

**Н. В. Артамонов**

**ВВЕДЕНИЕ  
В ЭКОНОМЕТРИКУ**

УДК 330.43

ББК 65.05

А86

Артамонов Н. В.

Введение в эконометрику

Электронное издание

М.: МЦНМО, 2019

251 с.

ISBN 978-5-4439-3381-8

Учебник знакомит читателя с базовыми понятиями и методами современной эконометрики, которая является неотъемлемой частью современного экономического образования. В первых двух главах подробно излагаются линейные и нелинейные регрессионные модели, их статистические свойства и возможности применения в экономике. В третьей главе рассматриваются возможные отклонения от стандартных предположений линейной модели регрессии, встречающиеся при моделировании экономических ситуаций и при анализе экономических данных. Обсуждаются корректировки регрессионной модели для описания таких ситуаций. Последняя глава посвящена регрессионным моделям временных рядов. Учебник основан на лекциях по курсу «Эконометрика-1», читаемых автором в МГИМО (У) МИД России на факультете Международных экономических отношений.

Книга предназначена студентам (бакалавриата и магистратуры), аспирантам и преподавателям, специалистам и исследователям, работающим в области прикладной экономики и финансов.

В настоящее издание добавлено новое приложение о линейной регрессии в языке R.

Подготовлено на основе книги: *Н. В. Артамонов. Введение в эконометрику. — 3-е изд., испр. и доп. — М.: МЦНМО, 2019. — ISBN 978-5-4439-1381-0.*

Учебное издание для вузов

12+

Издательство Московского центра

непрерывного математического образования

119002, Москва, Большой Власьевский пер., 11. Тел. (499) 241-08-04

<http://www.mcme.ru>

ISBN 978-5-4439-3381-8

© Артамонов Н. В., 2019.

© МЦНМО, 2019.

# Оглавление

<b>Введение</b>	<b>5</b>
Структура книги . . . . .	8
Статистические данные в эконометрике . . . . .	11
Список обозначений . . . . .	12
<b>Глава 1. Парная регрессия</b>	<b>15</b>
§ 1.1. Парный коэффициент корреляции . . . . .	15
1.1.1. Коэффициент корреляции . . . . .	15
1.1.2. Выборочный коэффициент корреляции . . . . .	17
§ 1.2. Подгонка прямой. Метод наименьших квадратов . . . . .	21
§ 1.3. Парная линейная модель регрессии . . . . .	23
1.3.1. Теорема Гаусса—Маркова . . . . .	25
1.3.2. Статистические свойства OLS-оценок коэффициентов . . . . .	30
1.3.3. Доверительные интервалы. Проверка гипотез . . . . .	32
1.3.4. Коэффициент $R^2$ и «качество подгонки» . . . . .	35
§ 1.4. Прогнозирование в модели парной регрессии . . . . .	38
§ 1.5. Парная регрессия без константы . . . . .	40
§ 1.6. Нелинейные модели . . . . .	45
§ 1.7. Стохастические регрессоры . . . . .	48
§ 1.8. Задачи . . . . .	52
<b>Глава 2. Многофакторная регрессия</b>	<b>63</b>
§ 2.1. Метод наименьших квадратов . . . . .	64
§ 2.2. Основные предположения. Теорема Гаусса—Маркова . . . . .	65
§ 2.3. Статистические свойства OLS-оценок. Доверительные интервалы и проверка гипотез . . . . .	69
§ 2.4. Коэффициент $R^2$ . Проверка сложных гипотез о коэффициентах регрессии . . . . .	72
§ 2.5. Прогнозирование в линейной модели регрессии . . . . .	79
§ 2.6. Множественная регрессия без константы . . . . .	81
§ 2.7. Нелинейные модели . . . . .	86
§ 2.8. Бинарные переменные . . . . .	89
§ 2.9. Стохастические регрессоры . . . . .	93
2.9.1. Асимптотические свойства OLS-оценок . . . . .	98
§ 2.10. Мультиколлинеарность . . . . .	100
§ 2.11. Задачи . . . . .	103

<b>Глава 3. Разные аспекты линейной регрессии</b>	<b>127</b>
§ 3.1. Спецификация модели регрессии . . . . .	127
3.1.1. Невключение в модель значимого фактора . . . . .	127
3.1.2. Включение в модель незначимого фактора . . . . .	129
3.1.3. Сравнение вложенных моделей . . . . .	130
3.1.4. Сравнение невложенных моделей . . . . .	131
3.1.5. Выбор функциональной формы зависимости . . . . .	132
§ 3.2. Гетероскедастичность ошибок регрессии. Взвешенный метод наименьших квадратов . . . . .	135
3.2.1. Тесты на гетероскедастичность . . . . .	136
3.2.2. Корректировка на гетероскедастичность . . . . .	142
§ 3.3. Корреляция во времени ошибок регрессии . . . . .	149
3.3.1. Автокорреляция первого порядка . . . . .	150
3.3.2. Автокорреляция произвольного порядка . . . . .	155
§ 3.4. Корректировка модели на гетероскедастичность и автокор- реляцию . . . . .	158
§ 3.5. Задачи . . . . .	161
<b>Глава 4. Модели временных рядов</b>	<b>175</b>
§ 4.1. Условия Гаусса — Маркова для регрессионных моделей вре- менных рядов . . . . .	175
§ 4.2. Модель тренда и сезонность . . . . .	177
§ 4.3. Модель распределенных лагов . . . . .	180
§ 4.4. Модель авторегрессии временных рядов . . . . .	181
4.4.1. Стационарные временные ряды . . . . .	182
4.4.2. Модель авторегрессии . . . . .	184
4.4.3. Прогнозирование авторегрессионных случайных процессов	188
4.4.4. Эконометрические методы исследования стационарных временных рядов . . . . .	190
§ 4.5. Динамические модели стационарных временных рядов . . . .	195
§ 4.6. Задачи . . . . .	197
<b>Приложение А. Статистические таблицы</b>	<b>201</b>
<b>Приложение В. Информационные критерии</b>	<b>217</b>
<b>Приложение С. Линейная регрессия в R (совм. с Д. В. Артамоновым)</b>	<b>219</b>
<b>Литература</b>	<b>248</b>

# Глава 1

## Парная регрессия

### § 1.1. Парный коэффициент корреляции

#### 1.1.1. Коэффициент корреляции

Пусть  $(X, Y)$  — двумерная нормально распределенная случайная величина. Тогда «степень зависимости» случайных величин  $X$  и  $Y$  характеризуется парным коэффициентом корреляции

$$\rho = \text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{Var}(X) \cdot \text{Var}(Y)}} = \frac{E(XY) - EX \cdot EY}{\sqrt{\text{Var}(X) \cdot \text{Var}(Y)}}.$$

Из определения коэффициента корреляции следует, что

- 1)  $-1 \leq \rho \leq 1$ ;
- 2) коэффициент корреляции не меняется при линейных преобразованиях величин, т. е.

$$\text{corr}(X, Y) = \text{corr}(a_0 + a_1X, b_0 + b_1Y), \quad a_1, b_1 \neq 0.$$

Коэффициент корреляции принимает крайние значения  $\pm 1$  в том и только том случае, когда между случайными величинами  $X$  и  $Y$  существует линейная функциональная зависимость, т. е.

$$\rho = \pm 1 \Leftrightarrow Y = \beta_0 + \beta_1X, \quad \beta_1 \neq 0,$$

причем

$$\beta_1 = \rho \sqrt{\frac{\text{Var}(Y)}{\text{Var}(X)}},$$

т. е. знак коэффициента  $\beta_1$  совпадает со знаком коэффициента корреляции.

В общем случае коэффициент корреляции возникает при решении следующей экстремальной задачи: подобрать линейную функцию  $l(x) = \beta_0 + \beta_1x$  так, чтобы случайная величина  $l(X)$  меньше всего отклонялась от  $Y$  в среднеквадратичном, т. е.

$$E(Y - \beta_0 - \beta_1X)^2 \xrightarrow{\beta_0, \beta_1} \min.$$

Решение этой задачи задается равенствами

$$\beta_1^* = \frac{\text{cov}(X, Y)}{\text{Var}(X)} = \rho \sqrt{\frac{\text{Var}(Y)}{\text{Var}(X)}}, \quad \beta_0^* = EY - \beta_1^* \cdot EX,$$

и наименьшее среднееквадратичное отклонение равно

$$E(Y - \beta_0^* - \beta_1^* X)^2 = (1 - \rho^2) \text{Var}(Y).$$

Кроме того, для всех  $x \in \mathbb{R}$  верно равенство

$$E(Y | X = x) = \beta_0^* + \beta_1^* x,$$

т. е. наилучший прогноз случайной величины  $Y$  при условии, что известно значение случайной величины  $X = x$ , равен  $\hat{Y} = \beta_0^* + \beta_1^* x$ . Рассмотрим три случая:

- 1)  $\rho > 0$ ; тогда  $\beta_1^* > 0$  и при увеличении  $x$  ожидаемое (среднее) значение  $E(Y | X = x)$  случайной величины  $Y$  также увеличивается; в этом случае говорят о *прямой линейной зависимости* между величинами;
- 2)  $\rho < 0$ ; тогда  $\beta_1^* < 0$  и при увеличении  $x$  ожидаемое (среднее) значение  $E(Y | X = x)$  случайной величины  $Y$  уменьшается; в этом случае говорят об *обратной линейной зависимости* между величинами;
- 3)  $\rho = 0$ ; тогда  $\beta_1^* = 0$ ,  $E(Y | X = x) = \beta_0^*$  и знание значения случайной величины  $X$  не улучшает прогноз  $Y$ .

Важное значение коэффициента корреляции обусловлено следующей теоремой.

**Теорема.** Пусть  $(X, Y)$  — двумерная нормально распределенная случайная величина. Тогда случайные величины  $X$  и  $Y$  независимы в том и только том случае, когда  $\text{corr}(X, Y) = 0$ .

Таким образом, парный коэффициент корреляции можно рассматривать как *меру зависимости* двух случайных величин (факторов), имеющих совместное нормальное распределение, причем

- $\rho = 0 \Leftrightarrow$  величины независимы;
- $\rho = \pm 1 \Leftrightarrow$  между величинами имеется линейная функциональная зависимость:  $y = \beta_0^* + \beta_1^* x$ .

**1.1.2. Выборочный коэффициент корреляции**

Пусть  $(x_i, y_i)_{i=1}^n$  — выборка из двумерной нормально распределенной случайной величины,  $n$  — объем выборки.

Напомним, что выборочные (неисправленные) оценки дисперсий случайных величин  $X$  и  $Y$  определяются как

$$\widehat{\text{Var}}(X) = \widehat{\sigma}_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = (\bar{x}^2) - (\bar{x})^2,$$

$$\widehat{\text{Var}}(Y) = \widehat{\sigma}_y^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 = (\bar{y}^2) - (\bar{y})^2,$$

где

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{x}^2 = \frac{1}{n} \sum_{i=1}^n x_i^2.$$

Напомним также, что  $\widehat{\text{Var}}(X)$  и  $\widehat{\text{Var}}(Y)$  — состоятельные, но смещенные оценки дисперсий  $\text{Var}(X)$  и  $\text{Var}(Y)$  соответственно.

Выборочный коэффициент ковариации определяется как<sup>1</sup>

$$\widehat{\text{cov}}(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \overline{xy} - \bar{x} \cdot \bar{y},$$

а выборочный коэффициент корреляции определяется равенством<sup>2</sup>

$$r = \widehat{\text{corr}}(X, Y) = \frac{\widehat{\text{cov}}(X, Y)}{\sqrt{\widehat{\text{Var}}(X) \cdot \widehat{\text{Var}}(Y)}}, \quad -1 \leq r \leq 1.$$

Выборочные коэффициенты ковариации и корреляции являются состоятельными оценками коэффициентов ковариации и корреляции в генеральной совокупности. Выборочный коэффициент корреляции может рассматриваться как выборочная «мера линейной зависимости» между случайными величинами.

**Проверка значимости коэффициента корреляции.** Проверка значимости подразумевает проверку статистической гипотезы

$$H_0: \rho = 0$$

<sup>1</sup> В MS Excel функция КОВАР(·, ·).

<sup>2</sup> В MS Excel функция КОРРЕЛ(·, ·).

против двусторонней альтернативы

$$H_1: \rho \neq 0.$$

Другими словами, проверяется статистическая гипотеза о том, что в генеральной совокупности случайные величины (факторы)  $X$  и  $Y$  не коррелируют. Так как двумерная случайная величина  $(X, Y)$  по предположению имеет совместное нормальное распределение, некоррелируемость означает независимость факторов. Проверка гипотезы о независимости факторов основана на следующем результате: при справедливости нулевой гипотезы  $t$ -статистика

$$t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} \underset{H_0}{\sim} t_{n-2}$$

имеет распределение Стьюдента с  $n - 2$  степенями свободы. Таким образом, получаем следующий статистический критерий проверки нулевой гипотезы:

при заданном уровне значимости  $\alpha$  гипотеза  $H_0$  отвергается в пользу альтернативы  $H_1$  при  $|t| > t_{кр}$ ,

где  $t_{кр} = t(\alpha; n - 2)$  есть *двустороннее* критическое значение распределения Стьюдента  $t_{n-2}$ . Напомним, что двустороннее критическое значение определяется как решение уравнения

$$P(|t_{n-2}| > t_{кр}) = \alpha.$$

При  $|t| < t_{кр}$  говорят, что данные *согласуются* с нулевой гипотезой или *не противоречат ей*, и  $H_0$  не отвергается.

**Пример.** На основе  $n = 62$  выборочных данных, был рассчитан выборочный коэффициент корреляции  $r = 0,68$  между дневными логарифмическими доходностями<sup>1</sup> биржевых индексов NASDAQ и FTSE. Проверим значимость коэффициента корреляции, т. е. проверим статистическую гипотезу  $H_0$  о *независимости* доходностей обоих биржевых индексов (в предположении их *нормальной распределенности!*). Вычислим значение  $t$ -статистики:

$$t = \frac{0,68 \cdot \sqrt{62-2}}{\sqrt{1-0,68^2}} \approx 7,1838.$$

---

<sup>1</sup> Логарифмическая доходность рассчитывается как  $h_t = \ln(S_t/S_{t-1})$ .



Критическое значение распределения Стьюдента при уровне значимости  $\alpha = 5\%$  равно  $t_{кр} = t(5\%; 62 - 2) \approx 2,003$ . Так как  $|t| > t_{кр}$ , гипотеза  $H_0$  о независимости доходностей *отвергается*, коэффициент корреляции значим.

**Доверительный интервал для коэффициента корреляции.** Задача о построении доверительного интервала для коэффициента корреляции связана с той проблемой, что в общем случае (при  $\rho \neq 0$ )  $t$ -статистика имеет неизвестное распределение. Однако Фишер заметил, что если взять  $z$ -преобразование Фишера<sup>1</sup> от выборочного коэффициента корреляции

$$z(r) = \frac{1}{2} \ln \frac{1+r}{1-r},$$

то эта статистика при больших объемах выборки (а фактически уже при  $n > 6$ ) имеет распределение, близкое к нормальному:

$$z(r) \approx \mathcal{N}\left(z(\rho), \frac{1}{n-3}\right).$$

Следовательно, при заданной доверительной вероятности  $\gamma$  приближенный (асимптотический) доверительный интервал для  $z$ -преобразования Фишера коэффициента корреляции определяется как

$$P\left(z(r) - \frac{z_\gamma}{\sqrt{n-3}} < z(\rho) < z(r) + \frac{z_\gamma}{\sqrt{n-3}}\right) \approx \gamma. \quad (1.1)$$

Здесь  $z_\gamma$  есть *двустороннее* критическое значение стандартного нормального распределения при уровне значимости  $1 - \gamma$  и находится как решение уравнения

$$\Phi(z_\gamma) = \frac{1+\gamma}{2},$$

где  $\Phi(x)$  — функция стандартного нормального распределения.

Доверительный интервал для коэффициента корреляции получается применением к интервалу (1.1) *обратного преобразования Фишера*<sup>2</sup>

$$z^{-1}(x) = \text{th}(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}.$$

<sup>1</sup> В MS Excel функция ФИШЕР(.).

<sup>2</sup> В MS Excel функция ФИШЕРОБР(.).

Таким образом, асимптотический доверительный интервал для коэффициента корреляции имеет вид

$$P\left(z^{-1}\left(z(r) - \frac{z_\gamma}{\sqrt{n-3}}\right) < \rho < z^{-1}\left(z(r) + \frac{z_\gamma}{\sqrt{n-3}}\right)\right) \approx \gamma.$$

**Замечание.** Зная доверительный интервал для коэффициента корреляции, можно проверить его значимость, т. е. статистическую гипотезу  $H_0: \rho = 0$  при уровне значимости  $\alpha = 1 - \gamma$ . Нулевая гипотеза отвергается тогда и только тогда, когда нуль не принадлежит доверительному интервалу.

**Пример.** На основе  $n = 100$  выборочных данных, был рассчитан выборочный коэффициент корреляции  $r = 0,68$  между дневными логарифмическими доходностями биржевых индексов NASDAQ и FTSE. Построим доверительный интервал для коэффициента корреляции с доверительной вероятностью  $\gamma = 0,95$ . Применим  $z$ -преобразование Фишера  $z = z(0,68) \approx 0,8291$ . Критическое значение  $z_\gamma$  определяет-ся как решение уравнения

$$\Phi(z_\gamma) = \frac{1 + 0,95}{2} = 0,975,$$

откуда получаем  $z_\gamma = 1,96$ . Доверительный интервал для  $z(\rho)$  равен

$$\left(0,8291 - \frac{1,96}{\sqrt{100-3}}; 0,8291 + \frac{1,96}{\sqrt{100-3}}\right) = (0,6301; 1,0281).$$

Применив обратное преобразование Фишера, получаем доверительный интервал для коэффициента корреляции

$$P(0,5581 < \rho < 0,7732) = 0,95$$

( $0,5581 = z^{-1}(0,6301)$  и  $0,7732 = z^{-1}(1,0281)$ ).

Проверим значимость коэффициента корреляции, т. е. проверим нулевую гипотезу о *независимости* доходностей обоих биржевых индексов (в предположении их *нормальной распределенности!*). Так как нуль не принадлежит доверительному интервалу, нулевая гипотеза отвергается при уровне значимости  $\alpha = 1 - 0,95 = 0,05$ .

## § 1.2. Подгонка прямой. Метод наименьших квадратов

Рассмотрим следующую вспомогательную задачу. Пусть на координатной плоскости заданы  $n$  точек с координатами  $(x_i, y_i)_{i=1}^n$ . Требуется найти прямую, «меньше всего отклоняющуюся от заданных точек». Так как прямая задается уравнением

$$y = f(x) = \beta_0 + \beta_1 x,$$

зависящим от двух параметров  $\beta_0$  и  $\beta_1$ , необходимо по заданным значениям  $\{x_i\}$  и  $\{y_i\}$  найти значения этих параметров «оптимальной» прямой. Основной вопрос: что понимать под «наименьшим отклонением прямой от точек» и, более общим образом, как определить «меру отклонения прямой от точек»? Приведем несколько возможных подходов к определению меры  $\mu$  отклонения прямой от заданных точек:

1) сумма модулей отклонений в каждой точке  $x_i$ :

$$\mu = \sum_{i=1}^n |y_i - f(x_i)| = \sum_{i=1}^n |y_i - (\beta_0 + \beta_1 x_i)|;$$

2) сумма квадратов отклонений в каждой точке  $x_i$ :

$$\mu = \sum_{i=1}^n (y_i - f(x_i))^2 = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2;$$

3) сумма отклонений в каждой точке  $x_i$  с заданной весовой функцией  $\omega(\cdot) > 0$ :

$$\mu = \sum_{i=1}^n \omega(y_i - f(x_i)) = \sum_{i=1}^n \omega(y_i - (\beta_0 + \beta_1 x_i)).$$

С вероятностной точки зрения в случае нормального распределения выборочных данных «наилучшими вероятностными и статистическими свойствами» обладают оценки параметров прямой, полученные минимизацией суммы квадратов отклонений (второй случай). Этот метод получения оценок параметров оптимальной прямой называется *методом наименьших квадратов* (сокращенно МНК) или *Ordinary Least Squares* (сокращенно OLS), а полученные оценки параметров называются МНК- или OLS-оценками.

Итак, в качестве меры отклонения прямой от заданных на плоскости точек  $(x_i, y_i)_{i=1}^n$  возьмем сумму квадратов отклонений в каждой точке <sup>1</sup>:

$$S = S(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2.$$

Тогда параметры прямой, для которой эта мера отклонения минимальна, находятся как решение экстремальной задачи без ограничений:

$$S(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 \rightarrow \min.$$

Согласно необходимым условиям существования экстремума параметры оптимальной прямой находятся как решение системы уравнений

$$\begin{cases} \frac{\partial S}{\partial \beta_0} = \sum_{i=1}^n 2(y_i - \beta_0 - \beta_1 x_i) \cdot (-1) = 0, \\ \frac{\partial S}{\partial \beta_1} = \sum_{i=1}^n 2(y_i - \beta_0 - \beta_1 x_i) \cdot (-x_i) = 0. \end{cases}$$

После простых преобразований приходим к системе линейных уравнений

$$\begin{cases} n\beta_0 + \beta_1 \sum_{i=1}^n x_i = \sum_{i=1}^n y_i, \\ \beta_0 \sum_{i=1}^n x_i + \beta_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i, \end{cases} \quad (1.2)$$

называемой системой *нормальных уравнений*. Найдем явные формулы для решения этой системы. Для удобства разделим каждое уравнение в системе (1.2) на  $n$ :

$$\begin{cases} \beta_0 + \beta_1 \bar{x} = \bar{y}, \\ \beta_0 \bar{x} + \beta_1 \bar{x}^2 = \overline{xy}. \end{cases}$$

Выразим  $\beta_0$  из первого уравнения:

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

---

<sup>1</sup> Очевидно,  $S(\beta_0, \beta_1)$  есть многочлен второго порядка от параметров  $\beta_0$  и  $\beta_1$ .

и подставим во второе уравнение:

$$(\bar{y} - \beta_1 \bar{x})\bar{x} + \beta_1(\bar{x}^2) = \bar{x}\bar{y}.$$

После преобразования получаем (формально), что решение системы имеет вид

$$\hat{\beta}_1 = \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{(\overline{x^2}) - (\bar{x})^2} = \frac{\widehat{\text{cov}}(x, y)}{\widehat{\text{Var}}(x)} = \widehat{\text{corr}}(x, y) \sqrt{\frac{\widehat{\text{Var}}(y)}{\widehat{\text{Var}}(x)}} = \widehat{\text{corr}}(x, y) \frac{\hat{\sigma}_y}{\hat{\sigma}_x}$$

и

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x},$$

где  $\hat{\sigma}_x = \sqrt{\widehat{\text{Var}}(x)}$  и  $\hat{\sigma}_y = \sqrt{\widehat{\text{Var}}(y)}$  — выборочные стандартные отклонения величин  $x$  и  $y$  соответственно.

Несложно показать, что функция  $S(\beta_0, \beta_1)$  выпукла. Следовательно, решение системы нормальных уравнений (1.2) будет глобальным минимумом функции  $S(\beta_0, \beta_1)$ . Таким образом, оптимальная прямая задается уравнением

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x.$$

**Замечание 1.** Из первого уравнения системы (1.2) следует, что

$$\bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x},$$

т. е. оптимальная прямая проходит через точку с координатами  $(\bar{x}, \bar{y})$ .

**Замечание 2.** Несложно заметить, что система нормальных уравнений (1.2) имеет единственное решение тогда и только тогда, когда  $\widehat{\text{Var}}(x) \neq 0$ , т. е. когда не все значения  $x_i$  совпадают.

**Замечание 3.** Метод наименьших квадратов может быть применен для нахождения параметров любой функции, меньше всего отклоняющейся от заданных точек. Эта задача корректно разрешима в случае, когда неизвестные параметры входят в функцию линейно. Тогда система нормальных уравнений будет системой линейных уравнений и в общем случае будет иметь единственное решение.

### § 1.3. Парная линейная модель регрессии

Перейдем теперь к задаче количественного описания зависимости между двумя экономическими факторами  $y$  и  $x$  (например,  $y$  — уровень зарплаты индивидуума, а  $x$  — уровень образования (в годах)).

Естественно ожидать, что значение фактора  $y$  не всегда однозначно определяется значением фактора  $x$ . Так, уровень зарплаты зависит не только от уровня образования, но и от множества других факторов (стажа работы, возраста, индивидуальных способностей, места работы и пр.). Кроме того, учесть *все* факторы, влияющие на  $y$  помимо  $x$ , просто не представляется возможным в силу недостаточного количества информации или невозможности ее получения (например, как оценить или измерить индивидуальные способности индивидуума, несомненно влияющие на уровень зарплаты?). Также для одного значения фактора  $x$  могут наблюдаться различные значения фактора  $y$ .

Обычно для описания ситуаций с недостаточной информацией используют различные вероятностные математические модели. Рассмотрим подробно модель зависимости между факторами, описываемую уравнением

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, \dots, n, \quad (1.3)$$

где  $y_i$  и  $\varepsilon_i$  суть случайные величины, а  $x_i$  — *неслучайная* (детерминированная) величина,  $i$  — номер наблюдения. Фактор  $y$  называется *зависимой переменной* (dependent variable), а фактор  $x$  называется *регрессором* или *объясняющей переменной* (explanatory variable). Параметр  $\beta_1$  называется параметром *наклона прямой* (slope), а  $\beta_0$  — *константой, свободным членом* или *параметром сдвига* (intercept).

Уравнение (1.3) называется *уравнением регрессии* или *регрессионным уравнением*, а случайные величины  $\varepsilon_i$  называются *ошибками* регрессии. Ошибки регрессии удобно представлять себе как «неучтенные факторы», влияющие на  $y$  помимо фактора  $x$ . Таким образом, уравнение (1.3) отражает наши представления о характере зависимости между факторами.

Относительно ошибок регрессии будем предполагать выполнение следующих условий, называемых иногда условиями Гаусса—Маркова:

- 1)  $E\varepsilon_i = 0$ ,  $i = 1, \dots, n$  (ошибки регрессии несистематические);
- 2)  $\text{Var}(\varepsilon_i) = \sigma^2$  не зависит от  $i$ ;
- 3)  $\text{cov}(\varepsilon_i, \varepsilon_j) = 0$  при  $i \neq j$  (некоррелируемость ошибок для разных наблюдений);
- 4)  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ ,  $i = 1, \dots, n$  (нормальная распределенность ошибок регрессии).

Из условия  $E\varepsilon_i = 0$  следует, что

$$E y_i = \beta_0 + \beta_1 x_i,$$

т. е. среднее значение фактора  $y$  при заданном значении  $x_i$  равно  $\beta_0 + \beta_1 x_i$  и не зависит от ошибок регрессии. Отсюда термин: несистематические ошибки.

Очевидно,  $\text{Var}(y_i) = \text{Var}(\varepsilon_i)$  (так как значения  $x_i$  детерминированы). Следовательно, условие постоянства дисперсий ошибок регрессии влечет за собой постоянство дисперсий случайных величин  $y_i$ . Следует напомнить, что дисперсию  $\text{Var}(y_i)$  можно рассматривать как «меру разброса» значений случайной  $y_i$  величины относительно своего среднего значения (математического ожидания)  $E y_i = \beta_0 + \beta_1 x_i$ . Если смотреть на ошибки регрессии как на «неучтенные факторы», условие постоянства дисперсий можно описательно трактовать следующим образом: «степень влияния» невключенных в модель факторы в разных наблюдениях постоянна. Условие постоянства дисперсий ошибок называется *гомоскедастичностью* (homoskedasticity), и говорят, что ошибки модели регрессии гомоскедастичны или однородны. При нарушении условия постоянства дисперсий ошибок регрессии говорят, что ошибки *гетероскедастичны* или неоднородны.

Условие некоррелируемости (независимости в случае нормального распределения) ошибок для разных наблюдений можно трактовать как «локальность» их влияния: не включенные в модель факторы, которые моделируются ошибками регрессии, влияют только на «свое» наблюдение и не влияют на другие. В случае пространственных выборок (cross-sectional data) это условие обычно считается выполненным. Оно, как правило, нарушается в случае построения регрессионных моделей для временных рядов.

### 1.3.1. Теорема Гаусса—Маркова

Итак, мы предполагаем, что зависимость между факторами  $y$  и  $x$  описывается уравнением регрессии (1.3), но параметры уравнения  $\beta_0$ ,  $\beta_1$  и  $\sigma^2$  нам неизвестны.

Основная задача — получить «наилучшие» оценки параметров регрессии на основе выборочных данных. Ограничимся рассмотрением только оценок параметров, линейных относительно  $y_i$ . Под «наилучшими» будем подразумевать несмещенные оценки с минималь-

ной дисперсией<sup>1</sup>. Такие оценки называются BLUE-оценками (BLUE = Best Linear Unbiased Estimators) или эффективными оценками.

Основным результатом является следующая теорема.

**Теорема** (Гаусс—Марков). Пусть для линейной модели парной регрессии

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, \dots, n,$$

выполнены условия 1—3 на ошибки регрессии  $\varepsilon_i$ . Тогда OLS-оценки  $\hat{\beta}_0$  и  $\hat{\beta}_1$  параметров  $\beta_0$  и  $\beta_1$  являются BLUE-оценками, т. е. среди несмещенных линейных (относительно  $y_i$ ) оценок имеют наименьшую дисперсию.

**Доказательство.** Докажем несмещенность OLS-оценок. Рассмотрим сначала оценку параметра  $\beta_1$ . Для нее имеем следующее выражение:

$$\hat{\beta}_1 = \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{\bar{x}^2 - (\bar{x})^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n((\bar{x})^2 - \bar{x}^2)} = \frac{\sum_{i=1}^n (x_i - \bar{x})y_i}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$

Так как величины  $x_i$  неслучайны и  $E y_i = \beta_0 + \beta_1 x_i$  (условие 1 на ошибки регрессии), мы получаем

$$\begin{aligned} E \hat{\beta}_1 &= \frac{\sum (x_i - \bar{x}) E y_i}{\sum (x_i - \bar{x})^2} = \frac{\sum (x_i - \bar{x}) (\beta_0 + \beta_1 x_i)}{\sum (x_i - \bar{x})^2} = \\ &= \beta_0 \frac{\sum (x_i - \bar{x})}{\sum (x_i - \bar{x})^2} + \beta_1 \frac{\sum (x_i - \bar{x}) x_i}{\sum (x_i - \bar{x})^2} = \beta_1. \end{aligned}$$

При выводе мы воспользовались равенствами

$$\sum (x_i - \bar{x}) = 0, \quad \sum (x_i - \bar{x})^2 = \sum (x_i - \bar{x}) x_i.$$

Далее, так как

$$E(\bar{y}) = E\left(\frac{1}{n} \sum_{i=1}^n y_i\right) = \frac{1}{n} \sum_{i=1}^n E y_i = \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_i) = \beta_0 + \beta_1 \bar{x},$$

<sup>1</sup> Напомним, что оценка параметров вероятностной модели в математической статистике рассматривается как случайная величина.



для оценки константы  $\hat{\beta}_0$  в уравнении регрессии получаем

$$\begin{aligned}\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} &\Rightarrow E(\hat{\beta}_0) = E(\bar{y}) - E(\hat{\beta}_1 \bar{x}) = \\ &= \beta_0 + \beta_1 \bar{x} - \bar{x} \cdot E(\hat{\beta}_1) = \beta_0 + \beta_1 \bar{x} - \beta_1 \bar{x} = \beta_0.\end{aligned}$$

Итак,  $\hat{\beta}_0$  и  $\hat{\beta}_1$  — несмещенные (unbiased) оценки параметров  $\beta_0$  и  $\beta_1$  уравнения регрессии.

Вычислим теперь дисперсии оценок  $\hat{\beta}_0$  и  $\hat{\beta}_1$ . Для этого воспользуемся тем фактом, что из условий 2 и 3 на ошибки регрессии следует, что  $\text{Var}(y_i) = \sigma^2$  и  $\text{cov}(y_i, y_j) = 0$  при  $i \neq j$ . Следовательно, используя свойства дисперсии, для оценки  $\hat{\beta}_1$  получаем

$$\begin{aligned}\text{Var}(\hat{\beta}_1) &= \text{Var} \left( \frac{\sum_{i=1}^n (x_i - \bar{x}) y_i}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) = \frac{\text{Var} \left( \sum_{i=1}^n (x_i - \bar{x}) y_i \right)}{\left( \sum_{i=1}^n (x_i - \bar{x})^2 \right)^2} = \\ &= \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \text{Var}(y_i)}{\left( \sum_{i=1}^n (x_i - \bar{x})^2 \right)^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \sigma^2}{\left( \sum_{i=1}^n (x_i - \bar{x})^2 \right)^2} = \sigma^2 \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\left( \sum_{i=1}^n (x_i - \bar{x})^2 \right)^2} = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}.\end{aligned}$$

Для нахождения дисперсии оценки  $\hat{\beta}_0$  сначала перепишем ее в виде

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = \sum_{i=1}^n \frac{1}{n} y_i - \bar{x} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) y_i = \sum_{i=1}^n \left( \frac{1}{n} - \frac{\bar{x}(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) y_i.$$

Следовательно,

$$\begin{aligned}\text{Var}(\hat{\beta}_0) &= \text{Var} \left( \sum_{i=1}^n \left( \frac{1}{n} - \frac{\bar{x}(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) y_i \right) = \\ &= \sum_{i=1}^n \left( \frac{1}{n} - \frac{\bar{x}(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)^2 \text{Var}(y_i) = \\ &= \sigma^2 \sum_{i=1}^n \left( \frac{1}{n^2} - 2 \frac{\bar{x}(x_i - \bar{x})}{n \sum_{i=1}^n (x_i - \bar{x})^2} + \frac{(\bar{x})^2 (x_i - \bar{x})^2}{\left( \sum_{i=1}^n (x_i - \bar{x})^2 \right)^2} \right) =\end{aligned}$$

$$\begin{aligned}
 &= \sigma^2 \left( \sum_{i=1}^n \frac{1}{n^2} - \frac{2\bar{x} \sum_{i=1}^n (x_i - \bar{x})}{n \sum_{i=1}^n (x_i - \bar{x})^2} + \frac{(\bar{x})^2 \sum_{i=1}^n (x_i - \bar{x})^2}{\left( \sum_{i=1}^n (x_i - \bar{x})^2 \right)^2} \right) = \\
 &= \sigma^2 \left( \frac{1}{n} + \frac{(\bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) = \frac{\sigma^2 \sum_{i=1}^n x_i^2}{n \sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sigma^2 \cdot \bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}.
 \end{aligned}$$

Покажем теперь, что любая другая линейная несмещенная оценка имеет большую дисперсию. Пусть  $\tilde{\beta}_1 = \sum c_i y_i$  — произвольная линейная (по  $y_i$ ) несмещенная оценка параметра наклона  $\beta_1$ . Представим ее коэффициенты  $c_i$  как  $c_i = \omega_i + \theta_i$ , где

$$\tilde{\beta}_1 = \sum \omega_i y_i, \quad \omega_i = (x_i - \bar{x}) / \sum (x_i - \bar{x})^2.$$

Так как  $E\tilde{\beta}_1 = E\hat{\beta}_1 = \beta_1$ , мы получаем

$$\begin{aligned}
 0 = E\tilde{\beta}_1 - E\hat{\beta}_1 &= E(\tilde{\beta}_1 - \hat{\beta}_1) = E\left(\sum \theta_i y_i\right) = \sum \theta_i E y_i = \\
 &= \sum \theta_i (\beta_0 + \beta_1 x_i) = \beta_0 \sum \theta_i + \beta_1 \sum \theta_i x_i.
 \end{aligned}$$

Так как это равенство должно быть выполнено для произвольных значений  $\beta_0$  и  $\beta_1$ , получаем, что

$$\sum \theta_i = 0, \quad \sum \theta_i x_i = 0.$$

Далее,

$$\begin{aligned}
 \text{Var}(\tilde{\beta}_1) &= \text{Var}\left(\sum c_i y_i\right) = \sum c_i^2 \text{Var}(y_i) = \\
 &= \sigma^2 \sum (\omega_i + \theta_i)^2 = \sigma^2 \left(\sum \omega_i^2 + 2 \sum \omega_i \theta_i + \sum \theta_i^2\right).
 \end{aligned}$$

По условию  $\omega_i = (x_i - \bar{x}) / \left(\sum (x_i - \bar{x})^2\right)$ , поэтому

$$\sum \omega_i \theta_i = \sum \frac{(x_i \theta_i - \bar{x} \theta_i)}{\sum (x_i - \bar{x})^2} = \frac{\sum x_i \theta_i - \bar{x} \sum \theta_i}{\sum (x_i - \bar{x})^2} = 0.$$

Так как  $\text{Var}(\hat{\beta}_1) = \sigma^2 \sum \omega_i^2$ , окончательно получаем

$$\begin{aligned}
 \text{Var}(\tilde{\beta}_1) &= \sigma^2 \left(\sum \omega_i^2 + \sum \theta_i^2\right) = \sigma^2 \sum \omega_i^2 + \sigma^2 \sum \theta_i^2 = \\
 &= \text{Var}(\hat{\beta}_1) + \sigma^2 \sum \theta_i^2 \geq \text{Var}(\hat{\beta}_1).
 \end{aligned}$$

Таким образом,  $\text{Var}(\tilde{\beta}_1) \geq \text{Var}(\hat{\beta}_1)$ .

Аналогично можно показать, что для произвольной несмещенной оценки  $\hat{\beta}_0$  параметра  $\beta_0$  всегда выполняется неравенство

$$\text{Var}(\tilde{\beta}_0) \geq \text{Var}(\hat{\beta}_0).$$

Теорема доказана. □

**Замечание 1.** Из доказательства видно, что для несмещенности OLS-оценок достаточно *только* условия 1 на ошибки регрессии.

**Замечание 2.** Можно показать, что

$$\text{cov}(\hat{\beta}_0, \hat{\beta}_1) = -\frac{\sigma^2 \bar{x}}{\sum (x_i - \bar{x})^2}.$$

**Замечание 3.** Из теоремы Гаусса—Маркова следует, что среди линейных по  $y$  несмещенных оценок параметров  $\beta_0$  и  $\beta_1$  наилучшими (т. е. с минимальной дисперсией) будут OLS-оценки. Однако могут существовать и нелинейные оценки параметров  $\beta_0$  и  $\beta_1$  с дисперсией, меньшей, чем у OLS-оценок.

Найдем теперь оценку третьего параметра уравнения регрессии — дисперсии ошибок  $\sigma^2$ . Обозначим через

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

прогноз фактора  $y$  при заданном значении  $x_i$ . Значения  $\hat{y}_i$  также называются *подогнанными* (fitted value) или *предсказанными* значениями зависимой переменной.

**Определение.** *Остатки* (residual) модели регрессии определяются равенством  $e_i = y_i - \hat{y}_i$ .

Важно различать в модели регрессии ошибки  $\varepsilon_i$  и остатки  $e_i$ . Остатки также являются случайными величинами, но, в отличие от ошибок (имеющих теоретический характер), они наблюдаемы. Кроме того, для остатков всегда выполнено соотношение  $\sum_{i=1}^n e_i = 0$ , следующее из первого уравнения системы (1.2), т. е. остатки *всегда зависят*, в отличие от ошибок регрессии  $\varepsilon_i$ . Но, тем не менее, можно считать, что остатки в некотором смысле «моделируют» ошибки регрессии и «наследуют» их свойства. На этом основаны методы исследования отклонений выборочных данных от предположений теоремы Гаусса—Маркова.

Введем следующее обозначение:

$$\text{RSS} = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2.$$

Величина RSS называется *остаточной суммой квадратов* (residual sum of squares) в модели регрессии. Можно показать, что

$$E(\text{RSS}) = (n - 2)\sigma^2.$$

Следовательно, статистика

$$s^2 = \frac{\text{RSS}}{n - 2} = \frac{1}{n - 2} \sum_{i=1}^n e_i^2$$

является несмещенной оценкой дисперсии ошибок регрессии. Выборочная *стандартная ошибка* регрессии SER (Standard Error of Regression) определяется как

$$\text{SER} = s = \sqrt{s^2} = \sqrt{\frac{\text{RSS}}{n - 2}}.$$

### 1.3.2. Статистические свойства OLS-оценок коэффициентов

При доказательстве теоремы Гаусса—Маркова мы нашли дисперсии оценок параметров регрессии  $\hat{\beta}_0$  и  $\hat{\beta}_1$ . В выражениях для дисперсий участвует дисперсия ошибок  $\sigma^2$ , значение которой в большинстве прикладных задач неизвестно. Поэтому в прикладных вычислениях используют *оценки дисперсий* величин  $\hat{\beta}_0$  и  $\hat{\beta}_1$ :

$$\begin{aligned}\widehat{\text{Var}}(\hat{\beta}_0) &= \frac{s^2 \cdot \bar{x}^2}{\sum (x_i - \bar{x})^2}, \\ \widehat{\text{Var}}(\hat{\beta}_1) &= \frac{s^2}{\sum (x_i - \bar{x})^2}, \\ \widehat{\text{cov}}(\hat{\beta}_0, \hat{\beta}_1) &= -\frac{s^2 \cdot \bar{x}}{\sum (x_i - \bar{x})^2},\end{aligned}$$

получаемые формальной заменой неизвестного параметра  $\sigma^2$  в выражениях для дисперсии и ковариации оценок коэффициентов на его несмещенную оценку  $s^2$ . Стандартные ошибки оценок коэффици-