лагов

$$y_t = a + b_0x_t + b_1x_{t-1} + \dots + b_kx_{t-k} + \xi_t;$$

## — с бесконечным числом лагов

$$y_t = a + b_0 x_t + b_1 x_{t-1} + b_2 x_{t-2} + \dots + \xi_t.$$

Практическое применение чаще имеют модели с конечным числом лагов, т.е. модели, в которых число лагов экспериментально определено.

Предположим, рассматривается модель, в которой  $k = 4$ , т.е.  $\hat{y}_t = a + b_0 x_t + b_1 x_{t-1} + b_2 x_{t-2} + b_3 x_{t-3} + b_4 x_{t-4}$ . Данная модель означает, что изменение во времени  $t$  объясняющей переменной  $x$  будет влиять на значения результативного признака  $y$  в течение 4-х следующих моментов времени.

Коэффициент  $b_0$  называют краткосрочным мультипликатором, так как он характеризует среднее изменение результата  $y$  при изменении  $x_t$  на 1 ед. своего измерения в фиксированный момент времени  $t$ .

В момент времени  $(t + 1)$  воздействие объясняющей переменной  $x$  на результат  $y$  составит  $(b_0 + b_1)$  единиц, а в момент времени  $(t + 2)$  общее изменение  $y$  составит  $(b_0 + b_1 + b_2)$  единиц.

Любую сумму коэффициентов  $\sum_{j=0}^h b_j$ , где  $h < k$  называют **промежуточным мультипликатором**, а сумму всех коэффициентов регрессии  $\sum_{j=0}^k b_j$  — **долгосрочным мультипликатором**, который характеризует общее изменение  $y$  через  $k$  интервалов времени под воздействием изменения  $x$  в момент  $t$  на 1 ед.

При  $k = 4$  долгосрочный мультипликатор составит  $(b_0 + b_1 + b_2 + b_3 + b_4)$ . Он характеризует общее среднее изменение  $y$  через четыре временных интервала при увеличении  $x$  в момент времени  $t$  на 1 ед., а промежуточные мультипликаторы

$b_0 + b_1$  — изменение  $y$  в момент времени  $(t + 1)$ ;

$b_0 + b_1 + b_2$  — изменение  $y$  в момент времени  $(t + 2)$ ;

$b_0 + b_1 + b_2 + b_3$  — изменение  $y$  в момент времени  $(t + 3)$ .

Если все коэффициенты регрессии имеют одинаковые знаки, т.е. характеризуются однонаправленным изменением  $y$

в исследуемые  $k$  моментов времени, то можно определять относительные коэффициенты модели  $\beta_j$ , т.е.  $\beta_j = \frac{b_j}{\sum b_j}$ , где  $0 < \beta_j < 1$ , а  $\sum \beta_j = 1$ . Иными словами,  $\beta_j$  характеризует долю общего изменения  $y$  в момент времени  $(t + j)$ .

Предположим, что регрессия основных производственных фондов ( $y$  — в млн руб.) в зависимости от размера инвестиций ( $x$  — в млн руб.) характеризуется уравнением  $\hat{y}_t = 0,8 + 0,7x_t + 1,0x_{t-1} + 1,5x_{t-2} + 0,6x_{t-3} + 0,2x_{t-4}$ , где  $t$  — года.

Анализ уравнения показывает, что рост инвестиций на 1 млн руб. в текущем периоде приводит к росту основных производственных фондов:

- в том же периоде — на 0,7 млн руб. (краткосрочный мультипликатор);
- через 1 год — на  $0,7 + 1 = 1,7$  млн руб.;
- через 2 года — на  $0,7 + 1 + 1,5 = 3,2$  млн руб.;
- через 3 года — на 3,8 млн руб. (промежуточный, как и предыдущие два, мультипликатор);
- через 4 года — на 4 млн руб. (долгосрочный мультипликатор).

Относительные коэффициенты модели составят

$$\beta_0 = 0,7/4 = 0,175;$$

$$\beta_1 = 1/4 = 0,25;$$

$$\beta_2 = 1,5/4 = 0,375;$$

$$\beta_3 = 0,6/4 = 0,15;$$

$$\beta_4 = 0,2/4 = 0,05.$$

Следовательно, в текущем году реализуется 17,5% воздействия увеличения инвестиций на рост основных производственных фондов, а через год — еще 25%. Через 2 года — еще 37,5%, через 3 года — еще 15% и через 4 года — еще 5%.

Относительные коэффициенты модели  $\beta_j$  можно использовать как весовые коэффициенты для расчета **средней величины лага** по средней арифметической:

$$\bar{j} = \sum_{j=0}^k j\beta_j,$$

где  $j$  — величина лага.



Величина  $\bar{j}$  показывает средний интервал времени, в течение которого будет происходить изменение зависимой переменной  $x$  под воздействием изменения объясняющей переменной  $x$  в момент времени  $t$ . Чем меньше величина среднего лага, тем быстрее реагирует результат  $y$  на изменение  $x$ . И наоборот, высокое значение среднего лага показывает, что воздействие объясняющей переменной на результат будет сказываться с течением длительного промежутка времени. В рассматриваемом примере величина среднего лага составит

$$\bar{j} = 0 \cdot 0,175 + 1 \cdot 0,25 + 2 \cdot 0,375 + 3 \cdot 0,15 + 4 \cdot 0,05 = 1,65 \text{ года.}$$

Следовательно, основная часть эффекта увеличения инвестиций проявляется через 1,65 года. Кроме среднего лага можно рассчитывать *медианный лаг*  $j_{M_e}$ , т.е. тот период времени, в течение которого с момента времени  $t$  будет реализована половина общего эффекта воздействия объясняющей переменной  $x$  на результат  $y$ . Для медианного лага справедливо равенство

$$\sum_{j=0}^{M_e} \beta_j = 0,5.$$

В нашем примере медианный лаг составляет 2 года, т.е. увеличение инвестиций в период времени  $t$  на 1 млн руб. приводит к росту размера основных производственных фондов через 2 года на величину, составляющую половину долгосрочного мультипликатора, т.е. на 2 млн руб. Наибольший аналитический интерес представляет расчет величины медианного лага для моделей с большим числом лаговых переменных.

### 7.2.2. Оценка параметров моделей с распределенными лагами

Модель с конечным числом лагов при правильной ее спецификации может быть оценена обычным МНК. В этом случае в уравнении

$$y_t = a + b_0 x_t + b_1 x_{t-1} + b_2 x_{t-2} + \dots + b_k x_{t-k} + \zeta_t$$

переменные  $x_t, x_{t-1}, \dots, x_{t-k}$  рассматриваются как объясняющие переменные обычной множественной регрессии.

Вместе с тем применение МНК к моделям с конечным числом лагов может быть реально затруднено ввиду следующих причин:

- 1) при наличии тенденции переменные  $x_t, x_{t-1}, \dots, x_{t-k}$  тесно связаны между собой, что вызывает мультиколлинеарность факторов, которая может привести к неинтерпретируемому знакам у коэффициентов регрессии и к снижению их точности;
- 2) возможна автокорреляция остатков, так как МНК применяется к временным рядам с тенденцией.

Поэтому нередко для оценки параметров модели с распределенным конечным числом лагов используются специальные методы преобразования, как и для модели с бесконечным числом лагов. Разработаны разные методы оценивания параметров моделей с распределенными лагами, которые учитывают характер распределения коэффициентов регрессии при лаговых объясняющих переменных. Иными словами, методы оценивания параметров модели с распределенными лагами основаны на изучении структуры лага. Так, предполагая полиномиальное распределение лаговых коэффициентов, используют *метод Алмон*, а при гипотезе геометрической прогрессии для лаговых коэффициентов применяется *преобразование Койка*.

### 7.2.2.1. Полиномиально распределенные лаги Алмон

В 1965 г. Ш. Алмон предложила способ оценки параметров модели с распределенными лагами на основе гипотезы о том, что лаговые коэффициенты регрессии аппроксимируются полиномом соответствующей степени от величины лага. Это значит, что в модели  $y_t = a + b_0x_t + b_1x_{t-1} + \dots + b_kx_{t-k} + \xi_t$  параметр  $b_j$  рассматривается как функция:  $b_j = c_0 + c_1j + c_2j^2 + \dots + c_mj^m$ . При этом априори выдвигается предположение о степени полинома. Как правило, используется многочлен невысокой степени ( $m \leq 4$ ).

Предположим, что  $b_j$  имеет распределение в виде параболы второй степени, т.е.  $b_j = c_0 + c_1j + c_2j^2$ . Тогда каждый из коэффициентов  $b_j$  можно представить в виде

$$b_0 = c_0;$$

$$b_1 = c_0 + c_1 + c_2;$$

$$\begin{aligned} b_2 &= c_0 + 2c_1 + 4c_2; \\ b_3 &= c_0 + 3c_1 + 9c_2; \\ &\dots \\ b_k &= c_0 + kc_1 + k^2c_2. \end{aligned}$$

Подставим эти соотношения для  $b_j$  в модель с распределенными лагами

$$y_t = a + c_0x_t + (c_0 + c_1 + c_2)x_{t-1} + (c_0 + 2c_1 + 4c_2)x_{t-2} + (c_0 + 3c_1 + 9c_2)x_{t-3} + \dots + (c_0 + kc_1 + k^2c_2)x_{t-k} + \xi_t.$$

Перегруппируем слагаемые с одинаковыми значениями  $s$ :

$$\begin{aligned} y_t &= a + c_0(x_t + x_{t-1} + x_{t-2} + \dots + x_{t-k}) + \\ &+ c_1(x_{t-1} + 2x_{t-2} + 3x_{t-3} + \dots + kx_{t-k}) + \\ &+ c_2(x_{t-1} + 4x_{t-2} + 9x_{t-3} + \dots + k^2x_{t-k}) + \xi_t. \end{aligned}$$

Будем рассматривать слагаемые в скобках при  $c_0, c_1$  и  $c_2$  как новые переменные  $z$ , т.е. модель с распределенными лагами примет вид

$$y_t = a + c_0z_0 + c_1z_1 + c_2z_2 + \xi_t,$$

где  $z_0, z_1$  и  $z_2$  определяются как

$$\begin{aligned} z_0 &= \sum_{j=0}^k x_{t-j}; & ; \\ z_1 &= \sum_{j=1}^k jx_{t-j}; \\ z_2 &= \sum_{j=1}^k j^2x_{t-j}. \end{aligned}$$

Оценка параметров при преобразованных переменных  $z$  дается традиционным МНК. При этом случайные отклонения  $\xi_t$  удовлетворяют предпосылкам МНК. Далее на основе параметров  $c_0, c_1$  и  $c_2$  переходим к оценке параметров  $b_j$ , используя выражения коэффициентов  $b_j$  через коэффициенты полинома:

$$b_j = c_0 + c_1j + c_2j^2.$$

В общем виде при степени полинома  $m$  модель регрессии с распределенными лагами примет вид

$$\begin{aligned}
 y_t = & a + c_0(x_t + x_{t-1} + \dots + x_{t-k}) + c_1(x_{t-1} + 2x_{t-2} + 3x_{t-3} + \dots + kx_{t-k}) + \\
 & + c_2(x_{t-1} + 4x_{t-2} + 9x_{t-3} + \dots + k^2x_{t-k}) + \dots \\
 & + c_m(x_{t-1} + 2^m x_{t-2} + 3^m x_{t-3} + \dots + k^m x_{t-k}) + \xi_t, \\
 \text{или } y_t = & a + c_0 z_0 + c_1 z_1 + c_2 z_2 + \dots + c_m z_m + \xi_t.
 \end{aligned}$$

Как видим, в данной модели переменные  $z_1, z_2, \dots, z_m$  представляют собой линейную комбинацию переменных  $x_t$  и  $k$  лаговых переменных, веса при которых подчиняются полиномиальному распределению (рис. 7.1).

В матричном виде можно записать, что  $\mathbf{b} = \mathbf{H}\mathbf{c}$ , где

$$\mathbf{H} = \begin{bmatrix} 1 & 0 & 0 & \dots & \dots & \dots & 0 \\ 1 & 1 & 1 & \dots & \dots & \dots & 1 \\ 1 & 2 & 4 & \dots & \dots & \dots & 2^m \\ 1 & 3 & 9 & \dots & \dots & \dots & 3^m \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & k & k^2 & \dots & \dots & \dots & k^m \end{bmatrix} \quad \text{— матрица весов при лаговых ко-}$$

эффициентах  $b_j$ ;  $\mathbf{c}$  — вектор коэффициентов при переменных  $z$ .

Тогда модель в целом принимает вид

$$\mathbf{y} = \mathbf{X}\mathbf{H}\mathbf{c} + \boldsymbol{\xi} = \mathbf{Z}\mathbf{c} + \boldsymbol{\xi}.$$

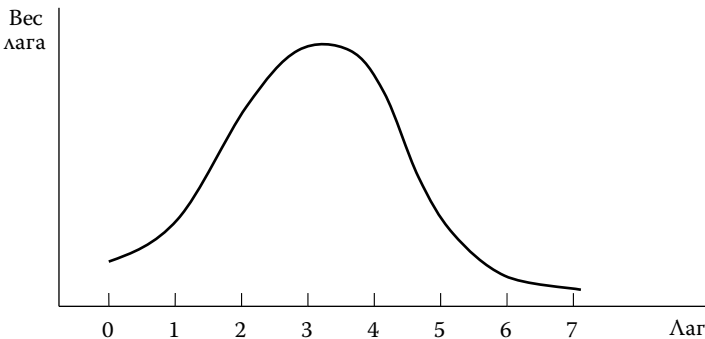


Рис. 7.1. Распределение лаговых переменных

Стандартная ошибка коэффициентов регрессии при лаговых переменных определится как  $\mu_{b_j} = S\sqrt{\mathbf{H}'[\mathbf{z}'\mathbf{z}]^{-1}\mathbf{H}}$ .

Далее через  $t$ -критерий Стьюдента оценивается значимость коэффициентов  $b_j$ .

Качество модели оценивается через коэффициент детерминации  $R^2$  для уравнения регрессии  $y_t$  от преобразованных переменных  $z$ , т.е. по модели

$$y = Zc + \xi.$$

Таким образом, применение метода Алмон включает в себя следующие этапы работы:

- 1) определение максимальной величины лага  $k$ ;
- 2) определение степени полинома  $m$ , описывающего распределение коэффициентов регрессии  $b_j$  в зависимости от величины лага;
- 3) расчет преобразованных переменных  $z_j$ ;
- 4) расчет параметров линейной регрессии  $y$  от преобразованных переменных  $z$ , т.е. оценка  $c_j$ ;
- 5) переход к исходным параметрам  $b_j$  модели с распределенными лагами.

Теоретически достаточно сложно определить максимальную величину лага  $k$ . В основном для этой цели используется экспериментальный путь: строится уравнение с большим числом последовательных лагов и с постепенным его уменьшением изучается значимость коэффициентов регрессии при лаговых объясняющих переменных. Останавливаются на варианте, для которого все коэффициенты регрессии статистически значимы.

Определение степени полинома  $m$  также связано с рядом трудностей. Формально можно изучать графически структуру лага (рис. 7.2).

Если с ростом величины лага  $j$  коэффициенты  $b_j$  описываются кривыми, представленными на рис. 7.2,  $a$  —  $z$ , то в расчетах могут быть использованы полиномы второй, третьей или четвертой степени. Рисунок 7.2,  $d$  предполагает линейную зависимость  $b_j$  от величины лага, а рис. 7.2,  $e$  показывает перевернутую  $V$ -образную структуру лага (например, при изучении капитальных вложений Де Люю в 1962 г. предложил подобную структуру лага). Однако учитывая, что оценки  $b_j$  по МНК часто

затруднены, исследователь, как правило, не располагает подобными графиками. Поэтому степень полинома задается исследователем, исходя из соответствующих теоретических соображений и результатов предыдущих исследований.

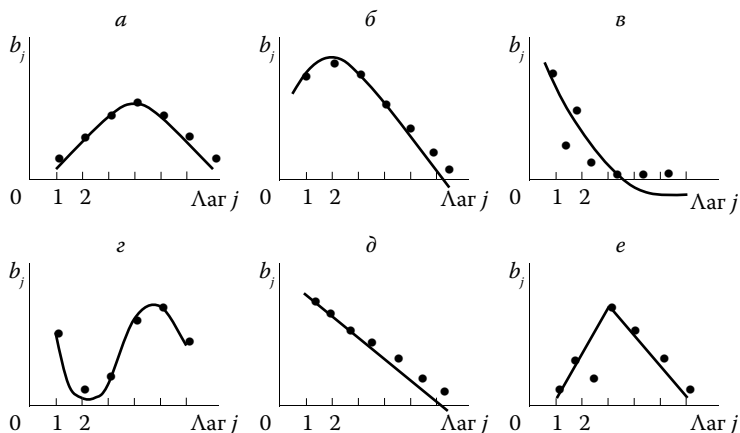


Рис. 7.2. Возможные распределения лаговых коэффициентов регрессии

### Пример 7.1

По данным за 32 квартала об объеме продукции ( $y$  — в млн руб.) и инвестициях в основной капитал ( $x$  — в млн руб.) строится модель с распределенными лагами

$$y_t = a + b_0 x_t + b_1 x_{t-1} + b_2 x_{t-2} + b_3 x_{t-3} + b_4 x_{t-4} + \xi_t.$$

Таблица 7.1. Объем продукции и инвестиции в основной капитал

№ квартала	$y_t$	$x_t$	$x_{t-1}$	$x_{t-2}$	$x_{t-3}$	$x_{t-4}$	$z_0$	$z_1$	$z_2$
1	5,2	0,87	—	—	—	—	—	—	—
2	5,6	0,9	0,87	—	—	—	—	—	—
3	6,5	1,05	0,9	0,87	—	—	—	—	—
4	6,4	1,04	1,05	0,9	0,87	—	—	—	—

Окончание табл. 7.1

№ квар-тала	$Y_t$	$x_t$	$x_{t-1}$	$x_{t-2}$	$x_{t-3}$	$x_{t-4}$	$z_0$	$z_1$	$z_2$
5	6,5	1,05	1,04	1,05	0,9	0,87	4,91	9,32	27,26
6	7,0	1,08	1,05	1,04	1,05	0,9	5,12	9,88	29,06
7	7,4	1,12	1,08	1,05	1,04	1,05	5,34	1,05	31,44
8	7,8	1,16	1,12	1,08	1,05	1,04	5,45	10,59	31,53
9	8,1	1,17	1,16	1,12	1,08	1,05	5,58	10,84	32,16
10	8,0	1,14	1,17	1,16	1,12	1,08	5,67	11,17	33,17
11	8,5	1,17	1,14	1,17	1,16	1,12	5,76	11,44	34,18
12	8,6	1,2	1,17	1,14	1,17	1,16	5,84	11,6	34,82
13	8,8	1,2	1,2	1,17	1,14	1,17	5,88	11,64	34,86
14	8,9	1,24	1,2	1,2	1,17	1,14	5,95	11,67	34,77
15	8,9	1,22	1,24	1,2	1,2	1,17	6,03	11,92	35,56
16	9,3	1,26	1,22	1,24	1,2	1,2	6,12	12,1	36,18
17	9,4	1,23	1,26	1,22	1,24	1,2	6,15	12,22	36,5
18	9,3	1,23	1,23	1,26	1,22	1,24	6,18	12,37	37,09
19	9,6	1,26	1,23	1,23	1,26	1,22	6,2	12,35	37,01
20	9,7	1,28	1,26	1,23	1,23	1,26	6,26	12,45	37,41
21	9,7	1,3	1,28	1,26	1,23	1,23	6,3	12,41	37,07
22	9,8	1,32	1,3	1,28	1,26	1,23	6,39	12,56	37,44
23	10,0	1,32	1,32	1,3	1,28	1,26	6,48	12,8	38,2
24	10,2	1,33	1,32	1,32	1,3	1,28	6,55	12,98	38,78
25	10,3	1,33	1,33	1,32	1,32	1,3	6,6	13,13	39,29
26	10,4	1,35	1,33	1,33	1,32	1,32	6,65	13,23	39,65
27	10,5	1,35	1,35	1,33	1,33	1,32	6,68	13,28	39,76
28	10,6	1,36	1,35	1,35	1,33	1,33	6,72	13,36	40
29	10,5	1,32	1,36	1,35	1,35	1,33	6,71	13,43	40,19
30	10,6	1,35	1,32	1,36	1,35	1,35	6,73	13,49	40,51
31	10,7	1,38	1,35	1,32	1,36	1,35	6,76	13,47	40,47
32	11	1,4	1,38	1,35	1,32	1,36	6,81	13,48	40,42

Предполагая квадратичную зависимость  $b_j$  от величины лага

$b_j = c_0 + c_1j + c_2j^2$ , имеем соотношения

$$b_0 = c_0;$$

$$b_1 = c_0 + c_1 + c_2;$$

$$b_2 = c_0 + 2c_1 + 4c_2;$$

$$b_3 = c_0 + 3c_1 + 9c_2;$$

$$b_4 = c_0 + 4c_1 + 16c_2.$$

Соответственно модель с распределенными лагами примет вид

$$Y_t = a + c_0 z_0 + c_1 z_1 + c_2 z_2 + \xi_t.$$

Расчет преобразованных переменных  $z_j$  представлен в табл. 7.1, где

$$z_0 = x_t + x_{t-1} + x_{t-2} + x_{t-3} + x_{t-4};$$

$$z_1 = x_{t-1} + 2x_{t-2} + 3x_{t-3} + 4x_{t-4};$$

$$z_2 = x_{t-1} + 4x_{t-2} + 9x_{t-3} + 16x_{t-4}.$$

Применяя к данным об  $y$ ,  $z_0$ ,  $z_1$  и  $z_2$  обычный МНК, получим следующее уравнение:

$$y_t = -4,7115 + 3,7713z_0 - 2,2668z_1 + 0,5065z_2 + \xi_t.$$

$t$	-15,0	5,8	-2,6	2,6
-----	-------	-----	------	-----

Все параметры уравнения статистически значимы ( $t_{0,05}$  при  $df = 242,06$ ).  $R^2 = 0,9955$  указывает на хорошее качество модели.

Далее найдем коэффициенты регрессии исходной модели, т.е.  $b_j$ , используя выражения  $b_j$  через коэффициенты  $c_0$ ,  $c_1$  и  $c_2$ :

$$b_0 = 3,7713;$$

$$b_1 = 3,7713 + (-2,2668) + 0,5065 = 2,011;$$

$$b_2 = 3,7713 - 2 \cdot 2,2668 + 4 \cdot 0,5065 = 1,2637;$$

$$b_3 = 3,7713 - 3 \cdot 2,2668 + 9 \cdot 0,5065 = 1,5294;$$

$$b_4 = 3,7713 - 4 \cdot 2,2668 + 16 \cdot 0,5065 = 2,8081.$$

Модель регрессии с распределенными лагами примет вид

$$\hat{y}_t = -4,711 + 3,771x_t + 2,011x_{t-1} + 1,264x_{t-2} + 1,529x_{t-3} + 2,808x_{t-4};$$

$$R^2 = 0,9955.$$

Стандартные ошибки коэффициентов регрессии по модели следующие:

$$m_{b_0} = 0,651; m_{b_1} = 0,296; m_{b_2} = 0,4045; m_{b_3} = 0,300; m_{b_4} = 0,494.$$

Для свободного члена  $a$  стандартная ошибка составила 0,313. Соответственно по  $t$ -критерию Стьюдента все параметры



оказались статистически значимыми:  $t_a = -15,0$ ;  $t_{b_0} = 5,8$ ;  $t_{b_1} = 6,8$ ;  $t_{b_2} = 3,1$ ;  $t_{b_3} = 5,1$ ;  $t_{b_4} = 5,7$ .

Модель показывает, что рост инвестиций в текущем периоде на 100 тыс. руб. способствует росту объема продукции в том же периоде в среднем на 377 тыс. руб., а через квартал — на 578 тыс. руб. В целом же через год прирост объема продукции за счет роста инвестиций на 100 тыс. руб. ожидается в размере 1,138 млн руб. ( $3,771 + 2,011 + 1,264 + 1,529 + 2,808 = 11,383$ ).

Определив относительные коэффициенты регрессии  $\beta_j$ , увидим, что половина воздействия фактора на результат реализуется с лагом в один квартал:

$$\beta_0 = \frac{3,771}{11,383} = 0,331;$$

$$\beta_1 = \frac{2,011}{11,383} = 0,177;$$

$$\beta_2 = \frac{1,264}{11,383} = 0,111;$$

$$\beta_3 = \frac{1,529}{11,383} = 0,134;$$

$$\beta_4 = \frac{2,808}{11,383} = 0,247.$$

На графике (рис. 7.3) рассматриваемые коэффициенты регрессии представляют собой параболу второй степени.

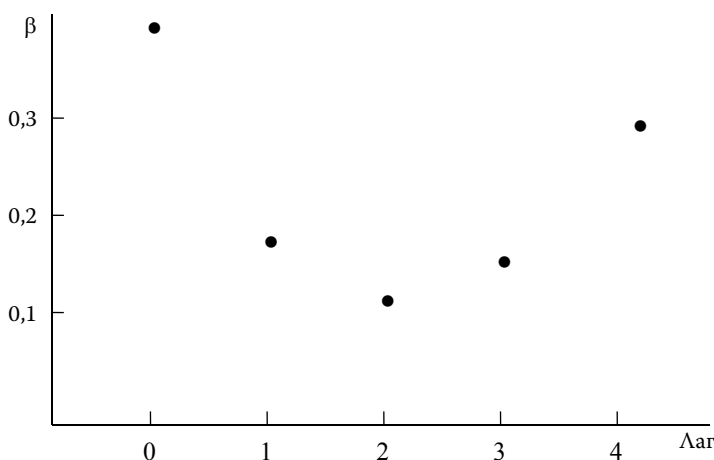
Если к исходным данным нашего примера применить традиционный МНК, то результаты окажутся следующими:

$$y_t = -4,782 + 4,390x_t + 1,193x_{t-1} + 1,372x_{t-2} + 2,001x_{t-3} + 2,479x_{t-4};$$

$$R^2 = 0,9957.$$

Хотя коэффициент детерминации здесь даже чуть-чуть выше, но коэффициенты регрессии при лаговых переменных  $x_{t-1}$  и  $x_{t-2}$  оказываются статистически незначимыми:

$$t_a = -14,7; \quad t_{b_0} = 5,1; \quad t_{b_1} = 1,2; \quad t_{b_2} = 1,7; \quad t_{b_3} = 3,1; \quad t_{b_4} = 4,1.$$



**Рис. 7.3.** Коэффициенты регрессии

Кроме того, применяя метод Алмон, получаем стандартные ошибки коэффициентов регрессии меньше, чем при традиционном МНК:

$$m_{b_0} = 0,868; \quad m_{b_1} = 0,983; \quad m_{b_2} = 0,811; \quad m_{b_3} = 0,654; \quad m_{b_4} = 0,610.$$

### 7.2.2.2. Метод Койка

Для модели с бесконечным числом лаговых значений объясняющей переменной

$$y_t = a + b_0x_t + b_1x_{t-1} + b_2x_{t-2} + \dots + \xi_t \quad (7.1)$$

оценка параметров не представляется возможной без какого-либо допущения относительно поведения коэффициентов при лаговых переменных. Одним из допущений является предположение о том, что после некоторой длины лага (например,  $k$ ) коэффициенты распределенного лага начнут убывать геометрически с одинаковым темпом  $\lambda$   $0 < \lambda < 1$ ). Тогда уравнение (7.1) может быть записано в виде

$$y_t = a + b_0x_t + b_1x_{t-1} + \dots + b_kx_{t-k} + b_k\lambda x_{t-k-1} + b_k\lambda^2 x_{t-k-2} + \dots + U_t. \quad (7.2)$$

В уравнении (7.2) первые  $k$  коэффициентов распределенного лага являются свободными (принимают любые значения), а остальные лаговые коэффициенты убывают в геометрической прогрессии.

Если в уравнении (7.2) предположить, что убывание лаговых коэффициентов в геометрической прогрессии происходит сразу же, а не через интервал времени  $k$ , то получим следующую модель:

$$y_t = a + b_0 x_t + b_0 \lambda x_{t-1} + b_0 \lambda^2 x_{t-2} + \dots + U_t. \quad (7.3)$$

Коэффициенты данной модели согласовываются с коэффициентами уравнения (7.1), а именно

$$b_j = b_0 \lambda^j; \quad j = 0, 1, 2, \dots \quad (7.4)$$

Это означает, что оценив три параметра уравнения (7.3), т.е.  $a$ ,  $b_0$  и  $\lambda$ , можно перейти к модели (7.1):  $a$  и  $b_0$  определены по модели (7.3),  $b_1 = b_0 \lambda$ ;  $b_2 = b_1 \lambda = b_0 \lambda^2$ ;  $b_3 = b_2 \lambda = b_0 \lambda^3$  и т.д.

Однако наличие в модели (7.3) бесконечного числа лаговых переменных затрудняет практическую ее реализацию, ибо исследователь имеет дело, как правило, с конечным числом лагов. Оценка параметров модели (7.3) возможна, если применить преобразование Койка.

Предполагая, что в модели (7.1) все лаговые коэффициенты имеют одинаковый знак и уменьшаются в геометрической прогрессии, Л. М. Койк предложил для оценки параметров модели (7.3) следующую процедуру:

- построить модель (7.3) для момента времени  $(t - 1)$ , т.е. получить уравнение

$$y_{t-1} = a + b_0 x_{t-1} + b_0 \lambda x_{t-2} + b_0 \lambda^2 x_{t-3} + \dots + \xi_{t-1}; \quad (7.5)$$

- умножить уравнение (7.5) на  $\lambda$ , т.е. получить уравнение

$$\lambda y_{t-1} = \lambda a + b_0 \lambda x_{t-1} + b_0 \lambda^2 x_{t-2} + b_0 \lambda^3 x_{t-3} + \dots + \lambda \xi_{t-1}; \quad (7.6)$$

- вычесть из уравнения (7.3) уравнение (7.6):

$$y_t - \lambda y_{t-1} = (1 - \lambda)a + b_0 x_t + (\xi_t - \lambda \xi_{t-1});$$

— после преобразования получить уравнение

$$y_t = (1 - \lambda)a + b_0x_t + \lambda y_{t-1} + U_t, \quad (7.7)$$

где  $U_t = \xi_t - \lambda \xi_{t-1}$ .

Уравнение (7.7) получило название **преобразование Койка**, так как Л. М. Койк впервые (1954 г.) предложил данный подход к оцениванию параметров модели с распределенными лагами.

Практически в модели (7.7) от уравнения с распределенными лагами с бесконечным их числом (7.1) Л. М. Койк перешел к модели авторегрессии, для которой требуется оценить всего три параметра:  $a$ ,  $b_0$  и  $\lambda$ . Далее из соотношения (7.4) наводятся параметры исходной модели (7.1).

Рассмотренный подход нашел широкое применение в исследовании кумулятивного эффекта рекламы на объем продаж, т.е. текущий объем продаж рассматривается в зависимости от расходов на рекламу текущего периода, объема продаж в предыдущий период времени и ошибки  $U_t$ <sup>1</sup>.

Преобразование Койка может быть использовано и при решении модели (7.2), когда несколько первых коэффициентов остаются свободными, а для оставшихся лагов реализуется данное преобразование. Например, считая, что  $b_0$  и  $b_1$  остаются свободными, а начиная с  $b_2$  все лаговые коэффициенты убывают с одинаковым темпом, можно записать

$$y_t = a + b_0x_t + b_1x_{t-1} + b_2x_{t-2} + \lambda b_2x_{t-3} + \lambda^2 b_2x_{t-4} + \dots + U_t.$$

Далее после применения преобразования Койка получается уравнение

$$y_t = a(1 - \lambda) + b_0x_t + (b_1 - \lambda b_0)x_{t-1} + (b_2 - \lambda b_1)x_{t-2} + \lambda y_{t-1} + (U_t - \lambda U_{t-1}),$$

т.е. происходит переход к модели авторегрессии с распределенными лагами.

Преобразование Койка приводит к существенным упрощениям, ибо вместе с уменьшением числа оцениваемых параметров

<sup>1</sup> См. подробное изложение в кн. : Берндта, Э. Практика эконометрики: классика и современность : пер. с англ. М. : ЮНИТИ, 2005. С. 457—467.

устраняется и проблема мультиколлинеарности факторов: теперь в модели (7.7) содержится две независимые переменные  $x_t$  и  $y_{t-1}$ .

Модель Койка позволяет анализировать краткосрочный и долгосрочный мультипликаторы. Краткосрочным мультипликатором является параметр  $b_0$ , а долгосрочным — сумма коэффициентов регрессии, представляющая собой сумму геометрической прогрессии

$$\sum_{j=0}^{\infty} b_j = b_0 + b_0\lambda + b_0\lambda^2 + b_0\lambda^3 + \dots = b_0(1 + \lambda + \lambda^2 + \lambda^3 + \dots) = b_0 \frac{1}{1 - \lambda}.$$

Например, по Великобритании для периода 1924—1938 гг. была построена модель  $C_t = 0,18R_t + 0,81C_{t-1}$ <sup>1</sup>, где  $C_t$  — потребление в период времени  $t$ ;  $R_t$  — доход в период времени  $t$ ;  $C_{t-1}$  — потребление в период времени  $(t - 1)$ .

Данное уравнение означает, что краткосрочная склонность к потреблению составляет 0,18, а долгосрочная склонность к потреблению равна  $0,18 / (1 - 0,81) = 0,95$ . Отсутствие в модели свободного члена не изменяет суть интерпретации краткосрочного и долгосрочного мультипликаторов, хотя естественно сказывается на величине параметров модели. Краткосрочный мультипликатор 0,18 показывает, что с ростом дохода на 1 ден. ед. потребление в тот же период времени увеличивается на 0,18 ден. ед. Долгосрочный мультипликатор 0,95 означает, что в долгосрочной перспективе увеличение дохода на 1 ден. ед. приведет к росту потребления на 0,95 ден. ед.

В модели Койка (7.7) случайная ошибка  $U_t = \xi_t - \lambda\xi_{t-1}$  коррелирована с переменной  $y_{t-1}$ . Поэтому оценивание параметров ее модели традиционным МНК дает смещенные и несостоятельные оценки. Вместо МНК могут быть применены инструментальные переменные (см. 7.3.2) или метод максимального правдоподобия.

Поскольку уравнение (7.7) является моделью авторегрессии, то остатки  $U_t$  могут быть автокоррелированы. Для их анализа не применим рассмотренный ранее критерий Дарбина — Уотсона ( $DW$ ). Вместо него необходимо использовать  $h$ -статистику Дарбина (см. далее формулу (7.15)).

<sup>1</sup> Маленво Э. Статистические методы эконометрии : пер. с фр. М. : Статистика, 1975. С. 137.

### 7.3. Модели авторегрессии

Преобразование Койка сворачивает модель с распределенными лагами к модели авторегрессии, т.е. к модели, в правой части которой используется лаговая зависимая переменная. Это не единственный вид авторегрессионных моделей. Но все же достаточно распространенный:

$$y_t = a + b_0 x_t + c_1 y_{t-1} + \xi_t. \quad (7.8)$$

Между тем интерпретация параметров данной модели имеет свою специфику, что и будет рассмотрено ниже.

#### 7.3.1. Интерпретация параметров модели авторегрессии

Для модели (7.8), как и в модели с распределенными лагами, параметр  $b_0$  характеризует краткосрочное изменение  $y_t$  под воздействием изменения  $x_t$  на одну единицу. Параметр  $c_1$  по существу представляет собой величину  $\lambda$  из преобразования Койка, т.е.  $|c_1| < 1$ , и показывает коэффициент снижения лаговых коэффициентов при увеличении величины лага в соответствии с концепцией их геометрического убывания. Следовательно, к моменту времени  $(t + 1)$  результат  $y$  изменится дополнительно на  $b_0 c_1$ , а к моменту времени  $(t + 2)$  дополнительное изменение  $y$  составит  $b_0 c_1^2$  единиц, к моменту времени  $(t + 3)$  —  $b_0 c_1^3$  единиц и т.д. Соответственно долгосрочный мультипликатор окажется равным

$$b = b_0 + b_0 c_1 + b_0 c_1^2 + b_0 c_1^3 + \dots$$

(в предположении бесконечного числа лагов).

С учетом геометрической прогрессии лаговых коэффициентов величина долгосрочного мультипликатора составит

$$b = b_0 (1 + c_1 + c_1^2 + c_1^3 + \dots) = \frac{b_0}{1 - c_1}.$$

Предположим, что по региону по данным временных рядов построена модель авторегрессии, описывающая зависимость сбережений на душу населения за год ( $S_t$  — в тыс. ден. ед.)

от среднедушевого совокупного годового дохода ( $y_t$  — в тыс. ден. ед.) и сбережений предшествующего года ( $S_{t-1}$ )

$$S_t = 52 + 0,24Y_t + 0,15S_{t-1} + \xi_t.$$

Уравнение показывает, что краткосрочное изменение размера сбережений с ростом дохода на 1 тыс. ден. ед. составляет в том же году 0,24 тыс. ден. ед. Через год рост дохода на 1 тыс. ден. ед. увеличит размер сбережений на 0,276 тыс. ден. ед. ( $0,24 + 0,24 \cdot 0,15$ ), т.е. дополнительно за год прирост составит 0,036 тыс. ден. ед. В дальнейшем величина дополнительного прироста будет убывать. Долгосрочный мультипликатор окажется равным 0,282 тыс. ден. ед. ( $0,24/0,85$ ). Его величина характеризует прирост сбережений в долгосрочной перспективе с ростом дохода на 1 тыс. ден. ед. Трактовка данного мультипликатора на примере зависимости потребления от доходов была показана в 7.2.2.2.

### 7.3.2. Инструментальные переменные как метод оценивания параметров модели авторегрессии

В силу того, что в модели авторегрессии в правой части содержатся лаговые эндогенные переменные, принято считать, что оценка параметров традиционным МНК дает неудовлетворительные результаты.

Предположим, что рассматривается модель авторегрессии вида (7.8).

Применение для оценивания параметров это уравнения традиционного МНК возможно, если выполняется предпосылка МНК относительно отсутствия автокорреляции остатков. Между тем при наличии в правой части лаговой зависимой переменной может иметь место автокорреляция остатков. Кроме того, может иметь место и зависимость объясняющей переменной  $Y_{t-1}$  с остатками  $\xi_t$ , т.е. нарушается предпосылка о гомоскедастичности остатков. В силу этого классический МНК в случае малых выборок даст смещенные оценки параметров.

Одним из возможных методов оценивания параметров модели (7.8) является *метод инструментальных переменных*. Суть метода состоит в том, что вместо лаговой зависимой переменной  $y_{t-1}$ , для которой нарушается предпосылка МНК,

используется другая переменная  $z$ , называемая **инструментальной**. При этом инструментальная переменная должна обладать двумя свойствами:

- она должна быть тесно коррелирована с лаговой переменной  $y_{t-1}$ ;
- она не должна коррелировать с остатками  $U_t$  (случайными ошибками).

Иными словами, от модели авторегрессии (7.8) необходимо перейти к модели вида

$$y_t = a + bx_t + cz_t + \xi_t. \quad (7.9)$$

Результаты регрессии по модели (7.9), естественно, зависят от того, насколько удачно подобрана инструментальная переменная. В качестве инструментальной переменной можно, например, взять оценку  $y_{t-1}$ , т.е.  $\hat{y}_{t-1}$ , полученную по регрессии  $y_{t-1}$  от  $x_{t-1}$ .

Поскольку в модели (7.9) предполагается наличие зависимости  $y_t$  от  $x_t$ , то можно предположить, что также имеет место зависимость  $y_{t-1}$  от  $x_{t-1}$ , т.е. найдем регрессию

$$\hat{y}_{t-1} = A + Bx_{t-1}. \quad (7.10)$$

Используя для оценки параметров уравнения (7.10) обычный МНК, что возможно ввиду отсутствия в правой части модели лаговой зависимой переменной, найдем теоретические значения  $\hat{y}_{t-1}$ , которые и будут рассматриваться как значения инструментальной переменной  $z$  в модели (7.9). Далее вновь применяем МНК уже к модели (7.9), т.е. по существу оценка параметров модели авторегрессии (7.8) будет найдена исходя из модели вида

$$y_t = a + b_0x_t + c_1\hat{y}_{t-1} + \xi_t. \quad (7.11)$$

Если вместо оценки  $\hat{y}_{t-1}$  подставить выражение (7.10), то получим следующую модель:

$$y_t = a + c_1A + b_0x_t + Bc_1x_{t-1} + \xi_t. \quad (7.12)$$



Она представляет собой модель с распределенным лагом, оценка параметров которой может быть дана МНК.

Таким образом, используя в качестве инструментальной переменной оценки  $\hat{y}_{t-1}$ , исходя из регрессии от  $x_{t-1}$  (7.10), модель авторегрессии (7.8) заменяют на модель с распределенным лагом (7.12).

Вместе с тем следует отметить, что применение рассмотренной инструментальной переменной может привести при практической реализации модели (7.8) к появлению коллинеарности факторов. Объясняется это тем, что в модель (7.8) одновременно вводятся в качестве объясняющих переменных линейно связанные и высокоррелируемые между собой  $\hat{y}_{t-1}$  и  $x_t$ , ибо  $\hat{y}_{t-1} = A + Bx_{t-1}$  и  $r_{\hat{y}_{t-1}x_{t-1}} = 1$ , а соответственно и  $r_{\hat{y}_{t-1}x_t}$  будет близок к единице. Однако если коллинеарность факторов не повлекла за собой неверные знаки у коэффициентов регрессии и не привела к большим стандартным ошибкам оценок, то применение инструментальной переменной можно считать возможным.

### Пример 7.2

Применим метод инструментальных переменных к модели авторегрессии (7.8) по данным фирмы об импорте сырья ( $y$  — в т) товара и величине производства ( $x$  — в тыс. ед.) за январь — декабрь 2008—2009 гг.

	Месяцы											
	1	2	3	4	5	6	7	8	9	10	11	12
	2008 г.											
$y$	164	162	165	168	172	177	182	186	187	191	196	201
$x$	78	81	89	76	105	101	93	94	107	103	116	170
	2009 г.											
$y$	213	211	219	228	232	239	244	249	255	264	265	267
$x$	101	110	138	145	180	165	144	130	155	142	151	305

Рассмотрим модель (7.8)

$$y_t = a + b_0x_t + c_1y_{t-1} + \xi_t.$$

Для оценивания параметров этой модели введем инструментальную переменную  $z = \hat{y}_{t-1} = A + Bx_{t-1}$ . Используя МНК, получим уравнение регрессии

$$\hat{y}_{t-1} = 102,4025 + 0,8705x_{t-1}; \quad R^2 = 0,6396;$$

$$t \quad 5,77 \quad 6,10; \quad F = 37,26.$$

Уравнение регрессии значимо, как и его параметры. Подставляя в это уравнение значения  $x_{t-1}$ , получим расчетные значения  $\hat{y}_{t-1}$ : 170,3; 172,9; 179,9; ...; 233,8. Далее вновь применяем МНК к модели (7.8), в которой вместо фактических значений  $y_{t-1}$  используются расчетные величины, т.е.  $\hat{y}_{t-1}$ . Результаты оказались следующими:

$$y_t = 12,6808 + 0,2636x_t + 0,7946y_{t-1} + \xi_t; \quad R^2 = 0,7819;$$

$$t \quad 0,4 \quad 2,9 \quad 4,86; \quad F = 35,8.$$

Уравнение авторегрессии в целом значимо, значимыми являются и коэффициенты регрессии.

Если к модели (7.8) сразу же применить МНК, т.е. без введения инструментальной переменной, то результаты окажутся следующими:

$$y_t = -0,4948 - 0,012x_t + 1,0315y_{t-1} + \xi_t; \quad R^2 = 0,9915;$$

$$t \quad -0,1 \quad -0,5 \quad 33,1.$$

Хотя коэффициент детерминации для модели, оцененной по обычному МНК, выше, чем для модели с инструментальной переменной, но коэффициент регрессии при  $x_t$  не только статистически не значим, но и имеет неверный знак, ибо увеличение объема продукции, на производство которой требуется ввоз сырья, ведет к росту величины импорта, что и показывает модель авторегрессии, оцененная с помощью метода инструментальных переменных.

### 7.3.3. Оценка автокорреляции остатков по модели авторегрессии

Рассмотренный ранее критерий Дарбина — Уотсона не применим для моделей авторегрессии, содержащих в составе

объясняющих переменных лаговые значения зависимой переменной. Связано это с тем, что критерий Дарбина — Уотсона для модели авторегрессии может принимать значение, близкое к двум, как при отсутствии, так и при наличии автокорреляции остатков.

Предположим, что в модели авторегрессии (7.8) имеет место автокорреляция остатков, т.е. случайное отклонение  $\xi_t$  можно рассматривать как авторегрессию вида

$$\xi_t = \rho \xi_{t-1} + U_t, \quad (7.13)$$

где  $\rho$  — коэффициент автокорреляции первого порядка;  $U_t$  — случайная составляющая.

Тогда уравнение (7.8) можно представить как

$$y_t = a + b_0 x_t + c_1 y_{t-1} + \rho \xi_{t-1} + U_t. \quad (7.14)$$

В уравнении (7.14)  $y_{t-1}$  связан с  $\xi_{t-1}$ , как и по уравнению (7.8)  $y_t$  связан с  $\xi_t$ . Таким образом, имеется систематическая связь лаговой зависимой переменной со случайной компонентой. Применение теста Дарбина — Уотсона к модели (7.14) может показать отсутствие автокорреляции в остатках  $U_t$  при наличии ее для остатков  $\xi_t$ . Как указано в работе Э. Маленво<sup>1</sup>, критерий Дарбина — Уотсона теряет мощность в авторегрессионных моделях. Дж. Дарбин предложил для моделей авторегрессии при оценке существенности автокорреляции остатков использовать другой критерий, который в литературе название ***h-статистика Дарбина***:

$$h = \rho \sqrt{\frac{n}{1-nV}}, \quad (7.15)$$

где  $\rho$  — коэффициент автокорреляции в остатках первого порядка, который практически используется при расчете критерия Дарбина — Уотсона, т.е.

$$\rho = \frac{\sum e_t e_{t-1}}{\sum e_t^2} = 1 - \frac{DW}{2},$$

<sup>1</sup> Маленво Э. Указ. соч. С. 148.

где  $n$  — число наблюдений в модели;  $V$  — выборочная дисперсия коэффициента при лаговой зависимой переменной  $y_{t-1}$ .

При большом числе наблюдений и при отсутствии в остатках автокорреляции первого порядка  $h$ -статистика Дарбина подчиняется стандартизированному нормальному распределению. Поэтому фактическое значение  $h$  сравнивается с табличным по заданному уровню значимости  $\alpha$ . Если  $|h|$  больше критического значения, то нулевая гипотеза об отсутствии автокорреляции ошибок отклоняется. При практических расчетах чаще всего  $\alpha$  берется как 0,05 и если  $|h| > 1,96$ , то гипотеза об отсутствии автокорреляции остатков отвергается.

Из уравнения (7.15) следует, что  $h$ -статистика не применима, если величина  $(nV) \geq 1$ . Кроме того, данный критерий предназначен для больших выборок (например, для  $n > 30$ ).  $h$ -статистика зависит от квадрата стандартной ошибки параметра только при лаговой зависимой переменной  $y_{t-1}(V)$  и не зависит от числа лагов, используемых в модели авторегрессии. Так, для модели  $y_t = a + b_0x_t + c_1y_{t-1} + c_2y_{t-2} + \xi_t$  оценка значимости автокорреляции остатков также проводится с помощью  $h$ -статистики Дарбина.

Аналогично данный критерий используется и для модели авторегрессии с несколькими экзогенными переменными

$$y_t = a + b_0x_{1t} + b_1x_{2t} + b_2x_{3t} + c_1y_{t-1} + \xi_t.$$

В рассматриваемом примере автокорреляция остатков не устранена, о чем свидетельствует  $h$ -статистики Дарбина: коэффициент автокорреляции в остатках  $\rho$  составил 0,440; стандартная ошибка коэффициента регрессии при переменной  $y_{t-1}$  оказалась равной 0,1635 (0,7946/4,86); соответственно  $V = 0,1635^2 = 0,026732$  и при  $n = 23$   $h = 0,440 \sqrt{\frac{23}{1 - 23 \cdot 0,026732}} = 3,4$ , что больше необходимого 1,96.

Автокорреляция в остатках по авторегрессионным моделям может быть устранена с помощью авторегрессионных преобразований с использованием моделей *ARMA* и *ARIMA*.

## 7.4. Авторегрессионные процессы и их моделирование (общая характеристика)

### 7.4.1. Авторегрессионные процессы

Рассмотренные ранее модели авторегрессии содержали в правой части наряду с лаговыми зависимыми переменными  $y_{t-1}$ ,  $y_{t-2}$  и т.п. независимые переменные  $x$ . Авторегрессионная модель, в которой отсутствуют независимые переменные и  $y_t$  рассматривается как линейная функция только предыдущих своих значений, представляет собой **авторегрессионный процесс**

$$y_t = a_0 + a_1 y_{t-1} + a_2 y_{t-2} + \dots + a_p y_{t-p} + \xi_t. \quad (7.16)$$

В зависимости от того, сколько предыдущих уровней временного ряда включено в уравнение (7.16), авторегрессионный процесс может быть разного порядка. Если текущее значение уровня динамического ряда  $y_t$  рассматривается как линейная функция от одного предыдущего значения, то имеем дело с авторегрессионным процессом первого порядка, что обычно в англоязычной литературе обозначается как  $AR(1)$ :

$$y_t = a_0 + a_1 y_{t-1} + \xi_t. \quad (7.17)$$

Увеличивая число лаговых переменных в модели (7.17), получим авторегрессионный процесс более высокого порядка. Например, процесс  $AR(3)$  сводится к уравнению

$$y_t = a_0 + a_1 y_{t-1} + a_2 y_{t-2} + a_3 y_{t-3} + \xi_t \quad (7.18)$$

и отражает авторегрессионный процесс третьего порядка.

Процессы  $AR$  могут быть стационарными и нестационарными. Чтобы процесс был стационарным, коэффициенты  $a_1, a_2, \dots, a_p$  в модели (7.16) должны образовывать сходящийся ряд и все корни характеристического уравнения  $1 - a_1 z - a_2 z^2 - \dots - a_p z^p = 0$  (вещественные и комплексные) должны лежать вне единичного круга, т.е.  $|z| > 1$ .

Рассмотренное условие стационарности для процесса  $AR(1)$  означает, что в уравнении (7.17) параметр  $a_1$  должен

соответствовать величине  $|a| < 1$ , так как характеристическое уравнение  $1 - a_1z = 0$  имеет корень  $z = \frac{1}{a_1}$  и  $|z| > 1$  при  $|a_1| < 1$ . Так, для ряда  $y_t = 3 + 0,8y_{t-1} + \xi_t$  при  $y_0 = 2$  (начальный уровень динамического ряда) характеристическое уравнение имеет вид  $1 - 0,8z = 0$ . Соответственно  $z = 1,25$  и рассматриваемый процесс является стационарным. Его асимптота  $\mu$  окажется равной  $\mu = \frac{a_0}{1 - a_1}$ , т.е. имеем  $\mu = \frac{3}{(1 - 0,8)} = 15$ , и траектория процесса флуктуирует и не превышает 15. Так, при  $n = 100$   $\hat{y}_t$  принимает значения 2; 4,6; 6,7; 8,3; 9,7; 10,7; 11,6; 12,3; 12,8; 13,3; 13,6; 13,8; 14,1 и далее возрастает до 15, а начиная с  $t = 27$  не превышает 15.

Предположим, что рассматривается процесс  $AR(2)$ , а именно

$$y_t = 5 + 0,9y_{t-1} + 0,2y_{t-2} + \xi_t.$$

Для него характеристическое уравнение составит

$$1 - 0,9z - 0,2z^2 = 0.$$

Корни этого уравнения составят  $z_1 = 2,5$  и  $z_2 = 2$ , что больше единицы и, следовательно, процесс является стационарным. Асимптота данного ряда окажется равной  $\mu = \frac{5}{1 - 0,9 + 2} = 16,6$ , т.е. начиная с  $t = 14$   $\hat{y}_t$  варьирует вокруг величины 16, (6): уровни ряда принимают значения 7; 6; 9; 11,9; 13,9; 15,1 и т. д. В рассмотренных примерах  $AR(1)$  и  $AR(2)$  динамические ряды обнаруживают вначале некоторую тенденцию, которая постепенно затухает и ряд становится стационарным (рис. 7.4).

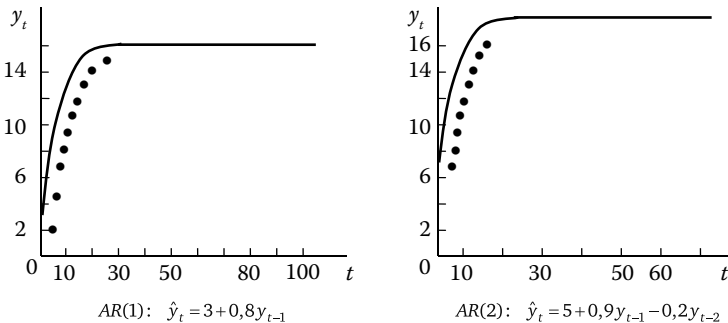


Рис. 7.4. Асимптотически стационарные временные ряды

Авторегрессионный процесс с большим числом лагов предполагает очень длинные динамические ряды, которые далеко не всегда имеются в эконометрических исследованиях. При наличии коротких временных рядов стационарные *AR*-процессы могут иметь место после удаления из уровней ряда тенденции и сезонных колебаний. Это означает, что исследователь должен вычленил эти компоненты динамического ряда и подвергнуть дальнейшей обработке остаточные величины. В этом случае авторегрессионный процесс первого порядка *AR*(1) примет вид

$$\xi_t = a\xi_{t-1} + V_t, \quad (7.19)$$

где  $\xi_t$  — остатки после устранения из уровней ряда  $y_t$  тенденции и периодической составляющей;  $V_t$  — белый шум.

#### 7.4.2. Модели скользящей средней

Среди моделей для стационарных временных рядов широкое распространение имеют *модели скользящей средней*.

Для стационарного ряда моделируемый уровень временного ряда можно представить как линейную функцию прошлых ошибок, т.е. разностей между прошлыми фактическими и теоретическими уровнями:

$$y_t = \mu + \xi_t - \theta_1\xi_{t-1} - \theta_2\xi_{t-2} - \dots - \theta_q\xi_{t-q}, \quad (7.20)$$

где  $\mu$  — константа;  $\xi_t, \xi_{t-1}, \dots, \xi_{t-q}$  — белый шум в текущий и предыдущий период времени:

$$\xi_t = y_t - \hat{y}_t.$$

Термин «скользящая средняя», используемый здесь, не синоним скользящей средней как методу сглаживания уровней динамического ряда.

В модели (7.20) уровень динамического ряда рассматривается как сумма константы  $\mu$  и скользящей средней между текущими и предыдущими значениями белого шума (случайных отклонений).

Обозначим скользящую среднюю модели (7.20) через  $x_t$ :

$$x_t = \xi_t - \theta_1\xi_{t-1} - \theta_2\xi_{t-2} - \dots - \theta_q\xi_{t-q}. \quad (7.21)$$

Уравнение (7.21) принято называть **процессом скользящего среднего порядка  $q$**  и обозначать как *МА* ( $q$ ) от английского *Moving Average*. Порядок скользящей средней определяется числом учитываемых в модели предыдущих значений случайных отклонений. Так, *МА* (2) можно записать как  $x_t = \xi_t - \theta_1 \xi_{t-1} - \theta_2 \xi_{t-2}$ , а модель уровня динамического ряда с использованием *МА* (2) будет иметь вид

$$y_t = \mu + \xi_t - \theta_1 \xi_{t-1} - \theta_2 \xi_{t-2}.$$

Соответственно модель уровня ряда с использованием *МА*(1) примет вид

$$y = \mu + \xi_t - \theta_1 \xi_{t-1}.$$

При  $q = 0$  и  $\mu = 0$  получаем процесс белого шума.

Временные ряды с использованием процесса скользящего среднего могут иметь место, когда уровни динамического ряда характеризуются случайной колеблемостью.

### 7.4.3. Модели *АРМА*

Соединение в одной модели авторегрессионного процесса *AR* и модели скользящего среднего *МА* приводит к модели **авторегрессионного процесса со скользящими средними в остатках** (*ARMA* — от английского *Auto Regressive — Moving Average*):

$$y_t = a_0 + a_1 y_{t-1} + a_2 y_{t-2} + \dots + a_p y_{t-p} + \xi_t - \theta_1 \xi_{t-1} - \theta_2 \xi_{t-2} - \dots - \theta_q \xi_{t-q} \quad (7.22)$$

В модели (7.22) в качестве объясняющих переменных рассматриваются лаговые значения зависимой переменной с  $p$  интервалами сдвига и скользящие средние порядка  $q$  для остатков авторегрессии. Иными словами, модель включает в себя *AR* ( $p$ ) и *МА* ( $q$ ). Ее принято обозначать *ARMA* ( $p, q$ ). Например, *ARMA* (3, 2) имеет вид

$$y_t = a_0 + a_1 y_{t-1} + a_2 y_{t-2} + a_3 y_{t-3} + \xi_t - \theta_1 \xi_{t-1} - \theta_2 \xi_{t-2}. \quad (7.23)$$

При практической реализации моделей *ARMA* наиболее сложным является выбор числа лагов  $p$  и  $q$ .

Инструментом идентификации модели *ARMA* является изучение частной автокорреляционной функции по моделям с разным числом лагов. **Частная автокорреляционная функция** (*PACF* — *Partial Autocorrelation Function*) представляет собой



серии частных коэффициентов автокорреляции (РАС), которые измеряют связь между текущим уровнем динамического ряда  $y_t$  и предыдущими значениями  $y_{t-1}, y_{t-2}, \dots, y_{t-k}$  в условиях, когда влияние других промежуточных временных лагов устранено. Так, частный коэффициент автокорреляции при лаге  $k$  будет представлять собой корреляцию  $y_t$  и  $y_{t-k}$ , очищенную от влияния  $y_{t-1}, y_{t-2}, \dots, y_{t-k-1}$ .

Обозначим частный коэффициент автокорреляции с лагом  $k$  через  $\rho(k)$ . При  $k = 0$   $\rho(0) = 1$  (уровни ряда коррелируют сами с собой); при  $k = 1$   $\rho(1) = \tau_{a_1}$ , где  $\tau_{a_1}$  — коэффициент автокорреляции первого порядка. Это равенство связано с тем, что при расчете  $\rho(1)$  отсутствуют промежуточные лаги. Вычисление  $\rho$  более высокого порядка можно производить по формулам

$$\rho(2) = \frac{\begin{vmatrix} 1 & \tau_{a_1} \\ \tau_{a_1} & \tau_{a_2} \end{vmatrix}}{\begin{vmatrix} 1 & \tau_{a_1} \\ \tau_{a_1} & 1 \end{vmatrix}} = \frac{(\tau_{a_2} - \tau_{a_1}^2)}{(1 - \tau_{a_1}^2)};$$

$$\rho(3) = \frac{\begin{vmatrix} 1 & \tau_{a_1} & \tau_{a_2} \\ \tau_{a_1} & 1 & \tau_{a_2} \\ \tau_{a_2} & \tau_{a_1} & \tau_{a_3} \end{vmatrix}}{\begin{vmatrix} 1 & \tau_{a_1} & \tau_{a_2} \\ \tau_{a_1} & 1 & \tau_{a_1} \\ \tau_{a_2} & \tau_{a_1} & 1 \end{vmatrix}};$$

...

$$\rho(k) = \frac{\begin{vmatrix} 1 & \tau_{a_1} & \tau_{a_2} & \dots & \tau_{a_1} \\ \tau_{a_1} & 1 & \tau_{a_1} & \dots & \tau_{a_2} \\ \tau_{a_2} & \tau_{a_1} & 1 & \dots & \tau_{a_3} \\ \dots & \dots & \dots & \dots & \dots \\ \tau_{a_{k-1}} & \tau_{a_{k-2}} & \tau_{a_{k-3}} & \dots & \tau_{a_k} \end{vmatrix}}{\begin{vmatrix} 1 & \tau_{a_1} & \tau_{a_2} & \dots & \tau_{a_{k-1}} \\ \tau_{a_1} & 1 & \tau_{a_1} & \dots & \tau_{a_{k-2}} \\ \tau_{a_2} & \tau_{a_1} & 1 & \dots & \tau_{a_{k-3}} \\ \dots & \dots & \dots & \dots & \dots \\ \tau_{a_{k-1}} & \tau_{a_{k-2}} & \tau_{a_{k-3}} & \dots & 1 \end{vmatrix}}.$$

В данных формулах  $\rho(k)$  определитель числителя отличается от определителя в знаменателе только заменой последнего столбца определителя в знаменателе столбцом из значений  $\rho(1), \rho(2), \dots, \rho(k)$ .

Для авторегрессионного процесса порядка  $p$  частная автокорреляционная функция отлична от нуля при  $k \leq p$  и равна нулю при  $k > p$ . Это и позволяет определять порядок  $p$  процесса  $AR$ . Так, для модели  $AR(1)$ :  $y_t = a_1 y_{t-1} + \xi_t$   $\rho(2)$  близко к нулю.

### Пример 7.3

За 50 мес. темпы прироста объема продукции  $K$  характеризовались авторегрессией вида  $y_t = 0,616y_{t-1} + \xi_t$   $R^2 = 0,662$  и  $F = 93,9$ . Автокорреляционная функция составила убывающие значения автокорреляции:

Лаг	1	2	3	4	5	6	7	...	13	14
$r_a$	0,813	0,747	0,617	0,474	0,408	0,457	0,353	...	-0,107	-0,051

Частная автокорреляционная функция, начиная с лага 2, достаточно близка к нулю:

Лаг	1	2	3	4	5
$\rho$	0,813	0,252	-0,134	-0,202	0,122

Для модели типа  $MA(q)$  порядок  $q$  определяется по поведению автокорреляционной функции: при  $k \geq p$   $r_a$  стремится к нулю. Для модели  $ARMA(p, q)$  автокорреляционная функция характеризуется убыванием, начинающимся с лага  $q$ , а частная автокорреляционная функция убывает, начиная с лага  $p$ . Так, для модели  $ARMA(1, 1)$  при  $a_1 > 0$   $ACF$  наблюдает экспоненциальное затухание с лага 1, а  $PACF$  — осциллирующее убывание с лага 1. При  $a_1 < 0$   $ACF$  для модели  $ARMA(1, 1)$  наблюдает осциллирующее убывание с лага 1, а  $PACF$  — экспоненциальное затухание с лага 1.

Выбор типа модели  $ARMA$  не ограничивается обычно исследованием автокорреляционных функций. С этой целью может использоваться, например, информационный критерий

Акайке<sup>1</sup>, рассмотрение которого не входит в задачу данного учебника.

#### 7.4.4. Модели *ARIMA*

Для получения стационарного ряда могут рассчитываться разности уровней временного ряда  $\Delta$  разного порядка  $d$ . Модель, в которой соединены нахождение последовательных разностей временного ряда порядка  $d$  и *ARMA*, — модель порядка  $(p, q)$ , получила название **авторегрессионной интегрированной модели скользящей средней** — *ARIMA* (*Autoregressive Integrated Moving Average*).

Модель *ARIMA* обладает тремя параметрами:  $p$  — порядок авторегрессии *AR*;  $d$  — порядок последовательных разностей уровней временных рядов, обеспечивающий стационарность ряда, и  $q$  — порядок скользящей средней *MA*.

В общем виде модель *ARIMA*  $(p, d, q)$  выражается формулой

$$\Delta^k y_t = a_1 \Delta^k y_{t-1} + \dots + a_p \Delta^k y_{t-p} + \xi_t - \theta_1 \xi_{t-1} - \dots - \theta_q \xi_{t-q}, \quad (7.24)$$

где  $\Delta^k y_t$  —  $k$ -я последовательная разность уровней  $y_t$ , т.е.  $\Delta^k y_t = y_t - y_{t-1} - \dots - y_{t-k}$ ;  $\xi_t, \xi_{t-1}, \dots, \xi_{t-q}$  — нормально распределенные случайные величины с нулевым математическим ожиданием и постоянной дисперсией.

Из модели (7.24) для  $\Delta^k y_t$  можно получить модель для исходного динамического ряда с помощью выражения

$$y_t = y_{t-1} + \Delta y_t. \quad (7.25)$$

Так, если модель *ARIMA*  $(1, 1, 1)$  имеет вид  $\Delta^1 y_t = 0,2 \Delta^1 y_{t-1} + \xi_t - 0,1 \xi_{t-1}$ , то динамический ряд описывается моделью

$$y_t = y_{t-1} + 0,2 y_{t-1} - 0,2 y_{t-2} + \xi_t - 0,1 \xi_{t-1},$$

так как  $\Delta^1 y_t = y_t - y_{t-1}$ ,  $\Delta^1 y_t = y_{t-1} - y_{t-2}$ .

<sup>1</sup> См., например, Носко В. П. Эконометрика. Элементарные методы и введение в регрессионный анализ временных рядов. М. : ИЭПП, 2004. С. 251.

Модель *ARIMA* практически пригодна для большинства временных рядов. При  $\rho = 1, 2, \dots, k$ ,  $d = 0$  и  $q = 0$  модель *ARIMA* превращается в процесс *AR*

$$y_t = a + b_1 y_{t-1} + b_2 y_{t-2} + \dots + b_k y_{t-k} + \xi_t.$$

Если  $\rho = 0$ ,  $d = 0$  и  $q = 1, 2, \dots, k$ , то имеем модель *MA*

$$y_t = \mu + \xi_t - \theta_1 \xi_{t-1} - \dots - \theta_k \xi_{t-k}.$$

Наиболее распространены модели *ARIMA* с параметрами  $\rho$ ,  $d$  и  $q$ , не превышающими двух. Современные компьютерные программы предлагают разные варианты оценивания параметров модели *ARIMA*, среди которых преобладает оценка методом максимального правдоподобия. Такой подход можно видеть при реализации модели *ARIMA* в системе *SPSS*<sup>1</sup>.

## Контрольные вопросы и задания

1. Что такое модель с распределенными лагами?
2. Как интерпретируются параметры модели с распределенными лагами?
3. В каких случаях оценка параметров модели с распределенными лагами может быть дана методом наименьших квадратов?
4. Объясните, что такое структура лага и как она используется при построении модели с распределенными лагами.
5. Опишите методику построения модели с использованием лагов Алмон.
6. В чем состоит суть преобразования Койка?
7. Что такое модели авторегрессии?
8. Как интерпретируются параметры модели авторегрессии?
9. Для чего используются инструментальные переменные?
10. Что такое авторегрессионные процессы?
11. Как строятся модели *ARMA*?
12. В чем состоят особенности модели *ARIMA*?

---

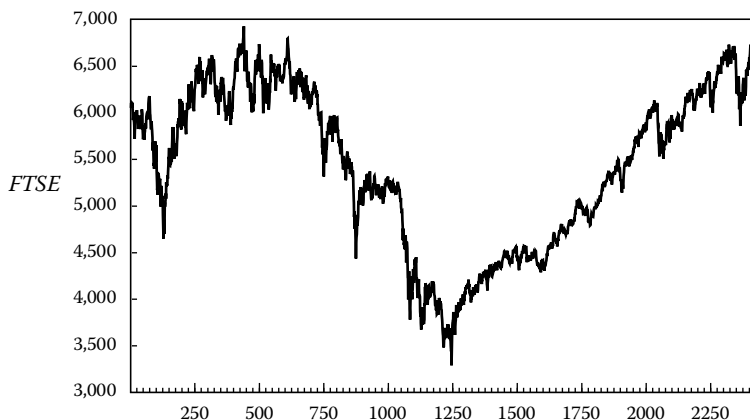
<sup>1</sup> Прогнозирование с помощью моделей *ARIMA* см. в кн. : Дубрава Т. А. Статистические методы прогнозирования. М. : ЮНИТИ, 2003. С. 178—184.

### 8.1. Стационарный ряд

При анализе пространственных данных, как правило, предпринимаются попытки объяснения изменений значений одного фактора в зависимости от изменений значений каких-то других факторов. При анализе временных рядов встречаются модели, базирующиеся на другом подходе. Такие модели называются одномерными временными моделями. В моделях этого класса предпринимается попытка смоделировать и спрогнозировать значения некоторых временных параметров, опираясь исключительно на информацию о прошлых значениях исследуемого параметра.

Обычно при использовании подобных моделей не строятся какие-то теоретические конструкции, объясняющие изменения исследуемого ряда. Априори полагается, что используемые данные уже содержат в себе такие модели.

Одна из причин использования временных моделей заключается в том, что они могут быть особенно полезны, когда невозможно найти объясняющие факторы, измеряемые с той же частотой, что и исследуемые переменные. Например, если исследуется дневные биржевые доходности, то в качестве объясняющих переменных могли бы выступать некоторые макроэкономические величины, которые измеряются не чаще, чем раз в месяц.



**Рис. 8.1.** Индекс *FTSE* со 2 апреля 1998 г. по 23 октября 2007 г.

При работе с временными моделями необходимо знать, являются ли рассматриваемые ряды стационарными. Этот вопрос важен потому, что стационарные и нестационарные ряды обладают различными статистическими характеристиками, поэтому должны оцениваться разными способами.

Итак, прежде всего необходимо определить, что такое стационарность временного ряда. Различают строго и слабо стационарные ряды.

**Строгая стационарность ряда** означает, что сдвиг по времени не меняет ни одну из функций плотности распределения ряда. То есть если  $F$  — функция распределения, то  $F_{y_{t_1}, y_{t_2}, \dots, y_{t_T}}(y_1, y_2, \dots, y_T) = F_{y_{t_1+k}, y_{t_2+k}, \dots, y_{t_T+k}}(y_1, y_2, \dots, y_T)$ .

У строго стационарных процессов есть полезные свойства. Так, например, если процесс является строго стационарным, то его математическое ожидание и дисперсия постоянны в любой момент времени.

**Слабо стационарным** называют такой процесс, у которого математическое ожидание и дисперсия существуют вне зависимости от времени и, кроме того, автоковариационная функция зависит только от разности значений  $(t_1 - t_2)$ . То есть  $y_t$  будет слабо стационарным, если для всех  $t = 1, 2, \dots, \infty$

$$E(y_t) = \mu;$$

$$E(y_t - \mu)(y_t - \mu) = \sigma^2,$$

где  $\sigma^2 < \infty$ ;

$$E(y_{t_1} - \mu)(y_{t_2} - \mu) = y_{t_2 - t_1} \text{ для } \forall t_1, t_2.$$

Последнее уравнение — автоковариационная функция. В общем виде она записывается как

$$E(y_t - E(y_t))(y_{t-s} - E(y_{t-s})) = y_s, \quad (8.1)$$

где  $s = 0, 1, 2, \dots$

При  $s = 0$  автоковариационная функция становится просто дисперсией ряда.

Поскольку автоковариационная функция показывает, как  $y$  зависит от предыдущих значений, можно заметить, что ковариация будет одинаковой для стационарного ряда, так как функция зависит только от того, насколько далеко моменты времени находятся друг от друга. Следовательно, при условии стационарности ковариация для  $y_t$  и  $y_{t-1}$  будет такой же, как, например, между  $y_{t-5}$  и  $y_{t-6}$ .

Часто используют автокорреляционную функцию (нормированную на величину дисперсии автоковариацией):

$$\tau_s = \frac{\gamma_s}{\gamma_0}, \text{ где } s = 0, 1, 2, \dots \quad (8.2)$$

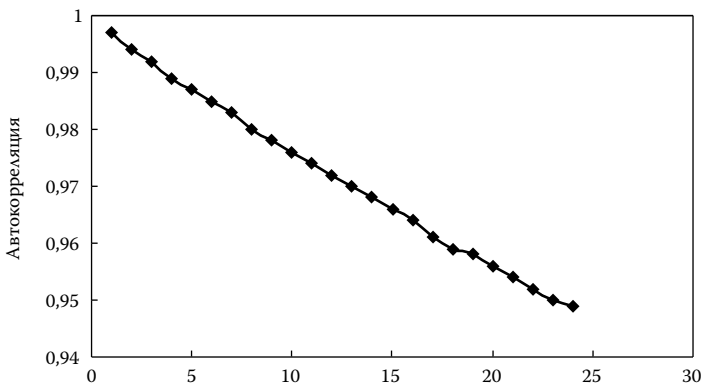


Рис. 8.2. Автокорреляционная функция индекса *FTSE*

Значения автокорреляционной функции лежат в интервале  $[-1, 1]$ . Если построить значения автокорреляционной функции для  $s = 0, 1, 2, \dots$ , то можно получить коррелограмму (график автокорреляционной функции).

Рассмотрим простейшие временные модели.

## 8.2. Базовые модели временных рядов

### Белый шум

Белым шумом называется процесс, имеющий постоянное математическое ожидание, постоянную дисперсию и нулевую, для всех, кроме нулевого лага, автоковариационную функцию. То есть  $y_t$  — белый шум, если

$$E(y_t) = \mu;$$

$$\text{var}(y_t) = \sigma^2;$$

$$\gamma_{t-s} = 0, \text{ если } t \neq s.$$

Белым шум, по определению, — слабо стационарный процесс. В самом деле, нулевая для всех, кроме нулевого лага, автоковариационная функция означает, что наблюдения не коррелированы между собой. График белого шума, смоделированно в пакете *Eviews*, изображен на рис. 8.3.

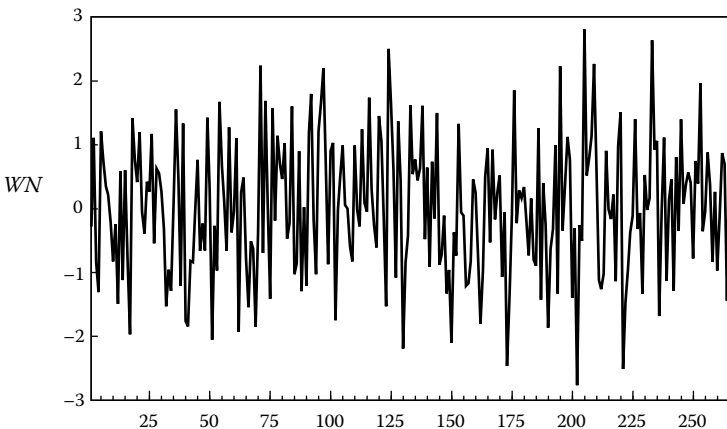


Рис. 8.3. Белый шум (модельный пример)



Если к условиям стандартного белого шума добавить условие нулевого математического ожидания, получится белый шум с нулевым средним.

Если выполняется предположение о нормальном распределении процесса  $y_t$ , то процесс будет строго стационарным. Кроме того, коэффициенты автокорреляции также будут определены согласно нормальному закону

$$\tau_s \sim N\left(0, \frac{1}{T}\right)$$

где  $\tau_s$  — коэффициент автокорреляции;  $T$  — размер выборки.

### Случайное блуждание

Случайным блужданием называется процесс, описываемый уравнением

$$X_t = X_{t-1} + \varepsilon_t, \quad (8.3)$$

где  $\varepsilon_t$  — белый шум.

Это просто авторегрессия первого порядка, с единичным коэффициентом. Модель случайного блуждания можно представить как

$$X_t = X_{t-1} + \varepsilon_t = X_{t-2} + \varepsilon_{t-1} + \varepsilon_t = X_{t-3} + \varepsilon_{t-2} + \varepsilon_{t-1} + \varepsilon_t = X_0 + \sum_{s=0}^t \varepsilon_{t-s}.$$

Таким образом, модель случайного блуждания выражается через начальное значение и сумму белых шумов (рис. 8.4). Математическое ожидание такого процесса равняется начальному значению:

$$E\{X_t\} = E\left\{X_0 + \sum_{s=0}^t \varepsilon_{t-s}\right\} = E\{X_0\} + E\left\{\sum_{s=0}^t \varepsilon_{t-s}\right\} = X_0 + 0 = X_0.$$

Дисперсия случайного блуждания при условии некоррелированности  $X_0$  и  $\varepsilon_t$  имеет вид

$$V(X_t) = E\left\{\left(\sum_{s=0}^t \varepsilon_{t-s}\right)^2\right\} = t \cdot \sigma^2.$$

Поскольку, как очевидно из уравнения, дисперсия зависит от времени, следовательно, условие постоянства дисперсии не выполняется и случайное блуждание не является стационарным процессом.



Рис. 8.4. Случайное блуждание (модельный пример)

Для того чтобы привести случайное блуждание к стационарному ряду, нужно просто взять первую разность  $X_t - X_{t-1} = \varepsilon_t$ . То есть  $\Delta X_t = \varepsilon_t$ , при этом  $\varepsilon_t$  — белый шум, слабо стационарный процесс. Прием взятия разностей очень часто встречается в эконометрике, на нем основывается один из принципов модели ARIMA, которая будет рассматриваться позже.

### Модель скользящей средней

Модель скользящей средней — одна из простейших моделей временных рядов. В модели в качестве объясняющих переменных выступает комбинация белых шумов (рис. 8.5). То есть ряд  $y_t$  описывается процессом MA ( $\tau$ ), если

$$y_t = \mu + u_t + \alpha_1 u_{t-1} + \alpha_2 u_{t-2} + \dots + \alpha_q u_{t-q} = \mu + u_t + \sum_{i=1}^q \alpha_i u_{t-i}, \quad (8.4)$$

где  $u_t$  (при  $t = 1, 2, 3, \dots$ ) — белый шум, с нулевым математическим ожиданием и дисперсией  $\sigma^2$ .

Уравнение может быть переписано с помощью оператора сдвига. В самом деле, если  $L^i y_t = y_{t-i}$ , то уравнение (8.4) можно представить как  $y_t = \mu + \sum_{i=1}^q \alpha_i L^i u_t + u_t = \mu + \alpha(L)u_t$ ,

где  $\alpha(L) = 1 + \alpha_1 L + \alpha_2 L^2 + \dots + \alpha_q L^q$ .

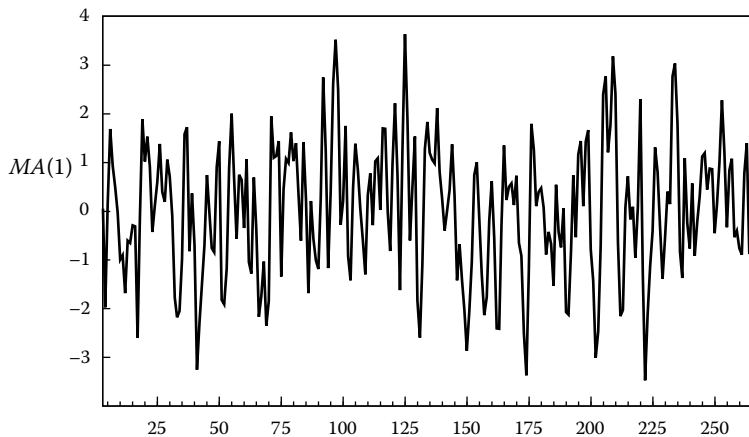


Рис. 8.5.  $MA(1) y_t = 0,8u_{t-1} + u_t$  (модельный пример)

Найдем характеристики процесса скользящей средней:

$$\begin{aligned}
 E(y_t) &= E(\mu + \alpha(L)u_t) = E(\mu) + E[\alpha(L)u_t] = \mu + 0 = \mu; \\
 V(y_t) &= V[\mu + \alpha(L)u_t] = V(\mu) + V[\alpha(L)u_t] = 0 + (1 + \alpha_1^2 + \alpha_2^2 + \dots + \alpha_q^2)\sigma^2 = \\
 &= (1 + \alpha_1^2 + \alpha_2^2 + \dots + \alpha_q^2)\sigma^2; \\
 \text{cov}(y_t, y_{t-s}) &= (\alpha_s + \alpha_{s+1}\alpha_1 + \alpha_{s+2}\alpha_2 + \dots + \alpha_q\alpha_{q-s})\sigma^2, \text{ при } s = 1, 2, \dots, q; \\
 \text{cov}(y_t, y_{t-s}) &= 0, \text{ при } s > q.
 \end{aligned}$$

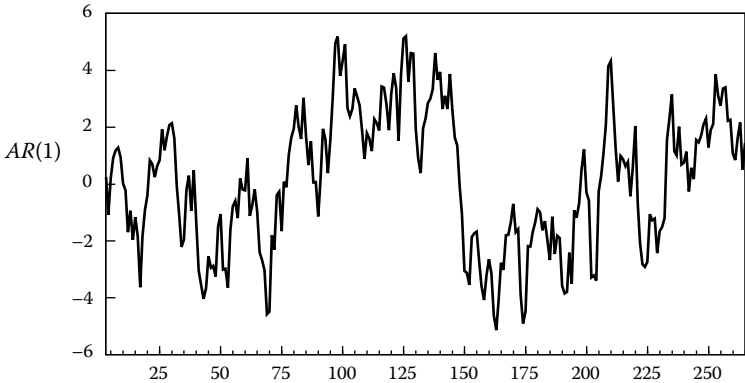
Таким образом, процесс скользящей средней имеет постоянные математическое ожидание и дисперсию и ненулевую до  $\tau$ -го лага автоковариацию.

### Авторегрессионная модель

В основе авторегрессионных моделей (рис. 8.6) лежит предположение о том, что последующие значения рассматриваемой переменной зависят только от ее значений в предыдущие периоды и ошибки. Таким образом, авторегрессионная модель порядка  $\tau$  —  $AR(\tau)$  записывается следующим образом:

$$y_t = \mu + \beta_1 y_{t-1} + \beta_2 y_{t-2} + \dots + \beta_p y_{t-p} + u_t = \mu + \sum_{i=1}^{\tau} \beta_i y_{t-i} + u_t, \quad (8.5)$$

где  $u_t$  — белый шум,  $y_{t-p}$  —  $(t-p)$ -й лаг рассматриваемой переменной.



**Рис. 8.6.**  $\mathcal{AR}(1) y_t = 0,8u_{t-1} + u_t$  (модельный пример)

Кроме того, уравнение авторегрессионной модели может быть записано с использованием оператора сдвига:

$$y_t = \mu + \sum_{i=1}^p \beta_i L^i y_t + u_t,$$

или

$$\beta(L)y_t = \mu + u_t,$$

где  $\beta(L) = (1 - \beta_1 L - \beta_2 L^2 - \dots - \beta_p L^p)$ .

При работе с процессом  $AR$  необходимо первым делом проверить его стационарность. Дело в том, что в моделях, где коэффициенты не удовлетворяют условиям стационарности, будет наблюдаться незатухающее влияние предыдущих значений ошибок на текущее значение параметра, что усложняет процесс оценивания подобной модели.

Для того чтобы вывести условие стационарности для авторегрессионных процессов, воспользуемся следующей формой записи модели  $AR$ :

$$\beta(L)y_t = \mu + u_t.$$

Пусть  $\mu = 0$ , тогда

$$\beta(L)y_t = u_t.$$

Процесс  $y_t$  будет стационарным, если его можно представить в виде

$$y_t = \beta(L)^{-1}u_t.$$

На самом деле,  $\beta(L)^{-1}u_t$  может быть представлено в виде  $MA(\infty)$  процесса  $\delta_1 u_{t-1} + \delta_2 u_{t-2} + \delta_3 u_{t-3} + \dots + u_t$ . Причем, если авторегрессионный процесс стационарен, то коэффициенты  $MA(\infty)$  будут уменьшаться для последующих лагов, а если процесс не стационарен, то коэффициенты  $MA(\infty)$  не будут сходиться к нулю по мере увеличения порядка лага.

Таким образом, условие стационарности  $AR(q)$  процесса заключается в том, чтобы корни характеристического уравнения

$$1 - \beta_1 z - \beta_2 z^2 - \dots - \beta_p z^p = 0$$

были по модулю больше единицы.

Рассмотрим условие стационарности авторегрессионного процесса на примере модели случайного блуждания

$$y_t = y_{t-1} + u_t.$$

Первым делом перепишем модель случайного блуждания с помощью оператора сдвига:

$$y_t = Ly_t + u_t;$$

$$y_t - Ly_t = u_t;$$

$$y_t(1 - L) = u_t.$$

Характеристическим уравнением для процесса случайного блуждания будет уравнение вида

$$1 - z = 0.$$

Единственный корень характеристического уравнения  $z = 1$  не является по модулю большим единицы. Следовательно, процесс случайного блуждания не является стационарным процессом.

### 8.3. Теорема декомпозиции Вольда

В 1938 г. Вольдом было доказано, что недетерминированный стационарный в широком смысле процесс может быть

представлен как сумма двух процессов: детерминированного и стохастического процесса  $MA(\infty)$ .

Применительно к процессам AR это означает, что любой стационарный авторегрессионный процесс порядка  $\tau$  может быть представлен с помощью модели скользящего среднего бесконечного порядка.

Интересно, что процессы скользящего среднего обладают похожим свойством. Если корни характеристического уравнения скользящего среднего  $\alpha(z) = 0$  по модулю больше единицы, то выполняется так называемое условие обратимости. Условие обратимости математически представляет то же самое, что условие стационарности авторегрессионных процессов, однако используется в других целях. Если выполняется условие обратимости, то конечный процесс  $MA(q)$  может быть представлен в виде бесконечного процесса  $AR(\infty)$ .

Таким образом, в процессах скользящего среднего и авторегрессионных процессах есть что-то общее. Однако есть принципиальное отличие. Процесс  $MA(\tau)$  всегда стационарен, условие обратимости просто обеспечивает его некоторым дополнительным полезным свойством. Для процесса  $AR(q)$  условие более жесткое: либо он стационарен и, следовательно, может быть представлен в виде скользящего среднего, либо он не стационарен.

Рассмотрим авторегрессионный процесс порядка  $p$ :

$$y_t = \varphi(L)u_t,$$

где  $\varphi(L) = \beta(L)^{-1} = (1 - \beta_1 L - \beta_2 L^2 - \dots - \beta_p L^p)^{-1}$ .

Математическое ожидание процесса  $y_t$  равно

$$E(y_t) = \frac{\mu}{1 - \beta_1 - \beta_2 - \dots - \beta_p}.$$

Автоковариационная и автокорреляционная функция может быть найдена путем решения системы уравнений Юла — Уолкера для процесса  $AR(p)$ :

$$\tau_1 = \beta_1 + \tau_1 \beta_2 + \dots + \tau_{p-1} \beta_p;$$

$$\tau_2 = \tau_1 \beta_1 + \beta_2 + \dots + \tau_{p-2} \beta_p;$$

...

$$\tau_p = \tau_{p-1}\beta_1 + \tau_{p-2}\beta_2 + \dots + \beta_p.$$

Для любого стационарного авторегрессионного процесса автокорреляционная функция будет уменьшаться по экспоненте. В дальнейшем мы будем использовать это свойство стационарных авторегрессионных функций.

## 8.4. Частная автокорреляционная функция

Как уже отмечалось ранее, частная автокорреляционная функция была введена с целью определения порядка авторегрессионного процесса. Дело в том, что в процессе скользящего среднего порядок модели достаточно просто определить, так как после него автокорреляционная функция резко стремится к нулю. Однако в авторегрессионном процессе все не так просто.

В этом случае на помощь приходит частная автокорреляционная функция, точно указывающая на порядок авторегрессионной модели. Частная автокорреляционная функция  $\tau_{kk}$  определяет корреляцию между текущим и произошедшим  $k$  периодов назад наблюдением после удаления косвенного влияния  $(y_{t-k+1}, y_{t-k+2}, \dots, y_{t-1})$ -х наблюдений. То есть  $\tau_{44}$  изменяет корреляцию между  $y_t$  и  $y_{t-4}$  без учета влияния  $y_{t-1}, y_{t-2}, y_{t-3}$ -го лага. Автокорреляционные функции для индекса *FTSE* представлены в табл. 8.1.

Очевидно, что для первого лага значения автокорреляционной и частной автокорреляционной функции совпадают:  $\tau_1 = \tau_{11}$ , поскольку отсутствует влияние промежуточных лагов, которое нивелировала бы частная автокорреляционная функция. Для второго лага частная автокорреляционная функция равна

$$\tau_{22} = \frac{\tau_2 - \tau_1^2}{1 - \tau_1^2},$$

где  $\tau_1$  и  $\tau_2$  — коэффициенты автокорреляции первого и второго порядка.

Таблица 8.1. График автокорреляционных функций по лагам для FTSE

Autocorrelation	Partial Correlation	Nº	AC	PAC	Q-STAT	Prob
		1	0,997	0,997	2401,0	0,000
		2	0,994	0,002	4789,0	0,000
		3	0,992	0,043	7165,4	0,000
		4	0,989	0,071	9532,1	0,000
		5	0,987	-0,011	11 889	0,000
		6	0,985	0,004	14 236	0,000
		7	0,983	0,025	16 573	0,000
		8	0,980	0,009	18 902	0,000
		9	0,978	-0,010	21 221	0,000
		10	0,976	-0,008	23 530	0,000
		11	0,974	0,021	25 831	0,000
		12	0,972	-0,013	28 122	0,000
		13	0,970	0,018	30 404	0,000
		14	0,968	-0,019	32 678	0,000
		15	0,966	0,041	34 943	0,000
		16	0,964	-0,027	37 199	0,000
		17	0,961	0,006	39 446	0,000
		18	0,959	0,006	41 658	0,000
		19	0,958	0,011	43 916	0,000
		20	0,956	0,022	46 139	0,000
		21	0,954	0,012	48 355	0,000
		22	0,952	0,024	50 564	0,000
		23	0,950	-0,020	52 766	0,000
		24	0,949	0,025	54 961	0,000

Для лагов высших порядков формула будет более сложной для вычислений. Однако заметим, что в случае авторегрессионного процесса определенного порядка есть более простой способ определения корреляция между лагами, порядком меньшими порядка авторегрессионного процесса.

Например, рассмотрим AR (4):

$$y_t = \beta_0 + \beta_1 y_{t-1} + \beta_2 y_{t-2} + \beta_3 y_{t-3} + \beta_4 y_{t-4} + u_t.$$

В модели представлено влияние на текущее значение параметра первых трех лагов рассматриваемого параметра. Следовательно, частная автокорреляционная функция принимает



ненулевые значения для лагов, меньших порядка модели, и нулевые значения для лагов, больших порядка модели.

В случае стационарного ряда значения выборочной частной автокорреляционной функции определяются как МНК-оценка последнего коэффициента в регрессии AR ( $p$ ).

Как говорилось ранее, при выполнении условия обратимости процесс скользящего среднего может быть представлен в виде авторегрессионного процесса. Следовательно, для скользящего среднего может быть использован аналогичный аппарат.

## 8.5. Модель ARMA

Комбинация авторегрессионного процесса и процесса скользящей средней образует еще один класс временных моделей — ARMA (авторегрессия — скользящая средняя). Данная модель базируется на предположении о том, что текущее значения исследуемого временного ряда зависит только от линейной комбинации предыдущих значений временного ряда и белого шума. Модель ARMA ( $p, q$ ) выглядит следующим образом:

$$y_t = \mu + \beta_1 y_{t-1} + \beta_2 y_{t-2} + \dots + \beta_p y_{t-p} + u_t + \alpha_1 u_{t-1} + \alpha_2 u_{t-2} + \dots + \alpha_q u_{t-q}, \quad (8.6)$$

где  $u_{t-1}, u_{t-2}, \dots, u_{t-p}$  и  $u_{t-1}, u_{t-2}, \dots, u_{t-q}$  — значения лагов ряда и белых шумов соответственно.

Процесс ARMA можно записать при помощи оператора сдвига

$$\beta(L)y_t = \mu + \alpha(L)u_t,$$

где  $\beta(L) = 1 - \beta_1 L - \beta_2 L^2 - \dots - \beta_p L^p$ ;  $\alpha(L) = 1 + \alpha_1 L + \alpha_2 L^2 + \dots + \alpha_q L^q$ .

Процесс ARMA является комбинацией процессов AR и MA. Следовательно, свойства и характеристики процесса ARMA также являются комбинациями свойств и характеристик используемых процессов AR и MA.

Возникает вопрос: является ли процесс ARMA ( $p, q$ ) стационарным? При условии нулевой средней процесс ARMA может быть представлен в виде  $\beta(L)y_t = \alpha(L)u_t$  при условии существования обратного оператора  $y_t = \beta(L)^{-1} \alpha(L)u_t$ . При этом обратный оператор может быть разложен в сумму элементарных

дробей, каждая из которых представима как бесконечно убывающая геометрическая прогрессия, т.е. в бесконечный операторный полином. При умножении на конечный полином снова получится бесконечный полином. Полученное выражение имеет смысл, только если корни характеристического уравнения  $\beta(z)=0$  по модулю меньше единицы. Но в таком случае полученное выражение есть не что иное, как разложение Вольда, и, следовательно, процесс стационарен. Таким образом, процесс ARMA будет стационарен, только если стационарным будет используемый процесс AR. Аналогично, процесс ARMA будет обратимым, если существует обратный оператор  $u_t = \alpha(L)^{-1}\beta(L)y_t$ . То есть процесс ARMA будет обратимым, если обратимым будет используемый процесс MA.

Найдем основные характеристики процесса ARMA. Очевидно, что математическое ожидание равняется нулю. Для того чтобы найти дисперсию процесса, представим процесс ARMA ( $p, q$ ) в виде бесконечного процесса MA ( $\infty$ ):

$$y_t = \frac{1+\alpha L}{1-\beta L} u_t = (1+\alpha L)(1+\beta L + \beta^2 L^2 + \dots)u_t = [1 + (\alpha + \beta)L + \alpha(\alpha + \beta)L^2 + \dots]u_t.$$

Тогда дисперсия  $y_t$  равна

$$V(y_t) = [1 + (\alpha + \beta)^2 + \alpha^2(\alpha + \beta)^2 + \dots]\sigma^2 = [1 + \frac{(\alpha + \beta)^2}{1 - \alpha^2}]\sigma^2 = \frac{1 + \beta^2 + 2\alpha\beta}{1 - \alpha^2}\sigma^2.$$



Автокорреляционные и частные автокорреляционные функции процессов AR, MA, ARMA обладают следующими свойствами.

- Автокорреляционная функция авторегрессионного процесса убывает по экспоненте, а количество ненулевых лагов частной автокорреляционной функции равно порядку авторегрессионного процесса.
- Частная авторегрессионная функция скользящей средней убывает по экспоненте, а число ненулевых лагов автокорреляционной функции равняется порядку процесса скользящей средней.
- Авторегрессионная функция и частная авторегрессионная функция процесса авторегрессия — скользящая средняя убывает по экспоненте.


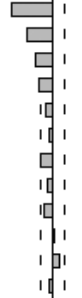
В табл. 8.2—8.12 приведены значения выборочных автокорреляционных и частных автокорреляционных функций стандартных процессов ARMA.

В табл. 8.2 представлена диаграмма автокорреляционных функций простейшего процесса скользящей средней. Как отмечалось выше, порядок скользящей средней может быть определен как число ненулевых лагов автокорреляционной функции. Существует значимая корреляция только с первым лагом. Следовательно, можно сделать вывод, что это скользящая средняя первого порядка.

**Таблица 8.2.** Выборочные автокорреляция и частная автокорреляция для  $MA(1)$ :  $y_t = 0,7u_{t-1} + u_t$

Autocorrelation	Partial Correlation	Nº	AC	PAC	Q-STAT	Prob
		1	0,481	0,481	61,667	0,000
		2	-0,003	-0,305	61,669	0,000
		3	-0,043	0,152	62,159	0,000
		4	-0,046	-0,135	62,735	0,000
		5	-0,036	0,063	63,087	0,000
		6	-0,069	-0,127	64,360	0,000
		7	-0,074	0,034	65,845	0,000
		8	0,001	0,015	65,845	0,000
		9	0,070	0,062	67,202	0,000
		10	0,121	0,074	71,280	0,000
		11	0,061	-0,055	72,236	0,000
		12	-0,008	0,026	72,253	0,000

**Таблица 8.3.** Выборочные автокорреляция и частная автокорреляция для  $MA(1)$  с отрицательным коэффициентом:  $y_t = -0,7u_{t-1} + u_t$




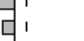
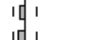
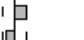
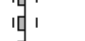
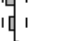



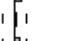

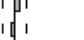
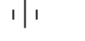
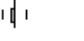



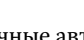
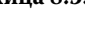
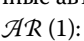
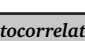
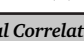
Autocorrelation	Partial Correlation	Nº	AC	PAC	Q-STAT	Prob
		1	-0,443	-0,443	52,110	0,000
		2	-0,016	-0,264	52,179	0,000
		3	-0,009	-0,180	52,199	0,000
		4	-0,013	-0,143	52,247	0,000
		5	0,024	-0,079	52,399	0,000
		6	0,002	-0,041	52,401	0,000
		7	-0,067	-0,115	53,628	0,000
		8	0,052	-0,058	54,359	0,000
		9	-0,043	-0,093	54,862	0,000
		10	0,077	0,011	56,489	0,000
		11	0,013	0,069	56,533	0,000
		12	-0,088	-0,034	58,694	0,000

В табл. 8.3 также представлена диаграмма автокорреляционных функций скользящей средней первого порядка. Коэффициент



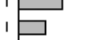



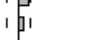



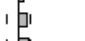



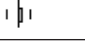









при первом лаге процесса по модулю равен коэффициенту предыдущего процесса, но противоположен по знаку. Этот факт выражается в отрицательной корреляции первого лага автокорреляционной функции. Частная автокорреляционная функция, как должно быть у процесса скользящей средней, убывает по экспоненте.

В табл. 8.4 представлены функции скользящей средней второго порядка. Как очевидно из диаграммы, значимая автокорреляция есть только с двумя первыми лагами. Более того, заметим, что поскольку оба коэффициента в модели положительны, корреляция с двумя первыми лагами также положительна.

**Таблица 8.4.** Выборочные автокорреляция и частная автокорреляция для  $MA(2)$ :  $y_t = 0,7u_{t-1} + 0,4u_{t-2} + u_t$

Autocorrelation	Partial Correlation	Nº	AC	PAC	Q-STAT	Prob
		1	0,602	0,602	96,296	0,000
		2	0,225	-0,215	109,80	0,000
		3	-0,040	-0,123	110,23	0,000
		4	-0,059	0,119	111,16	0,000
		5	-0,068	-0,093	112,42	0,000
		6	-0,080	-0,055	114,15	0,000
		7	-0,073	0,028	115,59	0,000
		8	0,008	0,076	115,60	0,000
		9	0,079	0,025	117,30	0,000
		10	0,122	0,049	121,40	0,000
		11	0,080	-0,030	123,18	0,000
		12	0,006	-0,045	123,19	0,000

**Таблица 8.5.** Выборочные автокорреляция и частная автокорреляция для  $AR(1)$ :  $y_t - 0,7y_{t-1} + u_t$

Autocorrelation	Partial Correlation	Nº	AC	PAC	Q-STAT	Prob
		1	0,696	0,696	128,91	0,000
		2	0,462	-0,043	186,00	0,000
		3	0,290	-0,031	208,48	0,000
		4	0,173	-0,011	216,52	0,000
		5	0,101	0,001	219,30	0,000
		6	0,049	-0,019	219,93	0,000
		7	0,027	0,015	220,13	0,000
		8	0,063	0,093	221,22	0,000
		9	0,087	0,017	223,30	0,000
		10	0,114	0,043	226,88	0,000
		11	0,080	-0,069	228,66	0,000
		12	0,032	-0,038	228,94	0,000

В табл. 8.5 можем увидеть выборочные автокорреляции для простейшего авторегрессионного процесса первого порядка. Как уже отмечалось, порядок авторегрессионной функции может быть найден по количеству ненулевых лагов частной автокорреляционной функции. В данном модельном примере это наглядно видно. Значимая корреляция есть только с первым лагом. При этом автокорреляционная функция, как и говорилось ранее, убывает по экспоненте.

В табл. 8.6 представлены диаграммы функций для авторегрессионного процесса первого порядка с отрицательным коэффициентом. Отрицательность коэффициента наглядно отражается в отрицательной значимой корреляции с первым лагом.

**Таблица 8.6.** Выборочные автокорреляция и частная автокорреляция для  $\mathcal{AR}(1)$  с отрицательным коэффициентом:

$$y_t = -0,7y_{t-1} + u_t$$












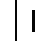












Autocorrelation	Partial Correlation	Nº	AC	PAC	Q-STAT	Prob
		1	-0,676	-0,676	121,51	0,000
		2	0,469	0,023	180,33	0,000
		3	-0,337	-0,022	210,80	0,000
		4	0,229	-0,022	224,96	0,000
		5	-0,172	-0,028	232,95	0,000
		6	0,129	0,003	237,48	0,000
		7	-0,133	-0,069	242,28	0,000
		8	0,102	-0,034	245,12	0,000
		9	-0,067	0,020	246,34	0,000
		10	0,057	0,020	247,24	0,000
		11	0,025	0,128	247,41	0,000
		12	-0,079	-0,043	249,16	0,000

Таблица 8.7. Выборочные автокорреляция и частная автокорреляция для  $AR(1)$ :  $y_t = -0,3y_{t-1} + u_t$

Autocorrelation	Partial Correlation	Nº	AC	PAC	Q-STAT	Prob
		1	0,321	0,321	27,495	0,000
		2	0,082	-0,024	29,301	0,000
		3	-0,009	-0,032	29,324	0,000
		4	-0,037	-0,026	29,700	0,000
		5	-0,032	-0,011	29,984	0,000
		6	-0,056	-0,046	30,839	0,000
		7	-0,071	-0,044	32,204	0,000
		8	0,016	0,060	32,273	0,000
		9	0,047	0,028	32,877	0,000
		10	0,111	0,089	36,292	0,000
		11	0,054	-0,015	37,110	0,000
		12	-0,018	-0,043	37,204	0,000

В табл. 8.7 представлены автокорреляционные диаграммы авторегрессионной модели первого порядка с коэффициентом 0,3. Порядок модели наглядно можно определить из частной автокорреляционной функции. Заметим, что если сравнить коррелограмму этой модели с коррелограммой модели, представленной на таблице 8.5, то можно увидеть, что меньшее по модулю значение коэффициента отразилось в меньшей корреляции частной автокорреляционной функции с первым лагом.

В табл. 8.8 представлены диаграммы выборочных функций корреляций авторегрессионного процесса второго уровня. Оба коэффициента в модели положительны, поэтому существует значимая положительная корреляция с двумя первыми

лагами. Заметим, что автокорреляционная функция не убывает по экспоненте. Этот факт наталкивает на мысль о нестационарности ряда.

**Таблица 8.8.** Выборочные автокорреляция и частная автокорреляция для  $\mathcal{AR}(2)$ :  $y_t = 0,6y_{t-1} + 0,4y_{t-2} + u_t$

Autocorrelation	Partial Correlation	№	AC	PAC	Q-STAT	Prob
		1	0,945	0,945	237,47	0,000
		2	0,931	0,358	468,98	0,000
		3	0,900	-0,033	686,06	0,000
		4	0,877	-0,007	892,85	0,000
		5	0,853	0,025	1089,6	0,000
		6	0,832	0,013	1277,2	0,000
		7	0,808	-0,017	1455,2	0,000
		8	0,794	0,070	1627,4	0,000
		9	0,772	-0,016	1791,0	0,000
		10	0,755	-0,012	1947,9	0,000
		11	0,729	-0,071	2094,9	0,000
		12	0,702	-0,086	2231,7	0,000

**Таблица 8.9.** Выборочные автокорреляция и частная автокорреляция для  $\mathcal{AR}(2)$ :  $y_t = 0,6y_{t-1} - 0,4y_{t-2} + u_t$

Autocorrelation	Partial Correlation	№	AC	PAC	Q-STAT	Prob
		1	0,441	0,441	51,829	0,000
		2	-0,144	-0,420	57,345	0,000
		3	-0,277	-0,002	77,927	0,000
		4	-0,117	-0,004	81,589	0,000
		5	0,021	-0,031	81,712	0,000
		6	-0,003	-0,080	81,715	0,000
		7	-0,068	-0,039	82,974	0,000
		8	-0,024	0,039	83,126	0,000
		9	0,062	0,018	84,184	0,000
		10	0,129	0,085	88,772	0,000
		11	0,075	-0,017	90,317	0,000
		12	-0,009	0,026	90,340	0,000

В табл. 8.9 можно увидеть диаграммы функций для авторегрессионного процесса второго порядка. Коэффициенты в модели

имеют разные знаки, что отражается в разных знаках значимой корреляции первых двух лагов частной автокорреляционной функции. Кроме того, модуль автокорреляционной функции, как было описано в теории, убывает по экспоненте.

В табл. 8.10 представлены данные стандартного процесса ARMA (1, 1). Автокорреляционная функция и модуль значений частной автокорреляционной функции процесса убывает по экспоненте. Заметим, что к 12-му лагу обе корреляционные функции убывают почти до нулевого уровня.

**Таблица 8.10.** Выборочные автокорреляция и частная автокорреляция для ARMA (1,1):  $y_t = 0,6y_{t-1} + 0,6u_{t-1} + u_t$

Autocorrelation	Partial Correlation	Nº	AC	PAC	Q-STAT	Prob
		1	0,835	0,835	185,55	0,000
		2	0,559	-0,458	268,96	0,000
		3	0,356	0,241	302,87	0,000
		4	0,216	-0,166	315,46	0,000
		5	0,125	0,098	319,69	0,000
		6	0,066	-0,076	320,89	0,000
		7	0,048	0,128	321,52	0,000
		8	0,070	0,019	322,85	0,000
		9	0,101	0,031	325,67	0,000
		10	0,115	-0,020	329,30	0,000
		11	0,088	-0,085	331,47	0,000
		12	0,046	0,043	332,05	0,000



























**Таблица 8.11.** Выборочные автокорреляция и частная автокорреляция для процесса случайного блуждания:  $y_t = y_{t-1} + u_t$

Autocorrelation	Partial Correlation	N <sup>o</sup>	AC	PAC	Q-STAT	Prob
		1	0,975	0,975	252,65	0,000
		2	0,948	-0,025	492,88	0,000
		3	0,922	-0,030	720,55	0,000
		4	0,896	0,019	936,73	0,000
		5	0,873	0,027	1142,7	0,000
		6	0,850	-0,009	1338,7	0,000
		7	0,828	0,007	1525,4	0,000
		8	0,810	0,064	1704,7	0,000
		9	0,790	-0,048	1876,0	0,000
		10	0,770	-0,018	2039,1	0,000
		11	0,745	-0,091	2192,6	0,000
		12	0,719	-0,041	2335,9	0,000

В табл. 8.11 представлены выборочные автокорреляционная и частная автокорреляционная функции процесса случайного блуждания. Модель случайного блуждания — классический пример нестационарного ряда. Как можно увидеть, нестационарность ряда отражается в очень медленном убывании автокорреляционной функции. Интересно заметить, что, по сути, процесс случайного блуждания является авторегрессионным процессом с единичным коэффициентом. Частная автокорреляционная функция и в этом примере хорошо определяет порядок авторегрессионной функции.

В табл. 8.12 представлены корреляционные данные модельного примера белого шума. Ни у автокорреляционной функции, ни у частной автокорреляционной функции не наблюдается значимой корреляции ни с одним из лагов.

**Таблица 8.12.** Выборочные автокорреляция и частная автокорреляция для процесса белого шума:  $y_t = u_t$

Autocorrelation	Partial Correlation	Nº	AC	PAC	Q-STAT	Prob
		1	0,026	0,026	0,1880	0,000
		2	-0,011	-0,012	0,2203	0,000
		3	-0,021	-0,020	0,3362	0,000
		4	-0,030	-0,029	0,5748	0,000
		5	-0,008	-0,007	0,5917	0,000
		6	-0,033	-0,034	0,8972	0,000
		7	-0,072	-0,072	2,3182	0,000
		8	0,024	0,026	2,4761	0,000
		9	0,016	0,012	2,5497	0,000
		10	0,095	0,091	5,0685	0,000
		11	0,040	0,033	5,5093	0,000
		12	-0,048	-0,048	6,1559	0,000

Чтобы ответить на вопрос, является ли ряд стационарным, можно воспользоваться **тестом Дики — Фуллера**.

Стандартный тест Дики — Фуллера заключается в проверке гипотезы стационарности временного ряда.

В рамках теста рассматривается модель

$$\Delta y_t = \alpha y_{t-1} + x_t \delta + \varepsilon_t, \quad (8.7)$$

где  $\alpha = p - 1$ .

Нулевая и альтернативная гипотеза записываются в следующем виде:

$$H_0: \alpha = 0;$$

$$H_1: \alpha < 0.$$

Оценка происходит на базе статистики Стьюдента

$$t_\alpha = \frac{\hat{\alpha}}{se(\hat{\alpha})},$$

где  $\hat{\alpha}$  — оценка  $\alpha$ ;  $se(\hat{\alpha})$  — стандартное отклонение.

В 1979 г. Д. Дики и У. Фуллер [Dickey, Fuller] показали, что при условии наличия единичного корня рассматриваемая статистика не подчиняется распределению Стьюдента. Оказалось, что простой тест Дики — Фуллера применим только для моделей AR (1).

В *расширенном тесте Дики — Фуллера (ADF)* учитываются корреляция лагов высших порядков путем предположения о том, что  $y_t$  описывается моделью AR ( $p$ ) и добавлением разностей порядка  $p$  в правую часть уравнений

$$\Delta y_t = \alpha y_{t-1} + x_t \delta + \beta_1 \Delta y_{t-1} + \beta_2 \Delta y_{t-2} + \dots + \beta_p \Delta y_{t-p} + v_t. \quad (8.8)$$

Фуллер доказал, что асимптотическое распределение  $t$ -статистики для  $\alpha$  не зависит от количества лагов разностей, включенных в ADF-тест. Более того, в 1984 г. Дики доказал, что если  $y_t$  описывается моделью AR ( $p$ ), то ADF-тест асимптотически эффективен при присутствии MA компонент.

### Подход Бокса — Дженкинса

В 1976 г. Д. Бокс и Г. Дженкинс предложили систематический подход практической оценки моделей ARMA. Подход включает три части:

- 1) идентификацию;
- 2) оценивание;
- 3) диагностику построенной модели.

На этапе идентификации Бокс и Дженкинс предложили определять порядок модели. Для того чтобы найти максимально подходящую спецификацию модели, можно провести визуальный анализ графика, построить автокорреляционную и частную автокорреляционную функции, провести расширенный тест Дики — Фуллера.

Авторы предложили оценивать параметры идентифицированной на первом шаге модели с помощью метода наименьших квадратов и метода максимального правдоподобия.

На третьем этапе проверяется адекватность полученной выше модели. Существуют различные критерии проверки модели.

При помощи *теста Стьюдента* можно проверить статистическую значимость отличия от нуля коэффициентов модели. С помощью *F-частного* можно определить, оптимальная ли комбинация регрессоров используется в модели.

Кроме того, согласно ARMA ошибки модели должны быть белым шумом. Можно проверить это условие. Как известно, белый шум имеет нулевую автокорреляцию, следовательно, выборочная автокорреляционная функция остатков также должна равняться нулю.

Выборочную автокорреляционную функцию можно найти как по формуле

$$r_s = \frac{\sum_{t=s+1}^T \varepsilon_t \varepsilon_{t-s}}{\sum_{t=1}^T \varepsilon_t^2}, s = 1, 2, \dots \quad (8.9)$$

Проверить одновременные статистически значимые отличия от нуля значений первых коэффициентов автокорреляции можно при помощи теста, разработанного в 1970 г. Д. Боксом и Д. Пирсом:

$$Q = T \sum_{s=1}^m r_s^2, \quad (8.10)$$

где  $T$  — размер выборки;  $m$  — порядок максимально лага.

При нулевой гипотезе  $Q$  подчиняется распределению  $\chi^2$  ( $m - p - q$ ). Нулевая гипотеза отвергается, если полученное  $Q$  превышает критическое значение.

К сожалению, статистика Бокса — Пирса часто приводит к ошибочным результатам при работе с высокочастотными маленькими выборками. В 1978 г. Г. Льюнг и Д. Бокс предложили **модифицированный тест Бокса — Пирса**:

$$Q = T(T+2) \sum_{s=1}^m \frac{r_s^2}{T-s} \sim \chi^2.$$

### Информационные критерии при идентификации модели ARMA

Процесс идентификации модели ARMA не всегда проводится с использованием только визуального анализа графиков автокорреляционных и частных автокорреляционных функций. При использовании реальных данных картинка не всегда будут такими очевидными, как описывается в теории. Следовательно,

не всегда можно будет определить порядок модели только по автокорреляционным функциям.

Альтернативный подход заключается в использовании так называемых информационных критериев. При добавлении в модель новой переменной или дополнительно лага в значении информационного критерия одновременно будут наблюдаться два эффекта: уменьшение суммы квадратов остатков и увеличение величины штрафа за введение дополнительной переменной.

Цель заключается в том, чтобы получить модель с минимальными значениями информационных критериев. Наиболее популярными информационными критериями являются критерий Акайке, критерий Шварца и информационный критерий Ханнан — Куина.

**Критерий Акайке (AIC)** был разработан в 1974 г. и вычисляется по формуле

$$AIC = \ln(\hat{\sigma}^2) + \frac{2k}{T}. \quad (8.11)$$

**Критерий Шварца (1978)** выглядит следующим образом:

$$SBIC = \ln(\hat{\sigma}^2) + \frac{k}{T} \ln T. \quad (8.12)$$

**Информационный критерий Ханнан — Куина** имеет вид

$$HQIC = \ln(\hat{\sigma}^2) + \frac{2k}{T} \ln[\ln(T)], \quad (8.13)$$

где  $\hat{\sigma}^2$  — дисперсия остатков (сумма квадратов остатков, деленная на число наблюдений  $T$ );  $k = p + q + 1$  — число оцениваемых параметров;  $T$  — размер выборки.

Заметим, что в критерии Шварца величина штрафа за количество параметров превышает штраф в критерии Акайке, а штраф в критерии Ханнан — Куина есть нечто среднее между штрафами в критериях Акайке и Шварца.

Внимательный читатель уже, наверно, обратил внимание, что по своей идеологии информационные критерии похожи на скорректированный коэффициент детерминации  $R^2$ .

К сожалению, нельзя утверждать, что какой-то критерий лучше других позволяет определить порядок авторегрессионной

скользящей средней модели. Как правило, при практической работе приходится пользоваться совокупностью описанных информационных критериев.

## 8.6. Модель *ARIMA*

Модель *ARIMA* в определенной степени является расширенной версией модели *ARMA*. Символ *I* (*Integrated*) отвечает за порядок оператора последовательной разности. Дело в том, что далеко не все ряды являются стационарными, но некоторые из них могут быть приведены к стационарным путем взятия последовательной разности. Если, например, временной ряд  $y_t$  стал стационарным после взятия последовательной разности порядка  $s$  и для описания уже стационарного ряда может быть использована модель *ARMA* ( $p, q$ ), то процесс  $y_t$  называется интегрированным процессом авторегрессии и скользящей средней (*ARIMA* ( $p, s, q$ )).

При этом для оценки модели можно пользоваться несколько модифицированным подходом Бокса — Дженкинса. Непосредственно перед первым этапом необходимо будет привести исследуемый ряд к стационарному путем взятия последовательных разностей. Остальные этапы подхода не изменятся, после этого стационарный ряд нужно будет идентифицировать, оценить, диагностировать и использовать.

### Примеры взятия последовательной разности

Рассмотрим примеры простейших нестационарных временных рядов, которые могут быть приведены к стационарным путем взятия последовательной разности.

#### 1. Уравнение с трендом

Временной ряд с трендом имеет вид

$$y_t = \alpha + \beta t + u_t,$$

где  $\alpha + \beta t$  — временной тренд;  $u_t$  — белый шум.

Таким образом, временной тренд состоит из детерминированной составляющей линейного тренда и случайной

составляющей белого шума (рис. 8.7). Найдем математическое ожидание временного тренда:

$$E(y_t) = E(\alpha + \beta t + u_t) = \alpha + \beta E(t).$$

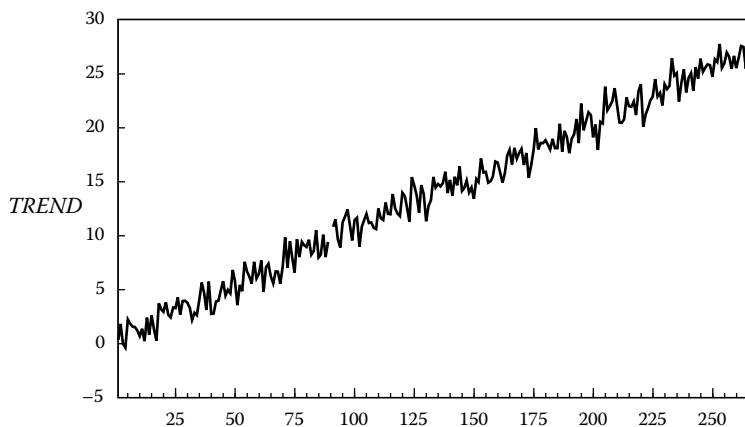


Рис. 8.7. Тренд (модельный пример)

Как очевидно, математическое ожидание зависит от времени, следовательно, ряд не является стационарным.

Для того чтобы привести ряд с линейный трендом к стационарному, нужно взять первую разность:

$$\Delta y_t = \beta + \varepsilon_t;$$

$$\varepsilon_t = u_t - u_{t-1}.$$

Возможен временной ряд с квадратичным временным трендом и временными трендами высших порядков. Для приведения их к стационарному ряду необходимо взятие последовательной разности тех же порядков. Так, например, чтобы привести к стационарному временной квадратичный тренд, нужно дважды взять последовательную разность:

$$y_t = \alpha + \beta t + \gamma t^2 + u_t;$$

$$\Delta^2 y_t = \Delta(\Delta y_t) = (y_t - y_{t-1}) - (y_{t-1} - y_{t-2});$$

$$\Delta^2 y_t = 2\gamma + \Delta^2 u_t,$$

где вторая последовательная разность  $\Delta^2 y_t$  — стационарный ряд.

## 2. Случайное блуждание

Модель случайного блуждания уже описывалась ранее. Уравнение модели имеет вид

$$y_t = y_{t-1} + u_t,$$

где  $u_t$  — белый шум.

Случайное блуждание — классический пример нестационарного ряда. Однако, чтобы привести его к стационарному, нужно просто взять первую разность:

$$z_t = y_t - y_{t-1} = y_{t-1} + u_t - y_{t-1} = u_t.$$

Первая разность случайного блуждания равняется белому шуму, следовательно, стационарна по определению.

## 3. Временной ряд с сезонностью

Сезонность часто встречается в статистических данных. В различных процессах сезонная составляющая может встречаться практически с любой частотой: месячная сезонная компонента, квартальная сезонная компонента, полугодовая сезонная компонента:

$$y_t = S(t) + u_t,$$

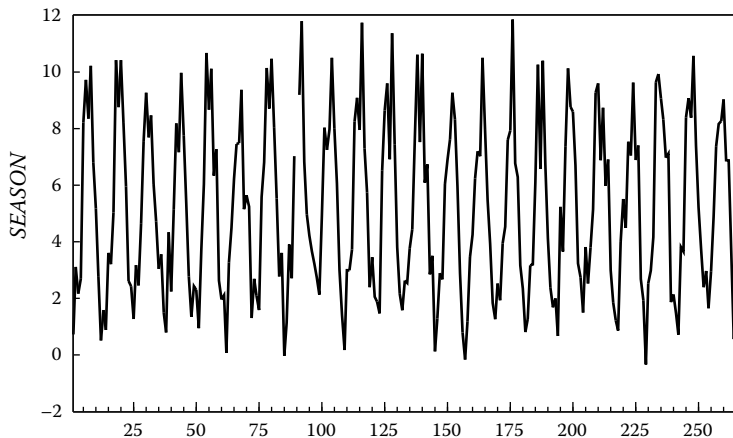
где  $S(t) = S(t+12)$ ,

или

$$y_t = S(t) + u_t,$$

где  $S(t) = S(t+12)$ .





**Рис. 8.8.** Сезонность (модельный пример)

Для того чтобы привести с сезонность к стационарному ряду, необходимо взять сезонную последовательную разность. Так, например, для квартальной сезонной компоненты

$$\Delta_4 y_t = y_t - y_{t-4} = S(t) + u_t - S(t-4) - u_{t-4} = u_t - u_{t-4}.$$

При этом сезонная последовательная разность  $\Delta_4 y_t$  будет стационарным временным рядом (рис. 8.8).

### Пример построения модели *ARIMA*

В качестве иллюстрирующего примера используем модель *ARIMA* для ряда значений индекса *DAX* со 2 апреля 1998 г. по 23 октября 2007 г. (рис. 8.9).

Первым шагом необходимо определить, является ли рассматриваемый ряд стационарным. Визуально анализ графика ряда не говорит о стационарности ряда. Необходимо провести дополнительный анализ.

Вторым шагом должно быть построение диаграммы выборочных автокорреляционной и частной автокорреляционной функций.

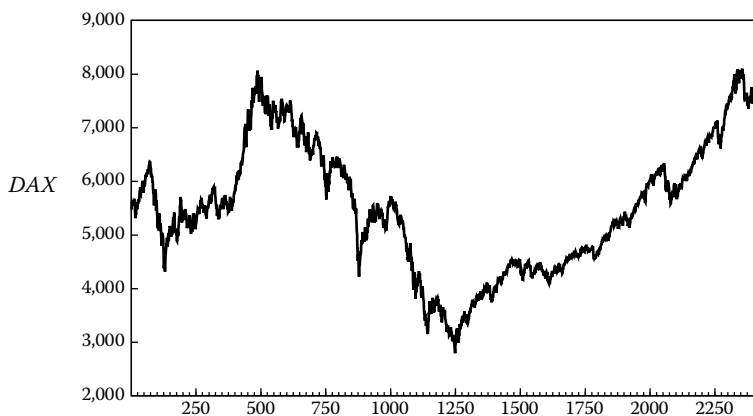


Рис. 8.9. Индекс *DAX* со 2 апреля 1998 г. по 23 октября 2007 г.

Таблица 8.13. Диаграмма автокорреляционных функций по лагам для *DAX*

Autocorrelation	Partial Correlation	№	AC	PAC	Q-STAT	Prob
		1	0,998	0,998	2403,2	0,000
		2	0,995	0,004	4795,6	0,000
		3	0,993	0,003	7177,5	0,000
		4	0,990	0,028	9549,5	0,000
		5	0,988	-0,020	11 911	0,000
		6	0,986	0,002	14 263	0,000
		7	0,984	0,028	16 604	0,000
		8	0,981	0,025	18 937	0,000
		9	0,979	-0,018	21 261	0,000
		10	0,977	0,029	23 576	0,000
		11	0,975	0,017	25 883	0,000
		12	0,973	-0,012	28 181	0,000
		13	0,971	0,013	30 471	0,000
		14	0,969	0,005	32 753	0,000
		15	0,967	-0,032	35 027	0,000
		16	0,965	-0,020	37 291	0,000
		17	0,963	-0,007	39 546	0,000
		18	0,961	-0,011	41 791	0,000
		19	0,959	0,001	44 027	0,000
		20	0,956	-0,002	46 253	0,000
		21	0,954	-0,024	48 469	0,000
		22	0,952	0,005	50 676	0,000
		23	0,949	0,021	52 873	0,000
		24	0,947	0,001	55 061	0,000

В случае стационарного ряда коррелограмма должна достаточно быстро убывать по мере увеличения лагов. Если же она убывает достаточно медленно, есть основания считать, что ряд нестационарен. Эта же логика применима и для частной автокорреляционной функции.

Как очевидно из коррелограммы, значения выборочной автокорреляционной функции практически не убывает по мере роста лагов, что дает еще большие основания предположить, что ряд нестационарен (рис. 8.10). Более того, коррелограмма ряда похожа на коррелограмму модели случайных блужданий.

Для того чтобы убедиться в нестационарности ряда, проведем **расширенный тест Дики — Фуллера (ADF)** (табл. 8.14).

Исходя из специфики работы биржи (пятидневная рабочая неделя), в модель включается пять лагов.

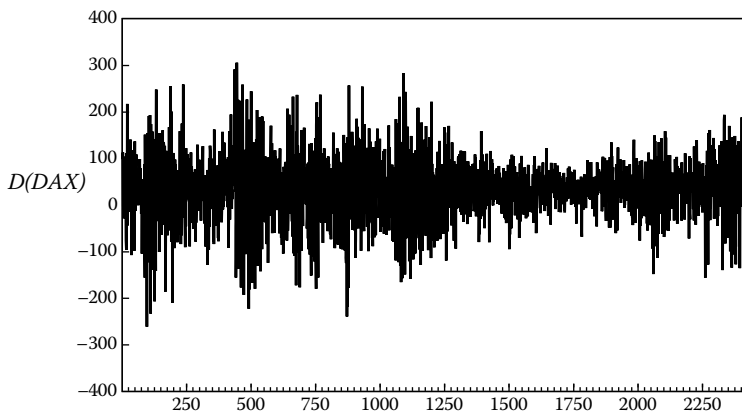
**Таблица 8.14.** Расширенный тест Дики — Фуллера для DAX

			t-Statistic	Prob*
Augmented Dickey — Fuller test statistic			0,527600	0,8298
Test critical values	1% level		-2,565923	
	5% level		-1,940955	
	10% level		-1,616611	
Variable	Coefficient	Std. Error	t-Statistic	Prob
DAX (-1)	0,000152	0,000289	0,527600	0,5978
D (DAX (-1))	0,004954	0,020418	0,242647	0,8083
D (DAX (-2))	-0,005194	0,020403	-0,254585	0,7991
D (DAX (-3))	-0,035070	0,020391	-1,719878	0,0856
D (DAX (-4))	0,036710	0,020398	1,799649	0,0720
D (DAX (-5))	-0,009063	0,020408	-0,444063	0,6570
R-squared	0,002591	Mean dependent var	1,084838	
Adjusted R-squared	0,000513	S. D. dependent var	75,84971	
S. E. of regression	75,83027	Akaike info criterion	11,49736	
Sum squared resid	13 800 550	Schwarz criterion	11,51179	
Log likelihood	-13 825,33	Hannan — Quinn criterion	11,50261	
Durbin — Watson stat	2,000697			

Null Hypothesis: DAX has a unit root.

Exogenous: None.

Lag Lengs: 5 (Fixed).



**Рис. 8.10.** График однодневных приращений индекса  $DAX$  со 2 апреля 1998 г. по 23 октября 2007 г.

Как можно увидеть из *статистики Дарбина — Уотсона*, проблемы с автокорреляцией нет. Часть коэффициентов при лагах оказались незначимыми. Тем не менее, есть основания считать, что единичный корень есть.

Построим *расширенный тест Дики — Фуллера* для первой разности ряда. Заметим, что первые разности индексов  $DAX$  также имеют экономический смысл. Этот ряд показывает однодневные абсолютные приращения индекса  $DAX$  (табл. 8.15).

**Таблица 8.15.** Расширенный тест Дики — Фуллера для однодневных приращений индекса  $DAX$

			<i>t</i> -Statistic	Prob*
<i>Augmented Dickey — Fuller test statistic</i>			-20,95237	0,0000
<i>Test critical values</i>	1% level		-2,565924	
	5% level		-1,940955	
	10% level		-1,616611	
<i>Variable</i>	<i>Coefficient</i>	<i>Std. Error</i>	<i>t</i> -Statistic	<i>Prob</i>
$DAX (-1)$	-1,052751	0,050245	-20,95237	0,0000
$D (DAX (-1), 2)$	0,057400	0,045858	1,251669	0,2108
$D (DAX (-2), 2)$	0,054322	0,041098	1,321766	0,1864
$D (DAX (-3), 2)$	0,018049	0,035259	0,511905	0,6088

Окончание табл. 8.15

$D(DAX(-4), 2)$	0,054415	0,028755	1,892346	0,0586
$D(DAX(-5), 2)$	0,045590	0,020385	2,236410	0,0254
<i>R-squared</i>	0,500551	<i>Mean dependent var</i>	0,002166	
<i>Adjusted R-squared</i>	0,499510	<i>S. D. dependent var</i>	107,0945	
<i>S. E. of regression</i>	75,76430	<i>Akaike info criterion</i>	11,49562	
<i>Sum squared resid</i>	13 770 809	<i>Schwarz criterion</i>	11,51006	
<i>Log likelihood</i>	-13 817,49	<i>Hannan — Quinn criterion</i>	11,50087	
<i>Durbin — Watson stat</i>	2,002092			

Null Hypothesis:  $D(DAX)$  has a unit root.

























Exogenous: None.

Lag Lengs: 5 (Fixed).

Ряд первых разностей индекса  $DAX$  стационарен. Значит, ряд индексов  $DAX$  имеет первый порядок интегрируемости.

Построим коррелограмму первых разностей ряда  $DAX$ .

**Таблица 8.16.** Диаграмма автокорреляционных функций по лагам для первых разностей  $DAX$

<i>Autocorrelation</i>	<i>Partial Correlation</i>	<i>Nº</i>	<i>AC</i>	<i>PAC</i>	<i>Q-STAT</i>	<i>Prob</i>
		1	0,004	0,004	0,0347	0,852
		2	-0,006	-0,006	0,1140	0,945
		3	-0,035	-0,035	3,0371	0,386
		4	0,037	0,037	6,2886	0,179
		5	-0,008	-0,009	6,4550	0,264
		6	0,045	-0,046	11,334	0,079
		7	-0,030	-0,027	13,473	0,061
		8	0,024	0,022	14,921	0,061
		9	-0,009	-0,012	15,122	0,088
		10	-0,008	-0,007	15,288	0,122
		11	0,010	0,013	15,551	0,159
		12	-0,023	-0,029	16,872	0,154

Как очевидно из коррелограммы, не наблюдается статистически значимой корреляции ни с одним из лагов. Более того, диаграмма автокорреляционной функции очень похожа на модельный пример коррелограммы белого шума.

Видимо, наиболее адекватной моделью значений ряда *DAX* является ARIMA (0, 1, 0).

В этом случае не имеет смысла пользоваться информационными критериями, поскольку добавление к процессу даже одного параметра качественным образом меняет суть модели.

Адекватность полученной модели можно проверить по остаткам. Ведь, как мы уже говорили раньше, остатки модели являются белым шумом. Следовательно, остатки регрессии тоже должны быть белым шумом. То есть остатки должны иметь нулевую автокорреляцию.

В табл. 8.16 в пятой графе представлены значения *статистики Льюнга — Бокса*. Напомним, что это тест, проверяющий гипотезу о равенстве нулю сразу нескольких значений автокорреляционной функции остатков. Как очевидно из таблицы, на 5%-ном уровне значимости нулевая гипотеза отвергается для всех 12 лагов. В остатках определенно присутствует автокорреляция. Для анализа ряда мы должны пользоваться более мощным аппаратом.

## 8.7. Коинтеграция

Как отмечалось ранее, одно из первых действий при построении регрессий временных рядов — проверка рядов на стационарность. Это связано с тем, что в случае регрессии, построенной по нестационарным рядам, ошибка также окажется нестационарным рядом. Следовательно, модель не будет удовлетворять условиям классической линейной регрессии и оценки параметров, полученные с помощью метода наименьших квадратов, могут оказаться несостоятельными.

Тем не менее, возможны случаи, когда МНК-оценки параметров регрессий, построенных на основе нестационарных временных рядов, будут состоятельными. Дело в том, что иногда нестационарность разных временных рядов связана с общим стохастическим трендом. Например, логично предположить, что цена и дивиденды от акции будут связаны некоторым

долгосрочным соотношением. В этом случае, даже если цена акции и дивиденды за некоторый период времени являются нестационарными временными рядами, существуют способы построения регрессии для них.

Р. Энгл и К. Грэнжер [Engle, Granger] предложили подход для построения регрессий по нестационарным рядам. Рассмотрим этот подход на следующем примере.

Пусть рассматривается некоторая регрессия по нестационарным временным рядам  $y_t$ ,  $x_t$  одного порядка интегрирования  $d$ :

$$y_t = \alpha + \beta x_t + \varepsilon_t.$$

Остатки модели будут равны

$$\varepsilon_t = y_t - \alpha - \beta x_t.$$

В случае если нестационарные временные ряды связаны долгосрочным соотношением, остатки  $\varepsilon_t$  будут стационарными. Временные ряды  $y_t$  и  $x_t$  в этом случае называются **коинтегрированными**. В случае коинтеграции оценки, полученные методом наименьших квадратов, будут состоятельными.

В общем случае ряды  $y_t$ ,  $x_{1t}$ ,  $x_{2t}$ , ...,  $x_{nt}$  будут коинтегрированными, если их линейная комбинация будет процессом нулевого порядка интегрирования. То есть ряды будут коинтегрированными, если ряд  $\varepsilon_t$  будет стационарным, где

$$\varepsilon_t = y_t - \beta_1 x_{1t} - \beta_2 x_{2t} - \dots - \beta_n x_{nt}.$$

Подход коинтегрированных рядов открыл возможность работать с нестационарными временными рядами. Однако при использовании подхода в первоначальном виде возникают некоторые сложности. Дело в том, что сама по себе коинтеграция не описывает долгосрочные соотношения и могут быть ситуации, когда коинтеграция возможна там, где долгосрочного соотношения нет.

Рассмотрим пример. Пусть исследуется два ряда  $y_t$  и  $x_t$  первого порядка интегрирования. Для них можно построить регрессию

$$y_t - y_{t-1} = \beta(x_t - x_{t-1}) + \varepsilon_t, \quad (8.14)$$

или

$$\Delta y_t = \beta \Delta x_t + \varepsilon_t.$$

Ряды  $y_t$  и  $x_t$  — некоторые экономические данные. Экономические данные могут не изменяться со временем, тогда  $y_t = y_{t-1} = \dots = y_0$  и  $x_t = x_{t-1} = \dots = x_0$ . В этом случае в уравнении (8.14) все члены модели будут равны нулю. Следовательно, остатки  $\varepsilon_t$  будут стационарными, но найти какое-то долгосрочное соотношение между  $y_t$  и  $x_t$  не представляется возможным.

С целью избежать проблем, связанных с коинтеграцией рядов, была предложена *модель корректировки ошибок* (*error correction model*). Для вышеописанного случая она будет выглядеть следующим образом:

$$\Delta y_t = \beta_1 \Delta x_t + \beta_2 (y_{t-1} - \lambda x_{t-1}) + \varepsilon_t.$$

Очевидно, что если  $y_t$  и  $x_t$  — коинтегрированные ряды, то линейная комбинация  $(y_{t-1} - \lambda x_{t-1})$  будет рядом нулевого порядка интегрирования.

В более общем случае, например, если рассматриваются три коинтегрированных ряда  $y_t$ ,  $x_t$  и  $z_t$ , модель корректировки ошибок будет иметь вид

$$\Delta y_t = \beta_1 \Delta x_t + \beta_2 \Delta z_t + \beta_3 (y_{t-1} - \lambda_1 x_{t-1} - \lambda_2 z_{t-1}) + \varepsilon_t.$$

Возникает вопрос: как можно узнать, есть ли коинтеграция в рассматриваемой модели? Если рассматриваемые ряды коинтегрированы, то остатки  $\varepsilon_t$  будут стационарными, и наоборот, если остатки модели будут стационарными, то ряды будут коинтегрированными. Следовательно, чтобы узнать, являются ли ряды коинтегрированными, нужно проверить остатки на стационарность.

Проверить, будут ли  $\varepsilon_t$  стационарными, можно с помощью тестов Энгла — Грэнжера, Дарбина — Уотсона и Филиппа — Перрона.

Существуют различные подходы при работе с нестационарными и возможно коинтегрируемыми рядами. Основные из них: подходы Энгла — Грэнжера и Йохансена. Как правило, используется подход Йохансена.

Дело в том, что в случае двух рассматриваемых рядов будет существовать только одно коинтеграционное соотношение. Однако если рассматривается много рядов, то будет существовать целый ряд коинтеграционных соотношений. Возникает



вопрос: как узнать, какое коинтеграционное соотношение будет лучшим для модели? Для этого необходим подход, при котором рассматривались бы все коинтеграционные соотношения. Именно такой подход был предложен С. Йохансеном.

Подробно прочитать про подход Йохансена можно, например, в [Brooks (2008)].

В качестве иллюстрирующего примера коинтеграции рассмотрим значения индекса газеты *Financial Times FTSE* со 2 апреля 1998 г. по 23 октября 2007 г. и индекс франкфуртской фондовой биржи *DAX* за тот же период. В обоих рядах взяты значения индексов при закрытии торгов (в конце дня).

Как было показано ранее, ряд значений индексов *DAX* имеет первый порядок интегрированности. Ряд значений индексов *FTSE* также имеет первый порядок интегрированности. Тест Дики — Фуллера для первых разностей *FTSE* приведен в табл. 8.17.

**Таблица 8.17.** Расширенный тест Дики — Фуллера для однодневных приращений индекса *FTSE*

1% Critical Value*	-3,4361	ADF Test Statistic -22,28069		
5% Critical Value	-2,8632			
10% Critical Value	-2,5677			
* Mac Kinnon critical values for rejection of hypothesis of a unit root.				
Variable	Coefficient	Std. Error	t-Statistic	Prob
<i>D (FTSE (-1))</i>	-1,226863	0,055064	-22,28069	0,0000
<i>D (FTSE (-1), 2)</i>	0,208974	0,049664	4,207718	0,0000
<i>D (FTSE (-2), 2)</i>	0,144011	0,043756	3,291202	0,0000
<i>D (FTSE (-3), 2)</i>	0,046163	0,036596	1,261405	0,2073
<i>D (FTSE (-4), 2)</i>	0,060264	0,029149	2,067461	0,0388
<i>D (FTSE (-5), 2)</i>	0,052314	0,020429	2,071232	0,0384
<i>C</i>	0,202829	1,231982	0,164636	0,8692
<i>R-squared</i>	0,513396	<i>Mean dependent var</i>	-0,027942	
<i>Adjusted R-squared</i>	0,512178	<i>S. D. dependent var</i>	86,49838	
<i>S. E. of regression</i>	60,41413	<i>Akaike info criterion</i>	11,04323	
<i>Sum squared resid</i>	8752381	<i>Schwarz criterion</i>	11,06007	
<i>Log likelihood</i>	-13 272,48	<i>F-statistic</i>	421,6717	
<i>Duubin — Watson stat.</i>	2 000 189	<i>Prob (F-statistic)</i>	0,000000	

Augmented Dickey — Fuller Test Equation.

Dependent Variable: *D (FTSE-2)*.

Method: Least Squares.

Date: 10/23/10. Time: 23:08.

Sample (adjusted): 8 2412.

Included observations: 2 405 after adjusting endpoints.

Для проверки на коинтегрируемость воспользуемся **процедурой Энгла — Грэнжера**. Для этого:

- 1) оценим регрессию *FTSE* на *DAX*;
- 2) с помощью теста не единичные корни проверим, являются ли остатки стационарными. Если да, то ряды коинтегрируемы.

В табл. 8.18 приведена оценка регрессии *DAX* на *FTSE*, выполненная при помощи эконометрического пакета *Eviews*.

**Таблица 8.18. Оцененная регрессия *DAX* на *FTSE***

<i>Variable</i>	<i>Coefficient</i>	<i>Std. Error</i>	<i>t-Statistic</i>	<i>Prob</i>
<i>DAX</i>	0,585042	0,005323	109,9024	0,0000
<i>C</i>	2384,582	28,56907	83,46727	0,0000
<i>R-squared</i>	0,833662	<i>Mean dependent var</i>	5424,903	
<i>Adjusted R-squared</i>	0,833593	<i>S. D. dependent var</i>	858,9729	
<i>S. E. of regression</i>	350,4014	<i>Akaike info criterion</i>	14,55686	
<i>Sum squared resid</i>	2,96 E +8	<i>Schwarz criterion</i>	14,56166	
<i>Log likelihood</i>	-17 553,58	<i>F-statistic</i>	12 078,53	
<i>Durbin — Watson stat</i>	0,046881	<i>Prob (F-statistic)</i>	0,000000	

*Dependent Variable: FTSE.*

*Method: Least Squares.*

*Date: 10/23/10. Time: 23:55.*

*Sample (adjusted): 1 2412.*

*Included observations: 2 512.*

Далее, согласно процедуре Энгла — Грэнжера, проведем **расширенный тест Дики — Фуллера** для остатков оцененной регрессии (табл. 8.19).

**Таблица 8.19. *ADF*-тест для остатков регрессии *DAX* на *FTSE***

<i>1% Critical Value*</i>	-3,4361	<i>ADF Test Statistic -3,455096</i>
<i>5% Critical Value</i>	-2,8632	
<i>10% Critical Value</i>	-2,5677	

*\* Mac Kinnon critical values for rejection of hypothesis of a unit root.*

<i>Variable</i>	<i>Coefficient</i>	<i>Std. Error</i>	<i>t-Statistic</i>	<i>Prob</i>
<i>A (-1)</i>	-0,015207	0,004401	-3,455096	0,0000
<i>D (A (-1))</i>	-0,121301	0,020505	-5,915740	0,0000
<i>D (A (-2))</i>	-0,141216	0,020613	-6,850871	0,0000
<i>D (A (-3))</i>	-0,144009	0,020571	-7,000536	0,0000
<i>D (A (-4))</i>	-0,054165	0,020533	-2,637995	0,0084
<i>D (A (-5))</i>	-0,078523	0,020371	-3,854646	0,0001
<i>C</i>	-0,716122	1,508182	-0,747825	0,6350

Окончание табл. 8.19

<i>R-squared</i>	0,053592	<i>Mean dependent var</i>	-0,487627
<i>Adjusted R-squared</i>	0,051225	<i>S. D. dependent var</i>	75,94046
<i>S. E. of regression</i>	73,96987	<i>Akaike info criterion</i>	11,44810
<i>Sum squared resid</i>	13 126 229	<i>Schwarz criterion</i>	11,46493
<i>Log likelihood</i>	-13 765,06	<i>F-statistic</i>	22,64122
<i>Durbin — Watson stat</i>	2,007991	<i>Prob (F-statistic)</i>	0,000000

*Augmented Dickey — Fuller Test Equation.*

*Dependent Variable: D (A).*

*Method: Least Squares.*

*Date: 10/23/10. Time: 23:57.*

*Sample (adjusted): 7 2412.*

*Included observations: 2 406 after adjusting endpoints.*

Уже на 1%-ном уровне значимости гипотеза о наличии единичного корня может быть отклонена. Следовательно, ряд остатков стационарен. И ряды значений индексов DAX и FTSE в период со 2 апреля 1998 г. по 23 октября 2007 г. коинтегрируемы.

## 8.8. Модели ARCH и GARCH

Как правило, гетероскедастичность связывают с пространственными данными, а при анализе временных рядов априори полагают, что выполняется условие гомоскедастичности. Однако предположение о гомоскедастичности временных рядов не всегда обосновано. В своей работе Р. Энгл [Engle (1982)], анализируя макроэкономические данные, заметил, что нарушения постоянства дисперсии во временных рядах случаются чаще, чем принято считать. Это связано с тем, что в некоторых временных рядах существует склонность к кластеризации отклонений от средних значений. Большие и малые отклонения как бы чередуются с течением времени.

Для анализа подобных рядов Энгл предложил пользоваться формой гетероскедастичности, при которой последующие значения отклонений будут зависеть от величин предыдущих. Подобный подход автор назвал *моделью ARCH* (моделью условной авторегрессионной гетероскедастичности). В дальнейшем эта модель была проверена и опробована исследователями. Так, Н. Колсон и Р. Робинс [Coulson, Robins (1985)] использовали эту модель для изучения инфляции; Р. Энгл, Д. Хендри, Д. Трамбал

[Engle, Hendry, Trumbull (1985)] при помощи ARCH-модели исследовали структуру процентных ставок; Р. Энгл, Д. Лильен, Р. Робинс [Engle, Lilien, Robins (1987)] исследовали изменения рыночных доходностей; поведение обменных курсов исследовалось Я. Домовитцом, К. Хаккио [Domowitz, Hakio (1985)].

### Модель ARCH

В моделях типа ARCH используемые ряды предполагаются стационарными. Простейшая форма модели ARCH, модель ARCH, (1) выглядит следующим образом:

$$y_t = \beta' x_t + \varepsilon_t;$$

$$\varepsilon_t = u_t \sqrt{\alpha_0 + \alpha_1 \varepsilon_{t-1}^2},$$

где  $u_t$  распределены по стандартному нормальному распределению. Следовательно,  $E(\varepsilon_t | x_t, \varepsilon_{t-1}) = 0$ , поэтому  $E(\varepsilon_t | x_t) = 0$ ,  $E(y_t | x_t) = \beta' x_t$ . То есть в контексте данной модели случайные остатки не зависят от значений факторов, входящих в модель регрессии.

Эта модель является классической регрессионной моделью. Однако  $\text{var}(\varepsilon_t | \varepsilon_{t-1}) = E(\varepsilon_t^2 | \varepsilon_{t-1}) = E(u_t^2)(\alpha_0 + \alpha_1 \varepsilon_{t-1}^2) = \alpha_0 + \alpha_1 \varepsilon_{t-1}^2$ .

Таким образом,  $\varepsilon_t$  — условная гетероскедастичность по отношению к  $\varepsilon_{t-1}$ . Безусловной дисперсией  $\varepsilon_t$  будет

$$\text{var}[\varepsilon_t] = \text{var}\{E[\varepsilon_t | \varepsilon_{t-1}]\} + E\{\text{var}[\varepsilon_t | \varepsilon_{t-1}]\} = \alpha_0 + \alpha_1 E[\varepsilon_{t-1}^2] = \alpha_0 + \alpha_1 \text{var}[\varepsilon_{t-1}].$$

В случае если «возмущения» в процессе проходят слабо стационарно, то безусловная дисперсия не меняется со временем и определяется по формуле

$$\text{var}[\varepsilon_t] = \text{var}[\varepsilon_{t-1}] = \alpha_0 + \alpha_1 \text{var}[\varepsilon_{t-1}] = \frac{\alpha_0}{1 - \alpha_1}.$$

Это отношение будет конечным и положительным, если  $|\alpha_1|$  будет меньше единицы. Тогда  $\varepsilon_t$  будет распределено с нулевым средним и дисперсией  $\sigma^2 = \frac{\alpha_0}{1 - \alpha_1}$ .

Вместо  $\sigma_t^2$  в литературе часто используют  $h_t$ . В таком случае модель ARCH (1) будет выглядеть следующим образом:

$$y_t = \beta' x_t + \varepsilon_t,$$

где  $\varepsilon_t \sim N(0, h_t)$ ;  $h_t = \alpha_0 + \alpha_1 u_{t-1}^2$ .

Как было показано выше, модель ARCH удовлетворяет классическим предположениям и обыкновенный метод наименьших квадратов дает лучшую линейную несмещенную оценку параметра  $\beta$ .

Однако может существовать лучшая нелинейная оценка. Логарифмическая функция правдоподобия для этой модели была задана Энглом:

$$\ln L = -\frac{1}{T} \sum_{t=1}^T \ln(\alpha_0 + \alpha_1 \varepsilon_{t-1}^2) - \frac{1}{2} \sum_{t=1}^T \frac{\varepsilon_t^2}{\alpha_0 + \alpha_1 \varepsilon_{t-1}^2}, \quad \varepsilon_t = y_t - \beta' x_t. \quad (8.14)$$

Если расширить простейшую модель условной гетероскедастичности до  $\tau$  лагов, получим модель ARCH ( $\tau$ ):

$$y_t = \beta' x_t + \varepsilon_t;$$

$$\sigma^2 = \alpha_0 + \alpha_1 \varepsilon_{t-1}^2 + \alpha_2 \varepsilon_{t-2}^2 + \dots + \alpha_q \varepsilon_{t-q}^2. \quad (8.15)$$

Существует другой способ описания стандартной модели ARCH. Для примера используем модель ARCH (1):

$$y_t = \beta' x_t + \varepsilon_t;$$

$$\varepsilon_t = v_t \sigma_t,$$

где  $v_t \sim N(0, 1)$ ;  $\sigma_t^2 = \alpha_0 + \alpha_1 \varepsilon_{t-1}^2$ .

Несмотря на то, что на первый взгляд два представленных способа описания модели ARCH отличаются, можно показать, что это одно и то же. В самом деле, если  $v_t$  подчиняется нормальному распределению с нулевым средним и единичной дисперсией, то  $u_t$  также будет подчиняться нормальному распределению с нулевым средним и дисперсией  $\sigma_t^2$ .

Позднее Р. Энглom, Д. Лильеном и Д. Робинсом [Engle, Lilien, Robins (1987)] была предложена модификация стандартной модели ARCH — модель ARCH-M:

$$y_t = \beta' x_t + \delta \sigma_t^2 + \varepsilon_t;$$

$$\text{Var}[\varepsilon_t | \psi_t] = \text{ARCH}(q). \quad (8.16)$$

Одной интересной особенностью модифицированной модели является то, что при выполнении основных предположений коэффициент  $\delta$  может трактоваться как коэффициент относительной склонности к риску.

Вскоре после своего появления модели ARCH-M были опробованы в ряде исследований. Так, например К. Френч с соавторами [French, Schwert, Stambaugh (1987)] использовали модели для изучения колебаний *Standart and Poor's Index*; Р. Чоу [Chou (1988)] опробовал модель при изучении доходностей Нью-Йоркской фондовой биржи.

Модель ARCH-M имеет несколько заслуживающих внимания статистических особенностей. В отличие от стандартной регрессионной модели ошибка спецификации дисперсионной функции оказывает влияние на состоятельность оценки параметров среднего. Напомним, что в классической регрессионной установке взвешенный метод наименьших квадратов дает состоятельную оценку до тех пор, пока веса не коррелируют с возмущениями. В этой модели все будет по-другому. Если в ARCH-части модели была сделана ошибка спецификации, то оценки  $\beta$  и  $\delta$  не будут состоятельны. Т. Боллерслев, Р. Чоу и К. Кронер [Bollerslev, Chou, Kroner (1992)] привели ряд исследований, в которых решался вопрос спецификации модели ARCH-M, и после переспецификации моделей ими были получены абсолютно другие результаты.

### Условие неотрицательности коэффициентов

Поскольку  $h_t^2$  — условная дисперсия, ее значение в любой момент времени должно быть сугубо положительно. Отрицательная дисперсия бессмысленна. Для того чтобы быть уверенными, что результат получен при положительной условной дисперсии, обычно вводят условие неотрицательности коэффициентов регрессии. Например, для модели ARCH ( $\tau$ ) все коэффициенты

должны быть неотрицательными:  $\alpha_i \geq 0$  для любых  $i = 0, 1, 2, \dots, q$ . Можно показать, что это достаточное, но не необходимое условие неотрицательности условной дисперсии.

Модели ARCH оказали серьезное влияние на развитие аппарата анализа временных рядов. Однако модель ARCH в первоначальном виде редко используется в последнее время. Это связано с тем, что при применении этих моделей возникает ряд проблем.

- Как должно находиться значение параметра  $q$ , определяющего количество лагов квадратов остатков в модели? Подход состоит в том, чтобы использовать метод максимального правдоподобия, хотя это не всегда является лучшим способом.
- Значение  $\tau$  количества лагов квадратов ошибок, которое необходимо для того, чтобы охватить все зависимости в условной дисперсии, может быть очень велико. Энгл обошел эту проблему тем, что использовал спецификацию модели, в которой искусственно ограничил длину лагов в модели ARCH (4)  $\sigma_t^2 = \gamma_0 + \gamma_1(0,4\varepsilon_{t-1}^2 + 0,3\varepsilon_{t-2}^2 + 0,2\varepsilon_{t-3}^2 + 0,1\varepsilon_{t-4}^2)$  таким образом, что в условной дисперсии есть только два параметра, в то время как в неограниченной модели ARCH (4) нужно было бы использовать пять параметров.
- Ограничение на неотрицательность может нарушаться. Чем больше параметров в уравнении условной дисперсии, тем больше вероятность того, что некоторые из них будут отрицательными.

Некоторых из этих проблем можно избежать при помощи модели GARCH, которая представляет собой естественную модификацию модели ARCH. В отличие от модели ARCH модели GARCH широко используются на практике.

Для того чтобы определить, являются ли ошибки в модели условно гетероскедастичными, можно провести следующую процедуру.

1. Применить МНК к уравнению регрессии и получить остатки  $\varepsilon_t$ .
2. С помощью МНК оценить регрессию  $\varepsilon_t^2 = b_0 + b_1\varepsilon_{t-1}^2 + \dots + b_p\varepsilon_{t-p}^2$ .
3. С помощью  $F$ -теста проверить гипотезу  $H_0: b_0 = b_1 = \dots = b_p = 0$ .

### Модель GARCH

Модель GARCH была предложена Т. Боллерслемом [Bollerslev (1986)]. В этой модели предполагается, что условная дисперсия будет зависеть также от собственных лагов. Простейшая форма модели GARCH выглядит следующим образом:

$$\sigma_t^2 = \alpha_0 + \alpha_1 \varepsilon_{t-1}^2 + \beta \sigma_{t-1}^2.$$

Это модель вида GARCH (1, 1) (поскольку используют первые лаги  $u_t^2$  и  $\sigma_t^2$ ). Заметим, что модель GARCH может быть представлена в виде модели ARMA для условной дисперсии. Для того чтобы убедиться в этом, проведем следующие математические преобразования:

$$\begin{aligned} u_t &= \varepsilon_t^2 - \sigma_t^2, \\ \text{или} \quad \sigma_t^2 &= \varepsilon_t^2 - u_t. \end{aligned}$$

Перепишем условную дисперсию в виде

$$\begin{aligned} \varepsilon_t^2 - u_t &= \alpha_0 + \alpha_1 \varepsilon_{t-1}^2 + \beta (\varepsilon_{t-1}^2 - u_{t-1}); \\ \varepsilon_t^2 &= \alpha_0 + \alpha_1 \varepsilon_{t-1}^2 + \beta (\varepsilon_{t-1}^2 - u_{t-1}) + u_t; \\ \varepsilon_t^2 &= \alpha_0 + (\alpha_1 + \beta) \varepsilon_{t-1}^2 - \beta u_{t-1} + u_t. \end{aligned}$$

Последнее уравнение есть не что иное, как процесс ARMA (1, 1) для квадрата ошибок.

В чем именно состоит преимущество моделей GARCH перед моделями ARCH? Основное преимущество моделей GARCH заключается в том, что для спецификации моделей GARCH требуется меньше параметров. Следовательно, модель в большей степени будет удовлетворять условиям неотрицательности.

Рассмотрим условную дисперсию модели GARCH (1, 1):

$$\sigma_t^2 = \alpha_0 + \alpha_1 \varepsilon_{t-1}^2 + \beta \sigma_{t-1}^2.$$

Для  $\tau = 1$  условной дисперсии будет выполняться уравнение

$$\sigma_{t-1}^2 = \alpha_0 + \alpha_1 \varepsilon_{t-2}^2 + \beta \sigma_{t-2}^2.$$

Для  $\tau = 2$  соответственно будет выполняться уравнение



$$\sigma_{t-2}^2 = \alpha_0 + \alpha_1 \varepsilon_{t-3}^2 + \beta \sigma_{t-3}^2.$$

Следовательно, условную дисперсию можно представить в виде

$$\sigma_t^2 = \alpha_0 + \alpha_1 \varepsilon_{t-1}^2 + \beta(\alpha_0 + \alpha_1 \varepsilon_{t-2}^2 + \beta \sigma_{t-2}^2);$$

$$\sigma_t^2 = \alpha_0 + \alpha_1 \varepsilon_{t-1}^2 + \beta \alpha_0 + \alpha_1 \beta \varepsilon_{t-2}^2 + \beta^2 \sigma_{t-2}^2.$$

Она в свою очередь равна

$$\sigma_t^2 = \alpha_0 + \alpha_1 \varepsilon_{t-1}^2 + \beta \alpha_0 + \alpha_1 \beta \varepsilon_{t-2}^2 + \beta^2 (\alpha_0 + \alpha_1 \varepsilon_{t-3}^2 + \beta \sigma_{t-3}^2);$$

$$\sigma_t^2 = \alpha_0 + \alpha_1 \varepsilon_{t-1}^2 + \beta \alpha_0 + \alpha_1 \beta \varepsilon_{t-2}^2 + \beta^2 \alpha_0 + \alpha_1 \beta^2 \varepsilon_{t-3}^2 + \beta^3 \sigma_{t-3}^2);$$

$$\sigma_t^2 = \alpha_0 (1 + \beta + \beta^2) + \alpha_1 \varepsilon_{t-1}^2 (1 + \beta L + \beta^2 L^2) + \beta^3 \sigma_{t-3}^2.$$

В итоге получим уравнение

$$\sigma_t^2 = \alpha_0 (1 + \beta + \beta^2 + \dots) + \alpha_1 \varepsilon_{t-1}^2 (1 + \beta L + \beta^2 L^2 + \dots) + \beta^\infty \sigma_t^2.$$

Первая скобка в этом уравнении — константа, причем при бесконечно большой выборке  $\beta^\infty$  будет стремиться к нулю. Следовательно, модель GARCH (1, 1) может быть представлена в виде

$$\sigma_t^2 = \gamma_0 + \alpha_1 \varepsilon_{t-1}^2 (1 + \beta L + \beta^2 L^2 + \dots),$$

или

$$\sigma_t^2 = \gamma_0 + \gamma_1 \varepsilon_{t-1}^2 + \gamma_2 \varepsilon_{t-2}^2 + \dots$$

Последнее уравнение есть не что иное, как модель ARMA. Таким образом, модель GARCH (1, 1), содержащая только три параметра в уравнении условной дисперсии, учитывает влияние на условную дисперсию бесконечно большого количества квадратов ошибок.

Модель GARCH (1, 1) может быть расширена до модели GARCH (p, q):

$$\sigma_t^2 = \alpha_0 + \alpha_1 \varepsilon_{t-1}^2 + \alpha_2 \varepsilon_{t-2}^2 + \dots + \alpha_q \varepsilon_{t-q}^2 + \beta_1 \sigma_{t-1}^2 + \beta_2 \sigma_{t-2}^2 + \dots + \beta_p \sigma_{t-p}^2;$$

$$\sigma_t^2 = \alpha_0 + \sum_{i=1}^q \alpha_i \varepsilon_{t-i}^2 + \sum_{j=1}^p \beta_j \sigma_{t-j}^2. \quad (8.17)$$

Необходимо отметить, что на практике возможностей модели *GARCH* (1, 1), как правило, хватает, и не всегда целесообразно пользоваться моделями *GARCH* более высоких порядков.

Несмотря на то, что условная дисперсия модели *GARCH* изменяется со временем, безусловная дисперсия будет постоянной при  $\alpha_1 + \beta < 1$ :

$$\text{var}(\varepsilon_t) = \frac{\alpha_0}{1 - (\alpha_1 + \beta)}.$$

В случае если  $\alpha_1 + \beta \geq 1$ , безусловная дисперсия не будет определена. Этот случай называется «нестационарностью дисперсии». В случае если  $\alpha_1 + \beta = 1$ , модель будет называться *IGARCH*. Нестационарность дисперсии не имеет строгой мотивации существования. Более того, модели *GARCH*, которых коэффициенты привели к нестационарности дисперсии, могут иметь некоторые более нежелательные свойства. Одним из них является невозможность сделать прогноз дисперсий исходя из модели. Для стационарных моделей *GARCH* прогнозы условной дисперсии сходились в долгосрочное среднее значение дисперсий. Для процесса *IGARCH* такой сходимости не будет. Прогноз условной дисперсии равняется бесконечности.

### Оценка модели

Метод наименьших квадратов не всегда может использоваться для оценки модели *GARCH*. Основная причина заключается в том, что фундамент МНК — минимизация суммы квадратов остатков, а эта сумма не всегда будет зависеть от условной дисперсии.

Для оценки модели *GARCH* нужно использовать метод максимального правдоподобия. Данный метод основывается на предположении о том, что вся информация о выборке содержится в функции максимального правдоподобия. Для оценки неизвестного параметра максимизируется функция максимального правдоподобия. Фактически метод заключается в том, чтобы найти наиболее близкие к реальным данным значения параметров модели. Максимальное правдоподобие

может использоваться для оценки параметров как линейных, так и нелинейных моделей.

Как уже говорилось ранее, логарифмическая функция правдоподобия для модели была предложена Энглom:

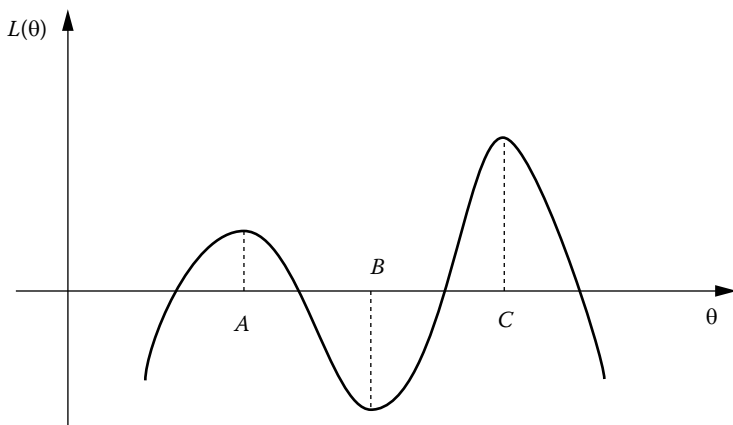
$$\ln L = -\frac{1}{T} \sum_{t=1}^T \ln(\alpha_0 + \alpha_1 \varepsilon_{t-1}^2) - \frac{1}{2} \sum_{t=1}^T \frac{\varepsilon_t^2}{\alpha_0 + \alpha_1 \varepsilon_{t-1}^2},$$

где  $\varepsilon_t = y_t - \beta' x_t$ .

К сожалению, процесс максимизации функции максимального правдоподобия для моделей с условной гетероскедастичностью более сложный, чем максимизация гомоскедастичных моделей. Метод нахождения производных в аналитической форме для функции максимального правдоподобия разработан только для простейших моделей GARCH. Более того, итоговые формулы сложны, и поэтому для максимизации функции максимального правдоподобия обычно используются числовые методы.

При использовании числовых методов также возникает ряд сложностей. Они связаны с максимизацией функций, имеющих несколько локальных точек экстремума. При работе с моделями GARCH это частое явление. Поэтому при использовании различных числовых алгоритмов могут быть найдены различные максимумы функции максимального правдоподобия.

На рис. 8.11 представлена зависимость значений функции максимального правдоподобия от параметра  $\theta$ . Как можно увидеть, функция достигает своего максимума при  $\theta = C$  и своего локального максимума при  $\theta = A$ . Как показали в своей работе К. Брукс с соавторами [Brooks, Burle, Persaud (2001)], разные оптимизационные методы могут приводить к разным оценкам коэффициентов модели и стандартных ошибок модели. Так, в случае, приведенном на рис. 8.11, в качестве оптимума могут выступать как точка A, так и точка C. На практике ситуация усложняется еще больше, поскольку приходится искать максимум функции максимального правдоподобия по отношению к нескольким параметрам. Кроме того, сложности при оптимизации функции максимального правдоподобия могут возникать, если функция не имеет в окрестности экстремума



**Рис. 8.11.** График зависимости функции максимального правдоподобия от параметра  $\theta$

ярко выраженного максимума. В этом случае также достаточно сложно найти оптимальное значение параметра.

Э. Берндт, Б. Холл, Р. Холл, Д. Хаусман [Berndt, B. Hall, R. Hall Hausman (1974)] предложили итерационный способ нахождения максимума подобных функций. Этот метод, получивший в честь авторов название *ВННН*, заключается в пошаговом нахождении первых и вторых производных функции максимального правдоподобия.

В качестве примера рассмотрим ряд значений однодневных приращений индекса *DAX* в период со 2 апреля 1998 г. по 23 октября 2003 г. Ранее было показано, что наиболее адекватным образом модель описывает процесс *ARIMA* (0, 1, 0). Однако **тест Льюнга — Бокса** показал наличие автокорреляции в остатках модели.

Наличие в модели *ARCH* эффекта может быть также определено с помощью теста *ARCH-LM*. В рамках данного теста оценивается уравнение

$$e_t = \alpha_0 + \alpha_1 e_{t-1}^2 + \alpha_2 e_{t-2}^2 + \dots + \alpha_p e_{t-p}^2, \quad (8.17)$$

где  $e_t$  — остатки модели. Проверяется гипотеза о том, что  $\alpha_1 = \alpha_2 = \dots = \alpha_p = 0$ .

Таким образом, принятие гипотезы означает, что в модели отсутствует автокорреляция остатков. Для того чтобы провести тест *ARCH-LM*, воспользуемся возможностями эконометрического пакета *Eviews* (табл. 8.20—8.22).

**Таблица 8.20.** Тест *ARCH-LM*: один лаг

<i>F</i> -statistic	41,19047	<i>Prob. F</i> (1,2408)	0,0000
<i>Obs*R-squared</i>	40,53137	<i>Prob. Chi-Square</i> (1)	0,0000

**Таблица 8.21.** Тест *ARCH-LM*: три лага

<i>F</i> -statistic	74,40652	<i>Prob. F</i> (3,2404)	0,0000
<i>Obs*R-squared</i>	204,5938	<i>Prob. Chi-Square</i> (3)	0,0000

**Таблица 8.22.** Тест *ARCH-LM*: шесть лагов

<i>F</i> -statistic	69,53690	<i>Prob. F</i> (3,2398)	0,0000
<i>Obs*R-squared</i>	356,4258	<i>Prob. Chi-Square</i> (6)	0,0000

Был построен тест *ARCH-LM* с одним, тремя и шестью лагами. Во всех случаях наблюдается значимая связь. В модели определенно присутствует автокорреляция.

Построим оценку однодневных приращений индекса *DAX* в рассматриваемый период, используя модели *ARCH* (1) (табл. 8.23).

**Таблица 8.23.** Оценка модели со свободным членом при помощи модели *ARCH* (1)

<i>Variable</i>	<i>Coefficient</i>	<i>Std. Error</i>	<i>t-Statistic</i>	<i>Prob</i>
C	2,534816	1,469567	1,724940	0,0845
<i>Variance Equation</i>				
C	4495,572	137,8380	32,61490	0,0000
<i>RESID</i> (-1) 2	0,237466	0,031614	7,511380	0,0000
<i>R-squared</i>	-0,000340	<i>Mean dependent var</i>		1,137876
<i>Adjusted R-squared</i>	-0,000340	<i>S. D. dependent var</i>		75,82268
<i>S. E. of regression</i>	75,83556	<i>Akaike info criterion</i>		11,46326
<i>Sum squared resid</i>	13 859 987	<i>Schwarz criterion</i>		11,47046
<i>Log likelihood</i>	-13 815,96	<i>Hannan — Quinn criterion</i>		11,46588
<i>Durbin — Watson stat</i>	1,991082			

Как очевидно из табл. 8.20, на 5%-ном уровне значимости свободный член модели статистически незначим. Исключим его из уравнения регрессии (табл. 8.24).

**Таблица 8.24.** Оценка модели без свободного члена при помощи модели ARCH (1)

Variable	Coefficient	Std. Error	t-Statistic	Prob
Variance Equation				
C	4527,292	138,6908	32,64305	0,0000
RESID (-1) 2	0,230330	0,031407	7,333714	0,0000
R-squared	-0,000225	Mean dependent var		1,137 876
Adjusted R-squared	0,000190	S. D. dependent var		75,82268
S. E. of regression	75,81550	Akaike info criterion		11,46368
Sum squared resid	13 858 403	Schwarz criterion		11,46848
Log likelihood	-13 817,46	Hannan — Quinn criterion		11,46542
Durbin — Watson stat	1,991310			

После исключения из модели свободного члена значения **информационных критериев Шварца и Ханнан — Куина** уменьшились. Следовательно, исключение повышает качество модели.

Проведем тест ARCH-LM на автокорреляцию остатков в этой модели (табл. 8.25).

**Таблица 8.25.** Тест ARCH-LM: три лага

F-statistic	46,57907	Prob. F (3,2404)	0,0000
Obs*R-squared	132,2807	Prob. Chi-Square (3)	0,0000

Уже при трех лагах в остатках присутствует автокорреляция.

Построим модель ARCH (2) для описания однодневных приращений индексов DAX и проведем тест ARCH-LM для остатков новой модели (табл. 8.26—8.28).

**Таблица 8.26.** Тест ARCH-LM: один лаг

F-statistic	3,422462	Prob. F (1,2408)	0,0644
Obs*R-squared	3,420443	Prob. Chi-Square (1)	0,0644

**Таблица 8.27.** Тест *ARCH-LM*: три лага

<i>F</i> -statistic	7,599411	<i>Prob. F</i> (3,2404)	0,0000
<i>Obs</i> * <i>R</i> -squared	22,62164	<i>Prob. Chi-Square</i> (3)	0,0000

**Таблица 8.28.** Тест *ARCH-LM*: шесть лагов

<i>F</i> -statistic	22,57308	<i>Prob. F</i> (6,2398)	0,0000
<i>Obs</i> * <i>R</i> -squared	128,5721	<i>Prob. Chi-Square</i> (6)	0,0000

Итак, начиная с трех лагов, в остатках модели присутствует статистически значимая связь. К сожалению, используя возможности модели *ARCH*, нам не удалось избавиться от автокорреляции остатков. Поэтому не остается ничего другого как использовать модель *GARCH* (1, 1) (см. табл. 8.26).

**Таблица 8.29.** Оценка модели при помощи модели *GARCH* (1, 1)

<i>Variable</i>	<i>Coefficient</i>	<i>Std. Error</i>	<i>t</i> -Statistic	<i>Prob</i>
<i>Variance Equation</i>				
<i>C</i>	25,2812	8,685661	2,910789	0,0036
<i>RESID</i> (-1) 2	0,070466	0,008637	8,158957	0,0000
<i>GARCH</i> (-1)	0,926667	0,008676	106,8090	0,0000
<i>R</i> -squared	-0,000225	<i>Mean dependent var</i>		1,137876
<i>Adjusted R</i> -squared	0,000190	<i>S. D. dependent var</i>		75,82268
<i>S. E. of regression</i>	75,81550	<i>Akaike info criterion</i>		11,22525
<i>Sum squared resid</i>	13 858 403	<i>Schwarz criterion</i>		11,23246
<i>Log likelihood</i>	-13 529,04	<i>Hannan — Quinn criterion</i>		11,22787
<i>Durbin — Watson stat</i>	1,991310			

*Dependent Variable*: DDAX.

*Method*: ARCH-ML (Marquardt) — Normal distribution.

*Sample (adjusted)*: 2 2412.

$GARCH = C(1) + C(2) RESID(-1)^2 + C(3) GARCH(-1)$ .

Все регрессоры в условной дисперсии значимы. Построим тест *ARCH-LM*, чтобы проверить, удалось ли нам избавиться от автокорреляции в остатках (табл. 8.30—8.33).

**Таблица 8.30.** Тест *ARCH-LM*: один лаг

<i>F</i> -statistic	3,334445	<i>Prob. F</i> (1,2408)	0,0680
<i>Obs</i> * <i>R</i> -squared	3,332600	<i>Prob. Chi-Square</i> (1)	0,0679

Таблица 8.31. Тест *ARCH-LM*: три лага

<i>F</i> -statistic	1,108968	Prob. <i>F</i> (3,2404)	0,3441
Obs* <i>R</i> -squared	3,327834	Prob. Chi-Square (3)	0,3438

Таблица 8.32. Тест *ARCH-LM*: шесть лагов

<i>F</i> -statistic	1,442166	Prob. <i>F</i> (6,2398)	0,1946
Obs* <i>R</i> -squared	8,647052	Prob. Chi-Square (6)	0,1944

Таблица 8.33. Тест *ARCH-LM*: 12 лагов

<i>F</i> -statistic	1,056489	Prob. <i>F</i> (12,2386)	0,3934
Obs* <i>R</i> -squared	12,67957	Prob. Chi-Square (12)	0,3928

Итак, тест *ARCH-LM* был построен для моделей с лагом, равным единице, трем, шести и 12. Во всех моделях уже на 5%-ном уровне значимости есть основания для отклонения гипотезы о значимости связи, т.е. о присутствии автокорреляции.

С помощью модели *GARCH* (1, 1) нам удалось корректно описать модель. Возникает вопрос: почему не хватило мощности моделей *ARCH* (1) и *ARCH* (2)? Дело в том, что как отмечалось раньше, модель *GARCH* (1, 1), содержащая только три параметра в уравнении условной дисперсии, учитывает влияние на условную дисперсию бесконечно большого количества квадратов ошибок.

## Контрольные вопросы и задания

1. Дайте определение строгой и слабой стационарностей. Почему при работе с временными рядами необходимо делать проверку стационарности?
2. Сформулируйте условие стационарности для авторегрессионных процессов. Будет ли стационарным процесс  $y_t = 2y_{t-1} - y_{t-2} + u_t$ ?
3. Какая гипотеза проверяется в тесте Дики — Фуллера? В чем преимущество расширенного теста Дики — Фуллера перед обычным тестом Дики — Фуллера?
4. Перечислите основные свойства автокорреляционных и частных автокорреляционных функций процессов *AR*, *MA* и *ARMA*.



5. Перечислите типичные нестационарные процессы. С помощью какой операции можно привести эти процессы к стационарному виду?
6. Опишите два способа идентификации модели *ARMA*.
7. Какие сложности возникают при использовании коинтеграции в первоначальном виде?
8. В чем заключается основное отличие моделей *ARCH* и *GARCH* от моделей *ARMA* и *ARIMA*?
9. Перечислите основные вопросы, возникающие у исследователя при использовании модели *ARCH* в первоначальном виде.
10. В чем состоит основное преимущество моделей *GARCH* перед моделями *ARCH*?

### 9.1. Панельные данные и их преимущества

Термин «**панельные данные**» (*panel data*) пришел из обследований индивидов, и в этом контексте «панель» представляла собой группу индивидов, за которыми регулярно осуществляли наблюдения в течение определенного периода времени. В настоящее время методы анализа панельных данных получили большое распространение, и понимание панельных данных стало намного шире. Наряду с термином «панельные данные» иногда также используется термин «**лонгитюдные данные**» (*longitudinal data*).

Панельные данные состоят из повторных наблюдений одних и тех же выборочных единиц, которые осуществляются в последовательные периоды времени. В качестве объектов наблюдения могут выступать индивиды, домашние хозяйства, фирмы, страны и т.д. Примером панельных данных могут быть ежегодные обследования одних и тех же домашних хозяйств или индивидов (например, для определения изменения их благосостояния), ежеквартальные данные об экономической деятельности отдельных компаний, ежегодные социально-экономические показатели для регионов одной страны или для группы стран и т.д.

Панельные данные совмещают в себе как пространственные данные, так и временные ряды и сочетают достоинства каждого из этих видов данных. Это позволяет строить более адекватные и содержательные модели для изучения истинной причинно-следственной связи между различными переменными, что представляется невозможным в рамках только временных или только пространственных данных. Выделяют следующие преимущества использования панельных данных.

1. Панельные данные позволяют учитывать индивидуальную неоднородность. Временные ряды или пространственные данные не всегда позволяют учесть неоднородность индивидов, фирм, регионов или стран, что может привести к смещенным оценкам. Так, например, в исследовании Б. Балтаги и Д. Левина<sup>1</sup> изучался спрос на сигареты в США. Спрос моделировался как функция от лагов потребления, цены и дохода. Эти переменные отличались по штатам и во времени. Однако существовало множество других факторов, различающихся по штатам или во времени, которые могли оказывать влияние на потребление, например такие факторы, как религия, образование и реклама на телевидении и радио. При измерении этих переменных для каждого штата и периода времени возникают определенные трудности и очень сложно достигнуть того, чтобы их можно было включить в уравнение потребления. Однако пропуск этих переменных приведет к смещению в оценках. Панельные данные способны учитывать переменные, отличающиеся по штатам и во времени, вне зависимости от того, измеряемы они или нет, в то время как временные ряды или пространственные данные не позволяют этого сделать. Таким образом, панельные данные дают возможность избежать ошибки спецификации, возникающей из-за того, что существенные переменные не включены в модель.
2. Панельные данные содержат большое число наблюдений и тем самым предоставляют исследователю большее количество информации, им свойственна большая вариация и меньшая коллинеарность объясняющих переменных, они дают большее число степеней свободы и обеспечивают

---

<sup>1</sup> Baltagi B. H. and Levin D. Cigarette taxation: Raising revenues and reducing consumption, *Structural Change and Economic Dynamics*. 1992. 3. Pp. 321—335.

большую эффективность оценок. При анализе временных рядов исследователи часто сталкиваются с мультиколлинеарностью факторов. Например, в рассмотренном выше случае со спросом на сигареты существовала высокая зависимость между ценой и доходом в агрегированных временных рядах для США. Высокая коллинеарность между этими факторами будет менее вероятна при использовании панельных данных по штатам, так как пространственное измерение немного увеличивает вариацию факторов и делает более информативными данные по цене и доходу. Действительно, вариация в данных может быть разложена на две составляющие: вариацию между штатами разных размеров и с различными характеристиками и вариацию внутри штатов, при этом последняя обычно всегда больше. К тому же более информативные данные могут привести к более надежным оценкам параметров.

3. Панельные данные предоставляют возможность изучать динамику изменений индивидуальных характеристик единиц совокупности. Панельные данные хорошо подходят для изучения перемены работы, периода безработицы, изменений в доходах, для исследования длительности пребывания в определенном экономическом состоянии, например в бедности или в качестве безработного, а также могут помочь изучить скорость приспособления индивидов к изменениям в экономической политике. Так, при измерении безработицы пространственные данные позволяют оценить, какую долю в совокупности составляют безработные в конкретный момент времени; временные данные могут показать, как эта доля менялась во времени; и только панельные данные позволяют оценить, какая доля тех, кто являлся безработным в один период, останется безработным в другой период. Панельные данные могут использоваться как для объяснения того, почему различные единицы совокупности ведут себя по-разному, так и для того, чтобы определить, почему конкретная единица совокупности ведет себя по-разному в различные периоды времени.
4. Панельные данные лучше способны идентифицировать и измерить эффекты, которые просто не определяемы только во временных рядах или только в пространственных данных. В качестве примера может выступать

исследование того, происходит ли увеличение или уменьшение заработной платы за счет членства в профсоюзе. На этот вопрос лучше всего ответить, если мы наблюдаем переход работника с работы с профсоюзом на работу без профсоюза или наоборот, а это могут отразить только панельные данные. Рассматривая индивидуальную характеристику работника в качестве константы, можно будет определить, оказывает ли влияние членство в профсоюзе на зарплату и насколько. Подобный анализ может также использоваться для оценки других типов дифференциации заработной платы, например, для оценки премии, выплачиваемой за опасную или неприятную работу.

5. Панельные данные позволяют конструировать и тестировать более сложные поведенческие модели, чем пространственные данные и временные ряды в отдельности. Например, техническая эффективность лучше изучается и моделируется с панельными данными. Также в панелях может быть наложено меньше ограничений на модели распределенного лага, которые обычно рассматриваются во временных рядах.
6. Панельные данные позволяют избежать смещения, связанного с агрегированием данных, так как панельные данные, собранные на микроуровне (по индивидам, фирмам или домашним хозяйствам), могут быть измерены более точно, чем аналогичные переменные, полученные на макроуровне. При этом во временных рядах рассматривается изменение во времени характеристик некоторой усредненной репрезентативной единицы совокупности, а в пространственных данных не учитываются ненаблюдаемые индивидуальные характеристики единиц совокупности.
7. Панельные данные макроуровня имеют более длинные временные ряды, и панельные тесты на единичный корень имеют стандартные асимптотические распределения в отличие от проблемы нестандартных распределений, типичной для теста на единичный корень в анализе временных рядов.

Однако у панельных данных есть и недостатки. Определенные проблемы связаны со сбором данных: проблема покрытия, т.е. неполный учет интересующей совокупности; отсутствие отклика, которое может быть связано как с отсутствием взаимодействия

с респондентом, так и с ошибкой интервьюера, искажения, связанные с ошибками измерения, которые могут возникнуть по причине неправильного ответа из-за неясной формулировки вопроса, ошибок памяти, намеренного искажения ответа (престижное смещение), неподходящих информантов, ошибочной записи ответов и эффектов интервьюера.

Отсутствие данных может быть вызвано различными причинами. Например, если индивиды выбирают, что им не стоит работать, потому что предлагаемая зарплата ниже минимального размера оплаты труда, то в этом случае для этих индивидов будут отсутствовать данные по зарплате, но будут иметься данные по другим характеристикам. Так как пропущена только их зарплата, то выборка будет цензурированной. Однако если мы не наблюдаем всех данных этих индивидов, то выборка будет уже усеченной, и в результате получатся смещенные оценки. В первой волне панели отсутствие отклика может быть связано с тем, что выбранный индивид или домохозяйство отказались участвовать в опросе или просто никого не оказалось дома. Частичное отсутствие отклика возникает, когда не дан ответ на один или несколько вопросов. В последующих волнах панели вследствие отсутствия отклика может возникнуть проблема истощения данных, когда респондент, ранее принимавший участие в опросах, может умереть, переехать на другое место жительства или обнаружить, что затраты на участие в опросе для него стали слишком большими, и отказаться от дальнейшего участия.

В настоящее время панельные обследования в разных формах проводятся во многих странах. Впервые панельные данные начали формироваться в США в 1960-х гг. Среди наиболее известных баз панельных данных США можно выделить *PSID* и *NLS*.

Панельное исследование динамики доходов (*The US Panel Study of Income Dynamics (PSID)*), <http://psidonline.isr.umich.edu/>— база панельных данных по американским домохозяйствам, собираемая Институтом социальных исследований Мичиганского университета. База *PSID* появилась в 1968 г. и включала данные по 4800 семьям. В настоящее время она охватывает около 9000 американских семей. Данные содержат более 5000 переменных по экономике, демографии, здоровью и социальному поведению.

Национальные лонгитюдные исследования (*National Longitudinal Surveys (NLS)*, <http://www.bls.gov/nls/home.htm>) — панельные исследования, спонсируемые Бюро трудовой статистики США, которые начали проводиться с 1966 г. На примере нескольких групп мужчин и женщин изучаются различные аспекты активности на рынке труда и другие значительные события в жизни этих людей.

Европейские панельные данные стали появляться только в 1980-х гг. Так, например первая волна панельного обследования Немецкой социально-экономической панели (*Sozio-oekonomisches Panel (SOEP)*, <http://www.diw.de/soep>), формируемой Немецким институтом экономических исследований (*Deutsches Institut für Wirtschaftsforschung (DIW), Berlin*), состоялась в 1984 г. и охватила более 5000 западногерманских домохозяйств. В настоящее время это обследование содержит данные около 11 000 домохозяйств, которые включают в себя демографические переменные, зарплату, доход, выплаты пособий, уровень удовлетворенности различными аспектами жизни, надежды и страхи, политическую активность и т.д. С 1991 г. Институтом социальных и экономических исследований Эссекского университета проводится панельное исследование британских домохозяйств (*The British Household Panel Survey (BHPS)*, <http://www.iser.essex.ac.uk/survey/bhps>). Это национальная репрезентативная выборка 5500 домохозяйств и 10 300 индивидов, выбранных из 250 районов Великобритании. Эти данные отражают демографические характеристики домохозяйств, рынок труда, здоровье, образование, жилищные условия, потребление, доход и т.д. В 1994—2001 гг. при содействии Евростата проводилось Европейское панельное обследование домохозяйств (*The European Community Household Panel (ECHP)*, <http://epunet.essex.ac.uk/echp.php>), в рамках которого в странах, являющихся членами Евросоюза, собирались данные домохозяйств по доходу, работе и безработице, бедности, жилью, здоровью и т.д.

В России панельные обследования стали проводиться в 1990-х г. Наиболее известной базой панельных данных является РМЭЗ — Российский мониторинг экономического положения и здоровья населения (*Russia Longitudinal Monitoring Survey (RLMS)*, <http://www.cpc.unc.edu/projects/rlms>). РМЭЗ представляет собой серию общенациональных репрезентативных

опросов домохозяйств и индивидов, проводившихся в России с 1992 г. Данные обследований содержат ответы на более чем 3000 вопросов, касающихся доходов и расходов, материального благосостояния, занятости, уровня образования, состояния здоровья и т.д.

## 9.2. Однонаправленные модели панельных данных

### 9.2.1. Модели панельных данных и основные обозначения

Введем обозначения, которые в дальнейшем будут использоваться в рамках данной главы:

$y_{it}$  — зависимая переменная для  $i$ -й единицы совокупности в момент времени  $t$ ,  $i=1, \dots, N$  и  $t=1, \dots, T$ ;

$x_{it}^j$  — значение  $j$ -й объясняющей переменной для  $i$ -й единицы совокупности в момент времени  $t$ ,  $j=1, \dots, k$ .

Рассмотрим сбалансированные панели, где для каждой пространственной единицы имеется одинаковое число наблюдений по всем периодам времени. Тогда общее число наблюдений будет  $NT$ . При  $N=1$  и достаточно большом  $T$  получаются временные ряды, а при  $T=1$  и достаточно большом  $N$  получаются пространственные данные. Метод оценивания панельных данных относится к случаю, когда  $N > 1$  и  $T > 1$ . Будем рассматривать панельные данные с короткими временными рядами, где  $N$  намного больше  $T$ , что очень часто встречается на практике, когда число наблюдаемых единиц достаточно велико (может достигать нескольких сотен или тысяч), а число моментов наблюдения ограничено.

Для  $i$ -й единицы совокупности данные можно представить в виде

$$y_i = \begin{bmatrix} y_{i1} \\ y_{i2} \\ \vdots \\ y_{iT} \end{bmatrix}; \quad X_i = \begin{bmatrix} x'_{i1} \\ x'_{i2} \\ \vdots \\ x'_{iT} \end{bmatrix} = \begin{bmatrix} x_{i1}^1 & x_{i1}^2 & \cdots & x_{i1}^k \\ x_{i2}^1 & x_{i2}^2 & \cdots & x_{i2}^k \\ \vdots & \vdots & \ddots & \vdots \\ x_{iT}^1 & x_{iT}^2 & \cdots & x_{iT}^k \end{bmatrix} \quad \varepsilon_i = \begin{bmatrix} \varepsilon_{i1} \\ \varepsilon_{i2} \\ \vdots \\ \varepsilon_{iT} \end{bmatrix}, \quad (9.1)$$

где  $\varepsilon_{it}$  — соответствующая ошибка. Тогда объединенные данные по всем единицам совокупности примут вид



$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}; \quad X = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_N \end{bmatrix}; \quad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_N \end{bmatrix}, \quad (9.2)$$

где  $y$  и  $\varepsilon$  — матрицы размерностью  $NT \times 1$ , а  $x$  имеет размерность  $NT \times k$ .

Стандартная линейная модель может быть записана в виде

$$y_{it} = x'_{it}\beta + z'_i\alpha + \varepsilon_{it}. \quad (9.3)$$

В этой модели  $x_{it}$  содержит  $k$  регрессоров, причем константа не входит в число этих регрессоров;  $z'_i\alpha$  отражает индивидуальный эффект, характеризующий неоднородность (гетерогенность) единиц совокупности. При этом  $z_i$  состоит из константы и набора индивидуальных или групповых переменных, которые могут быть как наблюдаемыми (например, пол, национальность, место жительства и т.п.), так и ненаблюдаемыми (например, профессиональные навыки индивидов, вкусы и предпочтения). Эти переменные можно рассматривать как постоянные во времени, которые не меняются с течением временем  $t$ . В зависимости от предположений относительно структуры  $z_i$  можно получить различные модели панельных данных.

1. **Объединенная модель** (*pooled model*). Если  $z_i$  состоит только из константы, то обычный метод наименьших квадратов дает состоятельные и эффективные оценки для общего коэффициента  $\alpha$  и вектора коэффициентов  $\beta$ .
2. **Модель с фиксированными эффектами** (*fixed effects model*). Если  $z_i$  являются ненаблюдаемыми переменными, которые коррелированы с  $x_{it}$ , то вследствие неправильной спецификации модели, связанной с пропущенными переменными, МНК-оценки  $\beta$ -коэффициентов будут смещенными и несостоятельными. Модель с фиксированными эффектами выглядит следующим образом:

$$y_{it} = x'_{it}\beta + \alpha_i + \varepsilon_{it}, \quad (9.4)$$

где  $\alpha_i = z'_i\alpha$ .

Получается, что неизвестные параметры  $\alpha_i$  для каждой отдельной единицы совокупности принимают некоторое постоянное значение, не изменяются с течением времени  $t$  и тем самым характеризуют так называемые фиксированные эффекты.

3. **Модель со случайными эффектами** (*random effects model*). Если ненаблюдаемые индивидуальные характеристики, формирующие неоднородность единиц совокупности, не коррелируют с включенными в модель объясняющими переменными  $x_{it}$ , то модель панельных данных может быть представлена в виде

$$y_{it} = x'_{it}\beta + E(z'_i\alpha) + [z'_i\alpha - E(z'_i\alpha)] + \varepsilon_{it} = x'_{it}\beta + \alpha + u_i + \varepsilon_{it}, \quad (9.5)$$

где  $E(z'_i\alpha)$  — математическое ожидание для  $z'_i\alpha$ .

Уравнение (9.5) представляет собой линейную регрессионную модель с составной ошибкой, при оценивании которой методом наименьших квадратов получатся состоятельные, но неэффективные оценки параметров.

### 9.2.2. Объединенная модель

Объединенная модель представляет обычную линейную модель регрессии и описывается уравнением

$$y_{it} = x'_{it}\beta + \alpha + \varepsilon_{it}. \quad (9.6)$$

При этом ошибки  $\varepsilon_{it} \sim iid(0, \sigma_\varepsilon^2)$ , т.е. ошибки  $\varepsilon_{it}$  являются независимыми, одинаково распределенными величинами (*iid* — *independent identically distributed*) с математическим ожиданием  $E(\varepsilon_{it}) = 0$  и дисперсией  $V(\varepsilon_{it}) = \sigma_\varepsilon^2$ . Также ошибки  $\varepsilon_{it}$  не коррелируют между собой по  $i$  и по  $t$  и не коррелируют со всеми независимыми переменными  $x_{it}$ .

Объединенная модель получается, если собрать все данные по  $i$  и по  $t$  в одну большую регрессию с  $NT$  наблюдениями. Объединенная модель фактически не учитывает панельную структуру данных и не позволяет реализовать потенциал панельных данных, в частности индивидуальные характеристики отдельных изучаемых единиц.

### 9.2.3. Модель с фиксированными эффектами

Модель с индивидуальными фиксированными эффектами имеет вид

$$y_{it} = x'_{it}\beta + \alpha_i + \varepsilon_{it}, \quad \varepsilon_{it} \sim iid(0, \sigma_\varepsilon^2). \quad (9.7)$$

В модели (9.6) индивидуальные эффекты  $\alpha_1, \alpha_2, \dots, \alpha_n$  измеряют ненаблюдаемую неоднородность,  $x_{it}$  и  $\beta$  представляют собой матрицы размерностью  $k \times 1$  и ошибки  $\varepsilon_{it}$  являются независимыми одинаково распределенными величинами с математическим ожиданием  $E(\varepsilon_{it}) = 0$  и дисперсией  $V(\varepsilon_{it}) = \sigma_\varepsilon^2$ .

Термин «фиксированные эффекты» означает, что константа в уравнении регрессии может различаться между индивидами, но для каждого конкретного индивида константа является постоянной во времени, т.е. не изменяется с течением времени  $t$ .

Для линейной модели с фиксированными эффектами существует несколько способов, для того чтобы получить состоятельные оценки  $\beta$ -коэффициентов, несмотря на присутствие мешающих параметров  $\alpha_i$ :

- оценивание обычным методом наименьших квадратов во внутригрупповой модели;
- прямое оценивание обычным методом наименьших квадратов модели с фиктивными переменными для каждого из  $N$  фиксированных эффектов;
- оценивание условным методом максимального правдоподобия в модели с фиксированными эффектами;
- оценивание обычным методом наименьших квадратов в модели первой разности.

#### 9.2.3.1. Внутригрупповые оценки, или оценки с фиксированными эффектами

Рассмотрим модель с фиксированными эффектами (9.7). Если из уравнения регрессии (9.7) вычесть уравнение, записанное в средних по времени значениях для каждого индивида

$$\bar{y}_i = \alpha_i + \bar{x}'_i\beta + \bar{\varepsilon}_i,$$

где  $\bar{y}_i = 1/T \sum_t y_{it}$ ;  $\bar{x}_i = 1/T \sum_t x_{it}$ ;  $\bar{\varepsilon}_i = 1/T \sum_t \varepsilon_{it}$ ,

то получим *внутригрупповую модель* (*within-group model*), построенную по отклонениям от групповых средних,

$$y_{it} - \bar{y}_i = (x_{it} - \bar{x}_i)\beta + (\varepsilon_{it} - \bar{\varepsilon}_i). \quad (9.8)$$

В отличие от первоначальной модели (9.7) внутригрупповая модель (9.8) уже не содержит фиксированные эффекты  $\alpha_i$ .

*Межгрупповая модель* (*between-group model*) представляет собой регрессию  $\bar{y}_i = \bar{x}_i\beta + \alpha_i + \bar{\varepsilon}_i$ , которая может быть переписана в виде

$$\bar{y}_i = \bar{x}_i'\beta + \alpha + (\alpha_i - \alpha + \bar{\varepsilon}_i), \quad (9.9)$$

где  $(\alpha_i - \alpha + \bar{\varepsilon}_i)$  — составные ошибки.

В межгрупповой модели используется вариация между различными индивидами, и межгрупповая модель является аналогом пространственной регрессии, которая представляет собой частный случай для  $T = 1$ .

Уравнения (9.7)—(9.9) — классические регрессионные модели, и их можно оценивать обычным методом наименьших квадратов, при помощи которого можно получить если не эффективные, то по крайней мере состоятельные оценки коэффициентов. Рассмотрим МНК-оценки  $\beta$ -коэффициентов для каждого из этих трех уравнений и определим, как они взаимосвязаны. Для нахождения оценок коэффициентов нам понадобятся суммы квадратов и смешанных моментов.

В уравнении (9.7) моменты отражают вариацию от общих средних  $\bar{y}$  и  $\bar{x}$ , суммы квадратов  $s_{xx}^{\text{общ}}$  и смешанных моментов  $s_{xy}^{\text{общ}}$  будут иметь вид

$$s_{xx}^{\text{общ}} = \sum_{i=1}^n \sum_{t=1}^T (x_{it} - \bar{x})(x_{it} - \bar{x}); \quad (9.10)$$

$$s_{xy}^{\text{общ}} = \sum_{i=1}^n \sum_{t=1}^T (x_{it}\bar{x})(y_{it} - \bar{y}). \quad (9.11)$$

В уравнении (9.8) данные представлены в виде отклонений от групповых средних, и суммы квадратов  $s_{xx}^{\text{внутри}}$  и смешанных моментов  $s_{xy}^{\text{внутри}}$  будут отражать внутригрупповые суммы квадратов

$$s_{xx}^{\text{внутр}} = \sum_{i=1}^n \sum_{t=1}^T (x_{it} - \bar{x}_i)(x_{it} - \bar{x}_i)'; \quad (9.12)$$

$$s_{xy}^{\text{внутр}} = \sum_{i=1}^n \sum_{t=1}^T (x_{it} - \bar{x}_i)(y_{it} - \bar{y}_i). \quad (9.13)$$

Для уравнения (9.9) суммы квадратов  $s_{xx}^{\text{меж}}$  и смешанных моментов  $s_{xy}^{\text{меж}}$  будут отражать межгрупповые суммы квадратов

$$s_{xx}^{\text{меж}} = \sum_{i=1}^n T(\bar{x}_i - \bar{x})(\bar{x}_i - \bar{x}); \quad (9.14)$$

$$s_{xy}^{\text{меж}} = \sum_{i=1}^n T(\bar{x}_i - \bar{x})(\bar{y}_i - \bar{y}). \quad (9.15)$$

При этом выполняются равенства  $s_{xx}^{\text{общ}} = s_{xx}^{\text{внутр}} + s_{xx}^{\text{меж}}$  и  $s_{xy}^{\text{общ}} = s_{xy}^{\text{внутр}} + s_{xy}^{\text{меж}}$ .

МНК-оценки  $\beta$ -коэффициентов для уравнения регрессии (9.7) будут определяться следующим образом:

$$\hat{\beta}^{\text{бщ}} = [s_{xx}^{\text{бщ}}]^{-1} s_{xy}^{\text{бщ}} = [s_{xx}^{\text{внутр}} + s_{xx}^{\text{меж}}]^{-1} [s_{xy}^{\text{внутр}} + s_{xy}^{\text{меж}}]. \quad (9.16)$$

Оценивая уравнение регрессии (9.8) для отклонений от групповых средних, получим МНК-оценки  $\beta$ -коэффициентов, которые называют **внутригрупповыми оценками**, или *within-groups estimators*):

$$\hat{\beta}^{\text{внутр}} = [s_{xx}^{\text{внутр}}]^{-1} s_{xy}^{\text{внутр}}. \quad (9.17)$$

Оценивая уравнение (9.9), получим межгрупповые МНК-оценки  $\beta$ -коэффициентов (*between-groups estimators*), которые также называют **оценками групповых средних**, или *between-оценками*:

$$\hat{\beta}^{\text{меж}} = [s_{xx}^{\text{меж}}]^{-1} s_{xy}^{\text{меж}}. \quad (9.18)$$

Из уравнений (9.17) и (9.18) следует, что  $s_{xx}^{\text{внутр}} \hat{\beta}^{\text{внутр}}$  и  $s_{xy}^{\text{меж}} = s_{xx}^{\text{меж}} \hat{\beta}^{\text{меж}}$ . Тогда можно показать, что

$$\hat{\beta}^{\text{бщ}} = ([s_{xx}^{\text{внутр}} + s_{xx}^{\text{меж}}]^{-1} s_{xx}^{\text{внутр}}) \hat{\beta}^{\text{внутр}} + ([s_{xx}^{\text{внутр}} + s_{xx}^{\text{меж}}]^{-1} s_{xx}^{\text{меж}}) \hat{\beta}^{\text{меж}}. \quad (9.19)$$

Получается, что общая МНК-оценка  $\beta$ -коэффициентов представляет собой средневзвешенное значение внутригрупповых и межгрупповых МНК-оценок  $\beta$ -коэффициентов.

Межгрупповые оценки  $\beta$ -коэффициентов будут состоятельными в том случае, если регрессоры  $\bar{x}_i$  не коррелируют с составными ошибками  $(\alpha_i - \alpha + \bar{\varepsilon}_i)$ . Однако в модели с фиксированными эффектами предполагается, что  $\alpha_i$  коррелируют с  $x_{it}$ , и следовательно, коррелируют с  $\bar{x}_i$ .

Внутригрупповая оценка — это состоятельная оценка  $\beta$ -коэффициентов в модели с фиксированными эффектами в отличие от оценки объединенной модели и межгрупповой оценки, которые являются несостоятельными. Внутригрупповая оценка измеряет связь между индивидуальными отклонениями регрессоров от их средних по времени значений и индивидуальными отклонениями зависимой переменной от ее среднего по времени значения.

Внутригрупповую оценку также называют **оценкой с фиксированными эффектами** ( $\hat{\beta} = \hat{\beta}^{\text{внутр}} = \hat{\beta}_{FE}$ ), так как она является эффективной оценкой  $\beta$ -коэффициентов в модели (9.7) в том случае, если  $\alpha_i$  представляют собой фиксированные эффекты и ошибки  $\varepsilon_{it}$  являются независимыми, одинаково распределенными величинами с нулевым математическим ожиданием и дисперсией  $\sigma_{\varepsilon}^2$ . Внутригрупповые оценки совпадают с МНК-оценками с фиктивными переменными (*LSDV*-оценками).

Главное ограничение внутригрупповых оценок состоит в том, что во внутригрупповой модели невозможно оценить коэффициенты при регрессорах, которые остаются постоянными для  $i$ -го индивида, т.е. не изменяются с течением времени  $t$ . Действительно,  $x_{it} - \bar{x}_i = 0$  в том случае, если  $x_{it} = x_i$ , для всех  $i$ . Однако многие исследователи заинтересованы в оценивании воздействия регрессоров, постоянных во времени, для каждого индивида. В качестве таких регрессоров, например, могут выступать пол, раса, вероисповедание индивида. Для этих целей исследователи предпочитают не использовать внутригрупповые оценки. Оценки объединенной модели или оценки со случайными эффектами позволяют оценить коэффициенты при объясняющих переменных, которые не меняются во времени и постоянны для каждого индивида, но если корректной моделью является модель с фиксированными эффектами, то эти оценки будут несостоятельными.

Из уравнения (9.12), (9.13) и (9.17) следует, что внутригрупповые МНК-оценки  $\beta$ -коэффициентов определяются по формуле

$$\hat{\beta}^{\text{внутри}} = \left( \sum_{i=1}^N \sum_{t=1}^T (x_{it} - \bar{x}_i)(x_{it} - \bar{x}_i)' \right)^{-1} \sum_{i=1}^N \sum_{t=1}^T (x_{it} - \bar{x}_i)(y_{it} - \bar{y}_i). \quad (9.20)$$

После того, как оценки  $\hat{\beta}^{\text{внутри}}$  стали известными, могут быть оценены индивидуальные эффекты  $\alpha_i$ :

$$\hat{\alpha}_i = \bar{y}_i - \bar{x}_i' \hat{\beta}^{\text{внутри}}, \quad i = 1, \dots, N. \quad (9.21)$$

Найденные таким образом оценки  $\hat{\alpha}_i$  будут несмещенными оценками индивидуальных эффектов  $\alpha_i$ . Состоятельность оценки  $\hat{\alpha}_i$  обеспечивается условием  $T \rightarrow \infty$ , так как  $\hat{\alpha}_i$  выступает в качестве среднего значения для  $T$  наблюдений. В коротких панелях оценки  $\hat{\alpha}_i$  являются несостоятельными, однако оценки  $\hat{\beta}^{\text{внутри}}$  будут состоятельными. Коэффициенты  $\alpha_i$  можно рассматривать как мешающие параметры или вспомогательные параметры, и нет необходимости в том, чтобы они были состоятельно оценены для получения состоятельных оценок наиболее важных для нас параметров  $\beta$ . Однако этот результат не стоит переносить на более сложные модели с фиксированными эффектами, например такие, как нелинейные модели.

Внутригрупповые оценки  $\hat{\beta}^{\text{внутри}}$  являются состоятельными, если

$$\text{plim} \frac{1}{NT} \sum_i \sum_t (x_{it} - \bar{x}_i)(\varepsilon_{it} - \bar{\varepsilon}_i) = 0. \quad (9.22)$$

Этот предел по вероятности может быть равным нулю в том случае, если или  $T \rightarrow \infty$ , или  $N \rightarrow \infty$  и  $E[(\varepsilon_{it} - \bar{\varepsilon}_i) | (x_{it} - \bar{x}_i)] = 0$ . (9.23)

Принимая во внимание то, что  $\bar{x}_i$  и  $\bar{\varepsilon}_i$  являются средними величинами, можно показать, что достаточным условием для уравнения (9.23) будет выполнение строгого экзогенного условия, состоящего в том, что  $E[\varepsilon_{it} | x_{i1}, x_{i2}, \dots, x_{iT}] = 0$ .

Так как во внутригрупповой модели ошибки  $(\varepsilon_{it} - \bar{\varepsilon}_i)$  коррелируют по  $t$  для данных  $i$ , то распределение  $\hat{\beta}^{\text{внутри}}$  получается потенциально сложным, но, несмотря на это, можно применять результаты обычного метода наименьших квадратов. При строгом предположении о том, что  $\varepsilon_{it}$  являются независимыми одинаково распределенными величинами с нулевым

математическим ожиданием и дисперсией  $\sigma_\varepsilon^2$ , дисперсия оценок  $\hat{\beta}^{\text{внутри}}$  будет иметь вид

$$V[\hat{\beta}^{\text{внутри}}] = \sigma_\varepsilon^2 \left[ \sum_{i=1}^N \sum_{t=1}^T (x_{it} - \bar{x}_i)(x_{it} - \bar{x}_i)' \right]^{-1}. \quad (9.24)$$

Для дисперсии ошибок  $\sigma_\varepsilon^2$  состоятельной и несмещенной оценкой будет оценка

$$\hat{\sigma}_\varepsilon^2 = \frac{1}{N(T-1) - k} \sum_i \sum_t \hat{\varepsilon}_{it}^2, \quad (9.25)$$

где число степеней свободы равно размеру выборки  $NT$  без числа оцененных параметров  $k$  и индивидуальных эффектов  $N$ .

Для коротких панелей достигается робастная оценка асимптотической дисперсии

$$V[\hat{\beta}^{\text{внутри}}] = s_{xx}^{-1} \left[ \sum_{i=1}^N \sum_{t=1}^T \sum_{s=1}^T (x_{it} - \bar{x}_i)(x_{it} - \bar{x}_i)' (\varepsilon_{it} - \bar{\varepsilon}_i)(\varepsilon_{is} - \bar{\varepsilon}_i)' \right] s_{xx}^{-1}. \quad (9.26)$$

Эта предпочтительная оценка допускает произвольную автокорреляцию для ошибок  $\varepsilon_{it}$  и произвольную гетероскедастичность.

Если  $N$  является слишком большим, то альтернативным и более простым способом вычислить внутригрупповые оценки будет способ посредством МНК-оценок с фиктивными переменными.

### 9.2.3.2. МНК-оценки с фиктивными переменными

Рассмотрим модель с фиксированными эффектами  $y_{it} = x_{it}'\beta + \alpha_i + \varepsilon_{it}$ . Для  $i$ -го индивида можно записать следующее уравнение в матричной форме:

$$y_i = \mathbf{i}_T \alpha_i + \mathbf{X}_i \beta + \boldsymbol{\varepsilon}_i, \quad (9.27)$$

где  $\mathbf{i}_T$  — единичный вектор размерностью  $T$ .

Если объединить данные по всем индивидам, то получим уравнение

$$\begin{array}{c} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \\ \mathbf{y} \\ (NT \times 1) \end{array} = \begin{array}{c} \begin{bmatrix} \mathbf{i}_T & 0 & \dots & 0 \\ 0 & \mathbf{i}_T & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \mathbf{i}_T \end{bmatrix} \\ \mathbf{D}_N \\ (NT \times N) \end{array} \begin{array}{c} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_N \end{bmatrix} \\ \boldsymbol{\alpha} \\ (N \times 1) \end{array} + \begin{array}{c} \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \\ \vdots \\ \mathbf{X}_N \end{bmatrix} \\ \mathbf{X} \\ (NT \times K) \end{array} \beta + \begin{array}{c} \begin{bmatrix} \boldsymbol{\varepsilon}_1 \\ \boldsymbol{\varepsilon}_2 \\ \vdots \\ \boldsymbol{\varepsilon}_N \end{bmatrix} \\ \boldsymbol{\varepsilon} \\ (NT \times 1) \end{array}. \quad (9.28)$$



Уравнение (9.28) можно также переписать в более простой форме

$$y = D_N \alpha + X\beta + \varepsilon, \quad (9.29)$$

где матрица  $D_N$  состоит из  $N$  индивидуальных фиктивных переменных, и ее можно представить через произведение Кронекера<sup>1</sup>  $D_N = I_N \otimes i_T$ , где  $I_N$  — единичная матрица порядка  $N$ , в которой все диагональные элементы равны единице, а все недиагональные элементы равны нулю.

Можно убедиться, что матрица  $D_N$  обладает следующими свойствами.

1.  $D_N i_T = i_N \otimes i_T = i_{NT}$ .
2.  $D'_N D_N = T I_N$ .
3.  $D_N D'_N = I_N \otimes i_T i'_T = I_N \otimes J_T$ , где  $J_T = i_T i'_T$  — матрица порядка  $T$ , состоящая из единиц.
4.  $\frac{1}{T} D'_N y = (\bar{y}_1, \dots, \bar{y}_N)'$ ,  $\frac{1}{T} D'_N X = (\bar{x}_1, \dots, \bar{x}_N)'$ , где  $\bar{y}_i = \frac{1}{T} \sum_{t=1}^T y_{it}$  и  $\bar{x}'_i = \frac{1}{T} \sum_{t=1}^T x'_{it}$  — соответствующие средние значения переменных для  $i$ -го индивида.

В этой модели помимо того, что ошибки являются независимыми одинаково распределенными величинами с нулевым математическим ожиданием и постоянной дисперсией, делается предположение о том, что объясняющие переменные не зависят от ошибок. Пусть  $Z = [D_N X]$  — матрица размером  $NT \times (N + K)$ , которая содержит все объясняющие переменные, включая  $N$  фиктивных переменных, и эта матрица должна иметь полный столбцовый ранг. Получается, что  $NT$  должно быть больше, чем  $(N + K)$ , а это выполняется в случае с большим  $N$  и любым  $T \geq 2$ . Кроме того, столбцы  $X$  должны быть линейно независимы от  $D_N$ . Это соответствует случаю, когда в матрице  $X_i$  нет константы, а также какого-либо другого столбца, пропорционального константе, что в итоге

<sup>1</sup> Произведением Кронекера матриц  $A_{m \times n}$  и  $B_{p \times q}$  называется матрица  $A \otimes B$

$$\text{размером } mp \times nq, \text{ имеющая вид } A \otimes B = \begin{pmatrix} a_{11}B & a_{12}B & \dots & a_{1n}B \\ a_{21}B & a_{22}B & \dots & a_{2n}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}B & a_{m2}B & \dots & a_{mn}B \end{pmatrix}.$$

исключает из регрессии любые переменные, которые являются постоянными во времени для конкретного индивида, но могут различаться между индивидами.

При выполнении описанных выше предположений МНК-оценки всех регрессионных коэффициентов в модели  $y = D_N \alpha + X\beta + \varepsilon$  будут лучшими линейными несмещенными оценками. Оценивание этой модели методом наименьших квадратов приносит МНК-оценки с фиктивными переменными, или *LSDV*-оценки (*least squares dummy variables estimators*). Соберем все регрессионные коэффициенты в вектор  $\gamma$ ,  $\gamma' = [\alpha' \beta']$ , тогда МНК-оценки будут иметь вид  $\hat{\gamma} = (Z'Z)^{-1}Z'y$ . Для нахождения оценки  $\hat{\gamma}$  необходима обратимость матрицы  $Z'Z$  размером  $(N + K) \times (N + K)$ . Однако при большом количестве коэффициентов получается слишком много регрессоров, и это вызывает определенные трудности, связанные с обратимостью матрицы регрессоров размерностью  $(N + K) \times (N + K)$ . Применяв преобразования матричной алгебры, можно сократить задачу до обратимости матрицы  $K \times K$ . Оценка  $\beta$ -коэффициентов в этой модели оказывается равной внутригрупповой оценке. Это случай так называемой **теоремы Фриша — Во** (*Frisch — Waugh*) для регрессии подмножества. В результате МНК-оценки с фиктивными переменными определяются по формулам

$$\hat{\beta}_{LSDV} = (X'W_N X)^{-1} X'W_N y; \quad (9.30)$$

$$\hat{\alpha}_{LSDV} = (D'_N D_N)^{-1} D'_N (y - X\hat{\beta}) = \frac{1}{T} D'_N (y - X\hat{\beta}), \quad (9.31)$$

где  $W_N = I_{NT} - D_N (D'_N D_N)^{-1} D'_N = I_{NT} - \frac{1}{T} D_N D'_N = I_{NT} - I_N \otimes \frac{1}{T} J_T$  является идемпотентной<sup>1</sup> матрицей порядка  $NT$ , которая имеет ранг  $N(T - 1)$ . Матрица  $W_N$  также называется **within-проектором**, так как проектор  $W_N$  позволяет осуществить внутригрупповое преобразование и вычислить отклонения от индивидуальных средних значений  $y_{it}^* = y_{it} - \bar{y}_i$ .

Сумма квадратов ошибок и дисперсия ошибок в модели с фиксированными эффектами будут соответственно иметь вид

$$SS = y'W_N y - \hat{\beta}'X'W_N y; \quad (9.32)$$

<sup>1</sup> матрица  $W_N$  будет идемпотентной в том случае, когда выполняется условие  $W_N W'_N = W_N$ .

$$\hat{\sigma}^2 = \frac{SS}{NT - N - K}. \quad (9.33)$$

Дисперсии оценок коэффициентов  $\beta$  и  $\alpha$  в модели с фиксированными эффектами будут соответственно определяться по формулам

$$V(\hat{\beta}) = \sigma^2 (X'W_N X)^{-1}; \quad (9.34)$$

$$V(\hat{\alpha}) = \frac{\sigma^2}{T} I_N + \frac{1}{T} D'_N X; \quad V(\hat{\beta}) X D_N \frac{1}{T}. \quad (9.35)$$

Оценки, заданные уравнениями (9.30) и (9.32), также могут быть получены методом наименьших квадратов из трансформированной модели

$$y^* = X^* \beta + \varepsilon^*, \quad (9.36)$$

где  $y^* = W_N y$  и  $X^* = W_N X$ .

Преобразование  $W_N$  является очень простым: трансформированные переменные — это просто первоначальные переменные, выраженные как отклонения от индивидуальных средних. Поэтому  $i$ -е уравнение, соответствующее регрессии (9.36), будет иметь вид

$$(y_{it} - \bar{y}_i) = \beta_1(x_{1it} - \bar{x}_{1i}) + \dots + \beta_K(x_{Kit} - \bar{x}_{Ki}) + \varepsilon_{it}^* = (x_{it} - \bar{x}_i)' \beta + \varepsilon_{it}^*. \quad (9.37)$$

Уравнение регрессии (9.37) представляет собой не что иное, как внутригрупповую модель (9.8). Выходит, что МНК-оценки с фиктивными переменными совпадают с внутригрупповыми оценками  $\alpha_i$  и  $\beta$ -коэффициентов или оценками с фиксированными эффектами

$$\hat{\beta}_{LSDV} = \hat{\beta}^{BHTP} = \left[ \sum_{i=1}^N \sum_{t=1}^T (x_{it} - \bar{x}_i)(x_{it} - \bar{x}_i)' \right]^{-1} \sum_{i=1}^N \sum_{t=1}^T (x_{it} - \bar{x}_i)(y_{it} - \bar{y}_i); \quad (9.38)$$

$$\hat{\alpha}_{LSDV} = \hat{\alpha}_{i, FE} = \bar{y}_i - \bar{x}_i' \hat{\beta}_{LSDV}, \quad i = 1, \dots, N. \quad (9.39)$$

Дисперсия оценки  $\hat{\alpha}_i$  в преобразованной модели может быть получена из уравнения (9.35) и будет иметь вид

$$V(\hat{\alpha}_i) = \frac{1}{T} \sigma^2 + \bar{x}_i' V(\hat{\beta}) \bar{x}_i. \quad (9.40)$$

При работе с преобразованными переменными необходимо помнить, что истинное число степеней свободы будет не  $(NT - K)$ , а  $(NT - N - K)$ , так как при преобразовании вычисляются  $N$  индивидуальных средних и в результате теряется  $N$  степеней свободы. Следовательно, необходимо обращать внимание на дисперсию, вычисляемую при применении компьютерной программой к преобразованным данным, и корректировать ее соответствующим образом.

Для коротких панелей очевидна потенциальная проблема, состоящая в том, что не гарантирована состоятельность оценок  $\beta$  и  $\alpha$ , если имеется  $(N + k)$  параметров, которые необходимо оценить, и  $N$  стремится к бесконечности. Тем не менее, состоятельные оценки  $\beta$ -коэффициентов возможны даже в том случае, если индивидуальные коэффициенты оценены несостоятельно и пока вдобавок  $T$  не будет стремиться к бесконечности.

### 9.2.3.3. Оценка условным методом максимального правдоподобия

Оценка условным методом максимального правдоподобия, или условная *ML*-оценка (*conditional ML estimator*), максимизирует совместное правдоподобие  $y_{11}, \dots, y_{nT}$ , условное по индивидуальным средним  $\bar{y}_1, \dots, \bar{y}_T$ . Этот метод привлекателен тем, что при выполнении условия нормальности, в линейных панельных моделях исключаются фиксированные эффекты  $\alpha_i$ , и в результате происходит максимизация функции условного правдоподобия только относительно параметра  $\beta$ .

Предположим, что  $y_{it}$ , условное по регрессорам  $x_{it}$ , и параметры  $\alpha_i$ ,  $\beta$  и  $\sigma^2$  являются одинаково распределенными величинами с нормальным распределением  $N[\alpha_i + x'_{it}\beta, \sigma^2]$ . Тогда функция условного правдоподобия будет иметь вид

$$L_{cond}(\beta, \sigma^2, \alpha) = \prod_{i=1}^N f(y_{i1}, \dots, y_{iT} | \bar{y}_i) = \prod_{i=1}^N \frac{f(y_{i1}, \dots, y_{iT}, \bar{y}_i)}{f(\bar{y}_i)} = \prod_{i=1}^N \frac{(2\pi\sigma^2)^{-T/2}}{(2\pi\sigma^2 / T)^{-1/2}} \exp \left[ \sum_{t=1}^T -\frac{(y_{it} - x'_{it}\beta)^2 + (\bar{y}_i - \bar{x}'_i\beta)}{2\sigma^2} \right]. \quad (9.41)$$

Первое равенство определяет условное правдоподобие при предположении о независимости по  $i$ . Второе

равенство всегда соблюдается, так как  $f(y_1, \dots, y_T | \bar{y}) = f(y_1, \dots, y_T, \bar{y}) / f(\bar{y})$  и  $f(y_1, \dots, y_T, \bar{y}) = f(y_1, \dots, y_T)$ , поскольку  $\bar{y}$  — это функция от  $y_1, \dots, y_T$ . Третье равенство получается после осуществления алгебраических преобразований при предположении о нормальном распределении. В последнем равенстве уже отсутствуют фиксированные эффекты  $\alpha$ , функция условного правдоподобия  $L_{cond}(\beta, \sigma^2, \alpha)$  на самом деле является функцией  $L_{cond}(\beta, \sigma^2)$ , и логарифм функции правдоподобия будет максимизироваться только по отношению к  $\beta$  и  $\sigma^2$ .

Условная  $ML$ -оценка  $\hat{\beta}_{CML}$  определяется на основе условия первого порядка

$$\frac{1}{\sigma^2} \sum_{t=1}^T \sum_{i=1}^N [(y_{it} - x'_{it}\beta)x_{it} - (\bar{y}_i - \bar{x}'_i\beta)\bar{x}_i] = 0. \quad (9.42)$$

Если произвести перегруппировку данных, то получим уравнение

$$\sum_{t=1}^T \sum_{i=1}^N [(y_{it} - \bar{y}_i) - (x_{it} - \bar{x}_i)'\beta] (x_{it} - \bar{x}_i) = 0. \quad (9.43)$$

Уравнение (9.43) есть не что иное, как условие первого порядка для МНК регрессии  $(y_{it} - \bar{y}_i)$  на  $(x_{it} - \bar{x}_i)$ , и поэтому условная  $ML$ -оценка  $\hat{\beta}_{CML}$  равна внутригрупповой оценке  $\hat{\beta}^{внутр}$ . Условная  $ML$ -оценка  $\hat{\beta}_{CML}$  является состоятельной оценкой  $\beta$ -коэффициентов в модели с фиксированными эффектами (9.7).

#### 9.2.3.4. МНК-оценки модели первых разностей

Если из модели с фиксированными эффектами (9.7) вычесть модель с фиксированными эффектами с лагом на один период  $y_{i,t-1} = x'_{i,t-1}\beta + \alpha_i + \varepsilon_{i,t-1}$ , то получим модель первых разностей (*first-differences model*)

$$y_{it} - y_{i,t-1} = (x_{it} - x_{i,t-1})'\beta + (\varepsilon_{it} - \varepsilon_{i,t-1}); \quad i=1, \dots, N, \quad t=2, \dots, T. \quad (9.44)$$

Модель первых разностей в отличие первоначальной модели (9.7) уже не содержит фиксированные эффекты  $\alpha_i$ . Следует обратить внимание на то, что в уравнении (9.44) используется только  $N(T-1)$  наблюдений.

Оценка первых разностей (*first-differences estimator*) представляет собой МНК-оценку  $\beta$ -коэффициентов в модели (9.44):

$$\hat{\beta}_{FD} = \left( \sum_{i=1}^n \sum_{t=2}^T (x_{it} - \bar{x}_{i,t-1})(x_{it} - \bar{x}_{i,t-1})' \right)^{-1} \sum_{i=1}^n \sum_{t=2}^T (x_{it} - \bar{x}_{i,t-1})(y_{it} - \bar{y}_{i,t-1}). \quad (9.45)$$

Оценка первых разностей использует специальные характеристики панельных данных. В короткой панели оценка первых разностей измеряет связь между индивидуальными изменениями в регрессорах за один период времени и индивидуальными изменениями в зависимой переменной за тот же период времени. Оценка первых разностей, так же как и внутригрупповая оценка, не позволяет оценить коэффициенты постоянных во времени регрессоров для каждого индивида в модели с фиксированными эффектами. Для  $T = 2$  оценка первых разностей совпадает с внутригрупповой оценкой, так как  $\bar{y} = (y_1 + y_2) / 2$ , то  $(y_1 - \bar{y}) = (y_1 - y_2) / 2$  и  $(y_2 - \bar{y}) = -(y_1 - y_2) / 2$ , и аналогично для  $x$ , но для  $T > 2$  эти оценки не совпадают. При предположении о том, что ошибки  $\varepsilon_{it}$  являются независимыми одинаково распределенными величинами, можно показать, что оценка обобщенным методом наименьших квадратов модели первых разностей (9.44) равна внутригрупповой оценке.

Оценка первых разностей является состоятельной оценкой  $\beta$ -коэффициентов в модели с фиксированными эффектами. Оценка первых разностей будет состоятельной, если  $E(\varepsilon_{it} - \varepsilon_{i,t-1} | x_{it} - x_{i,t-1}) = 0$ . Это условие является более строгим, чем  $E(\varepsilon_{it} | x_{it}) = 0$ , и менее слабым, чем строгое экзогенное условие, необходимое для состоятельности внутригрупповой оценки. Однако оценка первых разностей менее эффективна, чем внутригрупповая оценка для случая  $T > 2$  при условии, что ошибки  $\varepsilon_{it}$  являются независимыми, одинаково распределенными величинами с нулевым математическим ожиданием и дисперсией  $\sigma_\varepsilon^2$ .

#### 9.2.4. Модель со случайными эффектами

В модели с фиксированными эффектами ненаблюдаемые индивидуальные эффекты были коррелированы с включенными в модель объясняющими переменными. Однако иногда есть основания полагать, что индивидуальные эффекты могут быть не коррелированы с регрессорами. В этом случае

индивидуальные эффекты можно рассматривать как одну из составляющих ошибки.

**Модель со случайными эффектами** (иногда ее называют **моделью со случайными компонентами**) имеет вид

$$y_{it} = x'_{it}\beta + \alpha + u_i + \varepsilon_{it}, \quad (9.46)$$

где  $u_i + \varepsilon_{it}$  можно рассматривать как составную ошибку регрессии, которая содержит две компоненты: индивидуальную компоненту  $u_i$  и остаточный член  $\varepsilon_{it}$ . Компонента  $u_i$  представляет индивидуальную ошибку (случайное отклонение), которая является постоянной во времени для  $i$ -го объекта. В модели со случайными эффектами делаются следующие предположения относительно случайных компонент  $u_i$  и  $\varepsilon_{it}$ .

- 1)  $E(\varepsilon_{it} | u_i) = 0$ .
- 2)  $E(\varepsilon_{it}^2 | u_i) = \sigma_\varepsilon^2$ .
- 3)  $E(u_i | x_{it}) = 0$  для всех  $i, t$ .
- 4)  $E(u_i^2 | x_{it}) = \sigma_u^2$ .
- 5)  $E(\varepsilon_{it} u_j) = 0$  для всех  $i, t, j$ .
- 6)  $E(\varepsilon_{it} \varepsilon_{js}) = 0$  для всех  $i \neq j, t \neq s$ .
- 7)  $E(u_i u_j) = 0$  для всех  $i \neq j$ .

Рассмотрим характеристики составной ошибки  $\omega_{it} = u_i + \varepsilon_{it}$ . Ошибка  $\omega_{it}$  будет иметь нулевое математическое ожидание, т.е.  $E(\omega_{it}) = 0$ . Тогда будут иметь место следующие равенства:

$$E(\omega_{it}^2) = \sigma_u^2 + \sigma_\varepsilon^2 \text{ для всех } i, t;$$

$$E(\omega_{it} \omega_{is}) = \sigma_u^2 \text{ для всех } t \neq s;$$

$$E(\omega_{it} \omega_{js}) = 0 \text{ для всех } i \neq j, t \neq s.$$

Таким образом, если  $\omega_i = (\omega_{i1}, \omega_{i2}, \dots, \omega_{iT})'$  представляет вектор ошибок  $\omega_{it}$  для  $i$ -го индивида,  $t = 1, \dots, T$ , то  $E(\omega_i \omega_i') = \Omega$ , где

$$\Omega = \sigma_\varepsilon^2 I_T + \sigma_u^2 \mathbf{1}_T \mathbf{1}'_T = \begin{pmatrix} \sigma_\varepsilon^2 + \sigma_u^2 & \sigma_u^2 & \dots & \sigma_u^2 \\ \sigma_u^2 & \sigma_\varepsilon^2 + \sigma_u^2 & \dots & \sigma_u^2 \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_u^2 & \sigma_u^2 & \dots & \sigma_\varepsilon^2 + \sigma_u^2 \end{pmatrix}. \quad (9.47)$$

Для полных панелей наблюдений ковариационная матрица вектора размерностью  $NT$  ошибок  $\omega = (\omega_1, \omega_2, \dots, \omega_N)'$  может быть получена как

$$\mathbf{V}_{(NT \times NT)} = \begin{pmatrix} \Omega & 0 & \dots & 0 \\ 0 & \Omega & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \Omega \end{pmatrix} = \mathbf{I}_N \otimes \mathbf{I}, \quad (9.48)$$

где  $\mathbf{I}_N$  — единичная матрица размерностью  $N$ ;  $\otimes$  — произведение Кронекера.

Получается, что ошибки в модели со случайными эффектами являются гетероскедастичными, и если регрессию  $y_{it} = x'_{it} \beta + \alpha + \omega_{it}$  оценить обычным методом наименьших квадратов, то полученные оценки параметров  $\beta$  и  $\alpha$  будут несмещенными и состоятельными, но неэффективными оценками. Для получения эффективных оценок параметров  $\beta$  и  $\alpha$  необходимо применять обобщенный метод наименьших квадратов (*generalized least squares*).

В случае если известна структура ковариационной матрицы  $\mathbf{V}$ , то уравнение регрессии  $\mathbf{y} = \mathbf{X}\beta + \alpha \mathbf{I}_{NT} + \omega$  можно оценивать при помощи обобщенного метода наименьших квадратов. Для этого необходимо преобразовать исходную регрессию таким образом, чтобы убрать нестандартную структуру ковариационной матрицы  $E(\omega \omega') = \mathbf{V}$ .

Определим взвешенную матрицу  $\mathbf{P} = \mathbf{V}^{-1/2}$  и преобразуем уравнение регрессии

$$\mathbf{P}\mathbf{y} = \mathbf{P}\mathbf{X}\beta + \mathbf{P}\alpha \mathbf{i}_{NT} + \mathbf{P}\omega, \text{ или } \mathbf{y}^* = \mathbf{X}^* \beta + \alpha^* \mathbf{i}_{NT} + \omega^*. \quad (9.49)$$

Теперь обратим внимание на то, что

$$E(\omega^* \omega^{*\prime}) = E(\mathbf{P}\omega \omega' \mathbf{P}) = \mathbf{P}E(\omega \omega')\mathbf{P} = \mathbf{P}\mathbf{V}\mathbf{P} = \mathbf{I}_{NT}. \quad (9.50)$$

Для преобразованной модели остатки будут иметь одинаковую дисперсию по  $i$  и  $T$ . Таким образом, если известно  $\mathbf{P}$ , то оценки регрессии (9.46) обобщенным методом наименьших квадратов будут определяться следующим образом:

$$\begin{pmatrix} \alpha_{GLS} \\ \beta_{GLS} \end{pmatrix} = \left[ \begin{pmatrix} \mathbf{I}_{NT} \\ \mathbf{X} \end{pmatrix} \mathbf{V}^{-1} (\mathbf{I}_{NT} \quad \mathbf{X}) \right]^{-1} \begin{pmatrix} \mathbf{I}_{NT} \\ \mathbf{X} \end{pmatrix} \mathbf{V}^{-1} \mathbf{y}. \quad (9.51)$$



При этом нужно помнить, что должно выполняться предположение о том, что  $E(u_i | x_{it}) = 0$ . Если  $E(u_i | x_{it}) \neq 0$ , то оценки, полученные при помощи обобщенного метода наименьших квадратов, будут несостоятельными.

Для модели со случайными эффектами можно получить специфическую форму весовой матрицы  $P = V^{-1/2}$ . Так как  $V^{-1/2} = I_N \otimes |^{-1/2}$ , то мы можем переписать  $V$  как

$$V = I_N \otimes | = \sigma_u^2 I_N + \sigma_u^2 \mathbf{i}\mathbf{i}', \quad (9.52)$$

где  $\mathbf{i}$  представляет вектор порядка  $N$ , состоящий из единиц.

Это позволяет нам написать форму для  $\Omega^{-1/2}$  в виде

$$\Omega^{-1/2} = \frac{1}{\sigma_\varepsilon} \left( I_T - \frac{\theta}{T} \mathbf{i}_T \mathbf{i}_T' \right), \quad (9.53)$$

где  $\theta = 1 - \frac{\sigma_\varepsilon}{(T\sigma_u^2 + \sigma_\varepsilon^2)^{1/2}}$  и принимает значения между нулем и единицей, а оператор  $\frac{1}{T} \mathbf{i}_T \mathbf{i}_T'$  вычисляет среднее по времени значение для соответствующего вектора.

Таким образом, для получения  $\mathbf{y}^*$  в модели со случайными эффектами необходимо произвести соответствующее преобразование, состоящее в умножении каждого  $\mathbf{y}_i = (y_{i1}, \dots, y_{iT})'$  на  $|^{-1/2}$ :

$$\mathbf{y}^* = |^{-1/2} \mathbf{y}_i = \frac{1}{\sigma_\varepsilon} \begin{bmatrix} y_{i1} - \theta \bar{y}_i \\ y_{i2} - \theta \bar{y}_i \\ \vdots \\ y_{iT} - \theta \bar{y}_i \end{bmatrix}. \quad (9.54)$$

Для определения  $\mathbf{x}_i^*$  производится аналогичное преобразование

$$\mathbf{x}_i^* = \Omega^{-1/2} \mathbf{x}_i = \frac{1}{\sigma_\varepsilon} \begin{bmatrix} x_{i1} - \theta \bar{x}_i \\ x_{i2} - \theta \bar{x}_i \\ \vdots \\ x_{iT} - \theta \bar{x}_i \end{bmatrix}. \quad (9.55)$$

Тогда преобразованное уравнение регрессии будет иметь вид

$$(y_{it} - \theta \bar{y}_i) = (1 - \theta)\alpha + (x_{it} - \theta \bar{x}_i)\beta + v_{it}, \quad (9.56)$$

где  $v_{it} = (1 - \theta)u_i + (\varepsilon_{it} - \theta \bar{\varepsilon}_i)$ .

МНК-оценки параметров  $\alpha$  и  $\beta$  в уравнении (9.56) составят

$$\begin{pmatrix} \hat{\alpha}_{RE} \\ \hat{\beta}_{RE} \end{pmatrix} = \left[ \sum_{i=1}^N \sum_{t=1}^T (q_{it} - \theta \bar{q}_i)(q_{it} - \theta \bar{q}_i)' \right]^{-1} \sum_{i=1}^N \sum_{t=1}^T (q_{it} - \theta \bar{q}_i)(y_{it} - \theta \bar{y}_i), \quad (9.57)$$

где  $q_{it} = (1 \ x_{it})$  и  $\bar{q}_i = (1 \ \bar{x}_i)$ .

Оценка с фиксированными эффектами вычитает средние по времени значения из соответствующих переменных, а преобразование для модели со случайными эффектами вычитает из них долю этих средних по времени значений, где доля зависит от  $\sigma_u$ ,  $\sigma_\varepsilon$  и числа временных периодов  $T$ . Оценка обобщенным методом наименьших квадратов просто представляет собой МНК-оценку объединенной модели (9.56).

В результате преобразования модели (9.56) в уравнении регрессии остаются постоянные во времени объясняющие переменные, и это является преимуществом модели со случайными эффектами по сравнению с моделью с фиксированными эффектами и преобразованием первых разностей.

В случае если параметр  $\theta = 1$ , то модель со случайными эффектами совпадает с моделью с фиксированными эффектами. Если же параметр  $\theta = 0$ , то это означает, что в модели отсутствуют некоррелируемые индивидуальные компоненты дисперсии, т.е.  $\sigma_u^2 = 0$ , и оценка со случайными эффектами сократится до МНК-оценки объединенной модели.

Оценки обобщенного метода наименьших квадратов представляют собой средневзвешенное значение внутригрупповых и межгрупповых оценок:

$$\hat{\beta} = [(s_{xx}^{\text{внутр}} + \lambda s_{xx}^{\text{меж}})^{-1} s_{xx}^{\text{внутр}}] \hat{\beta}^{\text{внутр}} + [(s_{xx}^{\text{внутр}} + \lambda s_{xx}^{\text{меж}})^{-1} \lambda s_{xx}^{\text{меж}}] \cdot \hat{\beta}^{\text{меж}}, \quad (9.58)$$

где  $\lambda = \frac{\sigma_\varepsilon^2}{T\sigma_u^2 + \sigma_\varepsilon^2} = (1 - \theta)^2$ .

Если  $\lambda = 1$ , то  $\sigma_u^2 = 0$  и оценки обобщенного метода наименьших квадратов будут совпадать с оценками обычного метода наименьших квадратов. Если  $\lambda = 0$ , то оценки обобщенного метода наименьших квадратов будут совпадать с МНК-оценками с фиктивными переменными в модели с фиксированными эффектами. Существует две возможности.

Если  $\sigma_\varepsilon^2 = 0$ , тогда источником вариации между индивидами будут выступать различные значения  $u_i$  для  $i$ -го индивида. Так как  $u_i$  являются постоянными во времени, то они будут эквивалентны фиктивному переменным, которые применялись в модели с фиксированными эффектами. Однако вопрос о том, являются ли индивидуальные компоненты  $u_i$  на самом деле фиксированными или случайными, может стать спорным.

Другой случай будет при  $T \rightarrow \infty$ . Если  $T \rightarrow \infty$ , тогда ненаблюдаемые  $u_i$  становятся наблюдаемыми. Рассмотрим, почему это происходит. Так как оценки коэффициентов  $\alpha$  и  $\beta$  в модели со случайными эффектами являются состоятельными, и если взять  $T$  наблюдений для  $i$ -го индивида, то

$$y_{it} - x'_{it}\beta - \alpha = u_i + \varepsilon_{it} \quad (9.59)$$

становятся наблюдаемыми.

Если записать (9.59) в индивидуальных средних, то мы получим  $\bar{y}_{it} - \bar{x}'_{it}\hat{\beta} - \hat{\alpha} = u_i + \bar{\varepsilon}_{it}$ . Однако так как случайная компонента  $\bar{\varepsilon}_{it}$  будет стремиться к нулю при  $T$ , стремящемся к бесконечности, то в результате мы можем определить  $u_i = \bar{y}_{it} - \bar{x}'_{it}\hat{\beta} - \hat{\alpha}$ . Таким образом, при  $T \rightarrow \infty$  индивидуальные компоненты ошибки  $u_i$  ведут себя как фиктивные переменные, которые мы использовали ранее.

В случае если  $\lambda$  отличается от единицы, то в результате неэффективного взвешивания внутригрупповых и межгрупповых оценок получаются неэффективные оценки обычным методом наименьших квадратов, в этом случае межгрупповой сумме квадратов придается слишком большой вес по сравнению с обобщенным методом наименьших квадратов.

На практике параметр  $\theta$  никогда не известен, однако его всегда можно найти, если оценить дисперсии ошибок  $\sigma_u^2$  и  $\sigma_\varepsilon^2$ , т. е. реализовать так называемый **доступный обобщенный метод наименьших квадратов** (*feasible generalized least squares*) в том случае, если  $\Omega$  неизвестно. Оценить дисперсии

ошибок можно, например, на основе объединенной регрессии или регрессии с фиксированными эффектами.

Из регрессии внутригрупповой модели

$$y_{it} - \bar{y}_i = (x_{it} - \bar{x}_i)' \beta + (\varepsilon_{it} - \bar{\varepsilon}_i),$$

которая не содержит значений индивидуальных эффектов, можно получить оценку для дисперсии случайной компоненты ошибки  $\varepsilon_{it}$

$$\hat{\sigma}_\varepsilon^2 = \frac{1}{N(T-1) - K} \sum_i \sum_t [(y_{it} - \bar{y}_i) - (x_{it} - \bar{x}_i)' \hat{\beta}^{\text{внутри}}]^2. \quad (9.60)$$

Из регрессии межгрупповой модели  $\bar{y}_i$  на константу и  $\bar{x}_i$ , где ошибка имеет дисперсию  $(\sigma_u^2 + \sigma_\varepsilon^2 / T)$ , можно получить оценку для дисперсии индивидуальной компоненты ошибки  $u_i$

$$\hat{\sigma}_u^2 = \frac{1}{N - (K + 1)} \sum_i [\bar{y}_i - \hat{\alpha}^{\text{меж}} - \bar{x}_i' \hat{\beta}^{\text{меж}}]^2 - \frac{1}{T} \hat{\sigma}_\varepsilon^2. \quad (9.61)$$

Возможны и более эффективные оценки дисперсий ошибок  $\sigma_u^2$  и  $\sigma_\varepsilon^2$ , однако эти оценки не приведут к увеличению эффективности оценки  $\hat{\beta}_{RE}$ . Следует обратить внимание на то, что оценка дисперсии (9.61) может быть отрицательной. Но необходимо помнить, что для доступного обобщенного метода наименьших квадратов нет необходимости в несмещенных оценках дисперсии, а важна состоятельность этих оценок. Если они являются состоятельными, то можно пренебречь поправками на число степеней свободы, и обе оценки дисперсии получатся неотрицательными.

### Пример

Для объяснения индивидуальной заработной платы использовались данные Российского мониторинга экономического положения и здоровья населения (РМЭЗ). База данных РМЭЗ представляет собой результаты двух панельных опросов нескольких тысяч человек за 1992—1993 гг. (раунды 1—4) и 1994—2008 гг. (раунды 5—17). Часть данных РМЭЗ находится в закрытом доступе, но результаты нескольких раундов можно найти в открытом доступе в Интернете. Информация,

собранная в РМЭЗ, касается размеров, источников, структуры доходов и расходов домохозяйств, занятости, распределения времени, уровня образования, состояния здоровья и других характеристик (всего свыше 500 переменных).

В нашем примере используются данные только трех раундов: 7—9 раундов, проводившихся в 1996—2000 гг., в ходе которых были обследованы одни и те же домашние хозяйства. Численность опрошенных респондентов отличалась по раундам. Так в 7, 8 и 9-м раундах было опрошено 8342, 8699 и 9074 взрослых респондентов соответственно. Из этой совокупности первоначально были выбраны те респонденты, которые принимали участие во всех трех рассматриваемых раундах РМЭЗ (7—9 раунды). Поскольку основное внимание в нашем примере сосредоточено на изучении характеристик занятого населения, то в дальнейшую выборку были включены респонденты, имевшие работу и действительно получавшие в течение последнего месяца заработную плату по основному месту работы. Для того чтобы оценить уравнения заработной платы, из выборки, объединяющей занятых, были исключены те, у кого не было полного набора переменных, необходимых для дальнейшего исследования. Данные 7—9 раундов, объединенные по респондентам, были проверены на то, согласуются ли характеристики респондентов, заявленные ими в предыдущих раундах, с их собственными характеристиками, объявленными в последующих раундах. Респонденты, у которых были выявлены противоречивые данные по основным индивидуальным характеристикам, таким, как возраст, пол, уровень образования, были также исключены из выборки. В результате в итоговой выборке были оставлены 870 респондентов, имеющих работу.

Традиционно нормы отдачи от инвестиций в человеческий капитал рассчитываются на базе стандартного уравнения заработной платы Минцера

$$\ln W_i = \beta_0 + \beta_1 SCH_i + \beta_2 EXP_i + \beta_3 EXP_i^2 + \varepsilon_i.$$

В качестве зависимой переменной выступает логарифм заработной платы  $\ln W$ . Переменная  $SCH$  характеризует уровень образования и обозначает число лет обучения, скорректированных по достигнутому уровню образования. Каждому

уровню образования соответствует среднее число лет обучения: начальное и неполное среднее образование (8 лет и менее), полное среднее образование (10 лет), профессионально-техническое обучение со средним образованием (11,5 лет), среднее специальное (13 лет), законченное высшее (15 лет) и послевузовское (аспирантура) образование (18 лет). Переменная  $EXP$  характеризует опыт работы индивида на рынке труда, измеряемый в годах.

Коэффициент  $\beta_1$  при переменной  $SCH$  представляет оценку нормы отдачи от инвестиций в образование, которая предполагается постоянной в данной модели. Выпуклость наблюдаемых профилей заработной платы выражается через квадратичную форму профессионального опыта. В этом случае коэффициенты  $\beta_2$  и  $\beta_3$  при переменных  $EXP$  и  $EXP^2$  имеют положительный и отрицательный знак соответственно.

Следует отметить, что в 1996—2000 гг. в России, стране с переходной экономикой, заработная плата часто выдавалась в неденежной форме и имел место высокий уровень инфляции. В этих условиях рассматривались заработки индивидов, полученные как в денежной, так и неденежной формах по основному месту работы в течение последних 30 дней перед опросом, скорректированные и приведенные к сопоставимому виду при помощи соответствующих индексов роста цен.

Помимо указанных выше переменных, в модель в качестве объясняющих переменных также было включено несколько фиктивных переменных. Все переменные, включенные в модель индивидуальной заработной платы, представлены в табл. 9.1.

**Таблица 9.1.** Переменные, включенные в модель индивидуальной заработной платы

$LN\_WPH$	Логарифм почасовой заработной платы
$SCH$	Число лет обучения (уровень образования)
$EXP$	Опыт работы индивида на рынке труда
$EXP\_2$	Опыт работы индивида на рынке труда в квадрате
Фиктивные переменные	
$GEN$	Пол (1 — мужской; 0 — женский)
$RUS$	Принадлежность места работы индивида к частному бизнесу (1 — владельцами предприятия являются российские частные фирмы, российские частные лица или коллектив предприятия; 0 — в противном случае)

Окончание табл. 9.1

<i>UNSKILL</i>	Профессиональная квалификация (1 — неквалифицированные рабочие; 0 — в противном случае)
Региональные фиктивные переменные	
<i>NW</i>	Северный и Северо-Западный районы
<i>CENTR</i>	Центральный и Центрально-Черноземный районы
<i>VOL</i>	Волго-Вятский и Поволжский районы
<i>CAUCAS</i>	Северный Кавказ
<i>URAL</i>	Урал
<i>WSIB</i>	Западная Сибирь
<i>FAREAST</i>	Восточная Сибирь и Дальний Восток

Таким образом, в модели присутствуют фиктивные переменные по половому признаку, принадлежности к частному бизнесу, наличию профессиональной квалификации и принадлежности к определенному региону.

Региональные фиктивные переменные принимают значение 1 в том случае, если индивид принадлежит соответствующему региону, и 0 — в противном случае. В данных РМЭЗ использовалось территориальное деление России на следующие восемь районов: города-мегаполисы (Москва и Санкт-Петербург), Северный и Северо-Западный районы, Центральный и Центрально-Черноземный районы, Волго-Вятский и Поволжский районы, Северный Кавказ, Урал, Западная Сибирь, Восточная Сибирь и Дальний Восток. За базу сравнения приняты города-мегаполисы Москва и Санкт-Петербург, и коэффициенты при региональных фиктивных переменных будут показывать, насколько средний логарифм почасовой заработной платы соответствующего региона отличается от среднего логарифма почасовой заработной платы для Москвы и Санкт-Петербурга.

Объединенная модель индивидуальной заработной платы будет иметь вид

$$\ln WPH_{it} = \beta_0 + \beta_1 SCH_{it} + \beta_2 EXP_{it} + \beta_3 EXP_{it}^2 + \beta_4 GEN_{it} + \beta_5 RUS_{it} + \beta_6 UNSKILL_{it} + \\ + \beta_7 NW_{it} + \beta_8 CENTR_{it} + \beta_9 VOL_{it} + \beta_{10} CAUCAS_{it} + \beta_{11} URAL_{it} + \beta_{12} WSIB_{it} + \\ + \beta_{13} FAREAST_{it} + \varepsilon_{it},$$

где  $i = 1, \dots, 870$  и  $t = 1, 2, 3$ .

Уравнение регрессии объединенной модели содержит 13 регрессоров и будет оцениваться на основе данных из 2610 наблюдений.

Результаты оценивания уравнения заработной платы представлены в табл. 9.2. В первой графе таблицы представлены межгрупповые оценки параметров модели, основанные на индивидуальных средних, а во второй графе — внутригрупповые оценки параметров модели, основанные на отклонениях от индивидуальных средних. В третьей графе таблицы размещены результаты оценивания объединенной модели с помощью обычного метода наименьших квадратов, а в четвертой графе — результаты применения доступного обобщенного метода наименьших квадратов для оценивания случайных эффектов.

**Таблица 9.2.** Результаты оценивания уравнения заработной платы (в круглых скобках указаны стандартные ошибки)

Зависимая переменная <i>LN_WPH</i>				
переменная модели	межгрупповая оценка	оценка с фиксированными эффектами	МНК-оценка	оценка со случайными эффектами
Константа	1,3899 (0,1669)	—	1,4285 (0,1245)	1,4697 (0,1611)
<i>SCH</i>	0,0812 (0,0095)	—	0,0812 (0,0071)	0,0815 (0,0094)
<i>EXP</i>	0,0189 (0,0079)	-0,0057 (0,0190)	0,0183 (0,0058)	0,0170 (0,0073)
<i>EXP_2</i>	-0,0005 (0,0002)	-0,0006 (0,0004)	-0,0005 (0,0001)	-0,0005 (0,0001)
<i>GEN</i>	0,3112 (0,0444)	—	0,3153 (0,0337)	0,3186 (0,0444)
<i>RUS</i>	0,3232 (0,0575)	0,1324 (0,0488)	0,2561 (0,0356)	0,2072 (0,0373)
<i>UNSKILL</i>	-0,2950 (0,0909)	-0,1481 (0,1043)	-0,2639 (0,0611)	-0,2311 (0,0687)
<i>NW</i>	-0,3772 (0,1021)	—	-0,3856 (0,0775)	-0,3937 (0,1023)
<i>CENTR</i>	-0,6371 (0,0837)	—	-0,6396 (0,0635)	-0,6421 (0,0838)
<i>VOL</i>	-0,7204 (0,0875)	—	-0,7295 (0,0663)	-0,7365 (0,0872)
<i>CAUCAS</i>	-0,8617 (0,1062)	—	-0,8695 (0,0806)	-0,8771 (0,1064)



Окончание табл. 9.2

Зависимая переменная LN_WPH				
переменная модели	межгрупповая оценка	оценка с фиксированными эффектами	МНК-оценка	оценка со случайными эффектами
URAL	-0,4356 (0,0859)	—	-0,4396 (0,0652)	-0,4428 (0,0860)
WESTSIB	-0,2427 (0,1044)	—	-0,2551 (0,0792)	-0,2652 (0,1044)
FAREAST	-0,5058 (0,1088)	—	-0,5057 (0,0826)	-0,5070 (0,1091)

Все оцененные уравнения регрессии оказались в целом статистически значимыми на 1%-ном уровне значимости. В объединенной модели все коэффициенты при регрессорах получились статистически значимыми на 5%-ном уровне значимости, коэффициент детерминации  $R^2 = 0,1611$ . Обратим внимание на то, что небольшое значение коэффициента детерминации характерно для моделей такого рода.

Оценка нормы отдачи от инвестиций в образование в разных моделях принимает достаточно близкие значения: 0,0812 в межгрупповой и в объединенной модели и 0,0815 в модели со случайными эффектами.

Модель с фиксированными эффектами не позволяет оценить норму отдачи от инвестиций в образование, так как оценка модели с фиксированными эффектами (оценка внутригрупповой модели) исключает из модели любые переменные, независимые от времени. Это означает, что влияния числа лет обучения, полового признака, а также региональных фиктивных переменных не учитываются. Обратим внимание также на то, что оценки параметров внутригрупповой модели очень сильно отличаются от оценок параметров межгрупповой модели.

Во всех оцененных моделях подтвердилась квадратичная зависимость уровня заработной платы от профессионального опыта, коэффициенты перед регрессорами  $EXP$  и  $EXP^2$  имеют положительный и отрицательный знак соответственно.

Для внутригрупповой модели по формуле (9.60) оценим дисперсию ошибок  $\sigma_e^2 = 0,4340$ , для межгрупповой модели оценка дисперсия ошибок составит 0,3968. Используя формулу (9.61), оценим дисперсию ошибок  $\sigma_u^2 = 0,3968 - 0,4340 / 3 = 0,2521$ .

Зная оценки дисперсии ошибок  $\sigma_u^2$  и  $\sigma_\varepsilon^2$ , можно вычислить

$$\hat{\lambda} = \frac{\sigma_\varepsilon^2}{T\sigma_u^2 + \sigma_\varepsilon^2} = \frac{0,4340}{3 \cdot 0,2521 + 0,4340} = 0,3646.$$

Так как  $\hat{\lambda} = (1 - \hat{\theta})^2$ , то  $\hat{\theta} = 1 - \hat{\lambda}^{1/2} = 1 - 0,3646^{1/2} = 0,3962$ .

Найденная оценка параметра  $\theta$  принимает значение, лежащее в интервале от нуля до единицы. Напомним, что при  $\theta = 1$  модель со случайными эффектами совпадает с моделью с фиксированными эффектами, а при  $\theta = 0$  в модели будут отсутствовать некоррелируемые индивидуальные компоненты дисперсии и оценка со случайными эффектами сократится до МНК-оценки объединенной модели. Обратим внимание на то, что оценки со случайными эффектами и МНК-оценки принимают значения, находящиеся между оценками с фиксированными эффектами и межгрупповыми оценками.

Получается, что оценки доступного обобщенного метода наименьших квадратов можно получить путем оценивания обычным МНК уравнения преобразованной регрессии (9.56):

$$(y_{it} - 0,3962\bar{y}_i) = (1 - 0,3962)\alpha + (x_{it} - 0,3962\bar{x}_i)\beta + v_{it},$$

где  $v_{it} = (1 - 0,3962)u_i + (\varepsilon_{it} - 0,3962\bar{\varepsilon}_i)$ .

Если выполняются предположения модели со случайными эффектами, то оценки доступного обобщенного метода наименьших квадратов будут эффективными. Если же индивидуальные эффекты  $u_i$  коррелируют хотя бы с одним регрессором, то только оценка модели с фиксированными эффектами будет состоятельной.

### 9.3. Качество подгонки

Использование меры *качества подгонки* (*goodness-of-fit*) для панельных данных имеет некоторые особенности. В первую

очередь это связано с тем, что в панельных данных уделяется особое внимание внутригрупповой и межгрупповой вариациям. Другая особенность заключается в том, что обычный или скорректированный коэффициент детерминации  $R^2$  уместен в роли меры качества подгонки только в тех случаях, где модель оценивается при помощи обычного метода наименьших квадратов. Так, например, коэффициент детерминации не является адекватной мерой качества подгонки для модели со случайными эффектами, оцениваемой при помощи обобщенного метода наименьших квадратов.

Один из возможных подходов для меры качества подгонки связан с использованием квадрата коэффициента корреляции между фактическими значениями  $y_{it}$  и теоретическими (подогнанными) значениями  $\hat{y}_{it}$ , который совпадает с коэффициентом детерминации в случае стандартной линейной модели регрессии. Рассмотрим вариацию переменной  $y_{it}$ , которую можно представить в виде суммы внутригрупповой и межгрупповой вариаций:

$$\frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T (y_{it} - \bar{y})^2 = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T (y_{it} - \bar{y}_i)^2 + \frac{1}{N} \sum_{i=1}^N (\bar{y}_i - \bar{y})^2, \quad (9.62)$$

где  $\bar{y} = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T y_{it}$  — среднее значение переменной  $y$  по всей выборке;  $\bar{y}_i = \frac{1}{T} \sum_{t=1}^T y_{it}$  — среднее значение переменной  $y$  для  $i$ -го индивида.

Внутригрупповая модель (9.8), построенная по отклонению от групповых средних, используется для того, чтобы объяснить внутригрупповую вариацию, и мера качества подгонки для внутригрупповой модели будет выражаться в виде «внутригруппового  $R^2$ »:

$$R^2_{\text{внутри}}(\hat{\beta}_{\text{внутри}}) = r^2 \left( y_{it} - \bar{y}_i, \hat{y}_{it}^{FE} - \hat{y}_i^{FE} \right), \quad (9.63)$$

где  $r(\cdot, \cdot)$  — коэффициент корреляции;  $\hat{y}_{it}^{FE} - \hat{y}_i^{FE} = (x_{it} - \bar{x}_i) \hat{\beta}_{\text{внутри}}$ ;  
 $\hat{y}_i^{FE} = \frac{1}{T} \sum_{t=1}^T \hat{y}_{it}^{FE}$ .

Напомним, что оценки внутригрупповой модели совпадают с оценками с фиксированными эффектами.

Межгрупповая модель (9.9) на основе групповых средних используется для того, чтобы объяснить межгрупповую вариацию, т.е. вариацию между различными индивидами, и мера качества подгонки для межгрупповой модели будет выражаться в виде «межгруппового  $R^2$ »:

$$R_{\text{меж}}^2(\hat{\beta}_{\text{меж}}) = r^2 \left( \bar{y}_i, \hat{y}_i^B \right), \quad (9.64)$$

где  $\hat{y}_i^B = \bar{x}_i' \hat{\beta}_{\text{меж}}$ ;  $\hat{y}_i^B = \frac{1}{T} \sum_{t=1}^T \hat{y}_{it}^B$ .

Объединенная модель используется для того, чтобы объяснить общую вариацию, и для нее мера качества подгонки будет выражаться в виде «общего  $R^2$ »:

$$R_{\text{общ}}^2(\hat{\beta}_{OLS}) = r^2(y_{it}, \hat{y}_{it}), \quad (9.65)$$

где  $\hat{y}_{it} = x_{it}' \hat{\beta}_{OLS}$ .

Внутригрупповой, межгрупповой и общий  $R^2$  можно вычислить для любого вектора оценок  $\hat{\beta}$  для вектора параметров  $\beta$ . При этом в качестве теоретических (подогнанных) значений для  $y_{it}$  будут выступать  $\hat{y}_{it} = x_{it}' \hat{\beta}$ ,  $\hat{y}_i = \frac{1}{T} \sum_{t=1}^T \hat{y}_{it}$  и  $\hat{y} = \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \hat{y}_{it}$ . Необходимо отметить, что в модели с фиксированными эффектами эффекты  $\hat{\alpha}_i$  скорее не объясняют вариацию между индивидами, а только улавливают ее. К тому же при оценивании модели с фиксированными эффектами эффекты  $\hat{\alpha}_i$  часто даже не вычисляются. При вычислении подогнанных значений  $\hat{y}_{it}$  для оценок с фиксированными эффектами игнорируется вариация, улавливаемая индивидуальными эффектами  $\hat{\alpha}_i$ .

Рассмотренные выше меры качества подгонки можно вычислить для межгрупповой, внутригрупповой, объединенной моделей и модели со случайными эффектами. В случае если предположения о случайных эффектах действительно имеют место, то оценка со случайными эффектами будет (асимптотически) наиболее эффективной. Однако для модели со случайными

эффектами «внутригрупповой  $R^2$ » будет обязательно меньше, чем для модели с фиксированными эффектами, «межгрупповой  $R^2$ » будет обязательно меньше, чем для межгрупповой модели, и «общий  $R^2$ » будет обязательно меньше, чем для объединенной модели. Это свидетельствует о том, что мера качества подгонки не подходит для выбора соответствующего метода оценивания в моделях с панельными данными. Однако меры качества подгонки можно использовать для выбора между моделями, оцененными одним и тем же методом, но отличающимися набором объясняющих переменных.

Рассчитаем меры качества подгонки для примера с индивидуальной заработной платой.

**Таблица 9.3.** Меры качества подгонки

Меры качества подгонки	Межгрупповая модель	Модель с фиксированными эффектами	Объединенная модель	Модель со случайными эффектами
Внутригрупповой $R^2$	0,0060	0,0157	0,0066	0,0075
Межгрупповой $R^2$	0,2516	0,0089	0,2502	0,2468
Общий $R^2$	0,1600	0,0077	0,1611	0,1602

Меры качества подгонки, представленные в табл. 9.3, показывают, что оценка с фиксированными эффектами приводит к наибольшему внутригрупповому  $R^2$ , наибольший межгрупповой  $R^2$  соответствует межгрупповой модели и наибольший общий  $R^2$  относится к объединенной модели.

## 9.4. Выбор модели

При анализе панельных данных необходимо выбрать, какая из моделей панельных данных (объединенная модель, модель с фиксированными эффектами или модель со случайными эффектами) больше всего подходит для определенной ситуации. Объединенная модель предполагает, что у единиц совокупности отсутствуют индивидуальные различия. Модель с фиксированными эффектами предполагает, что каждая единица совокупности имеет свои специфические индивидуальные

характеристики, которые для каждого конкретного объекта являются постоянными во времени. Если же единицы совокупности различаются по своим индивидуальным характеристикам, но эти различия носят случайный характер, то в этом случае лучше рассматривать модель со случайными эффектами.

Важной особенностью модели со случайными эффектами является то, что ошибки  $u_i$  должны не коррелировать с объясняющими переменными. Однако это условие достаточно часто не выполняется, даже в тех случаях, когда осуществляется выборка из большой исходной совокупности. Иногда в качестве основного различия между моделью со случайными эффектами и моделью с фиксированными эффектами рассматривают именно присутствие или отсутствие корреляции между индивидуальными эффектами и объясняющими переменными.

Качество оценок  $\beta$ -коэффициентов зависит от того, насколько корректно делается выбор в пользу той или иной модели. Оценки модели с фиксированными эффектами будут несмещенными и состоятельными, однако их эффективность зависит от того, присутствует ли корреляция между индивидуальными эффектами и объясняющими переменными. Если имеет место такая корреляция, то оценки фиксированных эффектов будут эффективными, в противном случае эти оценки будут неэффективными. Оценки модели со случайными эффектами будут несостоятельными в том случае, если в панельных данных присутствуют фиксированные эффекты.

Помимо содержательных выводов существует несколько формальных тестов, которые можно проводить для того, чтобы сделать выбор в пользу той или иной модели панельных данных.

### 1. Объединенная модель против модели с фиксированными эффектами

При выполнении предположения о нормальности могут использоваться обычный  $t$ -критерий и  $F$ -критерий. Так, при помощи  $F$ -критерия можно проверить, присутствуют ли вообще какие-либо индивидуальные эффекты в панельных данных. В этом случае будет тестироваться нулевая гипотеза об отсутствии индивидуальных эффектов, которая заключается в том, что все коэффициенты  $\alpha_i$  равны между собой, против альтернативной гипотезы о наличии индивидуальных эффектов,

т.е. хотя бы два любых коэффициента  $\alpha_i$  принимают разные значения:

$$H_0: \alpha_1 = \alpha_2 = \dots = \alpha_N = \alpha_0;$$

$$H_1: \text{не } H_0.$$

Для проверки нулевой гипотезы будем использовать тест Чоу. В качестве модели с ограничениями будет выступать объединенная модель  $y_{it} = \alpha_0 + x'_{it}\beta + \varepsilon_{it}$ , а в качестве модели без ограничений — модель с фиксированными эффектами  $y = D_N\alpha + X\beta + \varepsilon$ . Нулевая гипотеза проверяется с помощью  $F$ -критерия

$$F = \frac{(SS_R - SS_{UR}) / (N - 1)}{SS_{UR} / (NT - N - K)}, \quad (9.66)$$

где  $SS_R$  — сумма квадратов остатков в модели с ограничениями;  $SS_{UR}$  — сумма квадратов остатков в модели без ограничений.

$F$ -критерий (9.66) имеет  $F$ -распределение с  $(N - 1)$  и  $(NT - N - K)$  степенями свободы. В числителе число степеней свободы составляет  $(N - 1)$ , а не  $N$ , так как проверка гипотезы о том, что все  $N$  коэффициентов  $\alpha_i$  равны между собой, равносильно проверке гипотезы о том, что все  $(N - 1)$  разности  $(\alpha_{i+1} - \alpha_i)$ ,  $i = 1, \dots, N - 1$ , равны нулю.

В случае если вычисленное значение  $F$ -критерия окажется больше критического значения,  $F_{\text{факт}} > F_{\text{кр}}$ , для заданного уровня значимости, то мы можем отклонить нулевую гипотезу и принять альтернативную гипотезу о присутствии индивидуальных эффектов, т.е. сделать выбор в пользу модели с фиксированными эффектами. В случае если  $F_{\text{факт}} < F_{\text{кр}}$  для заданного уровня значимости, то у нас нет оснований отвергать нулевую гипотезу об отсутствии индивидуальных эффектов и мы можем сделать выбор в пользу объединенной модели.

В нашем примере сумма квадратов остатков в объединенной модели (модель с ограничениями)  $SS_R = 1784,19$  и сумма квадратов остатков в модели с фиксированными эффектами (модель без ограничений)  $SS_{UR} = 75,4$ . По формуле (9.66) вычислим  $F$ -критерий

$$F_{\text{факт}} = \frac{(1784,19 - 753,46) / 869}{753,46 / 1736} = 2,7328.$$

Критическое значение  $F$ -статистики на 5%-ном уровне значимости составляет  $F_{кр}(869, 1736) = 1,1005$ . Так как  $F_{факт} > F_{кр}$ , то у нас есть основания для того, чтобы отклонить нулевую гипотезу и принять альтернативную гипотезу о том, что коэффициенты  $\alpha_i$  не равны между собой и в модели присутствуют индивидуальные эффекты. По результатам проведенного теста мы выбираем модель с фиксированными эффектами.

Помимо рассмотренного выше теста Чоу возможны также другие способы для того, чтобы осуществить выбор между объединенной моделью и моделью с фиксированными эффектами. Например, можно также проверить гипотезу о том, что значения индивидуальных характеристик совпадают у двух разных индивидов:

$$H_0: \alpha_i = \alpha_j \text{ для некоторых } i \text{ и } j, i \neq j;$$

$$H_1: \alpha_i \neq \alpha_j.$$

Для проверки этой нулевой гипотезы необходимо вычислить  $t$ -критерий

$$t = \frac{\hat{\alpha}_i - \hat{\alpha}_j}{\sqrt{\hat{V}(\hat{\alpha}_i - \hat{\alpha}_j)}}, \quad (9.67)$$

который будет иметь  $t$ -распределение с  $(NT - N - k)$  степенями свободы.

В случае если вычисленное значение  $t$ -критерия окажется больше критического значения ( $|t_{факт}| > t_{кр}$ ) для заданного уровня значимости, то мы можем отклонить нулевую гипотезу и принять альтернативную гипотезу о том, что два разных индивида имеют различные индивидуальные характеристики. В случае если  $|t_{факт}| < t_{кр}$  для заданного уровня значимости, то у нас нет оснований отвергнуть нулевую гипотезу о равенстве индивидуальных характеристик у двух разных индивидов.

## 2. Объединенная модель против модели со случайными эффектами

### Анализ дисперсии (тест Фишера)

Первый способ проверки отсутствия индивидуальных эффектов состоит в проверке нулевой гипотезы о равенстве нулю дисперсии  $\sigma_u^2$ :

$$H_0: \sigma_u^2 = 0;$$

$$H_1: \sigma_u^2 \neq 0.$$



Этот тест легко выполняется если оценить внутригрупповую и межгрупповую дисперсии, из которых мы можно получить оценки дисперсии остатков. Во внутригрупповой модели дисперсия остатков  $\hat{\sigma}_w^2$  будет равна  $\hat{\sigma}_u^2$ , а в межгрупповой модели дисперсия остатков составит  $\hat{\sigma}_b^2 = \hat{\sigma}_u^2 + \hat{\sigma}_\varepsilon^2 / T$ .

При выполнении предположения о нормальности величина

$$\xi_1 = \frac{\hat{\sigma}_w^2}{\hat{\sigma}_\varepsilon^2} [N(T-1) - k_w] \quad (9.68)$$

будет иметь  $\chi^2$ -распределение с  $[N(T-1) - k_w]$  степенями свободы, где  $k_w$  — число регрессоров во внутригрупповой модели, и величина

$$\xi_2 = (N - k_b) \frac{\hat{\sigma}_b^2}{\hat{\sigma}_u^2 + \hat{\sigma}_\varepsilon^2 / T} = (N - k_b) \frac{T \hat{\sigma}_b^2}{T \hat{\sigma}_u^2 + \hat{\sigma}_\varepsilon^2} \quad (9.69)$$

будет иметь  $\chi^2$ -распределение с  $(N - k_b)$  степенями свободы, где  $k_b$  — число регрессоров в межгрупповой модели. Как следствие величина  $\xi_2 / \xi_1$  будет иметь  $F$ -распределение с  $(N - k_b)$  и  $[N(T-1) - k_w]$  степенями свободы

$$\frac{\sigma_\varepsilon^2}{T \sigma_u^2 + \sigma_\varepsilon^2} \frac{T \hat{\sigma}_b^2}{\hat{\sigma}_w^2} \rightarrow F[N - k_b, N(T-1) - k_w]. \quad (9.70)$$

Тогда при нулевой гипотезе, состоящей в том, что  $\sigma_u^2 = 0$ , выражение (9.66) примет вид

$$\frac{T \hat{\sigma}_b^2}{\hat{\sigma}_\varepsilon^2} \rightarrow F[N - k_b, N(T-1) - k_w]. \quad (9.71)$$

В результате мы можем отклонить нулевую гипотезу и принять альтернативную в том случае, когда  $F$ -статистика (9.71) больше, чем  $F_{кр}$  со степенями свободы  $(N - k_b)$  и  $[N(T-1) - k_w]$  для соответствующего уровня значимости. Получается, что если индивидуальная дисперсия межгрупповой регрессии, которая вычисляется по  $N$  наблюдениям, умноженная на  $T$ , больше, чем остаточная дисперсия во внутригрупповой регрессии, то гипотеза об отсутствии индивидуальных эффектов отклоняется: такие эффекты присутствуют.

Проведем анализ дисперсий в модели индивидуальной заработной платы. Во внутригрупповой модели оценка дисперсии

ошибок  $\hat{\sigma}_w^2 = \hat{\sigma}_u^2 = 0,2521$ , а в межгрупповой модели оценка дисперсии ошибок  $\hat{\sigma}_b^2 = \hat{\sigma}_u^2 + \hat{\sigma}_\varepsilon^2 / T = 0,2521 + 0,4340 / 3 = 0,3968$ .

Для проверки нулевой гипотезы о равенстве нулю дисперсии  $\sigma_u^2$  вычислим  $F$ -статистику

$$F_{\text{факт}} = \frac{\sigma_\varepsilon^2}{T\sigma_u^2 + \sigma_\varepsilon^2} \frac{T\hat{\sigma}_b^2}{\hat{\sigma}_w^2} = \frac{0,4340}{3 \cdot 0,2521 + 0,4340} \frac{3 \cdot 0,3968}{0,2521} = 1,7217.$$

Критическое значение  $F$ -статистики для уровня значимости 0,05 составляет  $F_{\text{кр}}(857, 1736) = 1,101$ . Так как вычисленное значение  $F$ -критерия больше критического значения ( $F_{\text{факт}} > F_{\text{кр}}$ ), то мы можем отклонить нулевую гипотезу и принять альтернативную гипотезу о том, что дисперсия  $\sigma_u^2$  не равна нулю. Таким образом, можно сделать вывод о присутствии индивидуальных эффектов и сделать выбор в пользу модели со случайными эффектами.

### Тест множителей Лагранжа Бреуша — Пагана

Т. Бреуш и Э. Паган<sup>1</sup> в 1980 г. предложили тест множителей Лагранжа, который может использоваться для тестирования модели со случайными эффектами против объединенной модели. При помощи этого теста проверяются следующие гипотезы:

$$H_0: \sigma_u = 0;$$

$$H_1: \sigma_u \neq 0.$$

Из формулы (9.52) очевидно, что при справедливости нулевой гипотезы о равенстве дисперсии  $\sigma_u$  нулю ковариационная матрица  $V = \sigma_\varepsilon^2 I_N$ , и в результате модель со случайными эффектами сокращается до объединенной МНК-регрессии с независимыми одинаково распределенными ошибками.

Проверка этой гипотезы основана на МНК-остатках  $\hat{\varepsilon}_{it}$  из объединенной регрессии и требует вычисления  $LM$ -статистики

$$LM = \frac{NT}{2(T-1)} \left[ \frac{\sum_{i=1}^N \left( \sum_{t=1}^T \hat{\varepsilon}_{it} \right)^2}{\sum_{i=1}^N \sum_{t=1}^T \hat{\varepsilon}_{it}^2} - 1 \right]^2. \quad (9.72)$$

<sup>1</sup> Breusch T. S. and Pagan A. R. The Lagrange Multiplier Test and its Applications to Model Specification in Econometrics. Review of Economic Studies, Vol. 47 (1980). Pp. 239—253.

При нулевой гипотезе  $LM$ -статистика будет асимптотически иметь  $\chi^2$ -распределение с одной степенью свободы. В результате, если  $LM$ -статистика, вычисленная на основе остатков МНК-регрессии, больше, чем 3,84 (при проверке гипотезы на 5%-ном уровне значимости), то нулевая гипотеза об отсутствии индивидуальных случайных эффектов отвергается и принимается альтернативная гипотеза об их присутствии. В противном случае у нас нет оснований для того, чтобы отвергнуть нулевую гипотезу.

$LM$ -статистика для нашего примера, вычисленная по формуле (9.72) на основе МНК-остатков  $\hat{\varepsilon}_{it}$  из объединенной регрессии, составила  $LM = 253,15$ , что намного больше критического значения  $\chi_{кр}^2$  с одной степенью свободы на 5%-ном уровне значимости, т.е. мы можем отклонить нулевую гипотезу о равенстве нулю дисперсии  $\sigma_u$  и принять альтернативную гипотезу о том, что дисперсия  $\sigma_u$  не равна нулю. Таким образом, тест множителей Лагранжа Бреуша — Пагана свидетельствует о том, что в нашей модели присутствуют индивидуальные эффекты, и мы делаем выбор в пользу модели со случайными эффектами.

### Тест Хонды

Недостаток теста Бреуша — Пагана состоит в том, что он представляет собой двусторонний тест, тогда как дисперсия должна быть либо нулевой, либо положительной. Для того чтобы охватить эту проблему, Ю. Хонда<sup>1</sup> в 1985 г. предложил очень простой односторонний тест. Тест Хонды заключается в том, чтобы рассматривать квадратный корень  $LM$ -статистики Бреуша — Пагана. Тест показал, что при нулевой гипотезе об отсутствии индивидуальных эффектов, квадратный корень имеет нормальное распределение. Статистика Хонды определяется по формуле

$$g = \sqrt{\frac{NT}{2(T-1)}} \left[ \frac{\sum_{i=1}^N \left( \sum_{t=1}^T \hat{\varepsilon}_{it} \right)^2}{\sum_{i=1}^N \sum_{t=1}^T \hat{\varepsilon}_{it}^2} - 1 \right]. \quad (9.73)$$

<sup>1</sup> Honda Y. Testing the Error Components Model with Non-normal Disturbances. Review of Economic Studies. Vol. 52 (1985). Pp. 681—690.

Если  $g$ -статистика Хонды больше, чем 1,64, то нулевая гипотеза об отсутствии индивидуальных эффектов отклоняется. Тест Хонды приводит к отклонению нулевой гипотезы гораздо чаще, чем это делает  $LM$ -статистика Бреуша — Пагана.

В нашем примере статистика Хонды  $g = 15,91$  также свидетельствует о присутствии индивидуальных эффектов в модели и показывает, что модель со случайными эффектами больше подходит для наших данных, чем объединенная модель.

### 3. Модель с фиксированными эффектами против модели со случайными эффектами

#### Тест Хаусмана

В модели со случайными эффектами предполагается, что индивидуальные эффекты не коррелируют с регрессорами. Важно проверить, действительно ли выполняется предположение о такой корреляции, которая может привести к несостоятельности большинства оценок модели со случайными эффектами. Исключение в этом составляют лишь внутригрупповые оценки  $\beta$ -коэффициентов, так как они базируются на преобразовании, которое избавляется от индивидуального эффекта в модели. В результате предположение о том, что они не коррелируют с регрессорами, не связано с несмещенностью и несостоятельностью внутригрупповых оценок.

Дж. Хаусман<sup>1</sup> в 1978 г. предложил тест, который использует тот факт, что две оценки могут быть определены таким образом, что одна оценка  $\hat{\beta}(1)$  является состоятельной как при нулевой, так и при альтернативной гипотезе, а вторая оценка  $\hat{\beta}(2)$  является состоятельной и эффективной, только если справедлива нулевая гипотеза, и несостоятельной — в противном случае. Если полученные оценки  $\hat{\beta}(1)$  и  $\hat{\beta}(2)$  принимают достаточно близкие друг к другу значения, то это может свидетельствовать о том, что нулевая гипотеза справедлива. Если полученные оценки  $\hat{\beta}(1)$  и  $\hat{\beta}(2)$  отличаются друг от друга, то это может рассматриваться как индикатор того, что нулевая гипотеза не верна. Хаусман показал, что при нулевой гипотезе статистика  $Q_H$  определяемая по формуле

<sup>1</sup> Hausman J. A. Specification Tests in Econometrics. *Econometrica*. Vol. 46 (1978). Pp. 1251—1271.

$$Q_H = [\hat{\beta}(1) - \hat{\beta}(2)]' \{ \hat{V}[\hat{\beta}(1)] - \hat{V}[\hat{\beta}(2)] \}^{-1} [\hat{\beta}(1) - \hat{\beta}(2)], \quad (9.74)$$

где  $\hat{V}[\hat{\beta}(1)]$  и  $\hat{V}[\hat{\beta}(2)]$  — оценки соответствующих матриц ковариаций,

является асимптотически распределенной (когда  $N \rightarrow \infty$ ) как  $\chi^2$ -распределение со степенью свободы, равной размерности  $\beta$ . Если статистика  $Q_H$  больше, чем квантиль  $\chi^2$ -распределения с соответствующим числом степеней свободы, то мы можем отклонить нулевую гипотезу и принять альтернативную, в противном случае у нас нет оснований для отклонения нулевой гипотезы.

В нашем случае в качестве  $\hat{\beta}(1)$  мы можем выбрать внутригрупповую оценку ( $\hat{\beta}(1) = \hat{\beta}_{\text{внутри}}$ ), а в качестве  $\hat{\beta}(2)$  можно выбрать оценку доступного обобщенного метода наименьших квадратов ( $\hat{\beta}(2) = \hat{\beta}_{\text{FGLS}}$ ). Тогда статистика  $Q_H$  будет иметь вид

$$Q_H = (\hat{\beta}_{\text{внутри}} - \hat{\beta}_{\text{FGLS}})' [\hat{V}(\hat{\beta}_{\text{внутри}}) - \hat{V}(\hat{\beta}_{\text{FGLS}})]^{-1} (\hat{\beta}_{\text{внутри}} - \hat{\beta}_{\text{FGLS}}). \quad (9.75)$$

Если значение статистики  $Q_H$  больше, чем критическое значение  $\chi^2$ -распределения с  $k_w$  степенями свободы, где  $k_w$  — число регрессоров во внутригрупповой модели, то можно отклонить нулевую гипотезу и принять альтернативную, сделав тем самым выбор в пользу модели с фиксированными эффектами. Если же  $Q_H < \chi^2_{\text{кр}}(k_w)$ , то у нас нет оснований для отклонения нулевой гипотезы и различия между оценками не являются систематическими. Это означает, что можно выбрать модель со случайными эффектами. Получается, что когда внутригрупповая оценка является состоятельной, то оценки доступного обобщенного МНК, обычного МНК и межгрупповая оценка в этом случае не являются состоятельными.

Стоит обратить внимание на то, что этот тест можно провести альтернативным способом. Дж. Хаусман и У. Тейлор (*Taylor*) в 1981 г. показали, что для вычисления статистики можно также использовать межгрупповую оценку  $\hat{\beta}_{\text{меж}}$  и оценку доступного обобщенного МНК  $\hat{\beta}_{\text{FGLS}}$ :

$$Q_{\text{fglS}, b} = (\hat{\beta}_{\text{меж}} - \hat{\beta}_{\text{FGLS}})' [\hat{V}(\hat{\beta}_{\text{меж}}) - \hat{V}(\hat{\beta}_{\text{FGLS}})]^{-1} (\hat{\beta}_{\text{меж}} - \hat{\beta}_{\text{FGLS}}), \quad (9.76)$$

или межгрупповую оценку  $\hat{\beta}_{\text{меж}}$  и внутригрупповую оценку  $\hat{\beta}_{\text{внутр}}$ :

$$Q_{w,b} = (\hat{\beta}_{\text{меж}} - \hat{\beta}_{\text{внутр}})' \left[ \hat{V}(\hat{\beta}_{\text{меж}}) - \hat{V}(\hat{\beta}_{\text{внутр}}) \right]^{-1} (\hat{\beta}_{\text{меж}} - \hat{\beta}_{\text{внутр}}). \quad (9.77)$$

Если сравнить полученные значения статистики (9.76) и (9.77) с критическим значением  $\chi^2(k_w)$ , то можно проверить ту же самую нулевую гипотезу, что и при помощи статистики (9.75). Все эти три статистики (9.75)—(9.77) на самом деле являются численно идентичными, и можно использовать любую из них.

Таким образом, при помощи теста Хаусмана можно проверить гипотезу об отсутствии корреляции между индивидуальными эффектами и регрессорами и сделать выбор в пользу модели со случайными эффектами либо в пользу модели с фиксированными эффектами.

Для данных нашего примера по индивидуальной заработной плате статистика  $Q_H$  составила 22,09. Критическое значение  $\chi^2$ -статистики для четырех степеней свободы на 5%-ном уровне значимости равняется 9,49. Получается, что  $Q_H > \chi^2_{\text{кр}}$ , т.е. различия между оценками являются систематическими и индивидуальные регрессоры коррелируют с регрессорами. Таким образом, можно сделать выбор в пользу модели с фиксированными эффектами.

## 9.5. Двухнаправленная модель панельных данных с фиксированными эффектами

Двухнаправленная модель с фиксированными эффектами помимо индивидуальных эффектов  $\alpha_i$  включает также временные эффекты  $\gamma_t$ :

$$y_{it} = x'_{it} \beta + \alpha_i + \gamma_t + \varepsilon_{it}, \quad (9.78)$$

где  $i = 1, \dots, N$ ;  $t = 1, \dots, T$ ;  $\varepsilon_{it} \sim iid(0, \sigma_\varepsilon^2)$ .

Эту модель можно оценить путем добавления  $(T - 1)$  фиктивных переменных для соответствующих периодов времени

наряду с  $(N - 1)$  фиктивными переменными, которые соответствуют различным объектам наблюдения. Однако если  $N$  или  $T$  являются большими, то в регрессии будет слишком много фиктивных переменных  $(N + T - 2)$ , что приведет к огромной потере в степенях свободы и к необходимости обратить большую матрицу размерностью  $(N + T + K - 1)$ . Другой способ получения оценок  $\beta$ -коэффициентов в модели (9.78) заключается в использовании внутригруппового преобразования, предложенного в 1969 г. Т. Уолласем и Э. Хусейном<sup>1</sup>, которое исключает из модели эффекты  $\alpha_i$  и  $\gamma_t$ :

$$Q = (I_N - \bar{J}_N) \otimes (I_T - \bar{J}_T) = I_N \otimes I_T - I_N \otimes \bar{J}_T - \bar{J}_N \otimes I_T + \bar{J}_N \otimes \bar{J}_T, \quad (9.79)$$

где  $\bar{J}_N = \bar{J}_N / N$ ,  $\bar{J}_T = \bar{J}_T / T$ .

Преобразования  $\tilde{y} = Qy$  и  $\tilde{x} = Qx$  приведут соответственно к следующим значениям  $\tilde{y}_{it}$  и  $\tilde{x}_{it}$ :

$$\tilde{y}_{it} = (y_{it} - \bar{y}_i - \bar{y}_t + \bar{y}),$$

где  $\bar{y}_i = 1/T \sum_{t=1}^T y_{it}$ ;  $\bar{y}_t = 1/N \sum_{i=1}^N y_{it}$ ;  $\bar{y} = 1/NT \sum_{i=1}^N \sum_{t=1}^T y_{it}$ ;

$$\tilde{x}_{it} = (x_{it} - \bar{x}_i - \bar{x}_t + \bar{x}),$$

где  $\bar{x}_i = 1/T \sum_{t=1}^T x_{it}$ ;  $\bar{x}_t = 1/N \sum_{i=1}^N x_{it}$ ;  $\bar{x} = 1/NT \sum_{i=1}^N \sum_{t=1}^T x_{it}$ .

Построим регрессию  $\tilde{y}$  на  $\tilde{x}$ :

$$(y_{it} - \bar{y}_i - \bar{y}_t + \bar{y}) = (x_{it} - \bar{x}_i - \bar{x}_t + \bar{x}) \tilde{\beta} + (\varepsilon_{it} - \bar{\varepsilon}_i - \bar{\varepsilon}_t + \bar{\varepsilon}). \quad (9.80)$$

Оценивание методом наименьших квадратов уравнения регрессии (9.80) дает внутригрупповую оценку для двухнаправленной модели  $\tilde{\beta} = (X'QX)^{-1} X'Qy$ . Оценка константы может быть найдена как разница между  $\bar{y}$  и  $\tilde{\beta}\bar{x}$ . Внутригрупповые оценки эффектов  $\alpha_i$  и  $\gamma_t$  будут иметь вид

$$\tilde{\alpha}_i = (\bar{y}_i - \bar{y}) - \tilde{\beta}(\bar{x}_i - \bar{x}), \quad \text{где } i = 1, \dots, N - 1; \quad (9.81)$$

$$\tilde{\gamma}_t = (\bar{y}_t - \bar{y}) - \tilde{\beta}(\bar{x}_t - \bar{x}), \quad \text{где } t = 1, \dots, T - 1. \quad (9.82)$$

<sup>1</sup> Wallace T. D. and Hussain A. A. The Use of Error Components Models in Combining Cross-section and Time-series Data. *Econometrica*. Vol. 37 (1969). Pp. 55—72.

Результаты оценивания модели (9.80) методом наименьших квадратов представлены в табл. 9.4.

**Таблица 9.4.** МНК-оценки уравнения регрессии (9.80)

Переменные	<i>EXP</i>	<i>EXP_2</i>	<i>RUS</i>	<i>UNSKILL</i>
Коэффициенты	0,1439	-0,0008	0,1220	-0,1411
Стандартные ошибки	(0,1484)	(0,0003)	(0,0388)	(0,0828)

Следует обратить внимание на то, что внутригрупповая оценка не позволяет оценить постоянные во времени эффекты, а также переменные, не изменяющиеся по индивидам, потому что Q-преобразование исключает из уравнения все эти переменные. Если истинной моделью является двунаправленная модель с фиксированными эффектами, то оценивание методом наименьших квадратов объединенной модели принесет смещенные и несостоятельные оценки коэффициентов регрессии, так как объединенная модель не включает оба множества фиктивных переменных, тогда как оценки однонаправленных фиксированных эффектов игнорируют только временные фиктивные переменные. Если же эти временные фиктивные переменные являются статистически значимыми, то это будет приводить к смещению оценок однонаправленных фиксированных эффектов из-за пропущенных переменных в модели.

Для того чтобы определить, присутствуют ли в модели индивидуальные и временные эффекты одновременно или по отдельности, можно проверить следующие гипотезы.

### 1. Гипотеза об отсутствии индивидуальных и временных эффектов

Гипотеза  $H_0$ :  $\alpha_1 = \alpha_2 = \dots = \alpha_{N-1} = 0$  и  $\gamma_1 = \gamma_2 = \dots = \gamma_{T-1} = 0$ ;  
гипотеза  $H_1$ : не  $H_0$ .

Для проверки нулевой гипотезы вычисляется  $F$ -критерий

$$F_1 = \frac{(SS_c - SS) / (N + T - 2)}{SS / (NT - N - T + 1 - K)} \sim F[(N + T - 2), (NT - N - T + 1 - K)], (9.83)$$

где  $SS$  — сумма квадратов остатков в модели без ограничений  $y_{it} = x'_{it} \beta + \alpha_i + \gamma_t + \varepsilon_{it}$ ,  $SS_c$  — сумма квадратов остатков в объединенной модели  $y_{it} = x'_{it} \beta + \alpha + \varepsilon_{it}$ .



В случае если вычисленное значение  $F$ -критерия окажется больше критического значения для соответствующих степеней свободы и уровня значимости, то мы можем отклонить нулевую гипотезу и принять альтернативную, в противном же случае у нас не будет оснований для того чтобы отвергнуть нулевую гипотезу.

В нашем примере суммы квадратов остатков в модели без ограничений и в объединенной модели будут соответственно равны  $SS = 712,74$  и  $SS_C = 1784,19$ . Тогда  $F_1 = 2,99$  при критическом значении  $F$ -критерия на 5%-ном уровне значимости  $F_{кр}(872, 1734) = 1,10$ . Так как  $F_1 > F_{кр, 0,05}$ , то мы можем отклонить нулевую гипотезу об отсутствии индивидуальных и временных эффектов и принять альтернативную гипотезу. Мы можем сделать вывод о том, что индивидуальные эффекты и временные эффекты являются совместно статистически значимыми на 5%-ном уровне значимости.

## 2. Гипотеза об отсутствии временных эффектов

Гипотеза  $H_0$ :  $\alpha_1 = \alpha_2 = \dots = \alpha_{N-1} = 0$  и  $\gamma_t \neq 0$  для всех  $t = 1, \dots, T-1$ ; гипотеза  $H_1$ : не  $H_0$ .

Нулевая гипотеза об отсутствии временных эффектов проверяется при помощи  $F$ -критерия

$$F_2 = \frac{(SS_N - SS)/(T-1)}{SS/(NT - N - T + 1 - K)} \sim F[(T-1), (NT - N - T + 1 - K)], \quad (9.84)$$

где  $SS_N$  — сумма квадратов остатков в модели с индивидуальными фиктивными переменными  $y_{it} = x'_{it} \beta + \alpha_i + \varepsilon_{it}$ .

В нашем примере сумма квадратов остатков в модели с индивидуальными фиктивными переменными  $SS_N = 753,46$ . Вычисленное значение  $F$ -критерия составляет  $F_2 = 49,53$  при критическом значении на 5%-ном уровне значимости  $F_{кр}(2, 1734) = 3,00$ . Так как  $F_2 > F_{кр, 0,05}$ , то мы можем отклонить нулевую гипотезу о том, что временные эффекты одновременно равны нулю, и принять альтернативную гипотезу. Получается, что временные эффекты являются статистически значимыми на 5%-ном уровне значимости.

## 3. Гипотеза об отсутствии индивидуальных эффектов

$H_0$ :  $\gamma_1 = \gamma_2 = \dots = \gamma_{T-1} = 0$  и  $\alpha_i \neq 0$  для всех  $i = 1, \dots, N-1$ ; гипотеза  $H_1$ : не  $H_0$ .

Нулевую гипотезу можно проверить при помощи  $F$ -критерия

$$F_3 = \frac{(SS_T - SS)/(N-1)}{SS/(NT - N - T + 1 - K)} \sim F[(N-1), (NT - N - T + 1 - K)], \quad (9.85)$$

где  $SS_T$  — сумма квадратов остатков в модели с временными фиктивными переменными  $y_{it} = x'_{it}\beta + \gamma_t + \varepsilon_{it}$ .

В нашем примере сумма квадратов остатков в модели с временными фиктивными переменными  $SS_T = 1737,67$ . Тогда вычисленное значения  $F$ -критерия  $F_3 = 2,87$  при критическом значении  $F$ -критерия на 5%-ном уровне значимости  $F_{кр}(869, 1734) = 1,10$ . Так как  $F_3 > F_{кр, 0,05}$ , то мы можем отклонить нулевую гипотезу об отсутствии индивидуальных эффектов и принять альтернативную гипотезу о статистической значимости индивидуальных эффектов на 5%-ном уровне значимости.

Таким образом, мы можем сделать вывод о присутствии как индивидуальных, так и временных эффектов.

## Контрольные вопросы и задания

1. Какие данные называются панельными?
2. Назовите преимущества использования панельных данных.
3. Какие панельные данные называются сбалансированными?
4. Каким образом можно получить внутригрупповые и межгрупповые оценки?
5. Назовите недостатки применения модели с фиксированными эффектами.
6. Какими способами можно оценить модель с фиксированными эффектами?
7. Что общего у внутригрупповых оценок и МНК-оценок с фиктивными переменными?
8. Какие предположения об ошибках делаются в модели со случайными эффектами?
9. Чем модель со случайными эффектами отличается от модели с фиксированными эффектами?

10. Какое преобразование используется в модели со случайными эффектами, для того чтобы получить оценки доступного обобщенного МНК?
11. В каком случае оценки обобщенного МНК будут совпадать с оценками обычного МНК в модели со случайными эффектами?
12. В чем состоят особенности применения меры качества подгонки в моделях с панельными данными?
13. Каким образом можно проверить статистическую значимость индивидуальных фиксированных эффектов?
14. Как связаны между собой  $LM$ -статистика в тесте множителей Лагранжа Бреуша — Пагана и  $g$ -статистика в тесте Хонды?
15. Какая гипотеза проверяется при помощи теста Хаусмана?
16. Какие тесты можно провести для проверки статистической значимости эффектов в двунаправленной модели с фиксированными эффектами?

---

## Выдающиеся ученые

**Пирсон Карл** (1857—1936) — выдающийся английский математик, статистик, биолог, философ (позитивист). Родился 27 марта 1857 г. в Лондоне. Профессор прикладной математики и механики (с 1884), а затем евгеники (с 1911) Лондонского университета. Основатель знаменитой биометрической школы.

Внес значительный вклад в развитие математической статистики. Автор фундаментальной философской работы «Грамматика науки» (1892). Основная философская позиция Пирсона формулируется следующим образом: понятия науки — искусственные конструкции, средства описания и упорядочивания чувственного опыта; правила связи их в научные предложения вычленяются грамматикой науки, которая и является философией науки. Связать же разнородные понятия и явления позволяет универсальная дисциплина — прикладная статистика.

Многие построения К. Пирсона напрямую связаны или разрабатывались с использованием антропологических материалов. Им разработаны многочисленные способы нумерической классификации и статистические критерии, применяемые во всех областях науки.

В период 1891—1894 гг. К. Пирсон работал в Грэхэм-колледже, где им были прочитаны несколько курсов лекций, основанные как на идеях «Грамматики науки», так и связанные с развитием статистики. Два обстоятельства способствовали пробуждению его интереса к статистике: знакомство

в 1889 г. с книгой Ф. Гальтона «Природа наследственности» и появление в Университетском колледже В. Велдона (1860—1906), профессора зоологии, с которым он быстро подружился. Исследование данных о креветках и прибрежных крабах, собранных Велдоном, привели к написанию серии важных статей, появившихся начиная с 1893 г. Уже первая из них была связана с методом моментов при поиске кривой, наиболее пригодной для аппроксимации асимметричных распределений. Эти статьи впоследствии были собраны в книгу «Вклад в математическую теорию эволюции», в которой введено понятие «семейство кривых распределений Пирсона». Вместе с В. Велдоном и Ф. Гальтоном он создал биометрию как самостоятельное научное направление.

К. Пирсон — основатель и издатель журнала *Biometrika* (1901 г.). В 1903 г. Пирсон основал Биометрическую лабораторию. К нему стекалось множество посетителей из разных стран, включая математика В. Госсета (Стьюдента) из Дублина, экономистов Г. Мура, Ф. Миллса, В. Митчела, И. Фишера и др. из США, Кикучи Дайроку из Японии. Много сил Пирсон отдал построению статистических таблиц. В 1911 г., согласно пожеланию Ф. Гальтона, Пирсон возглавил департамент прикладной статистики, который включал лаборатории по евгенике и биометрике в Университетском колледже. Это был настоящий центр статистической мысли, в который стремились попасть талантливые студенты. После Первой мировой войны 1914—1918 гг. статистика продолжила свое победоносное развитие, и появились новые таланты, такие, как Дж. Нейман, Р. А. Фишер. К сожалению, Пирсон был нетерпим к критике, что привело к разрыву сотрудничества и уходу Р. А. Фишера из департамента.

К. Пирсон был удостоен высшей математической награды — золотой медали Гая (награда присуждается Королевским статистическим обществом), но отказался от нее, сказав, что пусть лучше награждают молодых. Карл Пирсон умер 27 апреля 1936 г. в Колдхарбуре (*Coldharbour*), Суррее (Юго-Восточная Англия). Возглавлявшийся им департамент был разделен на департамент евгеники, который возглавил Р. А. Фишер (1890—1962), и департамент статистики, который возглавил сын Карла Пирсона — Эгон Пирсон (по данным сайта <http://www.ekonominform.ru/>).

**Гаусс Карл Фридрих** (1777—1855) родился в Брауншвейге, Германия, где вырос в простой семье (отец — каменщик и садовник, мать не умела даже писать). Однако ему суждено было стать одним из величайших математиков. В раннем возрасте мальчик показал себя вундеркиндом. В три года он умел читать и писать, даже исправлял счетные ошибки отца. В школе быстрота и точность его расчетов в уме поражали учителей арифметики. Вскоре Гаусс изучил всевозможные доступные учебники и перешел к иным источникам.

Исключительный талант юного Гаусса оценил герцог Брауншвейга и стал его наставником, отправив в университет Геттингема. Будучи студентом, Гаусс доказал возможность построения с помощью циркуля и линейки правильного семнадцатиугольника, подобное не удавалось математикам на протяжении 2000 лет. Вскоре после этого, в 1801 г. Гаусс опубликовал «*Disquisitiones Arithmeticae*» («Арифметические исследования»), произведение, которое перевело теорию чисел на более высокий уровень и выдвинуло Гаусса как первоклассного математического гения наравне с Архимедом и Ньютоном.

За первой великой работой последовала вторая, написанная в 1809 г., — «*Theoria Motus Corporum Coelestium*» («Теория движения небесных тел»), в которой Гаусс разработал методы вычисления орбит небесных тел. Именно здесь он впервые применил два наиболее распространенных инструмента, используемых во всех науках в наше время, — метод наименьших квадратов и кривую нормального распределения. Несмотря на то, что Гаусс признавал приоритет Лапласа, вскоре эта кривая получила известность как кривая Гаусса и до сих пор так называется во многих странах.

К тридцати годам Гаусс стал руководителем Геттингенской обсерватории, где работал до конца своей жизни. Он внес большой вклад в развитие не только астрономии, но и математики и физики. Сотни публикаций свидетельствуют о том, что Гаусс был человеком, который удивительно сочетал чистый математический интерес, выраженный в абстрактных идеях, интерес теоретического физика в создании математических моделей, талант астронома для проведения наблюдений и навыки экспериментатора в создании и применении измерительных методов (по данным сайта: «Википедия» <http://ru.wikipedia.org/wiki>).

**Фишер Роналд Эйлмер** (1890—1962) — английский статистик и генетик, один из основателей математической статистики и математической популяционной генетики. Член Лондонского королевского общества (1929). Окончил колледж в Кембридже (1912). Работал статистиком в «Меркантайл энд дженерал инвестмент компани» (1913—1915). В 1919—1933 гг. работал в отделе статистики Ротемстедской экспериментальной станции. В 1933—1943 гг. — профессор евгеники Лондонского университета. В 1943—1957 гг. — профессор генетики Кембриджского университета, в 1956—1959 гг. руководил одним из его колледжей. Основные труды по теории статистики и генетической теории эволюции. Ввел понятие достаточной статистики, построил теорию точечных и интервальных статистических оценок, разработал методику планирования экспериментов и внес существенный вклад в создание современной теории статистической проверки гипотез; создатель дисперсионного анализа. Ввел основные понятия генетики количественных признаков, исследовал стохастические процессы в популяциях, предложил ряд моделей действия естественного отбора, первый рассмотрел случай сверждоминирования по приспособленности, предложил теорию эволюции доминантности. Сформулировал так называемую фундаментальную теорему естественного отбора, носящую его имя. Его именем Дж. Снедекор назвал статистический критерий (*F*-критерий).

Член Королевского статистического общества. Почетный член многих академий, английских и иностранных научных обществ; почетный доктор наук и доктор права многих университетов (Лондон, Гарвард, Чикаго, Калькутта, Глазго и др.). Награжден Дарвиновской медалью Лондонского королевского общества (по данным сайта [http://bse.sci-lib.com/article116\\_512.html](http://bse.sci-lib.com/article116_512.html)).

**Маленво Эдмон** — французский экономист, статистик. Родился 25 апреля 1923 г. в Лиможе (Франция). Получил образование в учебных заведениях Лиможа и Парижа. В 1948 г. получил должность статистика в Национальном институте статистики и экономических исследований (*Institut National de la Statistique et des Etudes Economiques*). Работал в Чикагском (1951) и Калифорнийском, г. Беркли (1961, 1967), университетах. С 1957 г. — ректор Школы прикладных исследований

(*Ecole Pratique des Hautes Etudes*), а с 1974 г. он также возглавил Национальный институт статистики (*Institut National de la Statistique*).

Принимал активное участие в работе французских и международных научных организаций, органов государственного управления Франции. Награжден высшей государственной наградой Франции — Орденом Почетного Легиона (1986). Является членом редколлегии журнала *De Economist*.

Область научных интересов — статистика, эконометрика, микро- и макроэкономика. Мировую известность приобрела его книга «Статистические методы в экономике» (*Statistical Methods of Economics*, 1964 г.), которая быстро стала наиболее цитируемой публикацией по эконометрике. Большое внимание Маленво уделял исследованию проблем макроэкономики, в частности рынка труда и безработицы.

Основные труды: «Статистические методы в эконометрике» (*Statistical Methods in Econometrics*, 1964 г.); «Лекции по микроэкономической теории» (*Leçons de Théorie Microéconomique*, 1969 г.); «Макроэкономическая теория» в 2 т. (*Théorie Macroéconomique*, 2 volumes, 1981—1982 гг.) (по данным сайтов: «Галерея экономистов» <http://gallery.economicus.ru>, «Википедия» <http://ru.wikipedia.org/wiki>).

**Фриш Рагнар Антон Киттил** родился в Осло (Норвегия) 3 марта 1895 г. Изучать экономику начал в Университете Осло, где получил степень бакалавра. Свое образование он продолжил в учебных заведениях Франции, Германии, Великобритании, Италии и США. В 1925 г. вернулся в Осло. В 1931 г. назначен на должность действительного профессора социальной экономики и статистики и одновременно директора Института социальной экономики при Университете Осло. На этом посту Р. Фриш оставался вплоть до своего ухода в отставку. Помимо научной деятельности, он принимал активное участие в работе различных международных организаций, был экспертом и советником правительства Норвегии и других стран, в частности Индии и Египта.

Р. Фриш считал, что описание развития экономических процессов должно опираться на математические модели, т.е. экономическая теория должна быть точной наукой. Эти взгляды он изложил в своей первой опубликованной работе — статье



«О проблеме чистой экономической теории» (*Sur un probleme d'economie pure*, 1926 г.).

Круг вопросов, которые исследовал Р. Фриш, очень широк. В частности, в 1920—1930 гг. он изучал проблему измерения эластичности в полной системе функций спроса, что в дальнейшем нашло отражение в статье «Полная схема для измерения всех прямых и перекрестных эластичностей спроса в многосекторной модели» (*A Complete Scheme for Computing all Direct and Cross Demand Elasticities in a Model with Many Sectors*, 1959 г.). Эти исследования легли в основу модели, использованной министерством финансов Норвегии для долгосрочного планирования развития экономики.

В 1932 г. Р. Фриш опубликовал работу «Анализ статистических пересечений в системе полной регрессии» (*Statistical Confluence Analysis by Means of Complete Regression Systems*), в которой исследовал мультиколлинеарность переменных, определенную им как тенденцию многих переменных к совместному движению из-за их подчиненности общим тенденциям, циклам и другим общим закономерностям.

Р. Фриш предложил новый подход к теории производства, заложив основу современной неоклассической теории; в начале 1930-х гг. активно изучал проблемы макроэкономического анализа, экономической динамики и экономических циклов. В его статье «Проблемы распространения и проблемы импульса в динамичной экономике» (*Propagation Problems and Impulse Problems in Dynamic*, 1933 г.) содержалась оригинальная и математически законченная модель экономического цикла, впервые были введены в оборот термины «микроэкономика» и «макроэкономика».

Р. Фриш предложил одну из первых методик расчета национального дохода; занимался разработкой теории экономической политики. Модели ученого, базирующиеся на описании экономики по типу «затраты-выпуск», были использованы в качестве аналитических инструментов при расчете национального дохода Норвегии.

Р. Фриш и Я. Тинберген стали первыми лауреатами Нобелевской премии по экономике, которая была им присуждена в 1969 г. «за развитие и применение динамических моделей к анализу экономических процессов».

Рагнар Фриш скончался в Осло 31 января 1973 г. В его честь с 1978 г. Эконометрическое общество присуждает одну из престижнейших экономических наград — медаль Фриша (по данным сайтов: «Галерея экономистов» <http://gallery.economicus.ru>, «Википедия» <http://ru.wikipedia.org/wiki>).

**Бокс Джордж Эдвард Пелхам** родился 18 октября 1919 г. в Гревсенде (Англия). Химик по первому образованию, он во время Второй мировой войны проводил исследования для британской армии. В ходе этих исследований Бокс искал статистика, который мог бы помочь ему в обработке данных химических экспериментов. Не найдя такого человека, Бокс вынужден был самостоятельно осваивать статистику. После войны получил степень бакалавра математики и статистики, а в 1953 г. в Лондонском университете получил степень *Ph.D.*

С 1948 по 1956 г. Бокс работал статистиком в Имперской химической индустрии, из которой перешел в Принстонский университет на должность директора исследовательской статистической группы.

В 1960 г. перешел в Университет Висконсин-Мэдисона для создания департамента статистики, где вместе с Биллом Хантером в 1984 г. основал центр качества и инновационных технологий. В 1992 г. ушел на пенсию, оставаясь Почетным профессором Университета.

Джордж Бокс избирался президентом Американского статистического общества в 1978 г. и Института математической статистики в 1979 г. Был избран членом Американской академии науки и искусств в 1974 г. и членом Королевского общества в 1979 г.

Джордж Бокс внес значительный вклад в развитие контроля качества, анализ временных рядов и методику проектирования экспериментов (по материалам сайта <http://en.wikipedia.org>).

**Дженкинс Вильям Мэрион** родился в 1933 г. в Говертауне (Уэльс). Получил степень магистра математики в 1953 г. и степень *Ph.D.* в 1956 г. в Лондонском университете. Работал приглашенным лектором и занимал позицию профессора в Имперском Колледже Лондона, Стэнфордском университете, Принстонском университете и университете Висконсин-Мэдисона. В 1956 г. стал профессором системного проектирования в Ланкастерском университете.

В 1974 г. Вильям Дженкинс основал собственную консалтинговую компанию.

Дженкинс служил исследователем и советником в Королевском статистическом обществе с 1960 г. Основал журнал системного проектирования в 1969 г.

Вильям Дженкинс был разносторонне развитым человеком. Кроме науки искренне интересовался музыкой. Был поклонником джаза и блюза, играл на пианино (по материалам сайта <http://en.wikipedia.org>).

**Йохансен Сорен** родился 6 ноября 1939 г. В 1964 г. окончил Университет Копенгагена по специальности математическая статистика. С 1964 г. работал в Институте математической статистики Университета Копенгагена. В 1974 г. за работу, посвященную цепям Маркова, получил звание доктора философии. С 1989 г. — профессор Университета Копенгагена. С 1996 по 2001 г. Сорен Йохансен занимал должность профессора эконометрики департамента экономики Университета во Флоренции.

В 1976 г. Сорен Йохансен был избран членом Международного статистического общества, членом Эконометрического общества.

С 1974 по 1981 г. Йохансен был редактором журнал «Летописи статистики» (*Annals of Statistics*) С 1985 по 1990 гг. — редактор Скандинавского статистического журнала. С 1990 г. — редактор журнала «Эконометрическая теория», с 1997 г. — редактор журнала «Эконометрика».

Сорен Йохансен внес значительный вклад в развитие математической статистики, теории вероятностей, эконометрики и медицинской статистики (по материалам автобиографии с сайта <http://math.ku.dk>).

**Энгл Роберт Фрай** — американский экономист, лауреат Нобелевской премии по экономике 2003 г. «за методы анализа экономических временных рядов с изменяющейся во времени волатильностью».

Родился 10 ноября 1942 г. в Сиракузах. Получил степень магистра физики в 1966 г. и степень *Ph.D* экономики в 1969 г. в Корнельском университете. С 1969 по 1977 гг. — профессор Массачусетского технологического института. С 1975 г. работал в Калифорнийском университете. Почетный профессор

Калифорнийского университета, преподает в Нью-Йоркском университете.

Энгл внес значительный вклад в разработку методов анализа изменений цен на финансовых рынках и процентных ставок. Им была разработана новая статистическая модель, учитывающая тенденции биржевых цен и других финансовых переменных чередовать высокую и низкую дисперсии («Авторегрессионная условная гетероскедастичность: ARCH»).

Является одним из основателей Финансового эконометрического общества (по материалам сайта <http://en.wikipedia.org>).

**Хаусман Джерри А.** родился 5 мая 1946 г. в Уэйртоне, Западная Вирджиния, США. В 1968 г. с отличием окончил Брауновский университет и получил степень бакалавра искусств. После двух лет службы в армии продолжил свое образование в Оксфордском университете, где получил степень бакалавра философии в 1972 г. и степень доктора философии (*Ph.D*) в 1973 г. Преподает на факультете экономики в Массачусетском технологическом институте с 1973 г. (профессор — с 1979 г., профессор им. Джона и Дженни МакДональд — с 1992 г.). Хаусман занимает пост директора программы исследования экономики телекоммуникаций в Массачусетском технологическом институте.

Хаусман имеет большое количество публикаций в ведущих экономических журналах. Его основные научные интересы лежат в области эконометрики, государственных финансов и регулирования, прикладных моделей в микроэкономике, включая работы по дифференцированным продуктам, телекоммуникации, налогообложению, возрастной экономике, экономике энергии и окружающей среды. В эконометрике широко известен тест спецификации Хаусмана, опубликованный им в 1978 г.

Хаусман был членом редакционных коллегий целого ряда ведущих научных журналов по экономике, таких, как *Bell Journal of Economics* (1974—1983), *Rand Journal of Economics* (1984—1988), *Econometrica* (1978—1987), *Review of Economic Studies* (1979—1982), *Journal of Public Economics* (1982—1998) и *Journal of Applied Econometrics* (1985—1993); в настоящее время является членом редакционных коллегий журналов *Economics Research Network and Social Science Research* (с 1998 г.),

*Journal of Sports Economics* (с 1999 г.), *Journal of Competition Law & Economics* (с 2004 г.) и *Journal of Applied Economics* (с 2005 г.).

Джерри Хаусман и Дэвид Уайз в 1980 г. получили медаль Фриша от Эконометрического общества за работу, посвященную ошибкам в экспериментальных и панельных данных, связанных с выбыванием или исключением объектов из исследования. В 1985 г. Хаусман получил награду Джона Бейтса Кларка от Американской экономической ассоциации за выдающийся вклад в экономику, внесенный экономистами в возрасте до 40 лет.

---

## Литература

1. *Айвазян, С. А.* Прикладная статистика и основы эконометрики/С. А. Айвазян, В. С. Мхитарян. — М. : ЮНИТИ, 1998.
2. *Берндт, Э. Р.* Практика эконометрики: классика и современность : учебник / Э. Р. Берндт ; пер. с англ. под ред. проф. С. А. Айвазяна. — М. : ЮНИТИ-ДАНА, 2005.
3. *Вербик, М.* Путеводитель по современной эконометрике/М. Вербер ; пер. с англ. Б. А. Банникова : научн. ред. и предисл. С. А. Айвазяна. — М. : Научная книга, 2008.
4. *Вучков, И.* Прикладной регрессионный анализ/И. Вучков, Л. Бояджиева, Е. Солаков: пер. с болг. — М. : Финансы и статистика, 1987.
5. *Дрейпер, Н.* Прикладной регрессионный анализ/Н. Дрейпер, Г. Смит ; пер. с англ. — 3-е изд. — М. : ИД «Вильямс», 2007.
6. *Дуброва, Т. А.* Статистические методы прогнозирования/Т. А. Дуброва. — М. : ЮНИТИ, 2003.
7. *Кремер, Н. Ш.* Эконометрика: учебник для вузов / Ш. Н. Кремер, Б. А. Путко ; под ред. проф. Н. Ш. Кремера. — М. : ЮНИТИ-ДАНА, 2003.
8. *Ланге, О.* Введение в эконометрику : пер. с пол. /О. Ленге. — М. : Прогресс.
9. *Льюис, К. Д.* Методы прогнозирования экономических показателей/К. Д. Льюис ; пер. с англ. — М. : Финансы и статистика, 1986.
10. *Магнус Я. Р.* Эконометрика: начальный курс : учебник /Я. Р. Магнус, П. К. Катышев, А. А. Пересецкий. — 6-е изд., перераб. и доп. — М. : Дело, 2004.
11. *Маленво, Э.* Статистические методы эконометрии / Э. Маленво ; пер. с фр. — М. : Статистика, 1975.

12. Носко, В. П. Эконометрика. Элементарные методы и введение в регрессионный анализ временных рядов/В. П. Носко. — М. : ИЭПП, 2004.
13. Справочник по прикладной статистике : в 2 т. /под ред. Э. Ллойда, У. Ледермана, Ю. Н. Тюрина ; пер. с англ. — М. : Финансы и статистика, 1989.
14. Тинтнер, Г. Введение в эконометрию /Г. Тинтнер ; пер. с нем. — М. : Статистика, 1965.
15. Тихомиров, Н. П. Эконометрика : учебник/Н. П. Тихомиров, Е. Ю. Дорохина. — М. : Экзамен, 2003.
16. Четыркин, Е. М. Статистические методы прогнозирования/Е. М. Четыркин. — М. : Статистика, 1975.
17. Эконометрика : учебник/под ред. д-ра экон. наук., проф. В. С. Мхитаряна. — М. : Проспект, 2008.
18. Baltagi B. H. *Econometric Analysis of Panel Data* / B. H. Baltagi. — 3rd Edition. — Chichester : John Wiley & Sons, Ltd, 2005.
19. Brooks C. *Introductory Econometrics for Finance* / C. Brooks. — Cambridge Academ, Cambridge, 2002.
20. Breusch T. S. A Simple Test for Heteroskedasticity and Random Coefficient Variation / T. S. Breusch, A. R. Pagan. — *Econometrica*. 1979. Vol. 47.
21. Cameron A. C. *Microeconometrics: Methods and Applications* / A. C. Cameron, P. K. Trivedi. — Cambridge : Cambridge University Press, 2005.
22. Fair R. C. The Estimation of Simultaneous Equation Models with Lagged Endogenous Variables and First Order Serially Correlated Errors / R. C. Fair. — *Econometrica*. 1970 May. Vol. 38. Pp. 507—516.
23. Greene W. H. *Econometric Analysis* / W. H. Greene. — 4th Edition. — New Jersey : Prentice Hall, 2002.
24. Hamilton J. D. *Time Series Analysis* / J. D. Hamilton. — Princeton. Princeton University Press, 1994.
25. Hsiao C. *Analysis of Panel Data* / C. Hsiao. — 2nd Edition. — Cambridge : Cambridge University Press, 2003.
26. Hausman J. A. Specification Tests in Econometrics / J. A. Hausman. — *Econometrica*. — 1978. Vol. 46.
27. Honda Y. Testing the Error Components Model with Non-normal Disturbances / Y. Honda. — *Review of Economic Studies*. 1985. Vol. 52.
28. Wooldridge J. M. *Econometric Analysis of Cross Section and Panel Data* / J. M. Wooldridge. — The MIT Press, 2002.

29. *Johnston J.* *Econometric Methods* / J. Johnston, J. DiNardo. — 4th Edition. — New York : McGraw-Hill, 1997.
30. *Mátyás L.* *The Econometrics of Panel Data. Fundamentals and Recent Developments in Theory and Practice* / L. Mátyás, P. Sevestre (eds.). — 3rd Edition. — Berlin : Springer, 2008.
31. *Verbeek M.* *A Guide to Modern Econometrics* / M. Verbeek. — 2nd Edition. — Chichester : John Wiley & Sons, Ltd, 2004.
32. *Wallace T. D.* *The Use of Error Components Models in Combining Cross-section and Time-series Data* / T. D. Wallace, A. Hussain. — *Econometrica*, 1969. Vol. 37.
33. *Wooldridge J. M.* *Econometric Analysis of Cross Section and Panel Data* / M. J. Wooldridge. — Cambridge, MA : MIT Press, 2002.



---

# Предметный указатель

Ссылки по главе 1: Возникновение и развитие эконометрики.

## Парная регрессия

- Эконометрика: история возникновения 9—15
- Непараметрическая эконометрика 14
- Специфика эконометрических исследований 15—16
- Эконометрические методы 16
- Критика и апологетика эконометрики 17—23
- Модель парной регрессии 23
- Регрессионные остатки 24
- Значимость параметров парной регрессии 24
- Коэффициент эластичности по парной регрессии 27
- Таблица дисперсионного анализа по парной регрессии 27
- Сумма квадратов объясненная (факторная) 28
- Сумма квадратов остаточная 28
- Свойства остатков уравнения регрессии 30—38
- Гетероскедастичность остатков 31
- Прогнозирование по уравнению парной регрессии 35—37

## Ссылки по главе 2: Множественная регрессия

- Регрессия множественная 39
- Оценки несмещенные 41
- Оценки эффективные 41
- Оценки состоятельные 41—42
- Предпосылки МНК 42
- Метод наименьших квадратов для множественной регрессии 43—44
- Условия Гаусса — Маркова 43

- Коэффициенты эластичности 48—49  
Стандартизованные коэффициенты регрессии 50  
Уравнение регрессии в стандартизованном масштабе 51  
Геометрическая интерпретация МНК 52—53  
Теорема Гаусса — Маркова 53  
Коэффициент множественной детерминации 58  
Оценка дисперсии случайных остатков по МНК 57  
Коэффициент множественной корреляции 58  
Оценка значимости параметров множественной регрессии 60—64  
Скорректированный коэффициент детерминации 59  
Оценка значимости уравнения множественной регрессии 65—68  
Таблица дисперсионного анализа 66  
Ограничения на параметры множественной регрессии 68—75  
Производственная функция Кобба — Дугласа 76—77  
Тест Бокса — Кокса 83—86  
Преобразование Зарембки 88—89  
Метод максимального правдоподобия 90—97  
Оценка дисперсии случайных остатков по ММП 97—99  
Прогнозирование по множественной регрессии 101  
Мультиколлинеарность — понятие 101  
Интеркорреляция 101  
Мультиколлинеарность — признаки 101  
Мультиколлинеарность — измерение 101—104  
Показатель (число) обусловленности матрицы 104  
Фактор инфляции (вздутия) дисперсии (*VIF*) 103  
Ридж-регрессия 107—109  
Метод главных компонент 109—111  
Гетероскедастичность 114  
Гомоскедастичность 114  
Тест Гольдфелда — Квандта 114  
Тест Глейзера 116  
Тест Парка 116  
Тест Уайта 116  
Тест Бреуша — Пагана 117  
Тест ранговой корреляции Спирмена 118  
Обобщенный метод наименьших квадратов 124—128  
Обобщенная линейная модель 124  
Теорема Айткена 125  
ОМНК, доступный 128  
Метод взвешенных наименьших квадратов 129

**Ссылки по главе 3: Фиктивные переменные**

- Фиктивные переменные — определение 142
- Модели с фиктивными переменными сдвига 145—150
- Модели с фиктивными переменными наклона 150—152
- Фиктивные переменные — общий вид модели 152—158
- ANOVA модель 157
- ANCOVA модель 157
- Тест Чоу 158—161

**Ссылки по главе 4: Системы эконометрических уравнений**

- Переменные эндогенные 162
- Переменные экзогенные 163
- Переменные лаговые 163
- Переменные predetermined 163
- Система независимых эконометрических уравнений 163
- Система рекурсивных уравнений 164
- Система одновременных уравнений 164
- Структурная форма модели 164
- Уравнение идентифицированное 165
- Уравнение неидентифицированное 165
- Уравнение сверхидентифицированное 166
- Модель идентифицированная 166—167
- Модель неидентифицированная 167
- Модель сверхидентифицированная 167
- Тест Хаусмана — Ву на экзогенность 173—174
- Уравнения, кажущиеся несвязанными 174—176

**Ссылки по главе 5: Моделирование изолированного динамического ряда**

- Динамические ряды 178
- Временные ряды 178
- Тренд 179, 180, 194—217
- Тенденция 179, 194—206
- Колебания периодические 179, 182
- Колебания случайные 179, 183
- Колебания циклические 183—184
- Колебания сезонные 184—185
- Модель аддитивная 186—193
- Модель мультипликативная 190
- Автокорреляция уровней ряда 186—193

Автокорреляционная функция (АКФ) 190  
Коррелограмма 191, 193  
Кривые S-образные 202—206  
Кривая логистическая 202—204, 212—214  
Кривая Гомперца 204—205, 214—217  
Кривая Перла — Рида 202—203, 211  
Автокорреляция остатков 43, 128, 219  
Коэффициент автокорреляции остатков 219  
Критерий Дарбина — Уотсона 221—226  
Ряд Фурье 231—242  
Ряд стационарный 179, 231  
Частота динамического ряда 232  
Амплитуда временного ряда 232  
Фаза 232  
Гармоника 233, 237  
Модель сезонности аддитивная 242—249  
Скользящие средние 246—251  
Фиктивные переменные 243, 247, 252  
Модель сезонности мультипликативная 250—253

#### Ссылки по главе 6: Модели регрессии по временным рядам

Ложная корреляция 255  
— регрессия 255  
Метод последовательных разностей 257, 258—261  
Метод отклонений от тренда 257, 261—264  
Фактор времени 264—270  
МНК обобщенный 270—281  
Поправка Прайса — Уинстена 272  
Двухшаговая процедура Дарбина 277  
Итеративная процедура Кохрейна — Оркатта 278  
Сезонность в модели регрессии 281—286

#### Ссылки по главе 7: Модели с лаговыми переменными

Переменная лаговая 287, 288  
Лаг 289  
Модель динамическая 288  
Модель с распределенными лагами 288, 290—306  
Модель авторегрессии 288, 306—313  
Модель с распределенными лагами авторегрессионная 288  
Мультипликатор промежуточный 291  
Мультипликатор долгосрочный 291

Мультипликатор краткосрочный 291  
Средняя величина лага 292  
Лаг медианный 293  
Лаги Алмон 294—302  
Преобразование Койка 294, 303  
Метод Койка 302—306  
Метод инструментальных переменных 307—311  
Критерий Дарбина — Уотсона 311  
 $h$ -статистика Дарбина 312—313  
Авторегрессионный процесс (AR) 313—315  
Стационарность 314, 322  
Модель скользящей средней (MA) 316—317  
Авторегрессионный процесс со скользящими средними в остатках (ARMA) 317—319  
Частная автокорреляционная функция 317  
Авторегрессионная интегрированная модель скользящей средней (ARIMA) 319—321

#### Ссылки по главе 8: Модели ARMA, ARIMA, ARCH, GARCH

Строгая стационарность ряда 323  
Слабо стационарный процесс 323  
Белый шум 325—326  
Случайное блуждание 326—327  
Модель скользящей средней 327—328  
Авторегрессионная модель 328—330  
Условие стационарности AR ( $q$ ) процесса 330  
Теорема декомпозиции Вольда 331—332  
Условие обратимости 331  
Частная автокорреляционная функция 332—334  
ARMA модель 334—347  
ARMA процесс: стационарность 335, 343  
ARMA процесс: обратимость 335  
Тест Дики — Фуллера: простой 343  
Тест Дики — Фуллера: расширенный (ADF) 344  
Подход Бокса — Дженкинса 344—345  
Информационные критерии при идентификации модели ARMA 345  
Критерий Акайке (AIC) 346  
Критерий Шварца 346  
Критерий: информационный критерий Ханнан — Куина 346

- ARIMA* модель 347—355  
Коинтеграция 355—360  
Энгла и Грэнжера подход 356  
—процедура 359  
Модель корректировки ошибок 357  
*ARCH* модель 360—364  
*ARCH-M* модель 362—363  
*GARCH* модель 360, 364—373  
*IGARCH* модель 367
- Ссылки по главе 9: Анализ панельных данных**  
Данные панельные, определение 374  
Данные лонгитюдные 374  
Модель объединенная 381, 382  
Модель с фиксированными эффектами 381, 382—395  
Модель со случайными эффектами 381, 395—406  
Эффекты фиксированные 383  
Модель внутригрупповая 383  
Модель межгрупповая 384  
Оценка внутригрупповая 385, 386—388  
Оценка межгрупповая 385  
Оценка с фиксированными эффектами 386  
Кронекера произведение 389  
Модель с фиксированными эффектами, МНК-оценки с фиктивными переменными 388—392  
Матрица идемпотентная 390  
Модель с фиксированными эффектами, оценка условным методом максимального правдоподобия 392—393  
Модель первых разностей 394  
Эффекты случайные 395  
Метод наименьших квадратов обобщенный 396, 398  
Метод наименьших квадратов обобщенный доступный 400  
Оценка со случайными эффектами 396—400  
Качество подгонки 406—409  
*R*-квадрат 407  
*R*-квадрат общий 408  
*R*-квадрат внутригрупповой 407  
*R*-квадрат межгрупповой 408  
Тест Чоу 411—412  
Тест Фишера 413—414  
Тест множителей Лагранжа Бреуша-Пагана 414—415  
Тест Хонды 415—416  
Тест Хаусмана 416—418  
Модель с фиксированными эффектами двунаправленная 419—423

---

# Приложения

### Критические значения $t$ -критерия на уровне значимости 0,01; 0,05; 0,1 (двусторонний)

Число степеней свободы	Уровень значимости		
	0,01	0,05	0,10
1	63,65674	12,7062	6,313752
2	9,924843	4,302653	2,919986
3	5,840909	3,182446	2,353363
4	4,604095	2,776445	2,131847
5	4,032143	2,570582	2,015048
6	3,707428	2,446912	1,94318
7	3,499483	2,364624	1,894579
8	3,355387	2,306004	1,859548
9	3,249836	2,262157	1,833113
10	3,169273	2,228139	1,812461
11	3,105807	2,200985	1,795885
12	3,05454	2,178813	1,782288
13	3,012276	2,160369	1,770933
14	2,976843	2,144787	1,76131
15	2,946713	2,13145	1,75305
16	2,920782	2,119905	1,745884
17	2,898231	2,109816	1,739607
18	2,87844	2,100922	1,734064
19	2,860935	2,093024	1,729133
20	2,84534	2,085963	1,724718
21	2,83136	2,079614	1,720743
22	2,818756	2,073873	1,717144
23	2,807336	2,068658	1,713872
24	2,796939	2,063899	1,710882
25	2,787436	2,059539	1,708141
26	2,778715	2,055529	1,705618
27	2,770683	2,05183	1,703288
28	2,763262	2,048407	1,701131
29	2,756386	2,04523	1,699127
30	2,749996	2,042272	1,697261
31	2,744042	2,039513	1,695519
32	2,738481	2,036933	1,693889

Число степеней свободы	Уровень значимости		
	0,01	0,05	0,10
33	2,733277	2,034515	1,69236
34	2,728394	2,032244	1,690924
35	2,723806	2,030108	1,689572
36	2,719485	2,028094	1,688298
37	2,715409	2,026192	1,687094
38	2,711558	2,024394	1,685954
39	2,707913	2,022691	1,684875
40	2,704459	2,021075	1,683851
41	2,701181	2,019541	1,682878
42	2,698066	2,018082	1,681952
43	2,695102	2,016692	1,681071
44	2,692278	2,015368	1,68023
45	2,689585	2,014103	1,679427
46	2,687013	2,012896	1,67866
47	2,684556	2,01174	1,677927
48	2,682204	2,010635	1,677224
49	2,679952	2,009575	1,676551
50	2,677793	2,008559	1,675905
60	2,660283	2,000298	1,670649
70	2,647905	1,994437	1,666914
80	2,638691	1,990063	1,664125
90	2,631565	1,986674	1,661961
100	2,625891	1,983971	1,660234
125	2,615733	1,979124	1,657135
150	2,609003	1,975905	1,655076
200	2,600634	1,971896	1,652508
250	2,595638	1,969498	1,650971
300	2,592316	1,967903	1,649949
400	2,588176	1,965912	1,648672
500	2,585698	1,96472	1,647907
600	2,584048	1,963926	1,647397
1000	2,580755	1,962339	1,646379



**Таблица значений  $F$ -критерия  
на уровне значимости 0,05**

$df_1$ $df_2$	1	2	3	4	5	6	8	12	24	$10^9$
1	161,45	199,50	215,71	224,58	230,16	233,99	238,88	243,91	249,05	254,31
2	18,51	19,00	19,16	19,25	19,30	19,33	19,37	19,41	19,45	19,50
3	10,13	9,55	9,28	9,12	9,01	8,94	8,85	8,74	8,64	8,53
4	7,71	6,94	6,59	6,39	6,26	6,16	6,04	5,91	5,77	5,63
5	6,61	5,79	5,41	5,19	5,05	4,95	4,82	4,68	4,53	4,36
6	5,99	5,14	4,76	4,53	4,39	4,28	4,15	4,00	3,84	3,67
7	5,59	4,74	4,35	4,12	3,97	3,87	3,73	3,57	3,41	3,23
8	5,32	4,46	4,07	3,84	3,69	3,58	3,44	3,28	3,12	2,93
9	5,12	4,26	3,86	3,63	3,48	3,37	3,23	3,07	2,90	2,71
10	4,96	4,10	3,71	3,48	3,33	3,22	3,07	2,91	2,74	2,54
11	4,84	3,98	3,59	3,36	3,20	3,09	2,95	2,79	2,61	2,40
12	4,75	3,89	3,49	3,26	3,11	3,00	2,85	2,69	2,51	2,30
13	4,67	3,81	3,41	3,18	3,03	2,92	2,77	2,60	2,42	2,21
14	4,60	3,74	3,34	3,11	2,96	2,85	2,70	2,53	2,35	2,13
15	4,54	3,68	3,29	3,06	2,90	2,79	2,64	2,48	2,29	2,07
16	4,49	3,63	3,24	3,01	2,85	2,74	2,59	2,42	2,24	2,01
17	4,45	3,59	3,20	2,96	2,81	2,70	2,55	2,38	2,19	1,96
18	4,41	3,55	3,16	2,93	2,77	2,66	2,51	2,34	2,15	1,92
19	4,38	3,52	3,13	2,90	2,74	2,63	2,48	2,31	2,11	1,88
20	4,35	3,49	3,10	2,87	2,71	2,60	2,45	2,28	2,08	1,84
21	4,32	3,47	3,07	2,84	2,68	2,57	2,42	2,25	2,05	1,81
22	4,30	3,44	3,05	2,82	2,66	2,55	2,40	2,23	2,03	1,78
23	4,28	3,42	3,03	2,80	2,64	2,53	2,37	2,20	2,01	1,76
24	4,26	3,40	3,01	2,78	2,62	2,51	2,36	2,18	1,98	1,73
25	4,24	3,39	2,99	2,76	2,60	2,49	2,34	2,16	1,96	1,71
26	4,23	3,37	2,98	2,74	2,59	2,47	2,32	2,15	1,95	1,69
27	4,21	3,35	2,96	2,73	2,57	2,46	2,31	2,13	1,93	1,67
28	4,20	3,34	2,95	2,71	2,56	2,45	2,29	2,12	1,91	1,65
29	4,18	3,33	2,93	2,70	2,55	2,43	2,28	2,10	1,90	1,64
30	4,17	3,32	2,92	2,69	2,53	2,42	2,27	2,09	1,89	1,62
31	4,16	3,30	2,91	2,68	2,52	2,41	2,25	2,08	1,88	1,61

## Окончание таблицы

$\frac{df_1}{df_2}$	1	2	3	4	5	6	8	12	24	10 <sup>9</sup>
32	4,15	3,29	2,90	2,67	2,51	2,40	2,24	2,07	1,86	1,59
33	4,14	3,28	2,89	2,66	2,50	2,39	2,23	2,06	1,85	1,58
34	4,13	3,28	2,88	2,65	2,49	2,38	2,23	2,05	1,84	1,57
35	4,12	3,27	2,87	2,64	2,49	2,37	2,22	2,04	1,83	1,56
36	4,11	3,26	2,87	2,63	2,48	2,36	2,21	2,03	1,82	1,55
37	4,11	3,25	2,86	2,63	2,47	2,36	2,20	2,02	1,82	1,54
38	4,10	3,24	2,85	2,62	2,46	2,35	2,19	2,02	1,81	1,53
39	4,09	3,24	2,85	2,61	2,46	2,34	2,19	2,01	1,80	1,52
40	4,08	3,23	2,84	2,61	2,45	2,34	2,18	2,00	1,79	1,51
41	4,08	3,23	2,83	2,60	2,44	2,33	2,17	2,00	1,79	1,50
42	4,07	3,22	2,83	2,59	2,44	2,32	2,17	1,99	1,78	1,49
43	4,07	3,21	2,82	2,59	2,43	2,32	2,16	1,99	1,77	1,48
44	4,06	3,21	2,82	2,58	2,43	2,31	2,16	1,98	1,77	1,48
45	4,06	3,20	2,81	2,58	2,42	2,31	2,15	1,97	1,76	1,47
46	4,05	3,20	2,81	2,57	2,42	2,30	2,15	1,97	1,76	1,46
47	4,05	3,20	2,80	2,57	2,41	2,30	2,14	1,96	1,75	1,46
48	4,04	3,19	2,80	2,57	2,41	2,29	2,14	1,96	1,75	1,45
49	4,04	3,19	2,79	2,56	2,40	2,29	2,13	1,96	1,74	1,44
50	4,03	3,18	2,79	2,56	2,40	2,29	2,13	1,95	1,74	1,44
60	4,00	3,15	2,76	2,53	2,37	2,25	2,10	1,92	1,70	1,39
70	3,98	3,13	2,74	2,50	2,35	2,23	2,07	1,89	1,67	1,35
80	3,96	3,11	2,72	2,49	2,33	2,21	2,06	1,88	1,65	1,32
90	3,95	3,10	2,71	2,47	2,32	2,20	2,04	1,86	1,64	1,30
100	3,94	3,09	2,70	2,46	2,31	2,19	2,03	1,85	1,63	1,28
125	3,92	3,07	2,68	2,44	2,29	2,17	2,01	1,83	1,60	1,25
150	3,90	3,06	2,66	2,43	2,27	2,16	2,00	1,82	1,59	1,22
200	3,89	3,04	2,65	2,42	2,26	2,14	1,98	1,80	1,57	1,19
250	3,88	3,03	2,64	2,41	2,25	2,13	1,98	1,79	1,56	1,17
300	3,87	3,03	2,63	2,40	2,24	2,13	1,97	1,78	1,55	1,15
400	3,86	3,02	2,63	2,39	2,24	2,12	1,96	1,78	1,54	1,13
500	3,86	3,01	2,62	2,39	2,23	2,12	1,96	1,77	1,54	1,11
600	3,86	3,01	2,62	2,39	2,23	2,11	1,95	1,77	1,54	1,10
1000	3,85	3,00	2,61	2,38	2,22	2,11	1,95	1,76	1,53	1,08

### Критические значения $\chi^2$ на уровне значимости 0,01; 0,05; 0,1

Число степеней свободы	Уровень значимости			Число степеней свободы	Уровень значимости		
	0,01	0,05	0,10		0,01	0,05	0,10
1	6,634897	3,841459	2,705544	33	54,77554	47,39988	43,74518
2	9,21034	5,991465	4,60517	34	56,06091	48,60237	44,90316
3	11,34487	7,814728	6,251388	35	57,34207	49,80185	46,05879
4	13,2767	9,487729	7,77944	36	58,61921	50,99846	47,21217
5	15,08627	11,0705	9,236357	37	59,8925	52,19232	48,36341
6	16,81189	12,59159	10,64464	38	61,16209	53,38354	49,51258
7	18,47531	14,06714	12,01704	39	62,42812	54,57223	50,65977
8	20,09024	15,50731	13,36157	40	63,69074	55,75848	51,80506
9	21,66599	16,91898	14,68366	41	64,95007	56,94239	52,94851
10	23,20925	18,30704	15,98718	42	66,20624	58,12404	54,0902
11	24,72497	19,67514	17,27501	43	67,45935	59,30351	55,23019
12	26,21697	21,02607	18,54935	44	68,70951	60,48089	56,36854
13	27,68825	22,36203	19,81193	45	69,95683	61,65623	57,5053
14	29,14124	23,68479	21,06414	46	71,2014	62,82962	58,64054
15	30,57791	24,99579	22,30713	47	72,44331	64,00111	59,77429
16	31,99993	26,29623	23,54183	48	73,68264	65,17077	60,90661
17	33,40866	27,58711	24,76904	49	74,91947	66,33865	62,03754
18	34,80531	28,8693	25,98942	50	76,15389	67,50481	63,16712
19	36,19087	30,14353	27,20357	60	88,37942	79,08194	74,39701
20	37,56623	31,41043	28,41198	70	100,4252	90,53123	85,52704
21	38,93217	32,67057	29,61509	80	112,3288	101,8795	96,5782
22	40,28936	33,92444	30,81328	90	124,1163	113,1453	107,565
23	41,6384	35,17246	32,0069	100	135,8067	124,3421	118,498
24	42,97982	36,41503	33,19624	125	164,694	152,0939	145,643
25	44,3141	37,65248	34,38159	150	193,2077	179,5806	172,5812
26	45,64168	38,88514	35,56317	200	249,4451	233,9943	226,021
27	46,96294	40,11327	36,74122	250	304,9396	287,8815	279,0504
28	48,27824	41,33714	37,91592	300	359,9064	341,3951	331,7885
29	49,58788	42,55697	39,08747	400	468,7245	447,6325	436,649
30	50,89218	43,77297	40,25602	500	576,4928	553,1268	540,9303
31	52,19139	44,98534	41,42174	600	683,5156	658,0936	644,8004
32	53,48577	46,19426	42,58475	1000	1106,969	1074,679	1057,724

**Значение интеграла вероятностей  $F(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{x^2}{2}} dx$**

Сотые доли $x$												
$x$	0	0,01	0,02	0,03	0,04	0,05	0,06	0,07	0,08	0,09		
0,0	0,000000	0,007979	0,015957	0,023933	0,031907	0,039878	0,047844	0,055806	0,063763	0,071713		
0,1	0,079656	0,087591	0,095517	0,103434	0,11134	0,119235	0,127119	0,13499	0,142847	0,150691		
0,2	0,158519	0,166332	0,174129	0,181908	0,18967	0,197413	0,205136	0,21284	0,220522	0,228184		
0,3	0,235823	0,243439	0,251032	0,2586	0,266143	0,273661	0,281153	0,288618	0,296055	0,303463		
0,4	0,310843	0,318194	0,325515	0,332804	0,340063	0,34729	0,354484	0,361645	0,368773	0,375866		
0,5	0,382925	0,389949	0,396936	0,403888	0,410803	0,417681	0,424521	0,431322	0,438085	0,444809		
0,6	0,451494	0,458138	0,464742	0,471305	0,477827	0,484308	0,490746	0,497142	0,503496	0,509806		
0,7	0,516073	0,522296	0,528475	0,53461	0,5407	0,546745	0,552745	0,5587	0,564609	0,570472		
0,8	0,576289	0,58206	0,587784	0,593461	0,599092	0,604675	0,610211	0,6157	0,621141	0,626534		
0,9	0,63188	0,637177	0,642427	0,647629	0,652782	0,657888	0,662945	0,667954	0,672914	0,677826		
1,0	0,682689	0,687505	0,692272	0,69699	0,70166	0,706282	0,710855	0,715381	0,719858	0,724287		
1,1	0,728668	0,733001	0,737286	0,741524	0,745714	0,749856	0,753951	0,757999	0,762	0,765954		
1,2	0,769861	0,773721	0,777535	0,781303	0,785025	0,7887	0,792331	0,795915	0,799455	0,802949		
1,3	0,806399	0,809804	0,813165	0,816482	0,819755	0,822984	0,82617	0,829313	0,832413	0,835471		

Сотые доли х										
х										
1,4	0,838487	0,84146	0,844392	0,847283	0,850133	0,852941	0,85571	0,858438	0,861127	0,863776
1,5	0,866386	0,868957	0,871489	0,873983	0,87644	0,878858	0,88124	0,883585	0,885893	0,888165
1,6	0,890401	0,892602	0,894768	0,896899	0,898995	0,901057	0,903086	0,905081	0,907043	0,908972
1,7	0,910869	0,912734	0,914568	0,91637	0,918141	0,919882	0,921592	0,923273	0,924924	0,926546
1,8	0,928139	0,929704	0,931241	0,93275	0,934232	0,935686	0,937114	0,938516	0,939982	0,941242
1,9	0,942567	0,943867	0,945142	0,946393	0,94762	0,948824	0,950004	0,951162	0,952296	0,953409
2,0	0,9545	0,955569	0,956617	0,957643	0,95865	0,959636	0,960601	0,961548	0,962474	0,963382
2,1	0,964271	0,965142	0,965994	0,966828	0,967645	0,968445	0,969227	0,969993	0,970743	0,971476
2,2	0,972193	0,972895	0,973581	0,974253	0,974909	0,975551	0,976179	0,976792	0,977392	0,977979
2,3	0,978552	0,979112	0,979659	0,980194	0,980716	0,981227	0,981725	0,982212	0,982687	0,983152
2,4	0,983605	0,984047	0,984479	0,984901	0,985313	0,985714	0,986106	0,986489	0,986862	0,987226
2,5	0,987581	0,987927	0,988265	0,988594	0,988915	0,989228	0,989533	0,98983	0,99012	0,990402
2,6	0,990678	0,990946	0,991207	0,991462	0,991709	0,991951	0,992186	0,992415	0,992638	0,992855
2,7	0,993066	0,993272	0,993472	0,993667	0,993856	0,99404	0,99422	0,994394	0,994564	0,994729
2,8	0,99489	0,995046	0,995198	0,995345	0,995489	0,995628	0,995764	0,995895	0,996023	0,996148
2,9	0,996268	0,996386	0,9965	0,99661	0,996718	0,996822	0,996924	0,997022	0,997118	0,99721
3,0	0,9973	0,997388	0,997472	0,997554	0,997634	0,997712	0,997787	0,997859	0,99793	0,997998

**d**-статистика Дарбина — Уотсона:  $d_L$  и  $d_U$ ,  
уровень значимости 5%

n	k = 1		k = 2		k = 3		k = 4		k = 5		k = 6		k = 7		k = 8		k = 9		k = 10		k = 15	
	$d_L$	$d_U$	$d_L$	$d_U$	$d_L$	$d_U$	$d_L$	$d_U$	$d_L$	$d_U$	$d_L$	$d_U$	$d_L$	$d_U$	$d_L$	$d_U$	$d_L$	$d_U$	$d_L$	$d_U$	$d_L$	$d_U$
6	0,610	1,400	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—
7	0,700	1,356	0,467	1,896	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—
8	0,763	1,332	0,559	1,777	0,368	2,287	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—
9	0,824	1,320	0,629	1,699	0,455	2,128	0,296	2,588	—	—	—	—	—	—	—	—	—	—	—	—	—	—
10	0,879	1,320	0,697	1,641	0,525	2,016	0,376	2,414	0,243	2,822	—	—	—	—	—	—	—	—	—	—	—	—
11	0,927	1,324	0,658	1,604	0,595	1,928	0,444	2,283	0,316	2,645	0,203	3,005	—	—	—	—	—	—	—	—	—	—
12	0,971	1,331	0,812	1,579	0,658	1,864	0,512	2,177	0,379	2,506	0,268	2,832	0,171	3,149	—	—	—	—	—	—	—	—
13	1,010	1,340	0,861	1,562	0,715	1,816	0,574	2,094	0,445	2,390	0,328	2,692	0,230	2,985	0,147	3,266	—	—	—	—	—	—
14	1,045	1,350	0,905	1,551	0,767	1,779	0,632	2,030	0,505	2,296	0,389	2,572	0,286	2,848	0,200	3,111	0,127	3,360	—	—	—	—
15	1,077	1,361	0,946	1,543	0,814	1,750	0,685	1,977	0,562	2,220	0,447	2,472	0,343	2,727	0,251	2,979	0,175	3,216	0,111	3,438	—	—
16	1,106	1,371	0,982	1,539	0,857	1,728	0,734	1,935	0,615	2,157	0,502	2,388	0,398	2,624	0,304	2,860	0,222	3,090	0,155	3,304	—	—
17	1,133	1,381	1,015	1,536	0,897	1,710	0,779	1,900	0,664	2,104	0,554	2,318	0,451	2,537	0,356	2,757	0,272	2,975	0,198	3,184	—	—
18	1,158	1,391	1,046	1,535	0,933	1,696	0,820	1,872	0,710	2,060	0,603	2,257	0,502	2,461	0,407	2,667	0,321	2,873	0,244	3,073	—	—
19	1,180	1,401	1,074	1,536	0,967	1,685	0,859	1,848	0,752	2,023	0,649	2,206	0,549	2,396	0,456	2,589	0,369	2,783	0,290	2,974	—	—
20	1,201	1,411	1,100	1,537	0,998	1,676	0,894	1,828	0,792	1,991	0,692	2,162	0,595	2,339	0,502	2,521	0,416	2,704	0,336	2,885	0,063	3,676
21	1,221	1,420	1,125	1,538	1,026	1,669	0,927	1,812	0,829	1,964	0,732	2,124	0,637	2,290	0,547	2,460	0,461	2,633	0,380	2,806	0,091	3,583
22	1,239	1,429	1,147	1,541	1,053	1,664	0,958	1,797	0,863	1,940	0,769	2,090	0,677	2,246	0,588	2,407	0,504	2,571	0,424	2,734	0,120	3,495
23	1,257	1,437	1,168	1,543	1,078	1,660	0,986	1,785	0,895	1,920	0,804	2,061	0,715	2,208	0,628	2,360	0,545	2,514	0,465	2,670	0,153	3,409
24	1,273	1,446	1,188	1,546	1,101	1,656	1,013	1,775	0,925	1,902	0,837	2,035	0,751	2,174	0,666	2,318	0,584	2,464	0,506	2,613	0,186	3,327
25	1,288	1,454	1,206	1,550	1,123	1,654	1,038	1,767	0,953	1,886	0,868	2,012	0,784	2,144	0,702	2,280	0,621	2,419	0,544	2,560	0,221	3,251
26	1,302	1,461	1,224	1,553	1,143	1,652	1,062	1,759	0,979	1,873	0,897	1,992	0,816	2,117	0,735	2,246	0,657	2,379	0,581	2,513	0,256	3,179
27	1,316	1,469	1,240	1,556	1,162	1,651	1,084	1,753	1,004	1,861	0,925	1,974	0,845	2,093	0,767	2,216	0,691	2,342	0,616	2,470	0,291	3,112
28	1,328	1,476	1,255	1,560	1,181	1,650	1,104	1,747	1,028	1,850	0,951	1,958	0,874	2,071	0,798	2,188	0,723	2,309	0,650	2,431	0,325	3,050

n	k = 1		k = 2		k = 3		k = 4		k = 5		k = 6		k = 7		k = 8		k = 9		k = 10		k = 15	
	$d_t$	$d_u$	$d_t$	$d_u$	$d_t$	$d_u$	$d_t$	$d_u$	$d_t$	$d_u$	$d_t$	$d_u$	$d_t$	$d_u$	$d_t$	$d_u$	$d_t$	$d_u$	$d_t$	$d_u$	$d_t$	$d_u$
29	1,341	1,483	1,270	1,563	1,198	1,650	1,124	1,743	1,050	1,841	0,975	1,944	0,900	2,052	0,826	2,164	0,753	2,278	0,682	2,396	0,359	2,992
30	1,352	1,489	1,284	1,567	1,214	1,650	1,143	1,739	1,071	1,833	0,998	1,931	0,926	2,034	0,854	2,141	0,782	2,251	0,712	2,363	0,392	2,937
31	1,363	1,496	1,297	1,570	1,229	1,650	1,160	1,735	1,090	1,825	1,020	1,920	0,950	2,018	0,879	2,120	0,810	2,226	0,741	2,333	0,425	2,887
32	1,373	1,502	1,309	1,574	1,244	1,650	1,177	1,732	1,109	1,819	1,041	1,909	0,972	2,004	0,904	2,102	0,836	2,203	0,769	2,306	0,457	2,840
33	1,383	1,508	1,321	1,577	1,258	1,651	1,193	1,730	1,127	1,813	1,061	1,900	0,994	1,991	0,927	2,085	0,861	2,181	0,795	2,281	0,488	2,796
34	1,393	1,514	1,333	1,580	1,271	1,652	1,208	1,728	1,144	1,808	1,080	1,891	1,015	1,979	0,950	2,069	0,885	2,162	0,821	2,257	0,518	2,754
35	1,402	1,519	1,343	1,584	1,283	1,653	1,222	1,726	1,160	1,803	1,097	1,884	1,034	1,967	0,971	2,054	0,908	2,144	0,845	2,236	0,547	2,716
36	1,411	1,525	1,354	1,587	1,295	1,654	1,236	1,724	1,175	1,799	1,114	1,877	1,053	1,957	0,991	2,041	0,930	2,127	0,868	2,216	0,575	2,680
37	1,419	1,530	1,364	1,590	1,307	1,655	1,249	1,723	1,190	1,795	1,131	1,870	1,071	1,948	1,011	2,029	0,951	2,112	0,891	2,198	0,602	2,646
38	1,427	1,535	1,373	1,594	1,318	1,656	1,261	1,722	1,204	1,792	1,146	1,864	1,088	1,939	1,029	2,017	0,970	2,098	0,912	2,180	0,628	2,614
39	1,435	1,540	1,382	1,597	1,328	1,658	1,273	1,722	1,218	1,789	1,161	1,859	1,104	1,932	1,047	2,007	0,990	2,085	0,932	2,164	0,653	2,585
40	1,442	1,544	1,391	1,600	1,338	1,659	1,285	1,721	1,230	1,786	1,175	1,854	1,120	1,924	1,064	1,997	1,008	2,072	0,952	2,149	0,678	2,557
45	1,475	1,566	1,430	1,615	1,383	1,666	1,336	1,720	1,287	1,776	1,238	1,835	1,189	1,895	1,139	1,958	1,089	2,022	1,038	2,088	0,788	2,439
50	1,503	1,585	1,462	1,628	1,421	1,674	1,378	1,721	1,335	1,771	1,291	1,822	1,246	1,875	1,201	1,930	1,156	1,986	1,110	2,044	0,882	2,350
55	1,528	1,601	1,490	1,641	1,452	1,681	1,414	1,724	1,374	1,768	1,334	1,814	1,294	1,861	1,253	1,909	1,212	1,959	1,170	2,010	0,961	2,281
60	1,549	1,616	1,514	1,652	1,480	1,689	1,444	1,727	1,408	1,767	1,372	1,808	1,335	1,850	1,298	1,894	1,260	1,939	1,222	1,984	1,029	2,227
65	1,567	1,629	1,536	1,662	1,503	1,696	1,471	1,731	1,438	1,767	1,404	1,805	1,370	1,843	1,336	1,882	1,301	1,923	1,266	1,964	1,088	2,183
70	1,583	1,641	1,554	1,672	1,525	1,703	1,494	1,735	1,464	1,768	1,433	1,802	1,401	1,837	1,369	1,873	1,337	1,910	1,305	1,948	1,139	2,148
75	1,598	1,652	1,571	1,680	1,543	1,709	1,515	1,739	1,487	1,770	1,458	1,801	1,428	1,834	1,399	1,867	1,369	1,901	1,339	1,935	1,184	2,118
80	1,611	1,662	1,586	1,688	1,560	1,715	1,534	1,743	1,507	1,772	1,480	1,801	1,453	1,831	1,425	1,861	1,397	1,893	1,369	1,925	1,224	2,093
85	1,624	1,671	1,600	1,696	1,575	1,721	1,550	1,747	1,525	1,774	1,500	1,801	1,474	1,829	1,448	1,857	1,422	1,886	1,396	1,916	1,260	2,073
90	1,635	1,679	1,612	1,703	1,589	1,726	1,566	1,751	1,542	1,776	1,518	1,801	1,494	1,854	1,469	1,881	1,445	1,881	1,420	1,909	1,292	2,055
95	1,645	1,687	1,623	1,709	1,602	1,732	1,579	1,755	1,557	1,778	1,535	1,802	1,512	1,827	1,489	1,852	1,465	1,877	1,442	1,903	1,321	2,040
100	1,654	1,694	1,634	1,715	1,613	1,736	1,592	1,758	1,571	1,780	1,550	1,803	1,528	1,826	1,506	1,850	1,484	1,874	1,462	1,898	1,347	2,026
150	1,720	1,746	1,706	1,760	1,693	1,774	1,679	1,788	1,665	1,802	1,651	1,817	1,637	1,832	1,622	1,847	1,608	1,862	1,594	1,877	1,519	1,956
200	1,758	1,778	1,748	1,789	1,738	1,799	1,728	1,810	1,718	1,820	1,707	1,831	1,697	1,841	1,686	1,852	1,675	1,863	1,665	1,874	1,610	1,931

Примечание: n — число наблюдений, k — число объясняющих переменных (без постоянного члена).

**d**-статистика Дарбина — Уогсона:  $d_L$  и  $d_U$ ,  
уровень значимости 1%

n	k = 1		k = 2		k = 3		k = 4		k = 5		k = 6		k = 7		k = 8		k = 9		k = 10		k = 15	
	$d_L$	$d_U$	$d_L$	$d_U$	$d_L$	$d_U$	$d_L$	$d_U$	$d_L$	$d_U$	$d_L$	$d_U$	$d_L$	$d_U$	$d_L$	$d_U$	$d_L$	$d_U$	$d_L$	$d_U$	$d_L$	$d_U$
6	0,390	1,142	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—
7	0,435	1,036	0,294	1,676	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—
8	0,497	1,003	0,345	1,489	0,229	2,102	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—
9	0,554	0,998	0,408	1,389	0,279	1,875	0,183	2,433	—	—	—	—	—	—	—	—	—	—	—	—	—	—
10	0,604	1,001	0,466	1,333	0,340	1,733	0,230	2,193	0,150	2,690	—	—	—	—	—	—	—	—	—	—	—	—
11	0,653	1,010	0,519	1,297	0,396	1,640	0,286	2,030	0,193	2,453	0,124	2,892	—	—	—	—	—	—	—	—	—	—
12	0,697	1,023	0,569	1,274	0,449	1,575	0,339	1,913	0,244	2,280	0,164	2,665	0,105	3,053	—	—	—	—	—	—	—	—
13	0,738	1,038	0,616	1,261	0,499	1,526	0,391	1,826	0,294	2,150	0,211	2,490	0,140	2,838	0,090	3,182	—	—	—	—	—	—
14	0,776	1,054	0,660	1,254	0,547	1,490	0,441	1,757	0,343	2,049	0,257	2,354	0,183	2,667	0,122	2,981	0,078	3,287	—	—	—	—
15	0,811	1,070	0,700	1,252	0,591	1,464	0,488	1,704	0,391	1,967	0,303	2,244	0,226	2,530	0,161	2,817	0,107	3,101	0,068	3,374	—	—
16	0,844	1,086	0,737	1,252	0,633	1,446	0,532	1,663	0,437	1,900	0,349	2,153	0,269	2,416	0,200	2,681	0,142	2,944	0,094	3,201	—	—
17	0,874	1,102	0,772	1,255	0,672	1,432	0,574	1,630	0,480	1,847	0,393	2,078	0,313	2,319	0,241	2,566	0,179	2,811	0,127	3,053	—	—
18	0,902	1,118	0,805	1,259	0,708	1,422	0,613	1,604	0,522	1,803	0,435	2,015	0,355	2,238	0,282	2,467	0,216	2,697	0,160	2,925	—	—
19	0,928	1,132	0,835	1,265	0,742	1,415	0,650	1,584	0,561	1,767	0,476	1,963	0,396	2,169	0,322	2,381	0,255	2,597	0,196	2,813	—	—
20	0,952	1,147	0,863	1,271	0,773	1,411	0,685	1,567	0,598	1,737	0,515	1,918	0,436	2,110	0,362	2,308	0,294	2,510	0,232	2,714	0,038	3,639
21	0,975	1,161	0,890	1,277	0,803	1,408	0,718	1,554	0,633	1,712	0,552	1,881	0,474	2,059	0,400	2,244	0,331	2,434	0,268	2,625	0,055	3,521
22	0,997	1,174	0,914	1,284	0,831	1,407	0,748	1,543	0,667	1,691	0,587	1,849	0,510	2,015	0,437	2,188	0,368	2,367	0,304	2,548	0,077	3,412
23	1,018	1,187	0,938	1,291	0,858	1,407	0,777	1,534	0,698	1,673	0,620	1,821	0,545	1,977	0,470	2,140	0,404	2,308	0,340	2,479	0,100	3,311
24	1,037	1,199	0,960	1,298	0,882	1,407	0,805	1,528	0,728	1,658	0,652	1,797	0,578	1,944	0,507	2,097	0,439	2,255	0,375	2,417	0,125	3,218
25	1,055	1,211	0,981	1,305	0,906	1,409	0,831	1,523	0,756	1,645	0,682	1,776	0,610	1,915	0,540	2,059	0,473	2,209	0,409	2,362	0,131	0,116
26	1,072	1,222	1,001	1,312	0,928	1,411	0,855	1,518	0,783	1,635	0,711	1,759	0,640	1,889	0,572	2,026	0,505	2,168	0,441	2,313	0,050	0,141
27	1,089	1,233	1,019	1,319	0,949	1,413	0,878	1,515	0,808	1,626	0,738	1,743	0,669	1,867	0,602	1,997	0,536	2,131	0,473	2,269	0,076	0,167
28	1,104	1,244	1,037	1,325	0,969	1,415	0,900	1,513	0,832	1,618	0,764	1,729	0,696	1,847	0,630	1,970	0,566	2,098	0,504	2,229	0,097	0,194



n	k = 1		k = 2		k = 3		k = 4		k = 5		k = 6		k = 7		k = 8		k = 9		k = 10		k = 15	
	d <sub>t</sub>	d <sub>u</sub>	d <sub>t</sub>	d <sub>u</sub>	d <sub>t</sub>	d <sub>u</sub>	d <sub>t</sub>	d <sub>u</sub>	d <sub>t</sub>	d <sub>u</sub>	d <sub>t</sub>	d <sub>u</sub>	d <sub>t</sub>	d <sub>u</sub>	d <sub>t</sub>	d <sub>u</sub>	d <sub>t</sub>	d <sub>u</sub>	d <sub>t</sub>	d <sub>u</sub>	d <sub>t</sub>	d <sub>u</sub>
29	1,119	1,254	1,054	1,332	0,988	1,418	0,921	1,512	0,855	1,611	0,788	1,718	0,723	1,830	0,658	1,947	0,595	2,068	0,533	2,193	2,843	0,222
30	1,133	1,263	1,070	1,339	1,006	1,421	0,941	1,511	0,877	1,606	0,812	1,707	0,748	1,814	0,684	1,925	0,622	2,041	0,562	2,160	2,785	0,249
31	1,147	1,273	1,085	1,345	1,023	1,425	0,960	1,510	0,897	1,601	0,834	1,698	0,772	1,800	0,710	1,906	0,649	2,017	0,589	2,131	2,730	0,277
32	1,160	1,282	1,100	1,352	1,040	1,428	0,979	1,510	0,917	1,597	0,856	1,690	0,794	1,788	0,734	1,889	0,674	1,995	0,615	2,104	2,680	0,304
33	1,172	1,291	1,114	1,358	1,055	1,432	0,996	1,510	0,936	1,594	0,863	1,683	0,816	1,776	0,754	1,876	0,698	1,975	0,641	2,080	2,633	0,331
34	1,184	1,299	1,128	1,364	1,070	1,435	1,012	1,511	0,954	1,591	0,896	1,677	0,837	1,766	0,779	1,860	0,722	1,957	0,665	2,057	2,590	0,357
35	1,195	1,307	1,140	1,370	1,085	1,439	1,028	1,512	0,971	1,589	0,914	1,671	0,857	1,757	0,800	1,847	0,744	1,940	0,689	2,037	2,550	0,383
36	1,206	1,315	1,153	1,376	1,098	1,442	1,043	1,513	0,988	1,588	0,932	1,666	0,877	1,749	0,821	1,836	0,766	1,925	0,711	2,018	2,512	0,409
37	1,217	1,323	1,165	1,382	1,112	1,446	1,058	1,514	1,004	1,586	0,950	1,662	0,895	1,742	0,841	1,825	0,787	1,911	0,733	2,001	2,477	0,434
38	1,227	1,330	1,176	1,388	1,124	1,449	1,072	1,515	1,019	1,585	0,966	1,658	0,913	1,735	0,860	1,816	0,807	1,899	0,754	1,985	2,445	0,458
39	1,237	1,337	1,187	1,393	1,137	1,453	1,085	1,517	1,034	1,584	0,982	1,655	0,930	1,729	0,878	1,807	0,826	1,887	0,774	1,970	2,414	0,482
40	1,246	1,344	1,198	1,398	1,148	1,457	1,098	1,518	1,048	1,584	0,997	1,652	0,946	1,724	0,895	1,799	0,844	1,876	0,749	1,956	2,386	0,505
45	1,288	1,376	1,245	1,423	1,201	1,474	1,156	1,528	1,111	1,584	1,065	1,643	1,019	1,704	0,974	1,768	0,927	1,834	0,881	1,902	2,269	0,612
50	1,324	1,403	1,285	1,446	1,245	1,491	1,205	1,538	1,164	1,587	1,123	1,639	1,081	1,692	1,039	1,748	0,997	1,805	0,955	1,864	2,182	0,705
55	1,356	1,427	1,320	1,466	1,284	1,506	1,247	1,548	1,209	1,592	1,172	1,638	1,134	1,685	1,095	1,734	1,057	1,785	1,018	1,837	2,117	0,786
60	1,383	1,449	1,350	1,484	1,317	1,520	1,283	1,558	1,249	1,598	1,214	1,639	1,179	1,682	1,144	1,726	1,108	1,771	1,120	1,802	2,027	0,919
65	1,407	1,468	1,377	1,500	1,346	1,534	1,315	1,568	1,283	1,604	1,251	1,642	1,218	1,680	1,186	1,720	1,153	1,761	1,162	1,792	1,995	0,974
70	1,429	1,485	1,400	1,515	1,372	1,546	1,343	1,578	1,313	1,611	1,283	1,645	1,253	1,680	1,223	1,716	1,192	1,754	1,162	1,792	1,995	1,023
75	1,448	1,501	1,422	1,529	1,395	1,557	1,368	1,587	1,340	1,617	1,313	1,649	1,284	1,682	1,256	1,714	1,227	1,748	1,199	1,783	1,970	1,066
80	1,466	1,515	1,441	1,541	1,416	1,568	1,390	1,595	1,364	1,630	1,332	1,653	1,317	1,683	1,285	1,714	1,259	1,745	1,232	1,777	1,949	1,106
85	1,482	1,528	1,458	1,553	1,435	1,578	1,411	1,603	1,386	1,630	1,362	1,657	1,337	1,685	1,312	1,714	1,287	1,743	1,262	1,773	1,931	1,141
90	1,496	1,540	1,474	1,563	1,452	1,587	1,429	1,611	1,406	1,636	1,383	1,661	1,360	1,687	1,336	1,714	1,312	1,741	1,288	1,769	1,917	1,174
95	1,510	1,552	1,489	1,573	1,468	1,596	1,446	1,618	1,425	1,642	1,403	1,666	1,381	1,690	1,358	1,715	1,336	1,741	1,313	1,767	1,905	1,174
100	1,522	1,562	1,503	1,583	1,482	1,604	1,462	1,625	1,441	1,647	1,421	1,670	1,400	1,693	1,378	1,717	1,357	1,741	1,335	1,765	1,895	1,203
150	1,611	1,637	1,598	1,651	1,584	1,665	1,571	1,679	1,557	1,693	1,543	1,708	1,530	1,722	1,515	1,737	1,501	1,752	1,486	1,767	1,847	1,400
200	1,664	1,684	1,653	1,693	1,643	1,704	1,633	1,715	1,623	1,725	1,613	1,735	1,603	1,746	1,592	1,757	1,582	1,768	1,571	1,779	1,836	1,507

Примечание: n — число наблюдений, k — число объясняющих переменных (без постоянного члена).

**Покупайте наши книги:**

**Оптом** в офисе книготорга «Юрайт»:  
140004, Московская обл., г. Люберцы, 1-й Панковский проезд, д. 1,  
тел.: (495) 744-00-12, e-mail: sales@urait.ru, www.urait.ru

**В розницу** в интернет-магазине: www.urait-book.ru,  
e-mail: order@urait-book.ru, тел.: (495) 742-72-12

**Для закупок у Единого поставщика**  
в соответствии с Федеральным законом  
от 21.07.2005 № 94-ФЗ обращаться по тел.: (495) 744-00-12,  
e-mail: sales@urait.ru, vuz@urait.ru

*Учебное издание*

# **Эконометрика**

Учебник для магистров

Под редакцией *И. И. Елисейевой*

Редактор *А. В. Желонкин*  
Корректор *М. Б. Быкова*  
Художественное оформление *А. И. Гиренко*  
Компьютерная верстка *Д. Е. Дегтярев, Е. Е. Савина*

Формат 84×108<sup>1</sup>/<sub>32</sub>.  
Гарнитура «CharterITC». Печать офсетная.  
Усл. печ. л. 23,78. Тираж 1500 экз. Заказ №

**ООО «Издательство Юрайт»**  
140004, Московская обл., г. Люберцы, 1-й Панковский проезд, д. 1.  
Тел.: (495) 744-00-12. E-mail: izdat@urait.ru, www.urait.ru