

Федеральное государственное бюджетное  
образовательное учреждение высшего  
профессионального образования  
«Московский государственный университет  
путей сообщения»

---

Кафедра “Математика”

А.С.МИЛЕВСКИЙ

# ЭКОНОМЕТРИКА

*Учебное пособие*

МОСКВА – 2015

Федеральное государственное бюджетное  
образовательное учреждение высшего  
профессионального образования  
«Московский государственный университет  
путей сообщения»

---

Кафедра “Математика”

А.С.МИЛЕВСКИЙ

# ЭКОНОМЕТРИКА

Рекомендовано редакционно-издательским  
советом университета в качестве учебного  
пособия для студентов ИЭФ

МОСКВА – 2015

УДК-330.43

М-60

Милевский А.С. Эконометрика: Учебное пособие. – М.: МГУПС (МИИТ), 2015. – 154 с.

Конспект лекций предназначен для подготовки бакалавров по направлениям “Экономика”, “Менеджмент”. Включает в себя материал по корреляционному анализу, классической линейной модели множественной регрессии, нарушениям предположений классической модели, временным рядам, динамическим моделям и системам эконометрических уравнений.

Рецензенты:

Соболева Е.С., к.ф.-м.н., доцент кафедры  
“Математический анализ” МГУ им. М.В.Ломоносова.

Деснянский В.Н., к.ф.-м.н., зав. кафедрой  
“Математический анализ” МИИТ.

© МГУПС (МИИТ), 2015

Св. план 2015 г., поз. 175

Милевский Александр Станиславович

ЭКОНОМЕТРИКА.

Учебное пособие

---

Подписано в печать

Формат 60x84 / 16

Заказ №

Усл. печ. л.

Тираж – 100 экз.

---

150048, г. Ярославль, Московский пр-т, д.151.  
Типография Ярославского филиала МИИТ.

# 1. Введение

## 1.1. Предмет эконометрики

**Эконометрика** – это наука, изучающая конкретные **количественные и качественные взаимосвязи** экономических объектов и процессов при помощи **математических** методов и моделей.

Норвежский ученый Р. Фриш определил эконометрику как «...единство экономической теории, статистики и математики».

Основные результаты *экономической теории* носят *качественный* характер, например, даются ответы на вопросы вида:

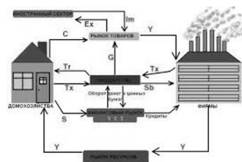
- Какие факторы влияют на совокупный объём инвестиций в частном бизнесе?
- Как зависят темпы экономического роста от размеров дефицита государственного бюджета, объёма инвестиций, внешнего долга, уровня инфляции?
- Каким образом изменит структуру потребления снижение ставок налогов?

Однако для принятия обоснованных решений часто требуется отвечать на вопросы, начинающиеся словами «**на сколько**». Например,

- Какого изменения объёма потребления электроэнергии можно ожидать, если предположить, что цены на неё возрастут на 30%?
- Насколько увеличится спрос на данную продукцию, если удвоить текущие расходы на её рекламу?

Это *количественное содержание* и является предметом изучения эконометрики.

Следует отметить, что зачастую не удаётся достичь желаемой степени точности предсказаний. Модели и уравнения, дающие хорошее согласование для одного периода времени, могут оказаться неудовлетворительными для другого. Это связано с большой сложностью экономических процессов.



Но в любом случае анализ в эконометрике производится на основе широкого применения самых разнообразных *математических методов и моделей*.

И, конечно, значимые результаты можно получить лишь путем обработки *реальных статистических данных*.

## 1.2. Модели и статистические данные

Закономерности в экономике проявляются как взаимосвязи между экономическими показателями. Как правило, их стремятся выразить в сравнительно простой математической форме.

Надо иметь в виду, что математическая модель является лишь упрощённым, концентрированным представлением реальности.

Искусство построения модели состоит в совмещении “желаемого с возможным”.

### Пример.

Изучая стоимость  $Y$  квадратного метра квартиры в Москве, естественно предполагать, что она зависит от

- количества комнат  $N$ ,
- жилой площади квартиры  $P$ ,
- нежилой площади  $Q$ ,
- расстояния от центра города  $R$ ,
- расстояния от ближайшей станции метро  $S$ ,



- вида дома (кирпичный, панельный и т.п.) Т,
- этажа U,..

$$Y = f(N, P, Q, R, S, T, U, \dots)$$

В любом случае влияние неучтённых (или неизвестных) факторов будет представлено как некая случайная добавка

$\varepsilon$  к  $Y$ .

Наличие *случайной составляющей*  $\varepsilon$  в модели может быть обусловлено причинами двоякой природы.

- Во-первых, она отражает влияние на формирование значения  $Y$  факторов, не учтённых в перечне объясняющих переменных  $X$
- Во-вторых, она может включать в себя погрешность в измерении  $Y$ .

В реальности возможны также ошибки в выборе вида функциональной зависимости и т.п.

В этом примере видна различная роль рассматриваемых переменных  $Y, N, P, Q, R, S, T, U$  (или, как их ещё называют, *факторов*).

В пределах данной модели одни переменные являются *независимыми* (или *объясняющими*, *экзогенными* –  $N, P, Q, R, S, T, U$ ), а другие *зависимыми* (или *эндогенными* –  $Y$ ).

В общем случае экономическая модель с одной *эндогенной*  $Y$  и  $m$  объясняющими переменными  $x_1, x_2, \dots, x_m$  имеет вид:

$$Y = f(x_1, x_2, \dots, x_m, \varepsilon)$$

Конкретный вид функции может быть различным.

Выбор вида функциональной зависимости называется **спецификацией модели**, а определение состава объясняющих переменных - **спецификацией переменных**.

### 1.3. **Статистические данные**

Основу эконометрического моделирования составляют статистические данные. Их различают по типам.

**Пространственные данные** собираются по какому-либо экономическому показателю для разных объектов в один момент времени или в разные моменты в случае, когда время несущественно.

Например, данные по курсам покупки/продажи валюты сегодня в различных обменных пунктах Москвы.

**Временные ряды** – данные для одного объекта в различные моменты времени.

Например, ежеквартальные данные об инфляции, данные о денежной эмиссии за последние годы, ежедневный курс доллара.

Промежуточное положение занимают **панельные данные**, которые отражают наблюдения по большому числу объектов за относительно небольшой период времени.

Например, данные об объёмах выпуска и капитальных вложениях российских предприятий топливно-энергетического комплекса за 1995–2000 г.



## 1.4. Основные этапы эконометрического моделирования

В процессе эконометрического исследования можно выделить следующие этапы:

*1-й этап (постановочный).* Определение конечных целей моделирования, формирование набора участвующих в модели факторов и их ролей.

*2-й этап (спецификация модели).* Выбор общего вида модели, состава и вида входящих в неё функциональных взаимосвязей.

*3-й этап (информационный).* Сбор необходимой статистической информации, т.е. значений участвующих в модели факторов.

*4-й этап (идентификация модели).* Анализ модели, статистическое оценивание параметров модели.

*5-й этап (проверка модели).*

Сопоставление реальных и оценённых данных, проверка адекватности модели.

В результате проверки модели нередко приходится неоднократно возвращаться к предыдущим этапам и вносить коррективы.

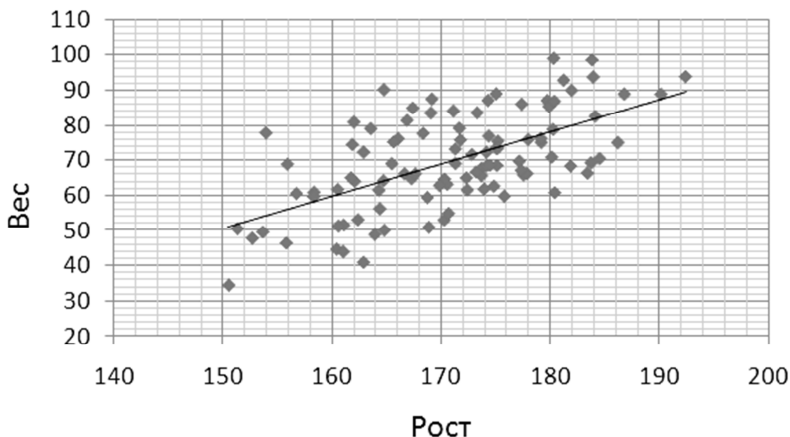


## 2. Корреляционный анализ количественных зависимостей

### 2.1. Выборочный коэффициент корреляции

При изучении связи между факторами сразу возникает

**Вопрос:** в какой степени две переменные **СОВМЕСТНО ИЗМЕНЯЮТСЯ?** (т.е., влечёт ли за собой увеличение одной переменной увеличение или уменьшение другой, или не влечёт).



Представляет интерес наклон (*направление связи*) и ширина (*сила связи*) воображаемого эллипса

Тесноту линейной связи между двумя количественными факторами обычно измеряют при помощи коэффициента корреляции. Он вычисляется по формуле:

$$r_{XY}^* = \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{\sqrt{D_x^*} \cdot \sqrt{D_y^*}}$$

Известно, что

$$-1 \leq r_{XY}^* \leq 1$$

и линейная зависимость между X и Y тем сильнее, чем ближе *модуль коэффициента корреляции* к 1.

**Пример.**

X\Y	-1	0	2
0-20	5	3	4
20-40	2	1	5

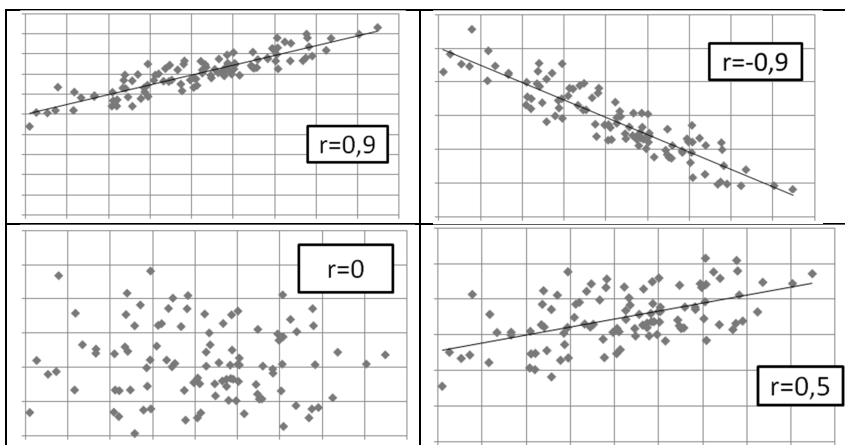
$$\bar{x} = \frac{10 \cdot 12 + 30 \cdot 8}{20} = \dots,$$

$$\bar{y} = \dots,$$

$$\overline{x^2} = \frac{10^2 \cdot 12 + \dots}{20} = \dots, D_x^* = \overline{x^2} - (\bar{x})^2 = \dots,$$

$$\overline{y^2} = \dots, D_y^* = \overline{y^2} - (\bar{y})^2 = \dots$$

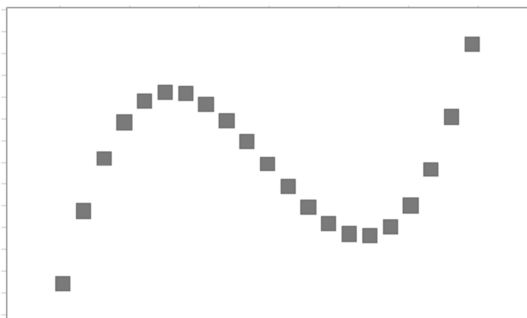
$$\overline{xy} = \frac{10 \cdot (-1) \cdot 5 + \dots}{20} = \dots; r_{xy}^* = \dots$$



## 2.2. **Замечания о коэффициенте корреляции**

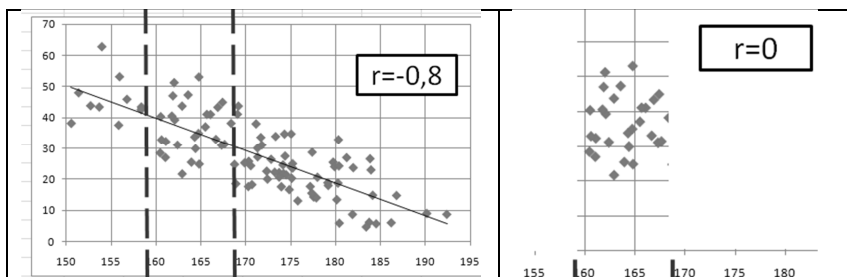
**А.** Коэффициент корреляции оценивает только линейную связь переменных!

Он не покажет наличие нелинейной связи.



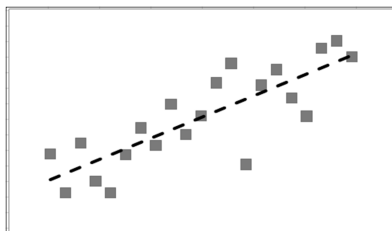
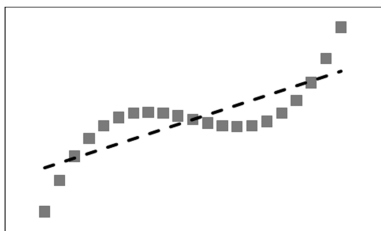
Здесь связь переменных есть, и она очень сильная, но  $r = 0$ .

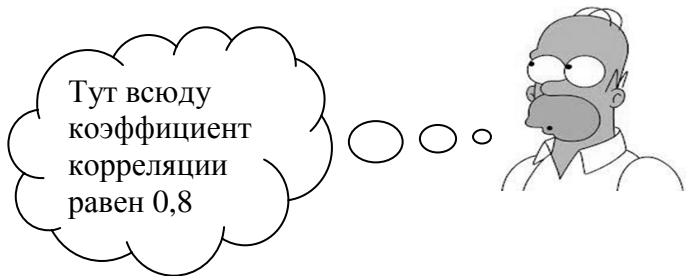
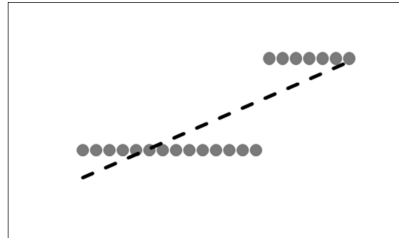
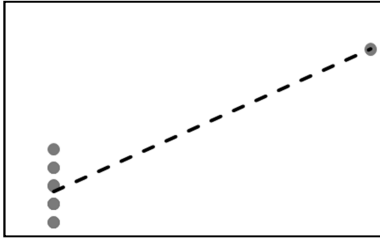
**Б.** Необходимо, чтобы у переменных была значительная изменчивость! Если сформировать выборку изначально однотипных по одному из рассматриваемых факторов, нечего надеяться выявить там корреляцию фактора с другими.



**В.** Корреляция совершенно не подразумевает наличие причинно-следственной связи!

**Г.** Зная только коэффициент корреляции, нельзя сколько-нибудь детально описать зависимость.





### 2.3. Матрица корреляций

Корреляционная матрица — это квадратная таблица, в которой на пересечении строки и столбца находится коэффициент корреляции между соответствующими параметрами.

**Пример.** Для следующих табличных данных составить корреляционную матрицу

$y$	$x_1$	$x_2$
7	2	3
14	2	4
11	5	7
12	7	10

**Решение.** 1) Составляем матрицу

$$M = \begin{pmatrix} 7 & 2 & 3 \\ 14 & 2 & 4 \\ 11 & 5 & 7 \\ 12 & 7 & 10 \end{pmatrix}$$

2) Из каждого столбца вычитаем его среднее значение (например, из второго столбца вычитаем “4”):

$$M_1 = \begin{pmatrix} -4 & -2 & -3 \\ 3 & -2 & -2 \\ 0 & 1 & 1 \\ 1 & 3 & 4 \end{pmatrix}$$

3) Умножаем  $M_1^T$  на  $M_1$  и делим на  $n$ , получаем *матрицу ковариаций*:

$$M_1^T \cdot M_1 = \begin{pmatrix} -4 & 3 & 0 & 1 \\ -2 & -2 & 1 & 3 \\ -3 & -2 & 1 & 4 \end{pmatrix} \cdot \begin{pmatrix} -4 & -2 & -3 \\ 3 & -2 & -2 \\ 0 & 1 & 1 \\ 1 & 3 & 4 \end{pmatrix} / 4 = \begin{pmatrix} 6,5 & 1,25 & 2,5 \\ 1,25 & 4,5 & 5,75 \\ 2,5 & 5,75 & 7,5 \end{pmatrix}$$

На диагонали в матрице стоят дисперсии, например  $Dy=6,5$ .

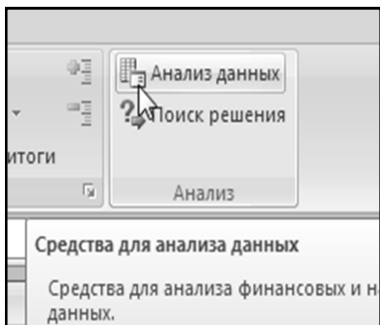
4) Делим каждое число на корень из произведения соответствующих дисперсий (например, 1,25 на корень из  $6,5 \cdot 4,5$ , 2,5 на корень из  $6,5 \cdot 7,5$ , получаем матрицу корреляций:

$$r_{YX} = \begin{pmatrix} 1 & 0,2311 & 0,3581 \\ 0,2311 & 1 & 0,9898 \\ 0,3581 & 0,9898 & 1 \end{pmatrix}$$

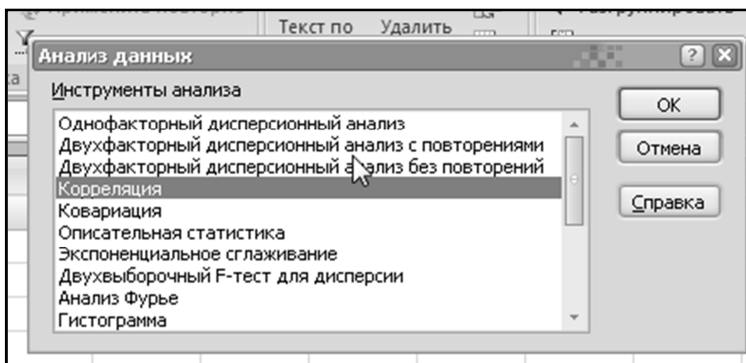
В MS Excel для вычисления корреляционных матриц используется инструмент **Корреляция** из надстройки **Анализ данных**.

	A	B	C
1	Y	X1	X2
2	7	2	3
3	14	2	4
4	11	5	7
5	12	7	10

Вызвать надстройку  
**Данные** → **Анализ**  
**данных**;

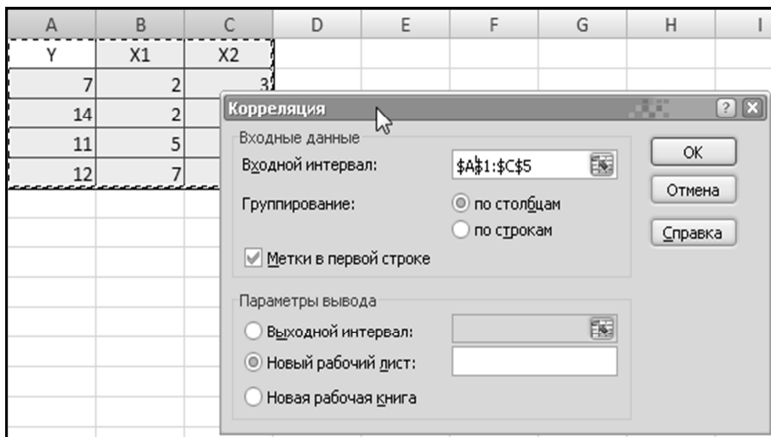


В появившемся списке **Инструменты анализа** выбрать строку **Корреляция** и нажать кнопку ОК;



В появившемся окне указать **Входной интервал**.

Установить переключатель в разделе Группировка (“по столбцам” или “по строкам”);



В выходной диапазон будет выведена корреляционная матрица.

	A	B	C	D
1		Y	X1	X2
2	Y	1		
3	X1	0,231125	1	
4	X2	0,358057	0,989762	1

Хотя в результате будет получена треугольная матрица, корреляционная матрица симметрична. Подразумевается, что в пустых клетках в правой верхней половине таблицы находятся те же коэффициенты корреляции, что и в нижней левой (симметрично расположенные относительно диагонали из единиц).

## 2.4. Частная корреляция

Корреляция между двумя переменными, вычисленная после устранения влияния всех других переменных, называется *частной корреляцией*.



Например, *длина волос* может коррелировать с *ростом* человека (чем выше человек, тем короче волосы), однако эта зависимость становится слабой или совсем исчезает, если устранить влияние *пола* людей, поскольку женщины обычно ниже ростом и чаще имеют более длинные волосы, чем мужчины.

**Частный коэффициент корреляции** между переменными X и Y при исключении влияния переменной Z вычисляется по формуле

$$r_{XY|Z} = \frac{r_{XY} - r_{XZ}r_{YZ}}{\sqrt{(1-r_{XZ}^2)(1-r_{YZ}^2)}}$$

**Пример.** Для предыдущего примера

$$r_{YX} = \begin{pmatrix} 1 & 0,2311 & 0,3581 \\ 0,2311 & 1 & 0,9898 \\ 0,3581 & 0,9898 & 1 \end{pmatrix}$$

$$r_{yx_1|x_2} = \frac{0,2311 - 0,3581 \cdot 0,9898}{\sqrt{(1 - 0,3581^2) \cdot (1 - 0,9898^2)}} = -0,925$$

$$r_{yx_2|x_1} = \frac{0,3581 - 0,2311 \cdot 0,9898}{\sqrt{(1 - 0,2311^2) \cdot (1 - 0,9898^2)}} = 0,9311$$

И так далее...

$$r_{YX}^{частн} = \begin{pmatrix} 1 & -0,925 & 0,9311 \\ -0,925 & 1 & 0,9984 \\ 0,9311 & 0,9984 & 1 \end{pmatrix}$$

### 3. Метод наименьших квадратов

#### 3.1. Постановка задачи

Нередко возникает задача: найти кривую заданного вида, наиболее точно приближающую экспериментальные данные.

Математически это формулируется так:

требуется подобрать такие значения параметров  $\beta$ , чтобы график функции

$$f(x) = \beta_1 \varphi_1(x) + \beta_2 \varphi_2(x) + \dots + \beta_m \varphi_m(x)$$

проходит как можно ближе к заданным точкам  $(x_i, y_i)$ ,  $i=1, 2, \dots, n$ . Здесь  $\beta_i$  – неизвестные коэффициенты,  $\varphi_i$  – известные функции.

Нужно как-то конкретизировать понятие *близости* точек к кривой. Это можно сделать многими способами, наиболее распространён следующий:

Для заданных точек  $(x_i, y_i)$ ,  $i=1, 2, \dots, n$ , найти такие значения  $\beta_i$ , чтобы минимизировать **остаточную сумму квадратов**:

$$ESS = \sum_{i=1}^n (y_i - f(x_i))^2 \rightarrow \min \quad (*)$$

### 3.2. Формулы МНК

Приравнивая производные ESS нулю, получаем систему линейных уравнений для нахождения  $\beta_i$

$$\begin{cases} A_{11}\beta_1 + A_{12}\beta_2 + \dots + A_{1m}\beta_m = B_1 \\ A_{21}\beta_1 + A_{22}\beta_2 + \dots + A_{2m}\beta_m = B_2 \\ \dots \quad \dots \quad \dots \quad \dots \\ A_{m1}\beta_1 + A_{m2}\beta_2 + \dots + A_{mm}\beta_m = B_m \end{cases}$$

где

$$A_{kl} = \sum_{i=1}^n \varphi_k(x_i)\varphi_l(x_i), \quad B_k = \sum_{i=1}^n y_i\varphi_k(x_i)$$

Можно сразу и решить эту систему, если использовать *матричные обозначения*

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad X = \begin{pmatrix} \varphi_1(x_1) & \varphi_2(x_1) & \dots & \varphi_m(x_1) \\ \varphi_1(x_2) & \varphi_2(x_2) & \dots & \varphi_m(x_2) \\ \vdots & \vdots & \vdots & \vdots \\ \varphi_1(x_n) & \varphi_2(x_n) & \dots & \varphi_m(x_n) \end{pmatrix}, \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_m \end{pmatrix}$$

Тогда

$$\beta = (X^T X)^{-1} X^T Y$$

**Пример.** Имеются следующие экспериментальные данные

№ п/п	1	2	3	4	5	6	7	8
X	2	3	2	2	4	3	5	5
Y	1	2	4	3	4	5	7	9

Для этих данных требуется подобрать наилучшую параболу

$y = \beta_1 + \beta_2 x + \beta_3 x^2$  методом наименьших квадратов

$$A_{11} = \sum \varphi_1 \varphi_1 = \sum (1 \cdot 1) = 8, \quad A_{12} = A_{21} = \sum \varphi_1 \varphi_2 = \sum (1 \cdot x) = 26,$$

$$A_{13} = A_{31} = \sum \varphi_1 \varphi_3 = \sum (1 \cdot x^2) = 96, \quad A_{22} = \sum \varphi_2 \varphi_2 = \sum x^2 = 96,$$

$$A_{23} = A_{32} = \sum \varphi_2 \varphi_3 = \dots, \quad A_{33} = \dots$$

$$B_1 = \sum y \varphi_1 = \sum y = 35, \quad B_2 = \sum (y \varphi_2) = \sum (yx) = \dots, \quad B_3 = \dots$$

Получаем систему уравнений

$$\begin{cases} 8\beta_1 + 26\beta_2 + 96\beta_3 = 35 \\ 26\beta_1 + 96\beta_2 + 392\beta_3 = 133, \\ 96\beta_1 + 392\beta_2 + 96\beta_3 = 559 \end{cases}$$

Откуда  $\beta_1 \approx 6,15$ ,  $\beta_2 \approx -3,06$ ,  $\beta_3 \approx 0,68$ .

Таким образом, наилучшее МНК-приближение в виде параболы для исходных данных имеет вид

$$y = 6,15 - 3,06x + 0,68x^2$$

Интересно сравнить экспериментальные и расчётные значения для этой формулы:

$x_i$	2	3	2	2	4	3	5	5
$y_i$	1	2	4	3	4	5	7	9
$y_i, \text{расч}$	2,76	3,1	2,76	2,76	4,8	3,1	7,87	7,87

## 4. Классическая модель множественной линейной регрессии

### 4.1. Описание модели

Модель множественной линейной регрессии имеет вид

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m + \varepsilon \quad (4.1)$$

где  $y$  – зависимая переменная,

$x_1, \dots, x_m$  – объясняющие переменные (также называемые регрессорами),

$\beta_0, \beta_1, \dots, \beta_m$  – коэффициенты регрессии,

$\varepsilon$  – случайный член (или “ошибки измерений”).

Коэффициенты  $\beta_0, \beta_1, \dots, \beta_m$  оцениваются на основе статистических данных.

Для нахождения оценок проводятся  $n$  наблюдений ( $n > m$ ), в результате чего получается  $n$  равенств

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_m x_{mi} + \varepsilon_i, \quad i = 1, 2, \dots, n$$

(индекс  $i$  показывает номер наблюдения).

Это удобно записать в матричной форме

$$Y = X\beta + \varepsilon, \quad (4.2)$$

где

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, X = \begin{pmatrix} 1 & x_{11} & \cdots & x_{m1} \\ 1 & x_{12} & \cdots & x_{m2} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{1n} & \cdots & x_{mn} \end{pmatrix}, \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_m \end{pmatrix}, \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix} \quad (4.3)$$

## 4.2. Оценки параметров регрессии

Будем искать оценки коэффициентов  $\beta$  в задаче (4.1) **методом наименьших квадратов**. Другими словами, оценки должны удовлетворять условию

$$ESS = \sum (y_i - \hat{y}_i)^2 \rightarrow \min; \quad \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i} + \dots$$

или, в матричной форме,

$$(Y - X\hat{\beta})^T \cdot (Y - X\hat{\beta}) \rightarrow \min$$

Если продифференцировать это равенство по  $\beta$  и приравнять производные нулю, получим

$$\boxed{X^T X \hat{\beta} = X^T Y}$$

Отсюда

$$\boxed{\hat{\beta} = (X^T X)^{-1} X^T Y, \quad (4.4)}$$

**Пример.** В таблице представлены объём импорта (Y) и ВВП ( $X_1$ ) и индекс потребительских цен ( $X_2$ ) в США

Годы	1964	1965	1966	1967	1968
$X_1$	636	689	753	796	868
$X_2$	93	95	97	100	104
Y	28	32	38	41	53

Требуется построить оценки параметров уравнения линейной регрессии Y на  $X_1, X_2$ :

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon.$$

**Решение.**

$$Y = \begin{pmatrix} 28 \\ 32 \\ 38 \\ 41 \\ 53 \end{pmatrix}, X = \begin{pmatrix} 1 & 636 & 93 \\ 1 & 689 & 95 \\ 1 & 753 & 97 \\ 1 & 796 & 100 \\ 1 & 868 & 104 \end{pmatrix}, \hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix},$$

$$X^T X = \begin{pmatrix} 1 & 1 & 1 & 1 & 1 \\ 636 & 689 & 753 & 796 & 868 \\ 93 & 95 & 97 & 100 & 104 \end{pmatrix} \cdot \begin{pmatrix} 1 & 636 & 93 \\ 1 & 689 & 95 \\ 1 & 753 & 97 \\ 1 & 796 & 100 \\ 1 & 868 & 104 \end{pmatrix} = \begin{pmatrix} 5 & 3742 & 489 \\ 3742 & 2833266 & 367516 \\ 489 & 367516 & 47899 \end{pmatrix};$$

$$(X^T X)^{-1} = \begin{pmatrix} 2452,8 & 1,822 & -39,02 \\ 1,822 & 0,0014 & -0,03 \\ -39,02 & -0,03 & 0,625 \end{pmatrix} \quad \hat{\beta} = (X^T X)^{-1} X^T Y = \begin{pmatrix} -150,99 \\ 0,02 \\ 1,78 \end{pmatrix}$$

$$y = -150,99 + 0,02x_1 + 1,78x_2$$

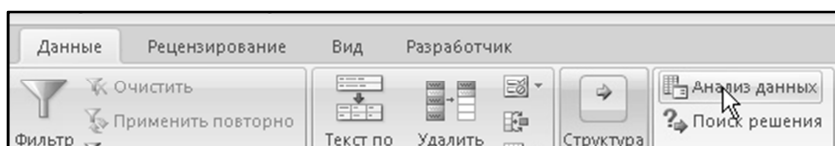
### 4.3. Нахождение параметров регрессии при помощи MS Excel

Для нахождения параметров регрессии при помощи MS Excel необходимо, чтобы была установлена надстройка «**Анализ данных**».

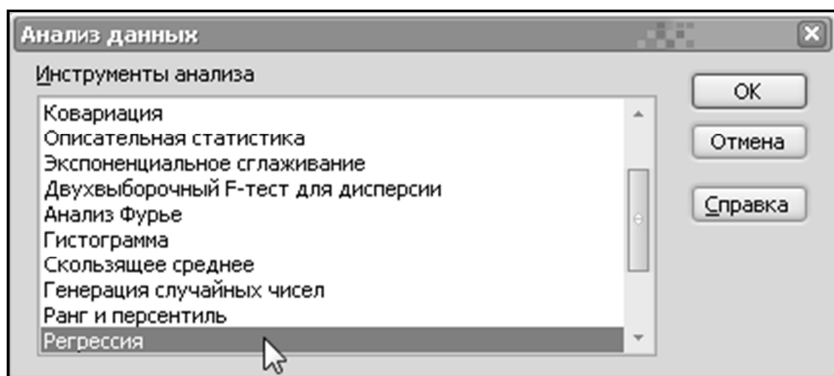
**Пример.** Возьмём тот же пример и внесём данные в таблицу MS Excel *по столбцам*:

	A	B	C	D
1	Годы	X1	X2	Y
2	1964	636	93	28
3	1965	689	95	32
4	1966	753	97	38
5	1967	796	100	41
6	1968	868	104	53

Выберем пункт меню *Данные – Анализ данных*



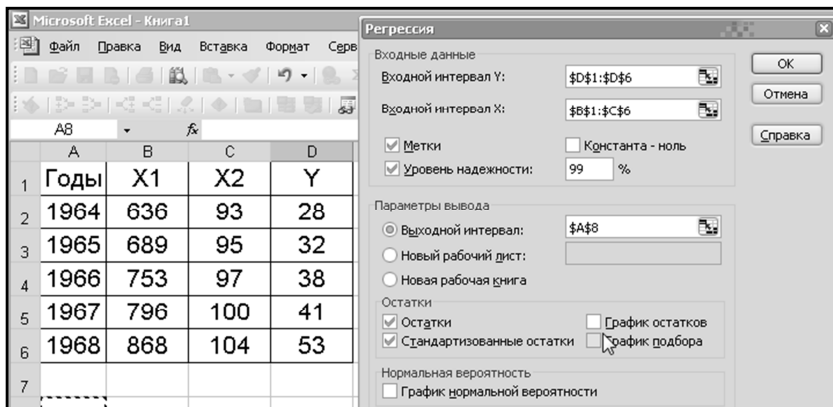
Укажем инструмент анализа *Регрессия*:



Появится диалог **“Регрессия”**.

Указываем области, где находятся значения X, Y (см. рис.), указываем, куда выводить результаты и нажимаем кнопку **ОК**.





Программа выдаст множество данных по регрессии, в том числе и оценки параметров регрессии:

<b>Дисперсионный анализ</b>			
	<i>df</i>	<i>SS</i>	<i>h</i>
Регрессия	2,0000	361,6094	
Остаток	2,0000	7,5906	
Итого	4,0000	369,2000	
	<i>Коэффициенты</i>	<i>Стандартная ошибка</i>	<i>t-стат.</i>
Y-пересечение	-150,9880	96,4838	
X1	0,0200	0,0736	
X2	1,7832	1,5403	

В частности, уравнение регрессии:

$$y = -150,99 + 0,02x_1 + 1,78x_2$$

(96,48)
(0,07)
(1,54)

(под коэффициентами в скобках написаны их стандартные отклонения – см. далее).

#### 4.4. Основные гипотезы

Гипотезы, лежащие в основе *классической модели множественной регрессии*, следующие:

- (A) “Истинная” зависимость  $Y$  от регрессоров  $X$  имеет вид  $y = X\beta + \varepsilon$ .
- (B) Величины  $X$  – детерминированные (не случайные).
- (C) Столбцы матрицы  $X$  линейно независимы (другими словами, *ранг матрицы  $X$  равен  $m+1$* ).
- (D)  $M\varepsilon_i = 0$ ,  $D\varepsilon_i = \sigma^2$  не зависит от  $i$
- (E)  $\text{cov}(\varepsilon_i, \varepsilon_j) = 0$  при  $i \neq j$   
(*некоррелированность ошибок для разных наблюдений*).

**Замечание.** Вместо (D) часто добавляют условие

- (F) Ошибки  $\varepsilon_i$  имеют *нормальное* распределение  $N(0, \sigma^2)$ .

В этом случае модель называется *нормальной*.

**Замечание.** Условия (D) – (F) удобно записать в *матричной форме*

- (G)  $M\varepsilon = 0$ ,  $M(\varepsilon\varepsilon^T) = \sigma^2 E$  ( $E$  – единичная матрица)
- (F)  $\varepsilon \sim N(0, \sigma^2 E)$

#### 4.5. Теорема Гаусса-Маркова

**Теорема.** В предположениях (A) – (E) оценки, полученные методом наименьших квадратов по формуле (4.4), являются *несмещёнными* и обладают *наименьшей дисперсией* среди всех линейных несмещённых оценок параметров  $\beta$ .

$$\hat{\beta} = (X^T X)^{-1} X^T Y, \quad (4.4)$$



#### 4.6. Коэффициент детерминации

Дисперсия зависимой переменной, помноженная на  $n$ , (или *полная сумма квадратов*) обозначается TSS:

$$TSS = \sum (y_i - \bar{y})^2 = nD_y^*$$

Сумма

$$RSS = \sum (\hat{y}_i - \bar{y})^2 = nD_{\hat{y}}^*$$

содержит “объяснённую часть дисперсии”.  
И, наконец, сумма квадратов “остатков

$$ESS = \sum (y_i - \hat{y}_i)^2 = \sum e_i^2 = nD_e^*$$

представляет собой “необъясненную регрессией часть дисперсии” зависимой переменной.

Нетрудно проверить, что

$$RSS + ESS = TSS$$

Величина



$$R^2 = \frac{RSS}{TSS} = 1 - \frac{ESS}{TSS}$$

называется **коэффициентом детерминации** и представляет собой долю дисперсии зависимой переменной  $Y$ , объясненную при помощи уравнения регрессии.

$$0 \leq R^2 \leq 1$$

Чем ближе коэффициент детерминации к 1, тем лучше приближаются исходные данные уравнением линейной регрессии.

**Замечание.** В модели *множественной линейной регрессии*, кроме  $\beta$ , есть ещё один параметр –  $\sigma^2$ . Его оценку можно найти по формуле

$$\hat{\sigma}^2 = S^2 = \frac{ESS}{n - m - 1}$$

**Пример.** В рассмотренном примере

$$Y = \begin{pmatrix} 28 \\ 32 \\ 38 \\ 41 \\ 53 \end{pmatrix}, \hat{Y} = X\hat{\beta} = \begin{pmatrix} 1 & 636 & 93 \\ 1 & 689 & 95 \\ 1 & 753 & 97 \\ 1 & 796 & 100 \\ 1 & 868 & 104 \end{pmatrix} \cdot \begin{pmatrix} 150,99 \\ 0,03 \\ 1,78 \end{pmatrix} = \begin{pmatrix} 27,59 \\ 32,22 \\ 37,07 \\ 43,28 \\ 51,85 \end{pmatrix},$$

$$\bar{Y} = \frac{28+32+\dots}{5} = 38,4; \quad RSS = (27,59 - 38,4)^2 + \dots = 361,609;$$

$$TSS = (28 - 38,4)^2 + \dots = 369,24; \quad ESS = (28 - 27,59)^2 + \dots = 7,5906;$$

$$S^2 = \frac{ESS}{5 - 2 - 1} = 3,7953; \quad R^2 = \frac{361,609}{369,24} = 0,9794$$

**Замечание.** На MS Excel эти параметры выглядят следующим образом:

ВЫВОД ОСТАТКА		
Наблюдение	Предсказанное Y	Остатки
1	27,58910468	0,4108953
2	32,21717357	-0,2171735
3	37,06559179	0,9344082
4	43,27653597	-2,2765359
5	51,85159399	1,148406

$$\hat{Y} = \begin{pmatrix} 27,59 \\ 32,22 \\ 37,07 \\ 43,28 \\ 51,85 \end{pmatrix},$$

Дисперсионный анализ			
	df	SS	MS
Регрессия	2	361,6094	180,8047
Остаток	2	7,5906	3,7953
Итого	4	369,2000	

$$RSS = 361,609;$$

$$ESS = 7,5906;$$

$$TSS = 369,24;$$

Регрессионная статистика	
Множественный R	0,9897
R-квадрат	0,9794
Нормированный R-кв	0,9589
Стандартная ошибка	1,9481
Наблюдения	5

$$R^2 = 0,9794$$

$$S^2 = 3,7953;$$

$$S = \sqrt{S^2} = 1,9481$$

**Замечание.** Важное свойство коэффициента детерминации  $R^2$  состоит в том, что при добавлении в модель новых объясняющих переменных  $x$  он не может уменьшиться. Поэтому при сравнении двух моделей с различным числом переменных не всегда ясно, за счёт чего возрос показатель детерминации: за счёт реального

влияния дополнительного фактора на результат или просто из-за возрастания числа переменных.

Для того чтобы можно было *сравнивать различные модели*, вводят так называемый *скорректированный* (или *нормированный*) *коэффициент детерминации*:

$$R_{корр}^2 = 1 - \frac{ESS/(n-m-1)}{TSS/(n-1)}$$

Это число также не превосходит 1, но в некоторых случаях может оказаться и отрицательным. В рассматриваемом примере  $R_{корр}^2 = 0,9589$

#### 4.7. Квантили

Формулы, описывающие статистические свойства оценок, содержат табличные значения, называемые *квантилями*.

Квантиль **нормального закона**  $N(0;1)$ :  $u_p$

p	0,9	0,95	0,975	0,99	0,995	0,999	0,9995
$u_p$	1,282	1,645	1,960	2,325	2,576	3,090	3,291

*Свойство:*

$$u_{1-\alpha} = -u_{\alpha}$$

Квантиль **закона распределения хи-квадрат:**  $\chi_p^2(k)$

$\begin{matrix} p \\ k \end{matrix}$	0,01	0,025	0,05	0,1	0,9	0,95	0,975	0,99
1	0,0002	0,001	0,004	0,0158	2,71	3,84	5,02	6,63
2	0,02	0,05	0,103	0,211	4,61	5,99	7,38	9,21

3	0,115	0,216	0,352	0,584	6,25	7,81	9,35	11,3
4	0,297	0,484	0,711	1,06	7,78	9,49	11,1	13,3
10	2,56	3,25	3,94	4,87	16	18,3	20,5	23,2
20	7,63	9,59	10,9	12,4	28,4	31,4	34,2	37,6
30	14,3	16,8	18,5	20,6	40,3	43,8	47,0	50,9
40	22,2	24,4	26,5	29,1	51,8	55,8	59,3	63,7
50	29,7	32,4	34,8	37,7	63,2	67,5	71,4	76,2
75	49,5	53	56,1	59,8	91,1	96,2	100,8	106,4

$$\chi_p^2(k) \approx k \cdot \left( 1 - \frac{2}{9k} + u_p \sqrt{\frac{2}{9k}} \right)^3 \text{ при больших } k.$$

Квантиль **закона распределения Стьюдента**:  $t_p(k)$

k \ p	0,9	0,95	0,975	0,99	0,995
1	3,078	6,314	12,706	31,821	63,657
2	1,886	2,92	4,303	6,965	9,925
3	1,638	2,353	3,182	4,541	5,841
4	1,533	2,132	2,776	3,747	4,604
5	1,476	2,015	2,571	3,365	4,032
10	1,372	1,812	2,228	2,764	3,169
20	1,325	1,725	2,086	2,528	2,845
30	1,312	1,697	2,042	2,457	2,750
40	1,303	1,684	2,021	2,423	2,704
120	1,289	1,658	2,98	2,358	2,617

**Свойство:**

$$t_{1-\alpha}(k) = -t_{\alpha}(k)$$

$$t_p(k) \approx u_p \cdot \left( 1 - \frac{u_p^2 + 1}{4k} \right)^{-1} \text{ при больших } k.$$

Квантиль закона распределения Фишера:  $F_p(k_1, k_2)$

Квантили  $F_{0,9}(n_1, n_2)$  закона распределения Фишера

$n_1 \backslash n_2$	1	2	3	4	5	10	15	20	30	120
1	40	8,53	5,54	4,54	4,06	3,29	3,07	2,97	2,88	2,75
2	49,5	9	5,46	4,32	3,78	2,92	2,7	2,59	2,49	2,35
3	53,6	9,16	5,39	4,19	3,62	2,73	2,49	2,38	2,28	2,13
4	55,8	9,24	5,34	4,11	3,52	2,61	2,36	2,25	2,14	1,99
30	62,2	9,46	5,17	3,82	3,17	2,16	1,87	1,74	1,61	1,41

Квантили  $F_{0,95}(n_1, n_2)$  закона распределения Фишера

$n_1 \backslash n_2$	1	2	3	5	10	15	20	30	120
1	161	18,5	10,13	6,61	4,96	4,54	4,35	4,24	3,92
2	199	19	9,55	5,79	4,1	3,68	3,49	3,39	3,07
3	216	19,16	9,28	5,41	3,71	3,29	3,1	2,99	2,68
4	225	19,25	9,12	5,19	3,48	3,05	2,87	2,76	2,45
30	250	19,46	8,62	4,5	2,7	2,25	2,04	1,84	1,55

**Свойство:**

$$F_{1-\alpha}(k_1, k_2) = \frac{1}{F_{\alpha}(k_2, k_1)}$$

В последних версиях MSExcel квантили можно найти так:

$$u_p = \text{НОРМ.СТ.ОБР}(p)$$

$$\chi_p^2(k) = \text{ХИ2.ОБР}(p, k)$$

$$t_p(k) = \text{СТЮДЕНТ.ОБР}(p, k)$$

$$F_p(k_1, k_2) = \text{F.ОБР}(p, k_1, k_2)$$



В более старых версиях MS Excel некоторые квантили вычисляются “задом наперёд”

$$\chi_{\alpha}^2(k) = \text{ХИ2ОБР}(1-\alpha, k)$$

$$\chi_{1-\alpha}^2(k) = \text{ХИ2ОБР}(\alpha, k)$$

$$t_{1-\frac{\alpha}{2}}(k) = \text{СТЬЮДРАСПОБР}(\alpha, k)$$

$$F_{1-\alpha}(k_1, k_2) = \text{ФРАСПОБР}(\alpha, k_1, k_2)$$

#### 4.8. Статистические свойства оценок параметров регрессии

**Теорема.** В предположениях (А) – (Е) и (F) теоремы Гаусса-Маркова оценки, полученные методом наименьших квадратов по формуле (2.4), обладают следующими свойствами:

1.  $\frac{(n-m-1) \cdot S^2}{\sigma^2} \sim \chi^2(n-m-1)$ ;
2.  $S^2$  и  $\hat{\beta}$  независимы.

#### 4.9. Доверительные интервалы

Будем считать, что условие нормальности (F) выполнено. Тогда можно построить **доверительные интервалы** для параметров регрессии, т.е. интервалы, **содержащие точные значения этих параметров с заданной вероятностью  $1-\alpha$ .**

Вот соответствующие формулы:

$$\hat{\beta}_i - t_{1-\alpha/2}(n-m-1) \cdot S_{\beta_i} < \beta_i < \hat{\beta}_i + t_{1-\alpha/2}(n-m-1) \cdot S_{\beta_i}$$

$$\frac{ESS}{\chi^2_{1-\alpha/2}(n-m-1)} < \sigma^2 < \frac{ESS}{\chi^2_{\alpha/2}(n-m-1)}$$

Здесь  $\alpha$  – заданный “уровень значимости”, т.е. *вероятность того, что построенный интервал не содержит истинного значения параметра*.

Далее,  $S_{\beta}$  – выборочные среднеквадратичные отклонения оценок параметров (или *стандартные ошибки параметров*), они находятся по формулам:

$$S^2_{\beta_i} = S^2 \cdot (X^T \cdot X)^{-1}_{ii}$$

Как правило, стандартные ошибки параметров вычисляются в эконометрических программах вместе с оценками самих параметров

**Пример.** В рассмотренном примере

$$(X^T X)^{-1} = \begin{pmatrix} 24528 & 1,822 & -39,02 \\ 1,822 & 0,0014 & -0,03 \\ -39,02 & -0,03 & 0,625 \end{pmatrix}, \quad \begin{matrix} y = -150,99 + 0,02x_1 + 1,78x_2 \\ S^2 = 3,7953 \end{matrix}$$

$$S^2_{\beta_0} = S^2 \cdot 2452,8 = 9309,16; \quad S_{\beta_0} = 96,48$$

$$S^2_{\beta_1} = S^2 \cdot 0,0014 = 0,0053; \quad S_{\beta_1} = 0,07$$

$$S^2_{\beta_2} = S^2 \cdot 0,625 = 2,372; \quad S_{\beta_2} = 1,54$$

Таким образом,

$$y = -150,99 + 0,02x_1 + 1,78x_2$$

$$\begin{matrix} (96,48) & (0,07) & (1,54) \end{matrix}$$

Далее, на уровне значимости  $\alpha=0,1$ :

$$t_{1-\alpha/2}(n-m-1) = t_{0,95}(2) = 2,92;$$

$$-150,99 - 2,92 \cdot 96,48 = -432,72 < \beta_0 < 130,74 = -150,99 + 2,92 \cdot 96,48$$

Аналогично,

$$-0,19 < \beta_1 < 0,23;$$

$$-2,71 < \beta_2 < 6,28;$$

Наконец,

$$\chi^2_{1-\alpha/2}(n-m-1) = \chi^2_{0,95}(2) = 5,99; \chi^2_{\alpha/2}(n-m-1) = \chi^2_{0,05}(2) = 0,1;$$

$$\frac{ESS}{5,99} = \frac{7,5906}{5,99} = 1,27 < \sigma^2 < 75,9 = \frac{7,5906}{0,1}$$

В MS Excel стандартные ошибки коэффициентов:

Дисперсионный анализ			
	<i>df</i>	<i>SS</i>	<i>h</i>
Регрессия	2,0000	361,6094	
Остаток	2,0000	7,5906	
Итого	4,0000	369,2000	
	<i>Коэффициенты</i>	<i>Стандартная ошибка</i>	<i>t-стат.</i>
Y-пересечение	-150,9880	96,4838	
X1	0,0200	0,0736	
X2	1,7832	1,5403	
			+

#### 4.10. Доверительный интервал для прогнозного значения

Рассмотрим модель

$$Y = X\beta + \varepsilon, \quad (4.2)$$

где

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, X = \begin{pmatrix} 1 & x_{11} & \cdots & x_{m1} \\ 1 & x_{12} & \cdots & x_{m2} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1n} & \cdots & x_{mn} \end{pmatrix}, \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_m \end{pmatrix}, \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix} \quad (4.3)$$

Напомним, что  $X$  – матрица, каждая из  $n$  строк которой соответствует какому-то набору значений независимых переменных.

Предположим, что у нас есть ещё один набор значений независимых переменных

$$x^{(n+1)} = (1, x_1^{(n+1)}, x_2^{(n+1)}, \dots, x_m^{(n+1)}).$$

Требуется найти доверительный интервал для соответствующего значения зависимой переменной  $y^{(n+1)}$ .

Тогда

$$\hat{y}^{(n+1)} - t_{1-\alpha/2}(n-m-1) \cdot S_y < y^{(n+1)} < \hat{y}^{(n+1)} + t_{1-\alpha/2}(n-m-1) \cdot S_y$$

где

$$\hat{y}^{(n+1)} = \hat{\beta}_0 + \hat{\beta}_1 x_1^{(n+1)} + \dots + \hat{\beta}_m x_m^{(n+1)} = x^{(n+1)} \cdot \hat{\beta}$$

$$S_y = S \cdot \sqrt{1 + x^{(n+1)} \cdot (X^T X)^{-1} \cdot (x^{(n+1)})^T}$$

**Пример.** В таблице приведены значения двух факторов. Построить уравнение линейной регрессии и найти доверительный интервал для  $Y$  при  $X=7$ . Взять  $\alpha=0,1$ .

X	1	2	3	6	8
Y	1	4	5	9	11

**Решение.**

$$Y = \begin{pmatrix} 1 \\ 4 \\ 5 \\ 9 \\ 11 \end{pmatrix}, X = \begin{pmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 6 \\ 1 & 8 \end{pmatrix}, \hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix},$$

Тогда

$$(X^T X)^{-1} = \begin{pmatrix} 0,671 & -0,118 \\ -0,118 & 0,029 \end{pmatrix}, y = 0,59 + 1,35x, S^2 = 0,5882$$

и

$$\hat{y}^{(n+1)} = \hat{\beta}_0 + \hat{\beta}_1 x_1^{(n+1)} = (1 \quad 7) \cdot \begin{pmatrix} 0,59 \\ 1,35 \end{pmatrix} = 10,04$$

$$S_y = \sqrt{0,5882} \cdot \sqrt{1 + (1 \quad 7) \cdot \begin{pmatrix} 0,671 & -0,118 \\ -0,118 & 0,029 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 7 \end{pmatrix}}$$

$$S_y \approx 0,921; \quad t_{1-\alpha/2}(n-m-1) = t_{0,95}(5-1-1) = 2,35$$

Окончательно,

$$10,04 - 2,35 \cdot 0,921 < y^{(n+1)} < 10,04 + 2,35 \cdot 0,921$$

#### 4.11. Проверка гипотез о параметрах регрессии

В этом параграфе также будем считать, что условие нормальности **(F)** выполнено:

**(F)** Ошибки  $\varepsilon_i$  имеют нормальное распределение  $N(0, \sigma^2)$ .

При исследовании модели приходится проверять гипотезы о равенстве истинных значений параметров регрессии заданным числам. Для этого **достаточно построить доверительный интервал и посмотреть, попадает ли в него это заданное число.**

Рассмотрим, например, гипотезу  $H_0: \beta_i = A$   
при альтернативной гипотезе  $H_1: \beta_i \neq A$ .

Тогда гипотеза  $H_0$  принимается в случае, *если число A лежит в интервале*

$$\hat{\beta}_i - t_{1-\alpha/2}(n-m-1) \cdot S_{\beta_i} < \beta_i < \hat{\beta}_i + t_{1-\alpha/2}(n-m-1) \cdot S_{\beta_i}$$

В противном случае принимается гипотеза  $H_1$

Особенно часто проверяют гипотезу  $H_0: \beta_i = 0$ .

Дело в том, что если принимается альтернативная гипотеза

$$H_1: \beta_i \neq 0$$

то говорят, что коэффициент  $\beta_i$  *значим*.

Таким образом, получаем следующий **критерий значимости параметра  $\beta_i$** :

$$\boxed{\frac{|\hat{\beta}_i|}{S_{\beta_i}} \geq t_{1-\alpha/2}(n-m-1) \quad (4.5)}$$

**Пример.** В рассмотренном примере было

$$-0,19 < \beta_1 < 0,23;$$

$$-2,71 < \beta_2 < 6,28;$$

Доверительные интервалы для коэффициентов содержат 0, так что оба коэффициента *не значимы*.

**Задача.** По выборке объема 8 получено следующее уравнение линейной регрессии:

$$y = 5,27 - 0,89x_1 + 1,69x_2 - 0,41x_3$$

$$(2,97) \quad (0,53) \quad (0,58) \quad (0,32)$$

Какие из его коэффициентов значимы на уровне значимости  $\alpha = 0,1$ ?

**Решение.**

$$n = 8, m = 3, t_{1-\alpha/2}(n-m-1) = t_{0,95}(4) = 2,13.$$

$$\frac{0,89}{0,53} \approx 1,68 < 2,13 \Rightarrow \dots$$

$$\frac{1,69}{0,59} \approx \dots \Rightarrow \dots$$

**Замечание.** Дробь в левой части неравенства 4.5 называют *t-статистикой*. Эти значения в MS Excel автоматически:

	Коэффициентная		<i>t-статистика</i>	<i>P-Значение</i>
Y-пересечение	1665,367	1442,75	1,154294475	0,264348886
Total investm	-68,6798	42,03165	-1,634003237	0,120639348
Inflation, aver	-1,3E-11	4,43E-11	-0,290871293	0,774668562
Unemployme	-34,9101	30,51556	-1,144010272	0,268472918
Population	7,296819	1,189727	6,133186538	1,10426E-05
General gover	20,79136	21,92427	0,948326344	0,356250483
Current accou	-0,51724	3,008832	-0,171907903	0,865539759

**Замечание.** Значимость коэффициентов удобнее всего оценивать, используя *P-значение*. Это минимальное значение уровня значимости, на котором коэффициент ещё значим. Таким образом, для значимости коэффициента необходимо и достаточно, чтобы *P-значение* было меньше  $\alpha$ .

	Коэффициентная		<i>t-статистика</i>	<i>P-Значение</i>
Y-пересечение	1665,367	1442,757	1,154294475	0,264348886
Total investm	-68,6798	42,03165	-1,634003237	0,120639348
Inflation, aver	-1,3E-11	4,43E-11	-0,290871293	0,774668562
Unemployme	-34,9101	30,51556	-1,144010272	0,268472918
Population	7,296819	1,189727	6,133186538	1,10426E-05
General gover	20,79136	21,92427	0,948326344	0,356250483
Current accou	-0,51724	3,008832	-0,171907903	0,865539759

На рисунке значимым на уровне 0,1 является только коэффициент при Population.



#### 4.12. Проверка гипотезы о значимости уравнения регрессии

При построении уравнения возникает вопрос: в какой мере можно ему доверять, насколько оно *значимо*?

*Уравнение считается значимым, если отвергается гипотеза о том, что все коэффициенты при переменных равны нулю.*

Для проверки значимости уравнения используется **F-статистика Фишера**:

$$F = \frac{RSS/m}{ESS/(n-m-1)}$$

Разделив числитель и знаменатель на TSS, получим ещё одну формулу для  $F$ :

$$F = \frac{R^2/m}{(1-R^2)/(n-m-1)}$$

Если вычисленное значение F-статистики оказывается больше  $F_{1-\alpha}(m, n-m-1)$ , то оцененное уравнение регрессии является значимым при выбранном  $\alpha$ .

**Задача.** По выборке получено уравнение линейной регрессии:  $y = 4,8 - 2,15x_1 + 2,18x_2$

Значимо ли это уравнение на уровне 0,1?

**Решение.** 1) Вычислим расчётные значения  $Y$ :

$y$	$x_1$	$x_2$
6	2	3
10	2	4
10	5	7
11	7	10
3	7	6

$\hat{y}$
7,04
9,22
9,31
11,55
2,83

2) Вычислим RSS и ESS

$$3) F = \frac{RSS/m}{ESS/(n-m-1)} = \frac{43,5/2}{2,5/(5-2-1)} = 17,4$$

Квантиль  $F_{0,9}(2;5-2-1)=F_{0,9}(2;2)$  по таблице равен 9.

Так как  $17,4 > 9$ , то уравнение значимо.

**Замечание.** В MS Excel *F-значение* вычисляется автоматически:

Дисперсионный анализ					
	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>Значимость F</i>
Регрессия	6	63713945	10618990,9	14,2521	8,42827E-06
Остаток	17	12666421	745083,611		
Итого	23	76380367			

**Замечание.** Значимость уравнения в MS Excel удобнее всего оценивать, используя величину “*значимость F*”. Это – минимальное значение уровня  $\alpha$ , на котором уравнение ещё значимо. Таким образом, для значимости уравнения необходимо и достаточно, чтобы число

“Значимость  $F$ ” было меньше  $\alpha$ .

Дисперсионный анализ					
	$df$	$SS$	$MS$	$F$	Значимость $F$
Регрессия	6	63713945	10618990,9	14,252	8,42827E-06
Остаток	17	12666421	745083,611		
Итого	23	76380367			

## 5. Спецификация модели

### 5.1. Основные этапы

Спецификация модели включает следующие этапы:

1. Определение списка объясняющих и зависимых переменных.
2. Выбор функциональной формы модели.

Рассмотрим подробнее первый пункт

### 5.2. Отбор объясняющих переменных

При определении состава факторов, включаемых в уравнение регрессии, руководствуются **теоретическими представлениями** о взаимосвязях этих факторов.

Однако часто встречается ситуация, когда имеется большое число факторов, **но нет априорной модели изучаемого явления**, так что неясно, какие именно переменные включать в модель. В этом случае используются различные **эвристические процедуры**.

Приведём простой пример такой процедуры

### Процедура пошагового отбора независимых переменных.

1. Из исходного набора переменных отбирается (включается в модель) переменная, имеющая **наибольший по модулю коэффициент корреляции с зависимой переменной  $Y$** .
2. На каждом последующем шаге в модель **добавляется та** из переменных, добавление которой **максимально увеличивает скорректированный коэффициент детерминации**.
3. Если добавление новых переменных не увеличивает этот коэффициент, процедура считается завершённой

**Пример.** Произвести пошаговый отбор переменных регрессии для выборки.

Y	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>
12,4	1,1	1,0	2,2
23,5	2,0	1,4	3,0
43,5	3,5	5,0	3,0
38,6	4,0	3,4	4,6
25,1	2,7	2,1	3,0
38,0	4,1	2,0	6,0
59,8	5,0	8,1	8,0
54,3	7,0	4,0	9,0
85,7	8,0	11,0	3,0
52,2	6,4	5,0	8,0

**Решение.** 1) Найдём матрицу корреляций.  
 Наибольший по модулю коэффициент корреляции с Y имеет переменная X<sub>2</sub>.

	Y	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>
Y	1,00			
X <sub>1</sub>	0,93	1,00		
X <sub>2</sub>	0,94	0,77	1,00	
X <sub>3</sub>	0,40	0,59	0,21	1,00

2) Строим регрессию Y на X<sub>2</sub>.

$$y = 16,48 + 6,24x_2, \quad R^2_{\text{корр}} = 0,868$$

3.1) Добавим к X<sub>2</sub> переменную X<sub>1</sub>.

$$y = 6,65 + 7,79x_1 + 3,65x_2, \quad R^2_{\text{корр}} = 0,986$$

3.2) Добавим к X<sub>2</sub> переменную X<sub>3</sub>.

$$y = 8,86 + 5,944x_2 + 1,786x_3, \quad R^2_{\text{корр}} = 0,9$$

$R^2_{\text{корр}}$  сильнее возрос в варианте 3.1. Оставляем этот вариант.

4.1) Добавим к оставленному варианту ( $X_1, X_2$ ) переменную  $X_3$ .

$$y = 6,9 + 4,96x_1 + 3,58x_2 - 0,145x_3, \quad R^2_{\text{корр}} = 0,980$$

Больше добавлять к ( $X_1, X_2$ ) нечего – других вариантов нет.

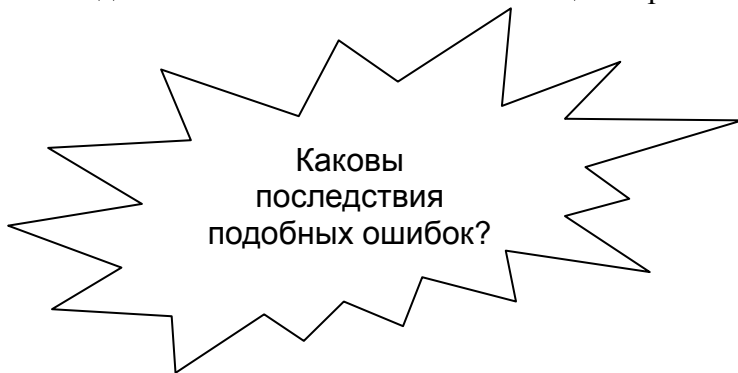
**Ответ.** Результат пошагового отбора переменных:

$$y = 6,65 + 7,79x_1 + 3,65x_2, \quad R^2_{\text{корр}} = 0,986$$

### **5.3. Влияние ошибок спецификации переменных**

Неправильная спецификация переменных может произойти из-за

- отбрасывания переменной, которая должна быть в составе объясняющих переменных
- или
- включения “лишней” переменной, которой не должно быть в составе объясняющих переменных.



#### **5.3.1. Пропуск существенных переменных**

Пусть истинная зависимость имеет вид

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m + \gamma_1 z_1 + \dots + \gamma_l z_l + \varepsilon, \quad (A)$$

или, в матричной форме,

$$Y = X\beta + Z\gamma + \varepsilon, \quad (A)$$

а мы, не зная об этом, ищем регрессию в виде:

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m + \varepsilon, \quad (B)$$

то есть,

$$Y = X\beta + \varepsilon \quad (B)$$

В случае такого пропуска существенных переменных:

- ❖ МНК-оценка параметров  $\beta$  в общем случае оказывается смещённой;
- ❖ Стандартные ошибки коэффициентов  $\beta$  имеют неотрицательное смещение;
- ❖ Оценка дисперсии  $\sigma^2$  имеет неотрицательное смещение

Таким образом, **последствия** достаточно **серьёзные**.

### 5.3.2. Включение несущественных переменных

Пусть истинная зависимость имеет вид

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m + \varepsilon, \quad (B)$$

или, в матричной форме,

$$Y = X\beta + \varepsilon \quad (B)$$

а мы, не зная об этом, ищем регрессию в виде:

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m + \gamma_1 z_1 + \dots + \gamma_l z_l + \varepsilon, \quad (A)$$

то есть,

$$Y = X\beta + Z\gamma + \varepsilon, \quad (A)$$

Тогда

- ❖ МНК-оценки параметров  $\beta$  несмещённые;
- ❖ Стандартные ошибки коэффициентов  $\beta$  имеют неотрицательное смещение;
- ❖ Оценка дисперсии  $\sigma^2$  имеет неотрицательное смещение

Последствия такой ошибки спецификации являются не столь серьёзными, как в предыдущем случае.

Однако увеличение стандартных ошибок может привести к неверным выводам при проверке гипотез и ухудшению интервальных оценок.

В частности, может возникнуть незначимость коэффициентов.

#### **5.4. Фиктивные переменные**

Как правило, независимые переменные принимают непрерывный ряд значений (цена товара, уровень инфляции и т.д.).

Однако нередко приходится учитывать *качественные* признаки.

Например, можно изучать влияние пола работника на уровень его зарплаты (переменная “пол” принимает всего



два значения), или спросом на прохладительные напитки в зависимости от времени года (сколько различных значений принимает эта переменная?).

В подобных ситуациях “качественному” значению условно присваивается числовая метка, как правило, это 0 или 1.

Если же качественный признак принимает  $k$  значений ( $k > 2$ ), то вводят  $k-1$  переменную, каждая из которых принимает значения 0 и 1.

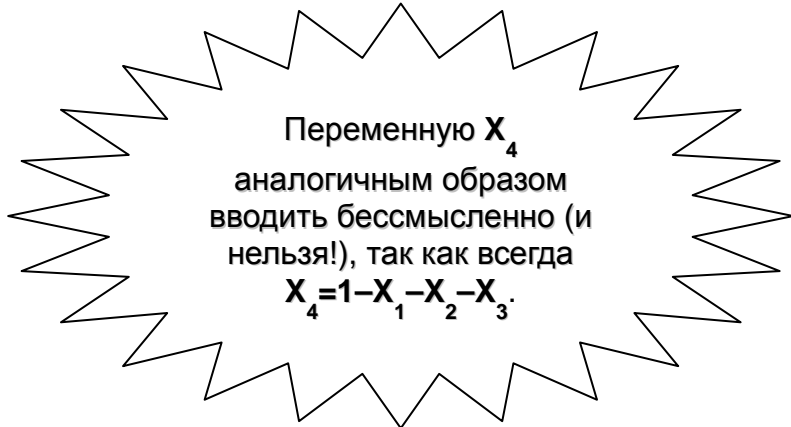
**Пример.** Качественный признак – время года, значения: зима, весна, лето, осень.

**Фиктивные переменные:**

$X_1=1$ , если зима, иначе  $X_1=0$ ,

$X_2=1$ , если весна, иначе  $X_2=0$ ,

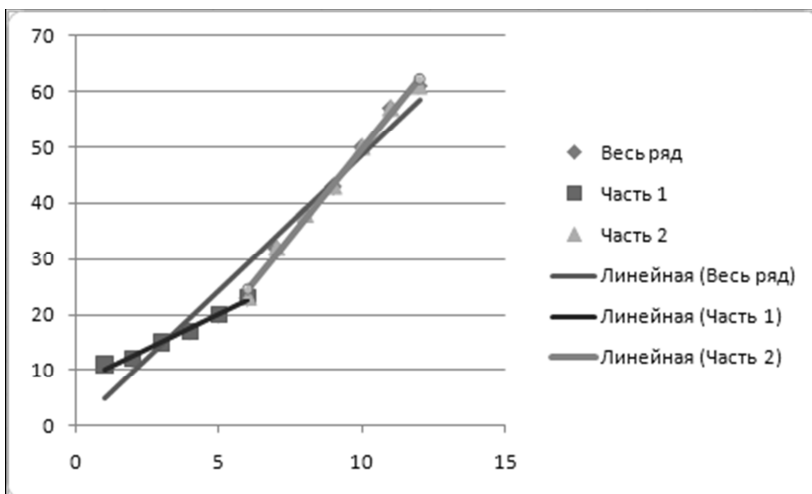
$X_3=1$ , если лето, иначе  $X_3=0$ .



### 5.5. Тест Чоу (Chow).

В некоторых случаях можно наблюдать изменения тренда, вызванные различными факторами.

В таком случае уравнения тренда на разных участках будут различаться



В ситуации на рисунке возникает альтернатива – описывать тренд одной синей линией или совокупностью из двух линий.

Другими словами, выборка состоит из нескольких частей и требуется выяснить, следует ли их объединить в одну, или рассматривать каждую подвыборку отдельно.

Обозначим регрессию по первой выборке (объёма  $n_1$ ) через  $A$ , по второй (объёма  $n_2$ ) – через  $B$ , а объединённую (объёма  $n=n_1+n_2$ ) – через  $R$ . Ясно, что

$$ESS_A + ESS_B \leq ESS_R.$$

Если левая часть неравенства близка к правой, то это говорит о том, что регрессия по объединённой выборке почти так же хороша, как “составная” их двух регрессий по двум выборкам. В этом случае объединение выборок в одну допустимо. Проверка этого составляет суть *критерия Чоу*.

1) Составляется статистика

$$F = \frac{(ESS_R - ESS_A - ESS_B)/(m+1)}{(ESS_A + ESS_B)/(n-2m-2)}$$

2) Если вычисленное значение оказывается меньше  $F_{1-\alpha}(m+1, n-2m-2)$ , то объединение выборок возможно при данном  $\alpha$ .

**Замечание.** Если одна из выборок слишком мала, то применяется другая формула. Пусть, например,  $n_2 \leq m+1$ , тогда

$$F = \frac{(ESS_R - ESS_A)/(n-n_1)}{ESS_A/(n_1-m-1)}$$

Если вычисленное значение оказывается меньше  $F_{1-\alpha}(n-n_1, n_1-m-1)$ , то объединение выборок возможно при данном  $\alpha$ .

**Замечание.** Другой способ проверки состоит в использовании *фиктивных переменных*.

Добавим в модель фиктивную переменную  $Z$ , равную 0 для первой выборки и 1 для второй, а также переменные  $ZX_1, ZX_2, \dots$ . Если **все** эти переменные окажутся *незначимыми*, то объединение выборок возможно.

**Пример.** Проверить возможность объединения выборок при  $\alpha=0,1$ .

$y$	$x_1$	$x_2$
3	2	3
5	3	4
5	3	6
6	4	1
7	5	7

$y$	$x_1$	$x_2$
9	4	7
9	6	8
10	6	9
12	7	9

### Решение. 1 способ (критерий Чоу).

Для первой выборки:

$$y = 0,86 + 1,26x_1 + 0,01x_2, ESS = 0,42$$

Для второй выборки:

$$y = 3 + 0,5x_1 + 0,5x_2, ESS = 2,5$$

Для объединённой выборки:

$$y = 0,11 + 1,32x_1 + 0,23x_2, ESS = 7,38$$

$$F = \frac{(7,38 - 0,42 - 2,5)/(2+1)}{(0,42 + 2,5)/(9 - 2 \cdot 2 - 2)} = 1,527$$

$$F_{0,9}(2+1; 9 - 2 \cdot 2 - 2) = F_{0,9}(3; 3) = 5,4$$

Так как  $1,527 < 5,4$ , то объединение возможно.

**2 способ (фиктивная переменная).** Добавляем новые переменные  $z, zx_1, zx_2$ .

Проверяем их значимость

$$\begin{aligned} y = & 0,86 + 1,26x_1 + 0,01x_2 \\ & (1,59) \quad (0,46) \quad (0,22) \\ & + 2,14z + 0,76zx_1 + 0,49zx_2 \\ & (6,63) \quad (1,13) \quad (0,38) \end{aligned}$$

$y$	$x_1$	$x_2$	$z$	$zx_1$	$zx_2$
3	2	3	0	0	0
5	3	4	0	0	0
5	3	6	0	0	0
6	4	1	0	0	0
7	5	7	0	0	0
9	4	7	1	4	7
9	6	8	1	6	8
10	6	9	1	6	9
12	7	9	1	7	9

Все добавленные переменные незначимы, следовательно, объединение выборок возможно.

## 5.6. Проверка гипотезы о линейной связи коэффициентов

Пусть рассматривается линейная модель:

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m + \varepsilon, \quad (U)$$

Мы хотим проверить гипотезу о том, не удовлетворяют ли коэффициенты модели  $Q$  линейным ограничениям:

$$H_0 : H\beta = r, \quad H : q \times (m+1)$$

Например, три ограничения ( $q=3$ ):

$$H_0 : \begin{cases} \beta_1 = 2\beta_2 \\ \beta_3 = 0 \\ \beta_4 = 5 \end{cases}$$

Чтобы проверить такую гипотезу, нужно

- 1) Вычислить  $ESS = ESS_U$  для исходной регрессии.
- 2) Вычислить  $ESS = ESS_R$  для регрессии, в которой выполнены условия гипотезы  $H_0$ .
- 3) Вычислить значение F-статистики

$$F = \frac{(ESS_R - ESS_U)/q}{ESS_U/(n - m - 1)}$$

Гипотеза  $H_0$  принимается, если **вычисленное значение F-статистики оказывается меньше  $F_{1-\alpha}(q; n - m - 1)$**

### 5.7. RESET-тест Рамсея (Ramsey)

RESET-тест Рамсея (сокр. от англ. *Regression Equation Specification Error Test*) отвечает на вопрос, надо ли включать

в регрессию дополнительные (в частности, нелинейные) члены.

Пусть рассматривается линейная модель:

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m + \varepsilon, \quad (*)$$

Мы хотим выяснить, не следует ли вместо (\*) рассмотреть более сложную функциональную модель, включающую степени  $x$ .

### Алгоритм теста Рамсея

- 1) Найти оценки параметров регрессии (\*).
- 2) Вычислить расчётные значения  $y^{\wedge}$  зависимой переменной  $y$ .
- 3) Оценить регрессию
 
$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m + \alpha_2 \hat{y}^2 + \dots + \alpha_k \hat{y}^k + \varepsilon, \quad (**)$$
- 4) Проверить гипотезу  $H_0$ : *все  $\alpha_j = 0$  (то есть что дополнительные слагаемые в (\*\*)) не нужны*.  
Для этого вычислить

$$F = \frac{(ESS_* - ESS_{**}) / (k - 1)}{ESS_{**} / (n - m - k)}$$

Гипотеза  $H_0$  принимается, если **вычисленное значение F-статистики оказывается меньше  $F_{1-\alpha}(k-1; n-m-1)$**

**Пример.** Заработная плата в Нидерландах.

Имеется 150 наблюдений, 75 мужчин и 75 женщин, работавших на полную ставку.

W - заработная плата (гульденов в час),  
 SEX - пол (1– мужчины, 2 – женщины),  
 EDU – уровень образования (1– нач. школа,... 5 – университет) ,  
 AGE возраст .

	W	SEX	EDU	AGE
1	10,44	1	1	19
2	13,52	1	1	20
3	19,12	1	1	21
4	20,28	1	1	25
5	14,63	1	1	26
...	...	...	...	...
22	20,32	1	2	29
23	15,96	1	2	32
24	17,07	1	2	37
25	13,04	1	2	39
...	...	...	...	...
148	26,09	2	5	27
149	24,95	2	5	37
150	41,05	2	5	52

Оценим линейную регрессию

$$W = 3,51 - 3,55 \cdot SEX + 3,24 \cdot EDU + 0,44 \cdot AGE,$$

$$ESS = 7631$$

Проверим при помощи **теста Рамсея** правильность функциональной формы. Для этого вычислим столбец  $W^2$  расчётных значений и добавим к исходным регрессорам его вторую степень (ограничимся второй степенью).

Получим уравнение регрессии

$$W = 12,09 + 0,29 \cdot SEX - 0,49 \cdot EDU - 0,09 \cdot AGE + 0,02 \cdot \hat{W}^2, \\ ESS = 7154$$

Вычислим F-статистику

$$F = \frac{(ESS_* - ESS_{**}) / (k - 1)}{ESS_{**} / (n - m - k)} = \frac{(7631 - 7154) / (2 - 1)}{7154 / (150 - 3 - 2)} = 9,67$$

Далее, на уровне значимости  $\alpha = 0,1$

$$F_{1-\alpha}(k - 1; n - m - k) = F_{0,9}(1; 145) = 2,74 < 9,67$$

В уравнение

$$W = \beta_0 + \beta_1 \cdot SEX + \beta_2 \cdot EDU + \beta_3 \cdot AGE + \varepsilon$$

надо попробовать включить нелинейные члены, например,  $AGE^2$  или другие переменные.

## 6. Некоторые дополнительные вопросы

### 6.1. Коэффициент эластичности

Зависимость изменения значения переменной  $y=f(x)$  от небольшого изменения значения переменной  $x$  выражается производной по  $x$ :

$$y' = \frac{dy}{dx} \approx \frac{\Delta y}{\Delta x}$$

Если, например, мы выражаем объём спроса  $Q$  через цену  $P$  товара (и, возможно, какие-то другие параметры), то (частная) производная

$$Q'_P \approx \frac{\Delta Q}{\Delta P}, \quad \text{где} \\ \Delta P - \text{изменение цены,} \\ \Delta Q - \text{вызванное им изменение объема спроса,}$$



покажет, **на сколько** единиц изменится спрос в расчете **на единичное изменение цены** при данных значениях прочих параметров.

Пусть повышение цены за 1 кг картофеля на 1 руб. снижает годовой объем спроса на 20 кг, т. е.  $\Delta P = 0,1$  руб./кг,  $\Delta Q = -20$  кг/год (знак "минус" соответствует уменьшению).

Считая эти изменения малыми, можно приближенно оценить производную:

$$Q'_p \approx -\frac{20}{1} = -20 \frac{\text{кг}^2}{\text{руб.} \cdot \text{год}}$$

Допустим, аналогичным образом мы установили, что для билетов в кино

$$Q'_p \approx -0,5 \frac{\text{шт}^2}{\text{руб.} \cdot \text{год}}$$

Какая из этих величин больше? Вопрос бессмыслен, и не только из-за того, что величины, измеренные в различных единицах, несопоставимы.

Эти трудности можно преодолеть, если в качестве основного показателя реакции спроса на изменение цены использовать не **производную**, а **эластичность** спроса по цене – предел отношения **относительного приращения** объема  $\delta Q = \Delta Q / Q$  к **относительному приращению цены**  $\delta P = \Delta P / P$  при условии, что последнее стремится к нулю

В общем случае эластичность вычисляется по формуле

$$E_x(y) = \frac{dy/y}{dx/x} = \frac{x}{y} \cdot y'$$

Эластичность - безразмерная величина; ее использование снимает сложности, связанные с единицами и масштабами рассматриваемых величин.

**Пример (степенная зависимость).** Вычислим эластичность для степенного вида зависимости.

$$y = a \cdot x^b$$

$$E_x(y) = y' \cdot \frac{x}{y} = a \cdot b \cdot x^{b-1} \cdot \frac{x}{a \cdot x^b} = b$$

То есть для степенного вида зависимости эластичность постоянна. При изменении  $x$  на 1 % величина  $y$  изменяется на  $b$  %.

**Пример (линейная зависимость).** Вычислим эластичность для линейного вида зависимости.

$$y = a + bx$$

$$E_x(y) = y' \cdot \frac{x}{y} = b \cdot \frac{x}{a + bx} = \frac{bx}{a + bx}$$

## 6.2. Парная линейная регрессия

В частном случае, когда  $m=1$  (одна независимая переменная  $x$ ) и ищется прямая линия  $y = \beta_0 + \beta_1 x$ , получаются легко запоминающиеся формулы:

$$\begin{aligned} y &= \beta_0 + \beta_1 x \\ \hat{\beta}_1 &= \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{x^2 - (\bar{x})^2}, \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \end{aligned}$$

### 6.3. *Нелинейные формы зависимости*

Использование лишь линейных зависимостей для описания экономических взаимосвязей часто оказывается недостаточным. Необходимо использовать и нелинейные соотношения.

Например, *производственная функция Кобба-Дугласа*:

$$Y = A \cdot K^\alpha \cdot L^\beta$$

где  $Y$  – выпуск,  $K$  – затраты капитала,  $L$  – затраты труда,  $A$ ,  $\alpha$ ,  $\beta$ , – параметры.

В некоторых случаях *нелинейные* зависимости путем замены переменных можно преобразовать к *линейному* виду.

#### **Пример (гиперболическая зависимость).**

По исходным данным оценить коэффициенты в формуле

$$y = a + \frac{b}{x} + \varepsilon$$

$y$	$x$
5	2
3	3
3	5
2	5
2	7
1	8

**Решение.** Заменой  $z = 1/x$  зависимость приводится к линейному виду  $y = a + bz + \varepsilon$ .

$$y = 0,43 + 8,94z + \varepsilon$$

*Ответ:*  $y = 0,43 + \frac{8,94}{x} + \varepsilon$

$y$	$z = 1/x$
5	0,5
3	0,3333
3	0,2
2	0,2
2	0,1428
1	0,125

**Пример (полиномиальная зависимость).**

По исходным данным оценить коэффициенты в формуле

$$y = a + bx + cx^2 + \varepsilon$$

**Решение.** Заменой  $z = x^2$  зависимость приводится к линейному виду  $y = a + bx + cz + \varepsilon$ .

$y$	$x$
5	2
3	3
3	5
2	5
2	7
1	8

$$y = 6,54 - 1,11x + 0,06z + \varepsilon$$

*Ответ:*  $y = 6,54 - 1,11x + 0,06x^2 + \varepsilon$

$y$	$x$	$z = x^2$
5	2	4
3	3	9
3	5	25
2	5	25
2	7	49
1	8	64

**Пример (логистическая кривая).**

По исходным данным оценить коэффициенты в формуле

$$y = \frac{1}{a + be^{-x} + \varepsilon}$$

**Решение.** Заменой  $z = 1/y$ ,  $t = e^{-x}$  зависимость приводится к линейному виду  $z = a + bt + \varepsilon$ .

$y$	$x$
1,61	1
1,78	2
1,92	2
1,90	3
2,06	4

$$z = 0,49 + 0,35t + \varepsilon$$

$$\text{Ответ: } y = \frac{1}{0,49 + 0,35e^{-x} + \varepsilon}$$

$z = 1/y$	$t = e^{-x}$
0,62	0,37
0,56	0,14
0,52	0,14
0,52	0,05
0,49	0,02

**Пример (степенная зависимость).** К линейному виду можно привести *степенную* зависимость:

$$y = a \cdot x_1^{b_1} x_2^{b_2} \cdot \dots \cdot x_m^{b_m} \cdot (1 + \varepsilon)$$

где  $a, b_1, \dots, b_m$  – параметры, а  $\varepsilon$  – случайный множитель. Прологарифмируем это соотношение

$$\ln y = \ln a + b_1 \ln x_1 + \dots + b_m \ln x_m + \ln(1 + \varepsilon)$$

Если теперь обозначить

$$\tilde{y} = \ln y; \beta_0 = \ln a; \tilde{x}_i = \ln x_i; \tilde{\varepsilon} = \ln(1 + \varepsilon)$$

то оно примет вид *линейной* регрессионной модели

$$\tilde{y} = \beta_0 + \beta_1 \tilde{x}_1 + \dots + \tilde{\varepsilon}$$

**Пример.** По исходным данным оценить коэффициенты  $A, \alpha, \beta$  в формуле Кобба-Дугласа

$$Y = A \cdot K^\alpha \cdot L^\beta \cdot (1 + \varepsilon)$$

**Решение.** Логарифмированием зависимость приводится к линейному виду

$$\ln Y = \ln A + \alpha \ln K + \beta \ln L + \ln(1 + \varepsilon).$$

$Y$	$K$	$L$
34	2	4
53	3	6
81	5	7
81	5	7
113	7	9
97	9	4

**Ответ.**

$$Y = 12 \cdot K^{0,69} \cdot L^{0,41} \cdot (1 + \varepsilon)$$

$\ln Y$	$\ln K$	$\ln L$
3,53	0,69	1,39
3,96	1,1	1,79
4,39	1,61	1,95
4,39	1,61	1,95

#### **6.4. Нелинейные модели, не приводимые к линейному виду**

Всё же далеко не всякую зависимость заменой переменных можно привести к линейному виду.

**Например,** если искать логистическую зависимость в виде

$$y = \frac{1}{a + b e^{-x}} + \varepsilon$$

вместо

$$y = \frac{1}{a + b e^{-x} + \varepsilon}$$

то никакими переобозначениями к линейному виду её не привести.

**Другие примеры:**

$$y = ax^b + \varepsilon, \quad y = a + bx^c + \varepsilon, \quad y = ae^{bx} + \varepsilon, \dots$$

В общем случае уравнение нелинейной регрессии с аддитивным случайным членом  $\varepsilon$  имеет вид

$$y = f(x_1, \dots, x_m, \beta_0, \dots, \beta_k) + \varepsilon,$$

Для нахождения оценок этих параметров можно использовать, как и в линейном случае, **метод наименьших квадратов**

$$ESS = \sum_1^n e_i^2 = \sum_1^n (y_i - \hat{y}_i)^2 \rightarrow \min$$

ESS является функцией  $k+1$  переменных  $\beta_0, \dots, \beta_k$ .  
Задача минимизации решается при помощи итеративных методов оптимизации, например, градиентным методом.

**Пример.** По исходным данным оценить коэффициенты  $a, b$  в формуле

$$y = a \cdot e^{bx} + \varepsilon$$

$y$	$x$
2	1
9	2
47	3

**Решение.** Задача не сводится к линейной. Выпишем формулу для ESS:

$$ESS = \sum_1^n (y_i - \hat{y}_i)^2 = (2 - ae^b)^2 + (9 - ae^{2b})^2 + (47 - ae^{3b})^2.$$

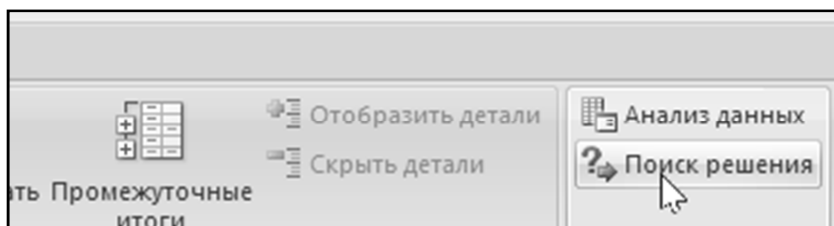
Требуется тем или иным способом найти **минимум** этой функции. Используем надстройку **Поиск решения** из **MS Excel**.

Внесём данные и формулы в таблицу:

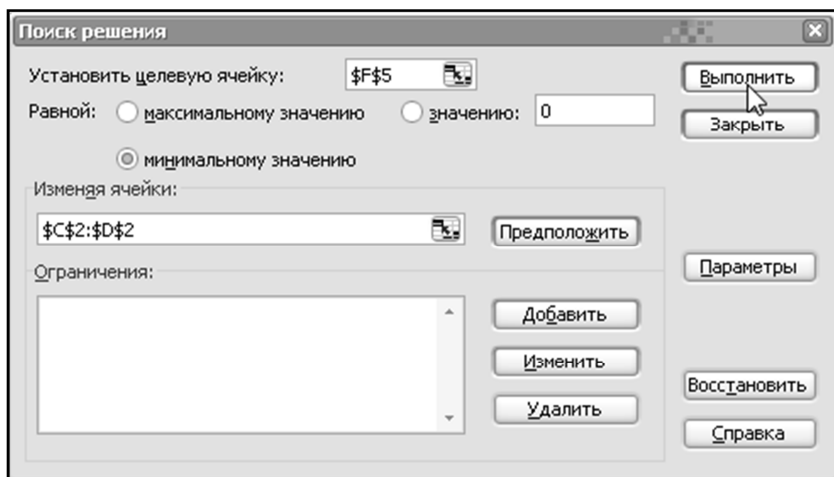
	A	B	C	D	E	F
1	Y	X	a	b	F(x)	
2	2	1	0	0	=C\$2*EXP(D\$2*B2)	=(E2-A2)^2
3	9	2			=C\$2*EXP(D\$2*B3)	=(E3-A3)^2
4	47	3			=C\$2*EXP(D\$2*B4)	=(E4-A4)^2
5					<b>ESS</b>	<b>=СУММ(F2:F4)</b>

Начальные значения a и b положены равными нулю.

Вызовем надстройку **Поиск решения**



и введём данные в диалоге



Получим результат:



	A	B	C	D	E	F
1	Y	X	a	b	F(x)	
2	2	1	0,341	1,642	1,76	0,06
3	9	2			9,09	0,01
4	47	3			46,99	0,00
5					ESS	0,07

Таким образом, оценки параметров модели:  $a=0,341$ ,  
 $b=1,642$ ,

$$y = 0,341 \cdot e^{1,642x} + \varepsilon$$

## 7. Нарушения допущений классической линейной модели

### 7.1. Мультиколлинеарность

#### 7.1.1. Определение

Одним из условий классической регрессионной модели является предположение (C) о линейной независимости столбцов матрицы X.

**Нарушение этого условия – линейная зависимость двух или более объясняющих переменных – называется (точной) мультиколлинеарностью.**

**В этой ситуации нельзя построить МНК-оценку параметров регрессии, так как тогда не существует**

$$(X^T X)^{-1},$$

в силу чего бесполезна формула

$$\hat{\beta} = (X^T X)^{-1} X^T Y, \quad (4.4)$$

**Точная** линейная зависимость переменных (например,  $x_1$  – расстояние в метрах,  $x_2$  – в километрах, так что  $x_1=1000 \cdot x_2$ )

– это явная ошибка. На практике более реальна ситуация, когда между объясняющими переменными существует не точная линейная, а сильная корреляционная зависимость.

Эту ситуацию мы и будем называть **мультиколлинеарностью**.

Например, если в состав объясняющих переменных  $X$  входят и доход, и потребление, то обе эти переменные будут сильно коррелированными.

### 7.1.2. Последствия мультиколлинеарности

- ❖ Оценки коэффициентов регрессии по-прежнему остаются наилучшими линейными несмещенными оценками.
- ❖ Стандартные ошибки коэффициентов при коррелированных регрессорах увеличиваются, так что коэффициенты могут стать незначимыми

### 7.1.3. Признаки мультиколлинеарности

- ❖ Небольшое изменение исходных данных приводит к существенному изменению оценок параметров модели.
- ❖ Оценки коэффициентов имеют малую значимость, в то время как модель в целом является значимой.
- ❖ Оценки коэффициентов имеют неправильные с точки зрения теории знаки.

**Пример.** Рассмотрим следующую выборку.

Матрица корреляций имеет вид

	Y	X <sub>1</sub>	X <sub>2</sub>
Y	1		
X <sub>1</sub>	0,78	1	
X <sub>2</sub>	-0,7	-0,99	1

Видна сильная линейная зависимость между X<sub>1</sub> и X<sub>2</sub>.

(На самом деле X<sub>2</sub>=8-0,5X<sub>1</sub>+ε).

Y	X <sub>1</sub>	X <sub>2</sub>
2,0	1,1	6,9
3,0	2	6,8
3,0	3,5	5,4
4,0	4	4,3
3,0	2,7	5,5
5,0	4,1	4,6
2,0	5	3,2
8,0	7	1,8
6,0	8	0,6
5,0	6,4	2,5

Оценим регрессию

$$Y = -12,54 + 2,22X_1 + 1,65X_2$$

F	Значимость F
6,961	0,022

Регрессия значима при  $\alpha > 0,022$ , т.е., скажем, при  $\alpha = 0,1$  или  $0,05$

	Коэффициенты	P-Значение
Y-пересечен	-12,54	0,33
X <sub>1</sub>	2,22	0,14
X <sub>2</sub>	1,65	0,28

Коэффициенты незначимы даже при  $\alpha = 0,1$

Изменим только одно значение X<sub>2</sub> с 6,9 на 7,1:

Получим регрессию

$$Y = -13,61 + 2,33X_1 + 1,77X_2$$

Видим, что совсем небольшое изменение исходных данных привело к заметному изменению коэффициентов.

Y	X <sub>1</sub>	X <sub>2</sub>
2,0	1,1	7,1
3,0	2	6,8
3,0	3,5	5,4
4,0	4	4,3
3,0	2,7	5,5
5,0	4,1	4,6
2,0	5	3,2
8,0	7	1,8
6,0	8	0,6
5,0	6,4	2,5

#### 7.1.4. Способы обнаружения

На практике о наличии МК первоначально судят по матрице парных корреляций между регрессорами. Большое по модулю значение коэффициента корреляции между двумя  $X$  говорит о наличии МК.

Близкое к нулю значение определителя этой матрицы также говорит о возможной МК.

**Ещё один часто используемый для обнаружения МК показатель называется VIF (“фактор инфляции вариации”).**

Пусть исследуется регрессия

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_m X_m + \varepsilon$$

Построим регрессию регрессора  $X_i$  на остальные  $X$  и вычислим соответствующий  $R^2$ :

Тогда

$$VIF(X_i) = \frac{1}{1 - R_i^2}$$

Если какие-то VIF оказываются больше 10, то это служит признаком МК.

В предыдущем примере  $VIF(X_1) = VIF(X_2) = 49,3$ .

#### 7.1.5. Методы устранения

Все кардинальные методы борьбы с МК имеют свои недостатки.

- Если задачей исследования является **прогноз** будущих значений зависимой переменной, то МК не является серьёзной проблемой, так как она не ухудшает значимость уравнения.
- Естественным методом устранения мультиколлинеарности является **исключение отдельных переменных** из модели.

Но при этом могут возникнуть новые трудности. Например, не всегда ясно, какие переменные являются “лишними”.

Удаление переменных может изменить содержательный смысл модели.

Кроме того, коэффициент при оставшейся переменной получает смещение, т.к. теперь измеряет не только свое «чистое» влияние на  $Y$ , но и влияние на  $Y$  всех отброшенных переменных

- Иногда можно **увеличить объём выборки или качество исходных данных.**

## **7.2. Гетероскедастичность**

### **7.2.1. Определение**

Условие (D) теоремы Гаусса-Маркова требует постоянства дисперсии случайного члена  $\varepsilon_i$  во всех наблюдениях:

$$D\varepsilon_i = \sigma^2 \text{ не зависит от } i$$

Это называется *гомоскедастичность*.

Если же дисперсия случайного члена *меняется* от наблюдения к наблюдению, то мы имеем дело с *гетероскедастичностью*.

## 7.2.2. Последствия гетероскедастичности

- ❖ Оценки коэффициентов регрессии несмещённые.
- ❖ Стандартные ошибки коэффициентов являются заниженными. Это приводит к завышению  $t$ -статистик и даёт неправильное (завышенное) представление о точности оценок.

## 7.2.3. Способы обнаружения

Существуют различные способы выявления гетероскедастичности.

### 7.2.3.1. Графический метод.

- При помощи метода наименьших квадратов оценивается регрессия.
- Вычисляются остатки  $e_j$ .
- Строится график зависимости остатков  $e_j$  от номера наблюдения или от значений какой-то из переменных.
- Далее анализируется внешний вид этого графика

### 7.2.3.2. Тест Уайта

- Для исходной модели (объём выборки = n) строится линейная регрессия и находятся остатки  $e_i$ .
- Строится регрессия квадратов этих остатков на все исходные переменные, их квадраты, попарные произведения и константу (пусть всего получится N переменных, не считая константы).
- Если для второй регрессии  $n \cdot R^2$  больше  $\chi^2_{1-\alpha}(N)$ , то гетероскедастичность есть.

**Пример.** Проверить на гетероскедастичность при  $\alpha=0,1$ , используя тест Уайта

**Решение.** 1) Построим регрессию и найдём расчётные значения и остатки  $e_i$ .

$$y = 0,77 + 1,15x_1 + 0,22x_2$$

2) Построим регрессию по выборке, описанной в тесте Уайта:

$e_i$	$e_i^2$	$x_1$	$x_2$	$x_1^2$	$x_2^2$	$x_1 \cdot x_2$
-0,72	0,52	2	3	4	9	6
0,13	0,017	3	3	9	9	9
-0,52	0,27	3	6	9	36	18
0,41	0,17	4	1	16	1	4
-1,03	1,06	5	7	25	49	35
2,12	4,5	4	7	16	49	28
-0,4	0,16	6	8	36	64	48

$\hat{y}$	$e = y - \hat{y}$	$y$	$x_1$	$x_2$
3,72	-0,72	3	2	3
4,87	0,13	5	3	4
5,52	-0,52	5	3	6
5,59	0,41	6	4	1
8,03	-1,03	7	5	7
6,99	2,12	9	4	7
9,4	-0,4	9	6	8

3) Для этой регрессии вычислим:

$$nR^2 = 7 \cdot 0,4835 = 3,3845,$$

$$\chi^2_{1-\alpha}(N) = \chi^2_{0,9}(5) = 9,236$$

$3,3845 < 9,236$ , так что гетероскедастичности нет.

### 7.2.3.3. Тест Голдфелда-Квандта

В этом тесте предполагается, что стандартное отклонение  $\sigma$  пропорционально значению некоторой независимой переменной  $x_i$ .

- 1) Наблюдения упорядочиваются по возрастанию значений этой переменной  $x_i$ .
- 2) Берутся первые  $n'$  и последние  $n'$  наблюдений, для них строятся регрессии, и оцениваются  $ESS_1$  и  $ESS_2$  соответственно.  
Число  $n'$  можно взять примерно равным  $n/3$ .
- 3) Составляется статистика

$$F = \frac{ESS_2}{ESS_1}$$

- 4) Значение этой статистики, большее

$$F_{1-\alpha}(n'-m-1, n'-m-1)$$

где  $m$  – число объясняющих переменных в исходном уравнении регрессии,  
означает **наличие гетероскедастичности**



**Пример.** Проверить на гетероскедастичность по переменной  $x_1$  при  $\alpha=0,1$ , используя тест Голдфелда-Квандта.

**Решение.** 1) Объем выборки=12, так что выберем  $n'=12/3=4$ . Отсортируем выборку по возрастанию  $x_1$  и построим две регрессии (по первым и последним четырём строкам для  $x_1$ ).

$$y = -3,37 + 2,79x_1 + 0,11x_2, \quad ESS_1 = 0,95$$

$$y = -0,67 - 13,67x_1 + 32x_2, \quad ESS_2 = 112,7$$

2)

$$F = \frac{ESS_2}{ESS_1} = \frac{112,7}{0,95} \approx 118$$

$$F_{1-\alpha}(n' - m - 1; n' - m - 1) = F_{0,9}(4 - 2 - 1; 4 - 2 - 1) = 40,$$

$118 > 40$ , так что гетероскедастичность в виде пропорциональности  $\sigma$  и  $x_1$  есть.

y	$x_1$	$x_2$
3	2	3
5	3	4
5	3	6
6	6	1
7	5	7
9	4	7
9	6	8
11	6	7
18	8	4
32	9	5
68	12	7
78	15	9

### 7.2.3.4. Тест Глейзера

- Оценивается регрессионная зависимость
 
$$|e_i| = a + bx_i^k + u_i$$
- Параметр  $k$  изменяется с некоторым шагом, и для каждого его значения строится регрессия.
- *Статистическая значимость* уравнения при каком-либо значении означает *наличие гетероскедастичности*.

### 7.2.4. Методы устранения

Рассмотрим только случай, когда **стандартное отклонение  $\sigma_i$  пропорционально значению некоторой независимой переменной  $x_k$** . Разделив уравнение регрессии

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m + \varepsilon$$

на  $x_k$ , запишем его в виде

$$\frac{y}{x_k} = \frac{\beta_0}{x_k} + \frac{\beta_1 x_1}{x_k} + \dots + \beta_k + \dots + \frac{\beta_m x_m}{x_k} + \frac{\varepsilon}{x_k}$$

Если теперь перейти к новым переменным

$$\tilde{y} = \frac{y}{x_k}; \tilde{\varepsilon} = \frac{\varepsilon}{x_k}; \tilde{x}_i = \dots; \tilde{\beta}_i = \dots$$

то оно примет вид

$$\tilde{y} = \tilde{\beta}_0 + \tilde{\beta}_1 \tilde{x}_1 + \dots + \tilde{\beta}_m \tilde{x}_m + \tilde{\varepsilon}$$

причем

$$D(\tilde{\varepsilon}) = const$$

**Пример.** В предыдущем примере мы выяснили наличие гетероскедастичности в виде пропорциональности  $\sigma$  переменной  $x_1$

Если использовать обычный МНК, то получится

$$y = -23,76 + 6x_1 + 0,92x_2$$

Но если учесть гетероскедастичность указанного вида, то более правильно, **разделив всё уравнение на  $x_1$** , ввести новые переменные

$$\tilde{y} = \frac{y}{x_1}; \tilde{x}_1 = \frac{1}{x_1}; \tilde{x}_2 = \frac{x_2}{x_1}$$

Получим таблицу

Из которой

$$\tilde{y} = 3,99 - 6,51\tilde{x}_1 - 0,16\tilde{x}_2,$$

так что

$$y = -6,51 + 3,99x_1 - 0,16x_2$$

y	x <sub>1</sub>	x <sub>2</sub>
3	2	3
5	3	4
5	3	6
6	6	1
7	5	7
9	4	7
9	6	8
11	6	7
18	8	4
32	9	5
68	12	7
78	15	9

$\tilde{y}$	$\tilde{x}_1$	$\tilde{x}_2$
3/2	1/2	3/2
5/3	1/3	4/3
5/3	1/3	6/3
6/6	1/6	1/6
...	...	...

## 7.3. Автокорреляция

### 7.3.1. Определение

Условие (E) теоремы Гаусса-Маркова требует отсутствия корреляции между значениями случайного фактора  $\varepsilon_i$  во всех наблюдениях:

$$\text{cov}(\varepsilon_i \varepsilon_j) = 0 \text{ при } i \neq j$$

Фактически это условие соответствует независимости случайных членов в разных наблюдениях.

При нарушении этого условия, т.е. при наличии связи между случайными членами для разных наблюдений, возникает явление *автокорреляции*.

Обычно автокорреляция рассматривается при изучении *временных рядов* (т.е. когда номер наблюдения соответствует *моменту времени*).



**Пример.** Авторегрессия первого порядка  $AR(1)$ .

$$\varepsilon_i = \rho \cdot \varepsilon_{i-1} + u_i, \quad -1 < \rho < 1$$

$\varepsilon$  – случайный член уравнения регрессии,

$\rho$  – коэффициент авторегрессии,

$u$  – случайный член, не подверженный автокорреляции

В этом случае  $\varepsilon$  в данном наблюдении прямо связан лишь с  $\varepsilon$  в предыдущем наблюдении.

Если  $\rho > 0$ , то автокорреляция *положительная*,  
если  $\rho < 0$ , то *отрицательная*.

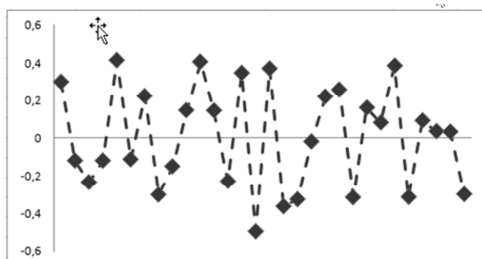
**Замечание.** Коэффициент  $\rho$  в  $AR(1)$  равен *коэффициенту корреляции соседних значений  $\varepsilon$* . Действительно:

$$r(\varepsilon_i, \varepsilon_{i-1}) = r(\rho \cdot \varepsilon_{i-1} + u_i, \varepsilon_{i-1}) = \rho \cdot r(\varepsilon_{i-1}, \varepsilon_{i-1}) = \rho$$

**Пример.** Авторегрессия второго порядка  $AR(2)$ .

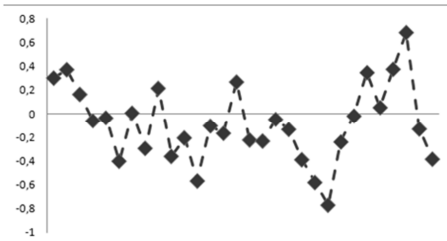
$$\varepsilon_i = \rho_1 \cdot \varepsilon_{i-1} + \rho_2 \cdot \varepsilon_{i-2} + u_i, \quad -1 < \rho_{1,2} < 1$$

**Пример.** Нулевая автокорреляция.



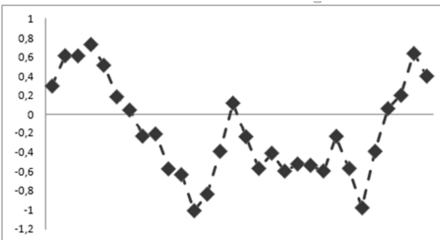
$$\varepsilon_i = 0 \cdot \varepsilon_{i-1} + u_i$$

**Пример.** Положительная автокорреляция.



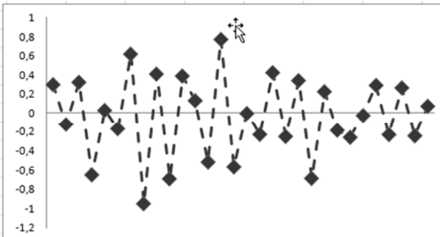
$$\varepsilon_i = 0,5 \cdot \varepsilon_{i-1} + u_i$$

**Пример.** Очень сильная положительная автокорреляция.



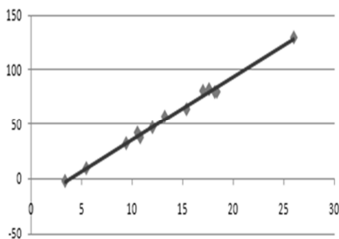
$$\varepsilon_i = 0,9 \cdot \varepsilon_{i-1} + u_i$$

**Пример.** Отрицательная автокорреляция.

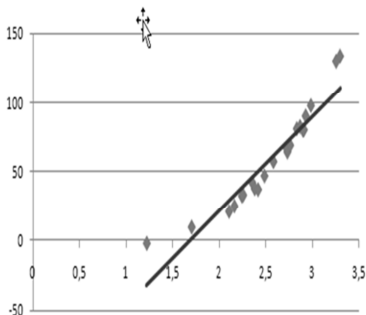


$$\varepsilon_i = -0,8 \cdot \varepsilon_{i-1} + u_i$$

**Пример.** Ложная автокорреляция, вызванная неправильной функциональной формой.



$$Y = \beta_0 + \beta_1 X_1 + \varepsilon$$



Те же данные, но другая формула

$$Y = \beta_0 + \beta_1 \ln X_1 + \varepsilon$$

### 7.3.2. Способы обнаружения

Существуют различные способы выявления автокорреляции.

#### 7.3.2.1. Графический метод.

- 1) При помощи метода наименьших квадратов оценивается регрессия.
- 2) Вычисляются остатки  $e_i$ .
- 3) Строится график зависимости остатков  $e_i$  от номера наблюдения или от значений какой-то из переменных.
- 4) Далее анализируется внешний вид этого графика

#### 7.3.2.2. Метод рядов

В этом методе рассматриваются знаки остатков  $e_i$

*“Ряд” это непрерывная последовательность одинаковых знаков.*

При большом числе наблюдений применяется следующая процедура.

1) Рассчитываются статистики.

$$M = \frac{2n_+n_-}{n_+ + n_-} + 1, \quad D = \frac{2n_+n_-(2n_+n_- - n_+ - n_-)}{(n_+ + n_-)^2(n_+ + n_- - 1)}$$

где  $n_+$  - количество положительных остатков,  $n_-$  - количество отрицательных.

2) Если

$$M - u_{1-\alpha/2} \cdot \sqrt{D} < k < M + u_{1-\alpha/2} \cdot \sqrt{D}$$

где  $k$  – количество рядов, то принимается гипотеза об отсутствии автокорреляции.

### 7.3.2.3. Критерий Дарбина-Уотсона.

Этот метод применяется для обнаружения авторегрессии AR(1)

$$\varepsilon_i = \rho \cdot \varepsilon_{i-1} + u_i$$

В критерии Дарбина-Уотсона проверяется гипотеза  $\rho=0$  при альтернативной  $\rho \neq 0$ .

При этом используется статистика Дарбина-Уотсона:

$$DW = \frac{\sum_{i=1}^n (e_i - e_{i-1})^2}{2 \sum_{i=1}^n e_i^2}, \quad e_i = y_i - \hat{y}_i$$

Для больших выборок  $DW \approx 2 \cdot (1 - \rho)$ .

При значительной положительной автокорреляции

$\rho \approx 1$  и  $DW \approx 0$ .

При отрицательной автокорреляции

$\rho \approx -1$  и  $DW \approx 4$ .

При отсутствии автокорреляции

$\rho \approx 0$  и  $DW \approx 2$ .

Критические значения статистики Дарбина-Уотсона зависят не только от числа переменных, но и от значений, которые они принимают в выборке. Поэтому **невозможно составить таблицы критических значений DW**.

**Но можно указать верхнюю  $d_U$  и нижнюю  $d_L$  границы для DW**. Они определяются в зависимости от  $n$  и числа оцениваемых параметров  $m$ .



Границы  $d_L$ ,  $d_U$  распределения Дарбина-Уотсона при  $\alpha = 0,05$



$m$ ↓	$n \rightarrow$	7	8	9	10	20	30
1	$d_L$	0,7	0,76	0,82	0,88	1,2	1,35
1	$d_U$	1,356	1,33	1,32	1,32	1,41	1,49
2	$d_L$	0,47	0,36	0,63	0,7	1,1	1,29
2	$d_U$	1,9	1,78	1,7	1,64	1,54	1,57
3	$d_L$		0,37	0,44	0,53	1	1,21
3	$d_U$		2,29	2,13	2,02	1,68	1,65
4	$d_L$			0,3	0,38	0,9	1,14
4	$d_U$			2,39	2,41	1,83	1,74

**Пример.** Построена регрессия с  $n=8$ ,  $m=2$ . Используя тест Дарбина-Уотсона, проверить гипотезу о наличии авторегрессии первого порядка при  $\alpha=0,05$ .

**Решение.** Заполним таблицу

$e_i = y_i - \hat{y}_i$	$e_{i-1}$	$(e_i - e_{i-1})^2$	$e_i^2$
-0,09	-	-	0,0081
-0,32	-0,09	0,17	0,1024
-0,79	-0,32	1,23	0,6241
-0,13	-0,79	0,44	0,0169
-1,43	-0,13	1,69	2,0449
1,61	-1,43	9,24	2,5921
-1,02	1,61	6,92	1,0404
1,54	-1,02	6,55	2,3716
$\Sigma$		= 26,2	= 8,8

$y$	$\hat{y}$
3	3,09
5	4,68
5	5,79
6	6,13
7	8,43
9	7,39
9	10,02
11	9,46

$$DW = \frac{26,2}{8,8} \approx 3$$

По таблице  $d_L=0,36$ ,  $d_U=1,78$ . Число  $DW=3$  попадает в зону  $(4-d_U; 4-d_L)=(2,22; 3,64)$ . Поэтому критерий Дарбина-Уотсона в этом примере ответа не даёт.

### 7.3.3. Методы устранения автокорреляции

#### 7.3.3.1. Исправление спецификации модели

Как уже отмечалось, важную роль играет правильный выбор функциональной формы зависимости. Например, при выборе линейной формы зависимости в ситуации, когда имеет место экспоненциальная, возникает положительная автокорреляция.

#### 7.3.3.2. Авторегрессионное преобразование

Рассмотрим случай AR(1). Если

$$\varepsilon_i = \rho \cdot \varepsilon_{i-1} + u_i,$$

и коэффициент  $\rho$  известен, то для исходного уравнения регрессии

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

можно выполнить *авторегрессионное преобразование*. Запишем это уравнение для предыдущего номера

$$y_{i-1} = \beta_0 + \beta_1 x_{i-1} + \varepsilon_{i-1},$$

умножим обе его части на  $\rho$  и вычтем из исходного уравнения

$$y_i - \rho y_{i-1} = \beta_0 \cdot (1 - \rho) + \beta_1 \cdot (x_i - \rho \cdot x_{i-1}) + u_i$$

Если ввести новые переменные

$$y_i^* = y_i - \rho \cdot y_{i-1}, \quad x_i^* = x_i - \rho \cdot x_{i-1}, \quad \beta_0^* = \beta_0 \cdot (1 - \rho),$$

то получим уравнение без автокорреляции

$$y_i^* = \beta_0^* + \beta_1 x_i^* + u_i.$$

Найдя оценки коэффициентов этого уравнения, мы тем самым сможем вычислить и оценки коэффициентов исходного.

Однако *проблема* заключается в том, что **величина  $\rho$  неизвестна**.

Поэтому используют разные методы получения **оценки для  $\rho$** , такие как *метод Хилдрета-Лу*, *метод*

*Кохрейна-Оркатта*, или просто  $\rho \approx 1 - \frac{DW}{2}$ .

#### **7.4. Стохастические объясняющие переменные**

До сих пор мы предполагали, что переменные  $x_1, x_2, \dots, x_m$  являются неслучайными, детерминированными (в отличие от  $y$ ).

Но иногда их приходится считать случайными. Выясним, к чему это может привести.

##### **7.4.1. Случай некоррелированности $X$ и $\varepsilon$ .**

Пусть, как и раньше,

$$Y = X \cdot \beta + \varepsilon$$

Будем считать элементы матрицы  $X$  случайными величинами. Потребуем, чтобы были выполнены условия

- (C<sub>1</sub>) Ранг матрицы  $X$  равен  $m+1$  с вероятностью единица  
 (D<sub>1</sub>)  $M(\varepsilon | X) = 0$ ,  $M(\varepsilon \varepsilon^T | X) = \sigma^2 E$ .

Здесь  $M(\dots | X)$  – условные математические ожидания.

В этих предположениях

- ✓ МНК-оценки параметров регрессии являются несмещёнными.
- ✓ Среди всех линейных условно несмещённых оценок
- ✓ МНК-оценка обладает наименьшей условной дисперсией.

Если в каждом наблюдении значения объясняющих переменных (регрессоров) выбираются из одной и той же генеральной совокупности, а случайные факторы независимы, одинаково распределены и не зависят от  $X$ , то МНК-оценки состоятельны.

#### **7.4.2. Коррелированность объясняющих переменных и случайного фактора.**

Если же объясняющие переменные  $X$  и ошибки  $\varepsilon$  коррелированы (т.е.  $M(\varepsilon|X) \neq 0$ ), то в общем случае

- ✓ МНК-оценки могут быть смещёнными и несостоятельными
- ✓ Содержательная интерпретация зависимостей ошибочна
- ✓ Рекомендации, сделанные на основе модели неверными.

## 8. Анализ временных рядов.

### 8.1. Определение

**Временной ряд** – это набор наблюдений какой-либо случайной величины, произведённых в последовательные моменты времени.



Это могут быть, например, цены на батон хлеба в соседнем магазине, курс обмена доллара на рубли в ближайшем обменном пункте или годовые объёмы добычи нефти странами ОПЕК.

### 8.2. Области практического применения временных рядов

Основное прикладное значение временных рядов состоит в прогнозировании значений экономических показателей.

Соответствующие задачи возникают в таких областях, как

- Планирование в производстве и торговле;
- Управление и оптимизация социально-экономических процессов в обществе;
- Частичное управление демографическими процессами;
- Принятие решений в бизнесе и т.п.

В предыдущих разделах значения результирующей переменной выводились из значений одной или нескольких объясняющих переменных.

Тем самым считалось, что объясняющая переменная может быть выведена из нескольких значимых факторов и у нас есть возможность выделить их и учесть влияние каждого.

Анализ временных рядов основан на другой идее: результирующая переменная складывается под влиянием большого числа факторов, многие из которых не поддаются непосредственному измерению.

Поэтому лучшим источником информации о совокупности этих факторов являются значения самой исследуемой переменной в предыдущие моменты времени.

### **8.3. Характер зависимости от времени.**

Чаще всего рассматривают временные ряды с равноотстоящими моментами наблюдений

$$t_2 - t_1 = t_3 - t_2 = \dots = t_N - t_{N-1} = \Delta.$$

где  $\Delta$  – временной интервал (минута, час, сутки, неделя, месяц, квартал, год и т.п.)

В таком случае бывает удобнее писать  $y(1)$ ,  $y(2)$  или даже  $y_1$ ,  $y_2$ , вместо  $y(t_1)$ ,  $y(t_2)$ ,...

### **8.4. Основные факторы, формирующие временной ряд**

Обычно выделяют четыре типа факторов, под влиянием которых формируются значения временного ряда.

1. Долговременные – формируют общую тенденцию изменений в длительной перспективе. Функция, описывающая эту тенденцию, называется трендом.
2. Сезонные – формируют периодические, повторяющиеся в определённое время года колебания изменяемого показателя.

3. Циклические – формируют изменения показателя под влиянием действия циклов экономической, демографической и другой природы.
4. Случайные – эти факторы не поддаются учёту.

Не обязательно, чтобы присутствовали все четыре составляющие временного ряда.

Но будем считать, что случайные факторы присутствуют всегда (иначе неинтересно).

Выводы о том, участвуют или нет факторы данного типа в формировании значений  $y(t)$ , могут базироваться как на анализе содержательной сущности задачи, так и на специальном статистическом анализе исследуемого временного ряда.

**Пример.** Рассмотрим данные о суммарных месячных расстояниях  $y(t)$  (в тысячах миль), пройденных британскими авиалайнерами за 96 месяцев с января 1963 г. по декабрь 1970 г. (т. е.  $t=1, 2, \dots, 96$ ), временной интервал  $\Delta$  равен месяцу



- ❖ Данные об авиалайнерах представляют собой образец сезонных колебаний, наслаивающихся на монотонно растущий тренд.

- ❖ Сезонный эффект легко расшифровывается. Например, в 1968 г. мы наблюдали три "всплеска" активности пассажирских авиаперевозок; все они объясняются перелетами в праздничный и отпускной периоды: один из них приходится на пасху, второй - на лето и третий - на рождественские праздники.
- ❖ Амплитуда колебаний меняется из года в год частично из-за возрастающего парка авиалайнеров, а частично из-за увеличения периода праздников.

### **8.5.      *Аддитивная и мультипликативная модели временного ряда***

Модель, в которой временной ряд представлен в виде суммы тренда, сезонной (или циклической) и случайной компонент, называется аддитивной.

$$y_t = f(t) + \varphi(t) + \varepsilon_t$$

Модель, в которой временной ряд представлен в виде произведения тренда, сезонной (или циклической) и случайной компонент, называется мультипликативной.

$$y_t = f(t) \cdot \varphi(t) \cdot \varepsilon_t$$

Мультипликативную модель часто можно свести к аддитивной логарифмированием.



## **8.6. Выделение основных компонентов временного ряда**

### **8.6.1. Порядок анализа модели**

Построение модели включает следующие шаги

1. Сглаживание исходного ряда.
2. Расчёт значений периодической компоненты.
3. Устранение периодической компоненты из исходного ряда и получение аналитического выражения для тренда.
4. Анализ остатков.

### **8.6.2. Сглаживание временного ряда методом скользящей средней**

Чтобы уменьшить влияние случайных и циклических факторов, упростив нахождение тренда, используется сглаживание временного ряда *методом скользящей средней*. Сглаживание представляет собой вычисление взвешенного среднего значений, наблюдаемых в окрестности рассматриваемой точки. Оно определяется для каждого момента времени, *за исключение нескольких первых и нескольких последних точек*.

В простейшем случае вычисление производится следующим образом: берется  $k$  последовательных значений временного ряда ( $k$  нечётно), со средним, соответствующим текущему моменту времени  $t$ :

$$Y_t, Y_{t-1}, Y_{t+1}, Y_{t-2}, Y_{t+2},$$

и находится их среднее арифметическое:

$$\tilde{y}_t = \frac{y_t + y_{t-1} + y_{t+1} + \dots}{k}$$

Затем точка  $t$  сдвигается вправо на один шаг, опять производится усреднение значений временного ряда и т.д.

Число  $k$  (т.е. временной интервал, по которому производится усреднение), называется **окном**.

Замечание. Иногда требуется усреднять за период, равный чётному количеству шагов. Например, за год=12 месяцев. В этом случае применяют взвешенные средние. Например, в случае года:

$$\tilde{y}_t = \frac{y_{t-6} + 2y_{t-5} + 2y_{t-4} + \dots + 2y_t + 2y_{t+1} + \dots + 2y_{t+5} + y_{t+6}}{24}$$

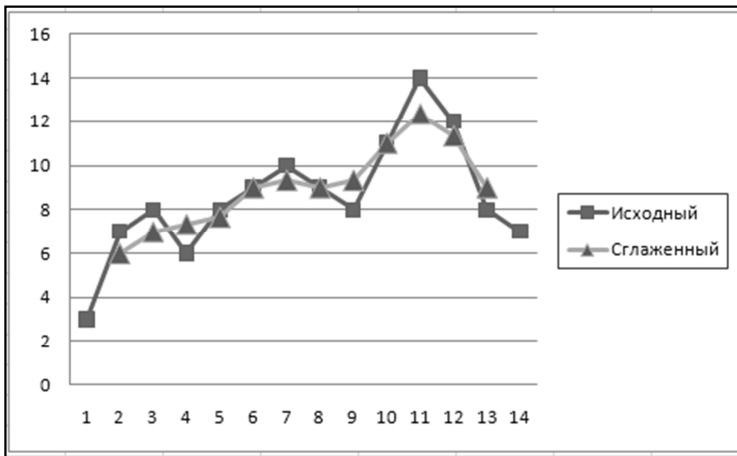
**Пример.** Сгладить временной ряд с окном 3.

$t$	1	2	3	4	5	6	7	8	9	10	11	12	13	14
$y_t$	3	7	8	6	8	9	10	9	8	11	14	12	8	7

**Решение.** Вычислим поочерёдно

$$\tilde{y}_2 = \frac{y_1 + y_2 + y_3}{3} = \frac{3 + 7 + 8}{3} = 6, \quad \tilde{y}_3 = \frac{y_2 + y_3 + y_4}{3} = \frac{7 + 8 + 6}{3} = 7, \dots$$

$t$	1	2	3	4	5	6	7	8	9	10	11	12	13	14
$y_t$	3	7	8	6	8	9	10	9	8	11	14	12	8	7
$\tilde{y}_t$	–	6	7	7,3	7,7	9	9,3	9	9,3	11	12,3	11,3	9	–



### 8.6.3. Экспоненциальное сглаживание

Экспоненциально сглаженный ряд  $S_t$  определяется формулами

$$s_1 = y_1, \quad s_t = s_{t-1} + \alpha(y_t - s_{t-1})$$

где  $\alpha$  – параметр сглаживания,  $0 \leq \alpha \leq 1$ . Чем меньше  $\alpha$ , тем сильнее сглаживание

Можно дать следующую интерпретацию экспоненциальной средней:

если  $S_{t-1}$  — прогноз значения ряда  $Y_t$ , то разность  $Y_t - S_{t-1}$  есть *погрешность* прогноза;

таким образом прогноз  $s_t$  для следующего момента времени  $t+1$  учитывает ставшую известной в момент  $t$  ошибку прогноза.

Метод экспоненциального сглаживания часто применяется для краткосрочного прогнозирования. При этом основной задачей является выбор параметра сглаживания  $\alpha$ .

**Пример.** Экспоненциально сгладить временной ряд с параметрами  $\alpha=0,5$  и  $\alpha=0,3$ .

$t$	1	2	3	4	5	6	7	8	9	10	11	12	13	14
$y_t$	3	7	8	6	8	9	10	9	8	11	14	12	8	7

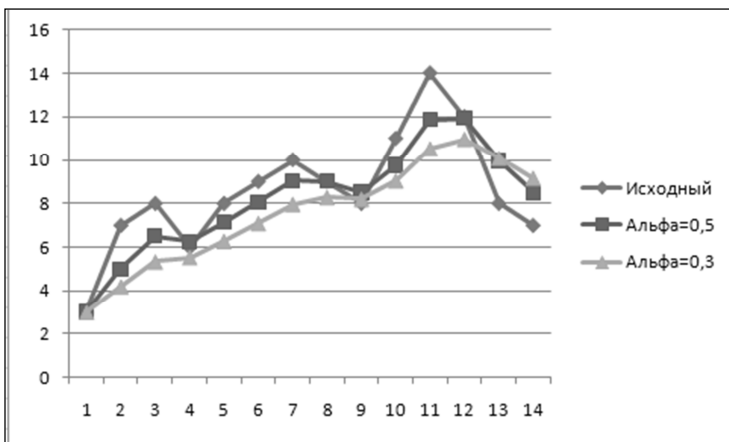
**Решение.** Пусть  $\alpha=0,5$ .

$$s_1 = y_1 = 3, \quad s_2 = s_1 + 0,5 \cdot (y_2 - s_1) = 3 + 0,5 \cdot (7 - 3) = 5, \dots$$

Пусть теперь  $\alpha=0,3$ .

$$s_1 = y_1 = 3, \quad s_2 = s_1 + 0,3 \cdot (y_2 - s_1) = 3 + 0,3 \cdot (7 - 3) = 4,2, \dots$$

$t$	1	2	3	4	5	6	7	8	9	10	11	12	13	14
$y_t$	3	7	8	6	8	9	10	9	8	11	14	12	8	7
$\alpha=0,5$	3	5	6,5	6,3	7,1	8,1	9	9	8,5	9,8	11,9	11,9	10	8,5
$\alpha=0,3$	3	4,2	5,3	5,5	6,3	7,1	8	8,3	8,2	9	10,5	11	10,1	9,2



#### 8.6.4. Выделение тренда временного ряда в аналитической форме

Трендом (или тенденцией) называют неслучайную медленно меняющуюся составляющую временного ряда, на которую могут накладываться циклические и случайные составляющие.

**Замечание.** Слова «медленно меняющаяся составляющая» в определении тренда носят относительный характер.

Например, медленное увеличение количества осадков в течение периода в сотню лет может быть понято как тренд.

Однако на самом деле рост осадков, характерный для этого столетия, может оказаться частью некоторого медленного колебательного процесса, происходящего в пределах нескольких тысячелетий.

При различении тренда и циклической компоненты невозможно полностью исключить из рассуждений элемент субъективности.

В качестве функции тренда чаще всего берут

- ❖ Гиперболический тренд:  $f(t)=a+b/t$ ;
- ❖ Линейный тренд:  $f(t)=a+bt$ ;
- ❖ Экспоненциальный тренд:  $f(t)=a \cdot b^t$ ;
- ❖ Степенной тренд:  $f(t)=a \cdot t^b$ ;
- ❖ Полиномиальный тренд:  $f(t)=a+bt+ct^2+\dots$ ;
- ❖ Логистическая кривая:  $f(t)=1/(a+be^{-x})$

Выбор вида функции обычно производится с учётом выводов экономической теории и визуального анализа графика ряда.

Вычисление коэффициентов тренда часто производится по МНК.

**Задача.** Вычислить коэффициенты линейного тренда  $f(t) = a + bt$  для временного ряда

$t$	1	2	3	4	5	6	7
$y_t$	82	87	99	104	107	121	118

**Решение.** При помощи обычного метода наименьших квадратов получаем.  $f(t) = 76,29 + 6,57 \cdot t$

**Задача.** Для того же временного ряда вычислить коэффициенты степенного тренда  $f(t)=a \cdot t^b$

**Решение.** Линеаризуем, взяв логарифм от обеих частей:

$\ln f = \ln a + b \ln t$	$\ln t$	0	0,7	1,1	1,4	1,6	1,8	1,9
	$\ln y_t$	4,4	4,5	4,6	4,6	4,7	4,8	4,8

$$\begin{pmatrix} \ln a \\ b \end{pmatrix} = \begin{pmatrix} 4,37 \\ 0,21 \end{pmatrix} \quad a = e^{4,37} = 79,13 \quad f(t) = 79,13 \cdot t^{0,21}$$

## Использование MS Excel для выделения тренда

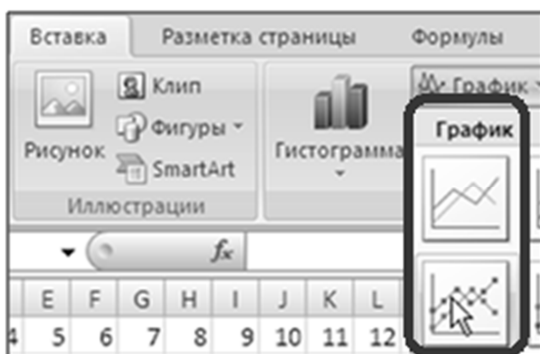
**Пример.** Найдём аналитическое выражение тренда в различных формах для временного ряда

$t$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
$y_i$	3	5	4	4	5	4	7	6	11	9	11	14	18	21	22	29

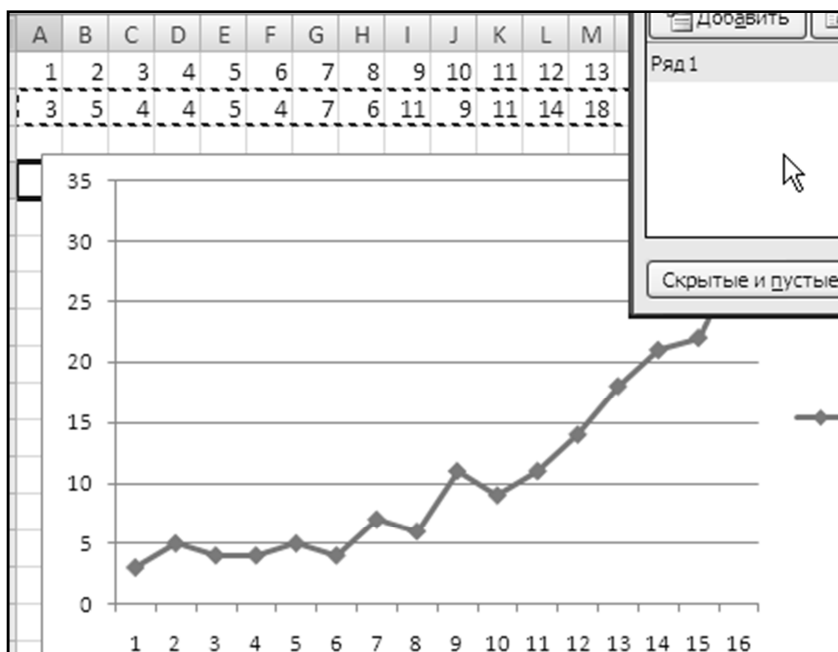
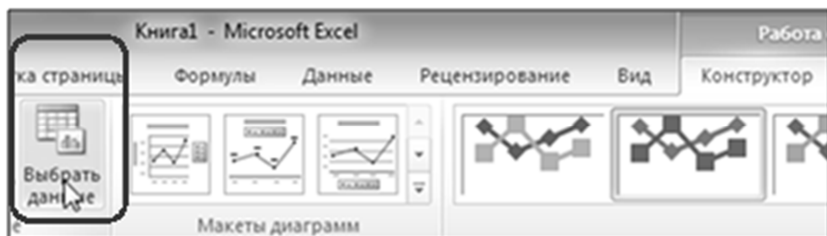
**Решение.** Занесём данные в таблицу MS Excel.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
2	3	5	4	4	5	4	7	6	11	9	11	14	18	21	22	29

Вставка → Диаграммы → График:

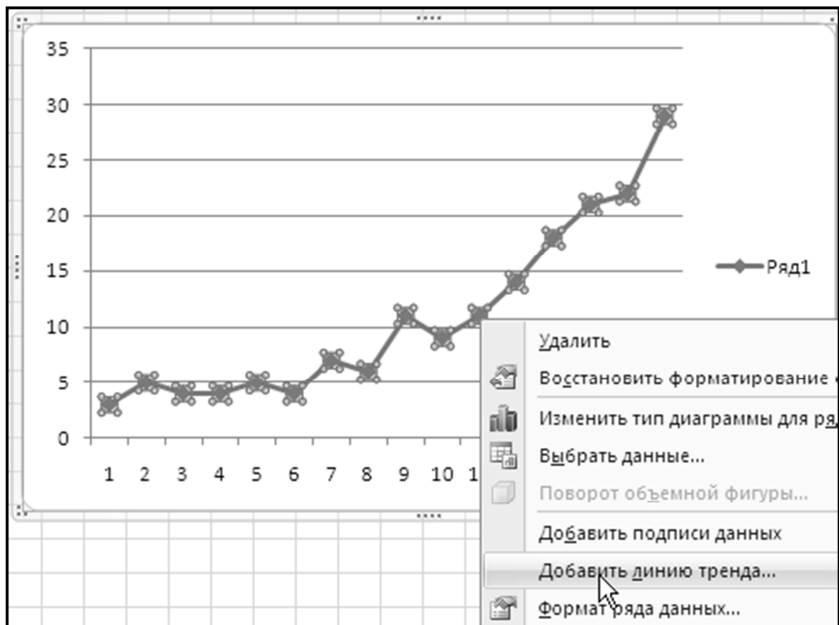


Работа с диаграммами →  
Конструктор→Выбрать данные





Выделим мышкой график и выберем пункт **меню**  
**Добавить линию тренда**



Появится диалог, в котором можно выбрать тип тренда

Параметры линии тренда

Цвет линии  
Тип линии  
Тень

Параметры линии тренда

Построение линии тренда (аппроксимация и сглаживание)

Экспоненциальная  
 Линейная  
 Логарифмическая  
 Полиномиальная    Степень: 2  
 Степенная  
 Линейная фильтрация    Точки: 2

Название аппроксимирующей (сглаженной) кривой

автоматическое:    Линейная (Ряд 1)  
 другое:    \_\_\_\_\_

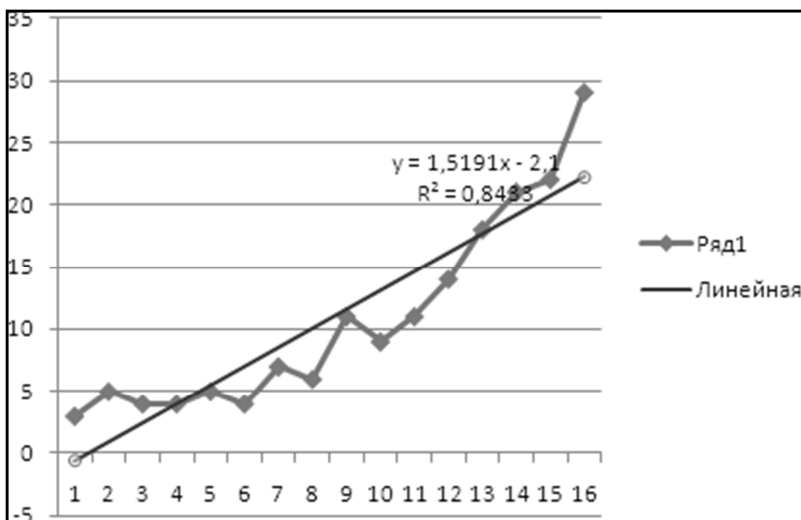
Прогноз

вперед на: 0,0    периодов  
назад на: 0,0    периодов

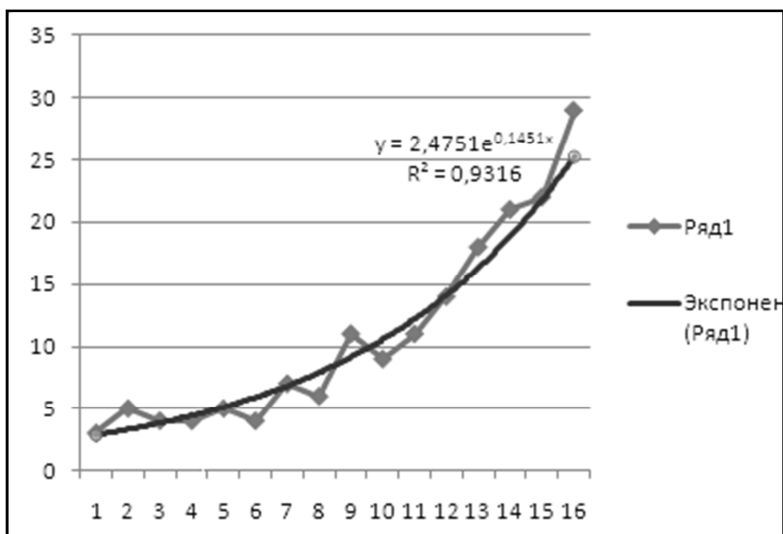
пересечение кривой с осью Y в точке: 0,0  
 показывать уравнение на диаграмме  
 поместить на диаграмму величину достоверности аппроксимации ( $R^2$ )

Заккрыть

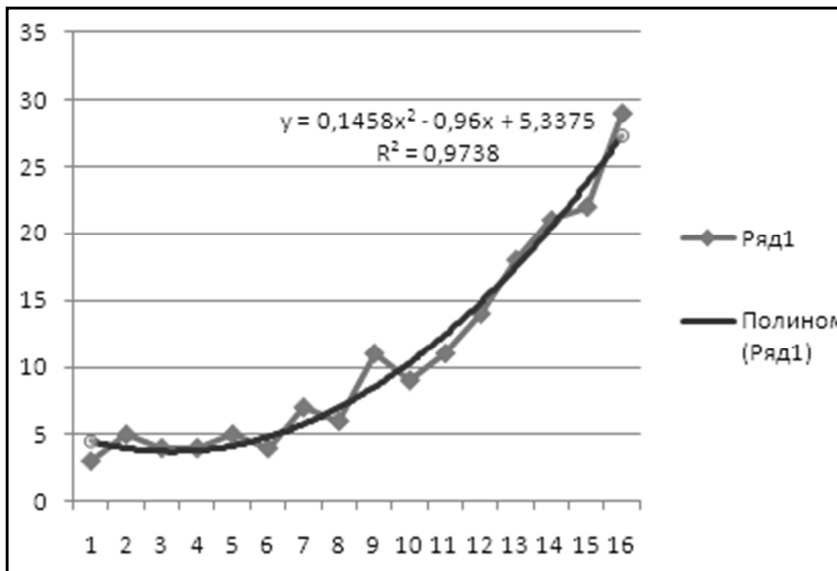
Выберем, например, линейный



Экспоненциальный:



Полиномиальный второй степени:



По критерию  $R^2$  наилучшее значение у последнего вида тренда.

### 8.6.5. Выделение периодической компоненты при помощи фиктивных переменных в аддитивной модели

Компоненты временного ряда можно выделять и другим способом, используя фиктивные переменные.

**Пример.** Выделить линейный тренд и циклическую компоненту временного ряда, считая, что циклическая имеет период 3.

$t$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
$y_t$	8	2	4	10	5	6	12	8	7	14	11	8	16	14	10

**Решение.** Для выделения циклической компоненты периода 3 заводим  $2=3-1$  фиктивные переменные  $D_1$  и  $D_2$ .

$t$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
$y_t$	8	2	4	10	5	6	12	8	7	14	11	8	16	14	10
$D_1$	1	0	0	1	0	0	1	0	0	1	0	0	1	0	0
$D_2$	0	1	0	0	1	0	0	1	0	0	1	0	0	1	0

Строим регрессию:

$$y_t = 0,6 + 0,711 \cdot t + 6,422 \cdot D_1 + 1,711 \cdot D_2$$

Среднее арифметическое коэффициентов при  $D$  (считая ещё один коэффициент нулём):

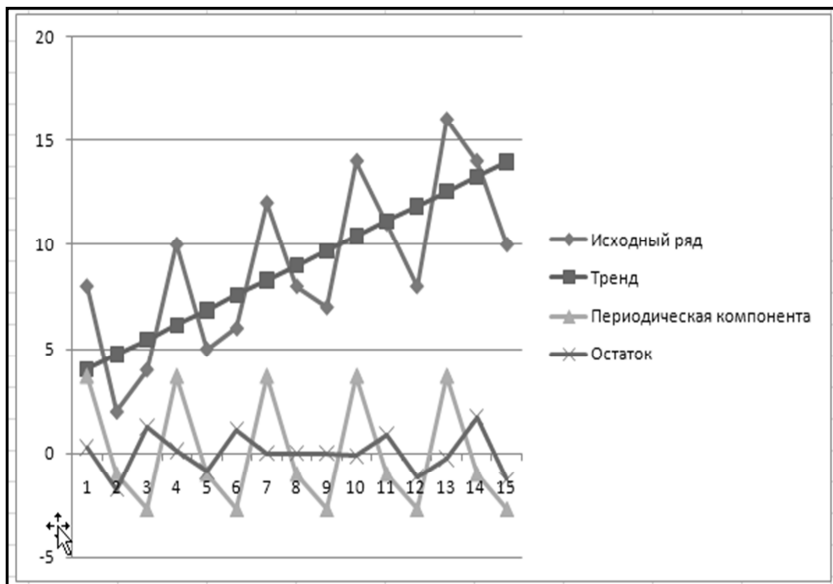
$$\frac{6,422 + 1,711 + 0}{3} = 2,711$$

Получаем циклическую компоненту

$\varphi(1)$	$\varphi(2)$	$\varphi(3)$
$6,422 - 2,711 = 3,711$	$1,711 - 2,711 = -1$	$0 - 2,711 = -2,711$

и линейный тренд:

$$f(t) = (0,6 + 2,711) + 0,711 \cdot t = 3,311 + 0,711 \cdot t$$



**Замечание.** *Сезонная* компонента имеет период 4, так что в этом случае заводятся  $3=4-1$  фиктивные переменные, примерно так:

$t$	1	2	3	4	5	6	7	8	9	10	11	12	...
$y_t$	...	...	...	...	...	...	...	...	...	...	...	...	...
$D_1$	1	0	0	0	1	0	0	0	1	0	0	0	...
$D_2$	0	1	0	0	0	1	0	0	0	1	0	0	...
$D_3$	0	0	1	0	0	0	1	0	0	0	1	0	...

И далее аналогично предыдущему примеру (только в среднем арифметическом деление на 4)

$$y_t = \beta_0 + \beta_1 \cdot t + \beta_2 \cdot D_1 + \beta_3 \cdot D_2 + \beta_4 \cdot D_3$$

### 8.6.6. Выделение компонентов временного ряда в случае мультипликативной модели

**Пример.** Выделить основные компоненты временного ряда, считая, что циклическая имеет период 4

$t$	1	2	3	4	5	6	7	8	9	...	16
$y_t$	72	100	90	64	70	92	80	58	62	...	30

Взятие логарифма превращает мультипликативный ряд в аддитивный:

$$y_t = f(t) \cdot \varphi(t) \cdot \varepsilon_t$$

$$\ln y_t = \ln f(t) + \ln \varphi(t) + \ln \varepsilon_t$$

Поэтому проще всего перейти от временного ряда к его логарифму, выделить компоненты, как в аддитивном, а затем вернуться к исходному.

$t$	1	2	3	4	5	6	...	16
$\ln y_t$	4,277	4,605	4,5	4,159	4,25	4,52	...	3,4

Выделяем периодическую компоненту и тренд, получаем:

$\ln \varphi(1)$	$\ln \varphi(2)$	$\ln \varphi(3)$	$\ln \varphi(4)$
-0,045	0,058	0,1	-0,113

$$\ln f(t) = 4,659 - 0,075 \cdot t$$

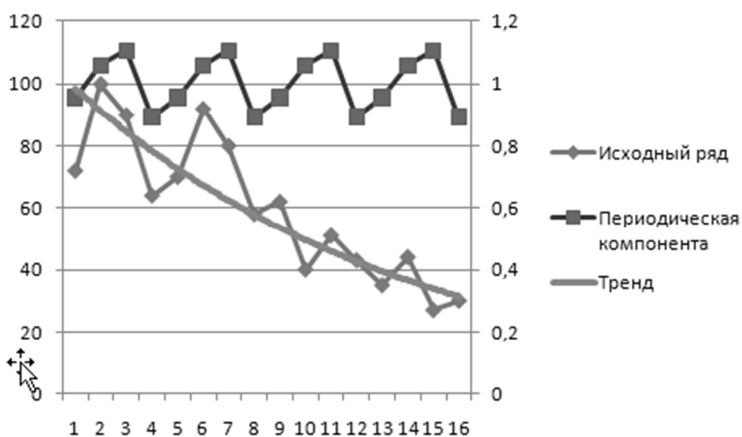
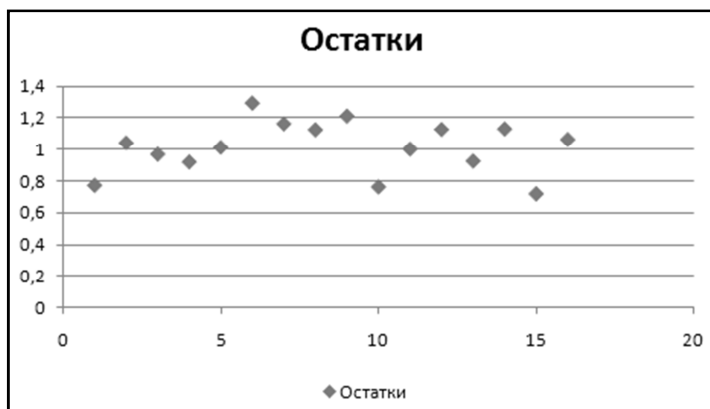
Отсюда

$\varphi(1)$	$\varphi(2)$	$\varphi(3)$	$\varphi(4)$
$e^{-0,045} = 0,956$	$e^{0,058} = 1,06$	$e^{0,1} = 1,106$	$e^{-0,113} = 0,893$

$$f(t) = e^{4,659 - 0,075 \cdot t} = 105,497 e^{-0,075 \cdot t}$$

Остаток находится делением исходного ряда на  $\varphi$  и  $f$

$t$	1	2	3	4	5	6	7	8	...	16
$y_t$	72	100	90	64	70	92	80	58	...	30
$\varepsilon$	0,77	1,04	0,967	0,918	1,012	1,293	1,162	1,124	...	1,062





## 8.7. Автокорреляционная функция. Коррелограмма.

Если временной ряд  $y(t)$  содержит тренд или периодическую составляющую, то, очевидно, имеется зависимость между предыдущими и последующими значениями ряда.

Коэффициент корреляции между уровнями исходного ряда  $y(t)$  и уровнями этого ряда, сдвинутыми на  $s$  шагов  $y(t-s)$  называется коэффициентом автокорреляции временного ряда порядка  $s$  и обозначается  $r(s)$ .

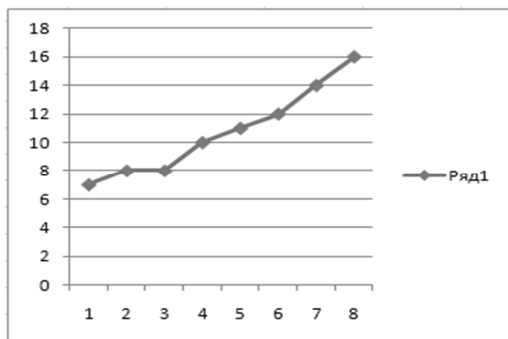
Число  $s$  ещё называют лагом.

Как известно, коэффициент корреляции показывает тесноту линейной связи между случайными величинами.

Поэтому  $r(s)$  показывает наличие (или отсутствие) линейной связи между значениями временного ряда, отстоящими друг от друга на время  $s$ .

**Задача.** Вычислить коэффициенты автокорреляции порядка 1 и 2 для временного ряда

$t$	1	2	3	4	5	6	7	8
$y_t$	7	8	8	10	11	12	14	16



**Решение.** Формула для выборочного коэффициента корреляции случайных величин  $x$  и  $y$  имеет вид

$$r_{xy} = \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{\sqrt{\overline{x^2} - (\bar{x})^2} \cdot \sqrt{\overline{y^2} - (\bar{y})^2}}$$

Найдём сначала  $r(1) = r_{y_t y_{t-1}}$ .

$t$	1	2	3	4	5	6	7	8
$y_t$	7	8	8	10	11	12	14	16
$y_{t-1}$	-	7	8	8	10	11	12	14

$$\bar{y}_t = \frac{8 + \dots + 16}{7}, \quad \bar{y}_{t-1} = \frac{7 + \dots + 14}{7},$$

$$\overline{y_t^2} = \frac{8^2 + \dots + 16^2}{7},$$

$$\overline{y_t y_{t-1}} = \frac{8 \cdot 7 + \dots + 16 \cdot 14}{7}$$

$$r(1) = r_{y_t y_{t-1}} = 0,976$$

Аналогично находим  $r(2) = r_{y_t y_{t-2}}$

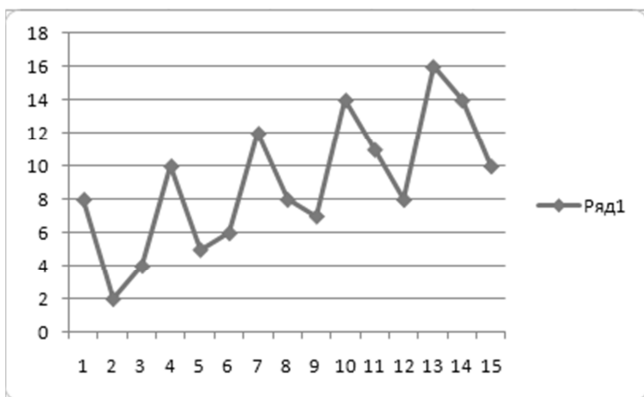
$t$	1	2	3	4	5	6	7	8
$y_t$	7	8	8	10	11	12	14	16
$y_{t-2}$	-	-	7	8	8	10	11	12

$$r(2) = r_{y_t y_{t-2}} = 0,973$$

- ✓ Большое по модулю значение коэффициента автокорреляции с единичным лагом  $r(1)$  говорит о наличии линейного тренда.
- ✓ Если заметно большим является коэффициент корреляции с лагом  $s$ , то это говорит о возможном наличии циклической компоненты периода  $s$
- ✓ Для более точного определения периода циклической компоненты следует перед вычислением *коррелограммы* удалить тренд.

График зависимости  $r(s)$  называется  
***коррелограммой***

**Задача.** Рассчитать коррелограмму для временного ряда до порядка 8 включительно.



t	y <sub>t</sub>
1	8
2	2
3	4
4	10
5	5
6	6
7	12
8	8
9	7
10	14
11	11
12	8
13	16
14	14
15	10

Аналогично предыдущему находим автокорреляции с лагом 1,2,3,...,8. Получаем

### Автокорреляционная функция для $y$

Лаг	ACF		PACF	
1	0,3000		0,3000	
2	0,0091		-0,0889	
3	0,5182	***	0,6002	***
4	0,1000		-0,4751	*
5	-0,2545		0,0720	
6	0,1273		-0,1892	
7	-0,0500		-0,0258	
8	-0,3818		-0,1456	

Столбец ACF содержит искомые коэффициенты автокорреляции. Видно, что максимальная по модулю автокорреляция (помечена звёздочками) имеет лаг 3, что на диаграмме примерно соответствует периоду циклической компоненты.

Более точно оценить зависимость уровней ряда с лагом позволяет коэффициент частной корреляции (последний столбец – PACF).

### **8.8. Стационарные и нестационарные временные ряды. Стационарность в узком и широком смысле**

Остаток временного ряда, освобождённый от периодической составляющей и тренда, должен быть в каком-то смысле стационарным.

Ряд  $y_t$  называется **стационарным в широком смысле**, если его среднее значение и  $\text{cov}(y_t, y_{t+k})$  (в частности, дисперсия) не зависят от  $t$ .

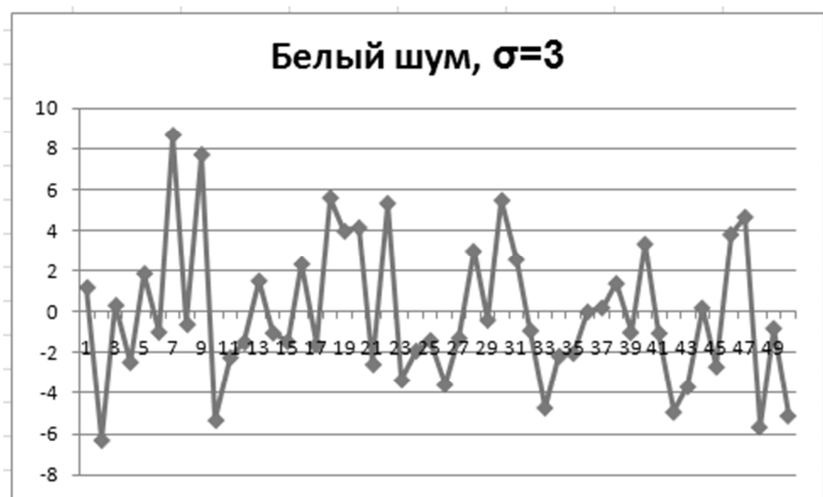
На интуитивном уровне стационарность временного ряда требует, чтобы он имел постоянное среднее значение и колебался вокруг этого среднего с постоянной дисперсией.

В частности, для стационарного ряда недопустимы тренд и циклическая компонента!

Ряд  $y(t)$  называется **стационарным в узком смысле**, если для любого числа  $T$  совместное распределение вероятностей  $m$  наблюдений  $y(t_1), y(t_2), \dots, y(t_m)$  такое же, как и для  $m$  наблюдений  $y(t_1 + T), y(t_2 + T), \dots, y(t_m + T)$ .

Другими словами, вероятностные свойства **стационарного в узком смысле** временного ряда **не меняются с течением времени**.

### 8.9. Пример: Белый шум



Белый шум (БШ), очевидно, является стационарным временным рядом.

### 8.10. Пример: Случайное блуждание

$$y_t = y_{t-1} + \varepsilon_t, \quad \varepsilon_t - \text{БШ}$$



Каждый шаг случаен



Случайное блуждание не является стационарным временным рядом!



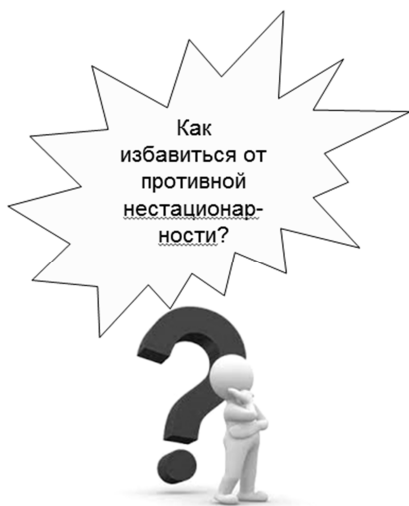
### 8.11. *Переход к разностям. Интегрированный ряд*

**Замечание.** Таким образом, мы видели два разных типа нестационарных временных рядов. Первый – ряд с трендом или циклической компонентой, например

$$y_t = \beta_0 + \beta_1 \cdot t + \varepsilon_t$$

Второй – случайное блуждание

$$y_t = y_{t-1} + \varepsilon_t$$



В первом случае можно просто удалить тренд, но во втором этот способ ничего не даст.

Есть другой метод – переход к первым разностям. По исходному ряду строится новый, состоящий из *первых разностей*:

$$\Delta y_t = y_t - y_{t-1}$$

В случае ряда с линейным трендом после взятия разности получается стационарный ряд:

$$\begin{aligned} \Delta y_t = y_t - y_{t-1} &= \beta_0 + \beta_1 \cdot t + \varepsilon_t - (\beta_0 + \beta_1 \cdot (t-1) + \varepsilon_{t-1}) = \\ &= \beta_1 + \varepsilon_t - \varepsilon_{t-1} \end{aligned}$$

Действительно, например

$$M(\Delta y_t) = M\beta_1 + M\varepsilon_t - M\varepsilon_{t-1} = \beta_1 + 0 - 0 = \beta_1$$

$$D(\Delta y_t) = D\beta_1 + D\varepsilon_t + D\varepsilon_{t-1} = 0 + \sigma^2 + \sigma^2 = 2\sigma^2$$

Если тренд представляется многочленом *второй* степени, то стационарным будет ряд, полученный *двукратным* взятием конечных разностей

$$y_t = \beta_0 + \beta_1 \cdot t + \beta_1 \cdot t^2 + \varepsilon_t$$



$$\begin{aligned} \Delta y_t &= \beta_0 + \beta_1 \cdot t + \beta_1 \cdot t^2 + \varepsilon_t - \\ &- (\beta_0 + \beta_1 \cdot (t-1) + \beta_2 \cdot (t-1)^2 + \varepsilon_{t-1}) = \\ &= \beta_1 + 2\beta_2 t + \varepsilon_t - \varepsilon_{t-1}, \\ \Delta^2 y_t &= \Delta y_t - \Delta y_{t-1} = 2\beta_2 + \varepsilon_t - 2\varepsilon_{t-1} + \varepsilon_{t-2} \end{aligned}$$

Если тренд показательный, то можно попробовать применить конечные разности к логарифму уровней ряда. Для случайного блуждания взятие первых разностей сразу приводит к стационарному ряду

$$\Delta y_t = y_t - y_{t-1} = \varepsilon_t$$

А вот для такого ряда  $y_t = 2y_{t-1} + \varepsilon_t$  взятие первых разностей не помогает совершенно.

Ряд  $y_t$  называется **интегрированным порядком k**, если его k-е разности представляют собой **стационарный ряд**.

## 8.12. Проверка гипотез о стационарности

### 8.12.1. Графический метод

Первым этапом при проверке стационарности является графический анализ. Наличие *видимого тренда* или *систематического изменения амплитуды колебаний* говорит о *нестационарности*.

*Корреллограмма* также иногда позволяет обнаружить нестационарность.

### 8.12.2. Критерий Кокса-Стюарта

Это один из быстрых критериев наличия тренда среднего и дисперсии.

1) Для проверки тренда среднего значения вычисляется

$$S_1 = \sum_{i=1}^{\lfloor n/2 \rfloor} (n - 2i + 1) \cdot h_{i, n-i+1}, \quad h_{i,j} = \begin{cases} 1, & \text{если } x_i > x_j \\ 0, & \text{если } x_i \leq x_j \end{cases}$$

Если

$$\frac{\left| S_1 - \frac{n^2}{8} \right|}{\sqrt{\frac{n^3 - n}{24}}} < u_{1-\alpha/2}$$

то гипотеза о наличии тренда среднего значения отклоняется

2) Для проверки тренда дисперсии выборка разбивается части по  $k$  элементов в каждой (если  $n$  не делится на  $k$ , то отбрасываются наблюдения в середине).

Для каждой части находят размах  $\omega_i$ . Далее полученные числа проверяются на тренд критерием  $S_1$ .

Рекомендуется выбирать  $k$  из соотношений.

$$k = \begin{cases} 5, & \text{если } n \geq 90 \\ 4, & \text{если } 64 \leq n < 90 \\ 3, & \text{если } 48 \leq n < 64 \\ 2, & \text{если } n < 48 \end{cases}$$

**Задача.** Проверить гипотезу о стационарности при  $\alpha=0,01$ .

**Решение.** 1)  $n=15$ ,  $[n/2]=7$ .

$$S_1 = \sum_{i=1}^7 (15-2i+1) \cdot h_{i,15-i+1}, \quad h_{i,j} = \begin{cases} 1, & \text{если } x_i > x_j \\ 0, & \text{если } x_i \leq x_j \end{cases}$$

$$h_{1,15} = 1, h_{2,14} = 0, h_{3,13} = 0, h_{4,12} = 1, h_{5,11} = 0,$$

$$h_{6,10} = 0, h_{7,9} = 1$$

$$S_1 = 14 \cdot 1 + 12 \cdot 0 + 10 \cdot 0 + 8 \cdot 1 + \dots + 2 \cdot 1 = 24$$

$$\frac{\left| S_1 - \frac{n^2}{8} \right|}{\sqrt{\frac{n^3 - n}{24}}} = 1,4258 < u_{1-\alpha/2} = u_{0,995} = 2,58$$

Гипотеза о наличии тренда среднего не подтверждается

t	y <sub>t</sub>
1	11
2	6
3	4
4	10
5	5
6	6
7	12
8	8
9	7
10	14
11	11
12	8
13	16
14	14
15	10

2)  $k=2$ . Разобьём на группы по 2.  
 $\omega_1 = 11 - 6 = 5$ ,  $\omega_2 = 6$ ,  $\omega_3 = 1$ ,  $\omega_4 = 1$ ,  
 $\omega_5 = 3$ ,  $\omega_6 = 8$ ,  $\omega_7 = 4$   
 $n=7$ ,  $[n/2]=3$ .  
 $h_{1,7} = 1$ ,  $h_{2,6} = 0$ ,  $h_{3,5} = 0$   
 $S_1 = 6 \cdot 1 + 4 \cdot 0 + 2 \cdot 0 = 6$

$$\frac{\left| S_1 - \frac{n^2}{8} \right|}{\sqrt{\frac{n^3 - n}{24}}} = 0,034 < u_{1-\alpha/2} = u_{0,995} = 2,58$$

Гипотеза о наличии тренда дисперсии не подтверждается.

t	y <sub>t</sub>
1	11
2	6
3	4
4	10
5	5
6	6
7	12
8	8
9	7
10	14
11	11
12	8
13	16
14	14
15	10

### 8.12.3. Критерий Фостера-Стюарта

Ещё один критерий наличия тренда среднего и дисперсии.

Вычисляются

$$S_{>} = \sum_{i=2}^n u_i, \quad u_i = \begin{cases} 1, & \text{если } x_i > \text{всех предыдущих} \\ 0, & \text{иначе} \end{cases}$$

$$S_{<} = \sum_{i=2}^n l_i, \quad l_i = \begin{cases} 1, & \text{если } x_i < \text{всех предыдущих} \\ 0, & \text{иначе} \end{cases}$$

$$t_1 = \frac{S_> + S_< - f^2}{l}, \quad t_2 = \frac{S_> - S_<}{f},$$

$$l \approx \sqrt{2 \ln n - 3,4253}, \quad f \approx \sqrt{2 \ln n - 0,8456},$$

Если  $|t_1| < t_{1-\alpha/2}(n)$ , то гипотеза о наличии тренда среднего отклоняется

Если  $|t_2| < t_{1-\alpha/2}(n)$ , то гипотеза о наличии тренда дисперсии отклоняется

**Задача.** Проверить гипотезу о стационарности при  $\alpha=0,01$ .

**Решение.**  $n=15$

$$S_> = \sum_{i=2}^n u_i = 3, \quad S_< = \sum_{i=2}^n l_i = 2$$

$$l \approx \sqrt{2 \ln 15 - 3,4253} \approx 2,14,$$

$$f \approx \sqrt{2 \ln 15 - 0,8456} \approx 1,5,$$

$$t_1 = \frac{3+2-1,5^2}{2,14} \approx 1,28, \quad t_2 = \frac{3-2}{1,5} \approx 0,67$$

$$t_{1-\alpha/2}(n) = t_{0,995}(15) \approx 2,94$$

$$|t_1| < t_{0,95}(15)$$

$$|t_2| < t_{0,95}(15)$$

Гипотеза о наличии трендов отвергается.

Стационарность возможна.

t	y <sub>t</sub>	u	l
1	11	—	—
2	6	0	1
3	4	0	1
4	10	0	0
5	5	0	0
6	6	0	0
7	12	1	0
8	8	0	0
9	7	0	0
10	14	1	0
11	11	0	0
12	8	0	0
13	16	1	0
14	14	0	0
15	10	0	0

t	y <sub>t</sub>
1	11
2	6
3	4
4	10
5	5
6	6
7	12
8	8
9	7
10	14
11	11
12	8
13	16
14	14
15	10

### 8.12.4. Критерий Льюнга-Бокса

Для белого шума коэффициенты корреляции с любым лагом должны равняться нулю. Отсюда получается такой критерий.

$H_0$ : временной ряд представляет собой белый шум.

Вычисляется

$$Q = n(n+2) \cdot \sum_{i=1}^p \frac{r^2(i)}{n-i}$$

где  $p$  — до какого лага проверяется автокорреляция.

Если  $Q > \chi_{1-\alpha}^2(P)$ , то гипотеза  $H_0$  отвергается, ряд не является белым шумом

Ещё один критерий стационарности – критерий Дики-Фуллера – будет позже.

### 8.13. Линейная регрессия для временных рядов

Рассмотрим уравнения регрессии в случае, когда переменные представляют собой временные ряды.

$$y_t = \beta_0 + \beta_1 x_{1t} + \beta_2 x_{2t} + \dots + \beta_m x_{mt} + \varepsilon_t,$$

В этом случае переменные в правой части являются случайными величинами (а не константами), что может сказаться на качестве оценок по МНК.



Вспомним основные особенности нахождения оценок параметров в этой ситуации. На следующие положения придётся не раз ссылаться.

- 1) Если переменные  $X_{it}$  в правой части **некоррелированы** со случайным фактором  $\varepsilon_{it}$ , то оценки параметров по МНК **несмещённые и состоятельные**.
- 2) Если переменные в правой части **коррелированы со случайным фактором, но связи при одинаковом  $t$  нет**, то оценки параметров по МНК **смещённые** (это плохо), но **состоятельные** (стремятся к правильным значениям при увеличении выборки).
- 3) Если переменные в правой части **коррелированы со случайным фактором даже при одинаковом  $t$** , то оценки параметров по МНК могут быть и **смещёнными**, и **несостоятельными**, то есть никуда не годными.

## 8.14. Модель авторегрессии AR(p)

### 8.14.1. Определение

Моделью авторегрессии AR(p) называется временной ряд, удовлетворяющий соотношению

$$y_t = c + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \phi_p y_{t-p} + \varepsilon_t,$$

где случайная составляющая  $\varepsilon_t$  – “белый шум”.

В частности, AR(1) – это временной ряд, удовлетворяющий соотношению

$$y_t = c + \phi \cdot y_{t-1} + \varepsilon_t$$

**Пример.** Случайное блуждание является процессом вида AR(1) с  $\alpha=1$  и  $c=0$ :

$$y_t = y_{t-1} + \varepsilon_t, \quad \varepsilon_t - \text{БШ}$$

### 8.14.2. Условия стационарности ряда AR(p)

Чтобы ряд AR(p) ,был стационарен, необходимо, чтобы все корни уравнения

$$Z^p = \phi_1 Z^{p-1} + \dots + \phi_p Z$$

по модулю были меньше 1



1) В частности, для AR(1)

$$y_t = c + \phi y_{t-1} + \varepsilon_t$$

$$Z^1 = \phi Z^0 \Leftrightarrow Z = \phi$$

так что,

Ряд AR(1) стационарен,  
если и только если  
 $|\phi| < 1$

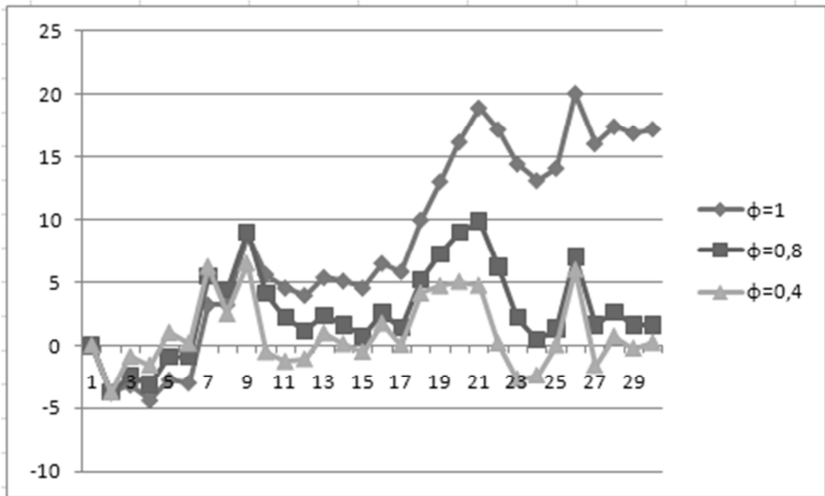
2) Для AR(2)

$$y_t = c + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \varepsilon_t$$

$$Z^2 = \phi_1 Z + \phi_2$$

Для стационарности ряда AR(2) корни этого уравнения должны по модулю быть меньше 1

Примеры реализации временных рядов AR(1) при различных  $\phi$ :



### 8.14.3. Тест Дики-Фуллера (Dickey-Fuller)

Перепишем уравнение AR(1) в виде

$$\Delta y_t = c + \rho y_{t-1} + \varepsilon_t, \rho = \phi - 1$$

Если  $\rho=0$ , что соответствует  $\phi = 1$ , то исходный ряд не стационарен.

Если  $\rho < 0$ , то  $\phi < 1$  и исходный ряд стационарен.

**Тест Дики-Фуллера** проверяет гипотезу  $H_0: \rho=0$  при альтернативе  $H_1: \rho < 0$ .

$$\text{Если } \tau = \frac{\hat{\rho}}{S_{\rho}} > \tau_{\text{крит}},$$

где  $\tau_{\text{крит}}$  – табличное значение, то гипотеза  $H_0$  принимается (ряд нестационарный), в противном случае отвергается (ряд считается стационарным).

Для использования теста Дики-Фуллера требуются соответствующие табличные значения. Они были получены эмпирически, с использованием метода Монте-Карло. За основу был взят процесс AR(1).

$\alpha \setminus n$	25	50	100	$\infty$
0,01	-3,75	-3,58	-3,51	-3,43
0,05	-3,33	-3,22	-3,17	-3,12
0,1	-3	-2,93	-2,89	-2,86

**Пример.** Проверить гипотезу о стационарности при помощи теста DF.  $\alpha=0,1$

**Решение.** Заполним таблицу.

t	$y_t$	$\Delta_t$	$y_{t-1}$
1	-7,1	—	—
2	6,0	13,1	-7,1
3	-7,2	-13,2	6,0
4	7,7	14,9	-7,2
5	-7,1	-14,8	7,7
6	14,4	21,5	-7,1
7	-12,1	-26,4	14,4
8	17,4	29,5	-12,1
9	-19,2	-36,6	17,4
10	13,2	32,4	-19,2
11	-12,0	-25,2	13,2
12	11,1	23,1	-12,0
13	-9,9	-21,0	11,1
14	6,5	16,4	-9,9
15	-2,8	-9,3	6,5

Построим регрессию обычным методом МНК

$$\Delta y_t = c + \rho y_{t-1} + \varepsilon_t$$

	Коэффициенты	Стандартная ошибка	t-статистика
Y-пересечен	0,517	0,939	0,551002
Переменная	-1,947	0,082	-23,8356

$$\Delta y_t = 0,517 - 1,947 y_{t-1}$$

$$\frac{\hat{\rho}}{S_{\rho}} = \frac{-1,947}{0,082} = -23,8356$$

Так как -23,836 гораздо меньше критического значения из таблицы, тест DF стационарности не отвергает.

t	$y_t$
1	-7,1
2	6,0
3	-7,2
4	7,7
5	-7,1
6	14,4
7	-12,1
8	17,4
9	-19,2
10	13,2
11	-12,0
12	11,1
13	-9,9
14	6,5
15	-2,8

**Модифицированный тест Дики-Фуллера** рассматривает авторегрессионные процессы более высокого порядка  $A(p)$ .

### 8.14.4. Оценка параметров модели AR(p)

Рассмотрим модель AR(p)

$$y_t = c + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \varepsilon_t,$$

Поскольку в правой части стоят лаговые значения переменной  $y_t$ , а они сами выражаются в виде такого же уравнения в предыдущие моменты времени, например

$$y_{t-1} = c + \phi_1 y_{t-2} + \dots + \varepsilon_{t-1},$$

то не равна нулю корреляция между переменными в правой части и случайным фактором:



$$\text{COV} \left( \begin{pmatrix} y_t \\ y_{t-1} \\ \dots \\ y_1 \end{pmatrix}, \begin{pmatrix} \varepsilon_t \\ \varepsilon_{t-1} \\ \dots \\ \varepsilon_1 \end{pmatrix} \right) \neq 0$$

Но при этом, если  $\varepsilon_t$  не зависит от  $\varepsilon_{t-1}, \varepsilon_{t-2}, \dots, \varepsilon_{t-p}$  (например, если  $\varepsilon_t$  — белый шум), лаговые переменные в предыдущие моменты времени не имеют связи с текущим  $\varepsilon_t$

$$y_t = c + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \phi_p y_{t-p} + \varepsilon_t,$$

Таким образом, если оценивать коэффициенты AR(p) по МНК в случае, когда  $\varepsilon_t$  не зависит от  $\varepsilon_{t-1}, \dots$ , то оценки будут смещённые (это плохо), но состоятельные.

Если же  $\varepsilon_t$  зависит от  $\varepsilon_{t-1}, \dots$ , то оценивать коэффициенты AR(p) по МНК нельзя, они получатся и смещённые и несостоятельные. В этом случае нужны другие методы.

### 8.14.5. Авторегрессия в остатках

Таким образом, для модели AR(p) чрезвычайно важно установить, зависит ли  $\varepsilon_t$  от своих значений в предыдущие моменты времени.

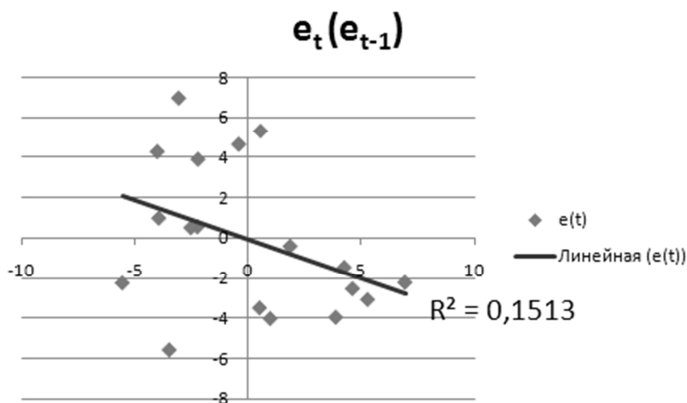
Например, наличие авторегрессии первого порядка

$$\varepsilon_t = \rho \varepsilon_{t-1} + u_t, \quad u_t - \text{БШ},$$

(или более высокого порядка) приводит к невозможности использования МНК.

Итак, при оценке модели  $AR(p)$  необходимо тщательно проверить по остаткам, является ли случайный фактор белым шумом, в частности, не имеет ли место авторегрессия случайного фактора.

Один из “быстрых” способов – построить точечный график зависимости остатков  $e_t$  от  $e_{t-1}$ . Не должно быть заметного тренда.



Если возникает подозрение на авторегрессию в остатках, следует проверить соответствующую гипотезу.

Для обнаружения авторегрессии первого порядка обычно применяют критерий Дарбина-Уотсона, но для моделей  $AR(p)$  правильнее использовать так называемый ***h-критерий Дарбина***.

Также имеет смысл использовать критерий Льюнга-Бокса, так как “правильные” остатки должны быть устроены как белый шум.

## 8.15. Модель скользящего среднего $MA(q)$

### 8.15.1. Определение.

Моделью скользящего среднего  $MA(q)$  называется временной ряд, удовлетворяющий соотношению

$$y_t = \varepsilon_t + \theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + \dots + \theta_q \varepsilon_{t-q},$$

где случайная составляющая  $\varepsilon_t$  – “белый шум”.

В частности,  $MA(1)$  – это временной ряд, удовлетворяющий соотношению

$$y_t = \varepsilon + \theta \cdot \varepsilon_{t-1}$$

### 8.15.2. Условия стационарности ряда $MA(q)$

Ряд  $MA(q)$  стационарен

В частности,

$$My_t = M\varepsilon_t + M\theta_1 \varepsilon_{t-1} + \dots + M\theta_q \varepsilon_{t-q} = 0,$$

$$Dy_t = D\varepsilon_t + D\theta_1 \varepsilon_{t-1} + \dots + D\theta_q \varepsilon_{t-q} = (1 + \theta_1^2 + \dots + \theta_q^2) \sigma^2$$

### 8.15.3. Оценка параметров $MA(q)$

В  $MA$  невозможно выразить остатки через значения  $Y_t$  (ведь  $\varepsilon_t$  нам неизвестны). Поэтому обычный МНК здесь не работает. Можно применить, например, **нелинейный МНК**.

## 8.16. Модель $ARMA(p,q)$

Эта модель включает в себя две предыдущие.

**Моделью ARMA(p,q)** называется временной ряд, удовлетворяющий соотношению

$$y_t = c + \phi_1 y_{t-1} + \dots + \phi_p y_{t-p} + \varepsilon_t + \theta_1 \varepsilon_{t-1} + \dots + \theta_q \varepsilon_{t-q},$$

где случайная составляющая  $\varepsilon_t$  – “белый шум”.

Если  $q=0$ , то получается процесс AR(p), если  $p=0$ , то получается процесс MA(q)

### 8.17. Модель ARIMA(p,k,q)

**ARIMA** — модель и методология анализа временных рядов, иногда называемая моделью Бокса-Дженкинса. Является расширением моделей ARMA для нестационарных временных рядов, которые можно сделать стационарными взятием разностей некоторого порядка.

**Моделью ARIMA(p,k,q)** называется временной ряд, для которого разности порядка  $k$  образуют ряд ARMA(p,q).

В подходе Бокса-Дженкинса исходный ряд заменяют конечными разностями, пока не получится стационарный ряд, а затем оценивают, как ряд ARMA.

## 9. Взаимосвязь временных рядов

### 9.1. Специфика оценки взаимосвязи временных рядов

Важную роль играет изучение взаимосвязи двух или большего числа временных рядов.

Оценка зависимостей традиционными методами может приводить к неправильным результатам из-за наличия в рядах тренда или периодической составляющей.

Так, большой коэффициент корреляции двух временных рядов может быть результатом того, что оба содержат ярко выраженный временной тренд.

Чтобы получить модель, характеризующие истинную причинно-следственную связь между изучаемыми рядами, следует избавиться от этой ложной корреляции.

Поэтому периодические составляющие и тренд обычно стремятся исключить перед проведением дальнейших исследований.

Для этой цели применяются следующие методы:

- Метод отклонений от трендов;
- Метод последовательных разностей;
- Включение фактора времени в модель регрессии.

### **9.2.        *Метод отклонений от трендов***

Для исходных временных рядов  $y_t$  и  $x_t$ , каждый из которых содержит трендовую компоненту, определяются расчетные по тренду уровни.

Затем из фактических уровней каждого ряда вычитают расчетные значения:

$$y'_t = y_t - \hat{y}_t, \quad x'_t = x_t - \hat{x}_t$$

Дальнейший анализ проводят с использованием не исходных уровней, а отклонений от трендов.

### **9.3.        *Метод последовательных разностей***

Если временной ряд содержит ярко выраженную линейную тенденцию, ее можно устранить путем замены исходных уровней ряда первыми разностями:



$$\Delta y_t = y_t - y_{t-1}, \quad \Delta x_t = x_t - x_{t-1}$$

Для тренда в виде многочлена более высокой степени применяют разности более высокого порядка.

#### 9.4. Включение фактора времени

Еще один путь учёта влияния фактора времени на причинно-следственную связь рядов – включение в модель фактора времени в качестве независимой переменной:

$$y_t = a + bx_t + ct + \varepsilon_t,$$

#### 9.5. Коинтеграция временных рядов

Общий недостаток методов исключения тренда заключается в том, что происходит модификация исходной модели

$$y_t = a + b \cdot x_t + e_t$$

Но большая часть соотношений экономической теории сформулирована для исходной модели. А в экономике временные ряды, как правило, содержат тренды.

Другой подход основывается на изучении таких экономических данных, которые, будучи нестационарными, могут быть скомбинированы в один ряд, который будет уже стационарным. Ряды, обладающие такой особенностью, называются **коинтегрированными** рядами.

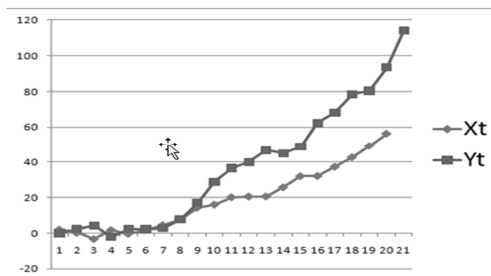
Напомним, что

Временной ряд  $Y_t$  ряд называется **интегрированным порядка 1**, (обозначение **I(1)**), если его первая разность  $\Delta Y_t$  стационарна.

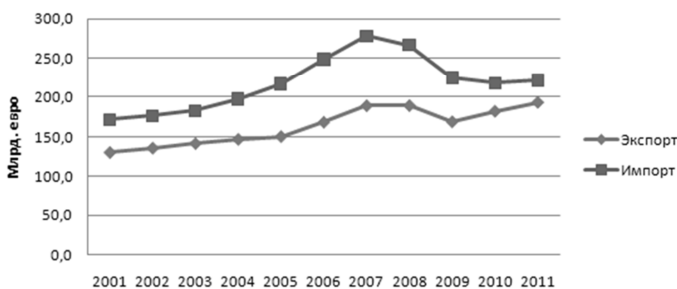
Два временных ряда, интегрированные порядка 1, называются **коинтегрированными**, если некоторая их линейная комбинация есть стационарный временной ряд.

**Пример.**

Здесь ряд  $Y_t - 2X_t$  стационарен.



**Пример** предположительно коинтегрированных рядов.  
**Объёмы экспорта и импорта Испании**



## Критерий Энгеля-Гранжера (Engel-Granger)

1. Выдвигается гипотеза  $H_0$  об отсутствии коинтеграции между рядами  $X_t$  и  $Y_t$ .
2. Строится уравнение регрессии  $Y_t$  на  $X_t$  и вычисляются остатки  $e_t$ . Нужно проверить стационарность этого ряда.
3. Рассчитывают параметры уравнения регрессии

$$\Delta e_t = a + \rho \cdot e_{t-1},$$

где  $\Delta e_t$  - первые разности остатков.

4. Если

$$\tau = \frac{\hat{\rho}}{s_{\rho}} < \tau_{крит} ,$$

то принимается гипотеза о коинтеграции.



$\alpha$	$\tau_{крит}$
0,01	-3,9
0,05	-3,34
0,1	-3,04

### 9.6. Модель распределённых лагов DL(q)

Моделью распределённых лагов DL(q) называется связь двух стационарных временных рядов вида

$$y_t = c + \beta_0 x_t + \beta_1 x_{t-1} + \dots + \beta_q x_{t-q} + \varepsilon_t,$$

где случайная составляющая  $\varepsilon_t$  – белый шум.

В частности, DL(1) – это связь между временными рядами вида

$$y_t = c + \beta_0 x_t + \beta_1 x_{t-1} + \varepsilon_t,$$

Оценивание коэффициентов  $c, \beta_j$  (в отличие от модели MA(q)), здесь, может, быть произведено обычным МНК. Но может оказаться, что количество неизвестных параметров слишком велико.

Идея состоит в том, чтобы, отталкиваясь от смысла коэффициентов  $\beta_j$ , выразить их через небольшое количество параметров.

### 9.6.1. Модель полиномиальных лагов Алмон

В этой модели предполагают, что зависимость  $\beta_j$  от  $j$  является полиномом некоторой степени  $r < q$ .

$$\beta_j = \gamma_0 + \gamma_1 j + \gamma_2 j^2 + \dots + \gamma_r j^r$$

**Пример.** Построить модель с распределёнными лагами

$$y_t = c + \beta_0 x_t + \beta_1 x_{t-1} + \beta_2 x_{t-2} + \beta_3 x_{t-3} + \varepsilon_t,$$

используя полиномиальные лаги степени 2.

**Решение.**

$$\beta_j = \gamma_0 + \gamma_1 j + \gamma_2 j^2$$

t	$y_t$	$x_t$
1	3,0	3,0
2	5,0	2,0
3	2,0	2,0
4	3,0	3,0
5	2,7	4,0
6	3,6	5,0
7	4,6	4,0
8	4,7	2,0
9	3,3	5,0
10	3,8	6,0
11	5,8	4,0
12	5,1	3,0

1) Заведём в MS Excel ячейки для  $C, \gamma_0, \gamma_1, \gamma_2$  и запишем туда для начала нули.

К	L
C	0
$\gamma_0$	0
$\gamma_1$	0
$\gamma_2$	0

2) Заведём в MS Excel ячейки для  $\beta_0, \beta_1, \beta_2, \beta_3$  и запишем туда соответствующие формулы.

N	O	P	Q	R
$\beta_0$	0			
$\beta_1$	0			
$\beta_2$				
$\beta_3$		=L2+3*L3+9*L4		

3) Заполним столбец  $Y_{расч}$  по формуле

$$y_{t,расч} = c + \beta_0 x_t + \beta_1 x_{t-1} + \beta_2 x_{t-2} + \beta_3 x_{t-3},$$

D	E	F	G	H	I	J	K	L	M	N	O
$Y_{расч}$	(y-урасч)^2			ESS	157		C	0		$\beta_0$	0
---	---						$\gamma_0$	0		$\beta_1$	0
---	---						$\gamma_1$	0		$\beta_2$	0
$=\$L\$1+\$O\$1*C5+\$O\$2*C4+\$O\$3*C3+\$O\$4*C2$									$\beta_3$	0	

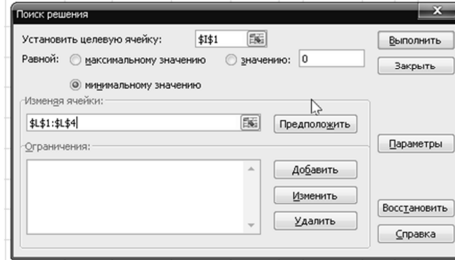
4) Сформируем ячейку с ESS

E	F	G	H	I	J	K
урасч)^2			ESS	=СУММ(E5:E13)		

B	C	D	E
$y_t$	$x_t$	Урасч	$(y - \text{Урасч})^2$
3	3	---	
5	2	---	
2	2	---	
3	3	0	$=(B5-D5)^2$
2,7	4	0	7,29
3,6	5	0	12,96
4,6	4	0	21,16
4,7	2	0	22,09

5) Поиск решения

G	H	I	J	K	L	M	N	O	P
	ESS	157		C	0		$\beta_0$	0	
				$\gamma_0$	0		$\beta_1$	0	
				$\gamma_1$	0		$\beta_2$	0	
				$\gamma_2$	0		$\beta_3$	0	



5) Ответ

ESS	0,4	C	1,17	$\beta_0$	-0,06
		$\gamma_0$	-0,06	$\beta_1$	0,416
		$\gamma_1$	0,712	$\beta_2$	0,432
		$\gamma_2$	-0,23	$\beta_3$	-0,02

$$y_t = 1,17 - 0,06x_t + 0,416x_{t-1} + 0,432x_{t-2} - 0,02x_{t-3} + \varepsilon_t,$$

## 9.6.2. Модель геометрических лагов Койка

В этой модели предполагают, что зависимость  $\beta_j$  от  $j$  имеет вид

$$\beta_j = \beta_0 \cdot \lambda^j, \quad j = 0, 1, 2, \dots, \quad 0 < \lambda < 1$$

Здесь всего два неизвестных параметра,  $\beta_0$  и  $\lambda$ , но оценить их нелегко, так как зависимость от них нелинейная.

Один из подходов – нелинейный МНК (как в предыдущем примере).

### 9.7. Модель $ADL(p,q)$

Эта модель содержит как частные случаи модели  $AR(p)$  и  $DL(q)$

Моделью  $ADL(p,q)$  называется связь двух стационарных временных рядов вида

$$y_t = c + \alpha_1 y_{t-1} + \dots + \alpha_p y_{t-p} + \beta_0 x_t + \beta_1 x_{t-1} + \dots + \beta_q x_{t-q} + \varepsilon_t,$$

где  $\varepsilon_t$  – белый шум.

### 9.8. Оценка параметров модели $ADL(p,q)$

В случае  $p=0$  получается модель  $DL(q)$ , которую, как говорилось ранее, можно оценивать по МНК.

Если  $p>0$ , то в правой части равенства есть переменная  $y_{t-1}$ , которая выражается через  $\varepsilon_{t-1}$ . Поэтому имеется корреляция между матрицами

$$\begin{pmatrix} 1 & y_{t-1} & \dots & x_{t-q} \\ 1 & y_{t-2} & \dots & x_{t-q-1} \\ \dots & \dots & \dots & \dots \end{pmatrix} \text{ и } \begin{pmatrix} \varepsilon_t \\ \varepsilon_{t-1} \\ \dots \end{pmatrix},$$

так что оценки по МНК будут смещёнными.

Можно использовать нелинейный МНК.

Ещё один из методов оценивания (метод инструментальных переменных), состоит в замене переменных  $Y$  в правой части на их оценки в регрессии на переменные  $X$ :

$$\hat{Y} = d + e \cdot X$$

**Пример.** Построить модель ADL(1,0)

$$y_t = c + \alpha_1 y_{t-1} + \beta_0 x_t + \varepsilon_t,$$

**Решение.** Решим задачу при помощи инструментальной переменной

1 шаг. Построим регрессию

$$y_t = d + ex_t$$

Обычным образом находим уравнение

$$y_t = 11,9029 - 1,2621x_t$$

Из него вычисляем значения переменной  $Z$ , которая будет заменять  $Y$  в нужном нам уравнении

$$z_1 = 11,9029 - 1,2621 \cdot 3 = 8,12$$

$$z_2 = 11,9029 - 1,2621 \cdot 4 = 6,85, \dots$$

t	$y_t$	$x_t$
1	9	3
2	7	4
3	6	5
4	6	4
5	5	6
6	6	4
7	4	6
8	3	7

t	$y_t$	$x_t$	$z_t$
1	9	3	8,12
2	7	4	6,85
3	6	5	5,59
4	6	4	6,85
5	5	6	4,33
6	6	4	6,85
7	4	6	4,33
8	3	7	3,07

2 шаг. Построим регрессию

t	$y_t$	$x_t$	$z_{t-1}$
1	9	3	-
2	7	4	8,12
3	6	5	6,85
4	6	4	5,59
5	5	6	6,85
6	6	4	4,33
7	4	6	6,85
8	3	7	4,33

$$y_t = c + \alpha_1 z_{t-1} + \beta_0 x_t + \varepsilon_t,$$

$$y_t = 8,67 + 0,28z_{t-1} - 0,99x_t$$

Отсюда находим уравнение модели ADL(1,0):

$$y_t = 8,67 + 0,28y_{t-1} - 0,99x_t + \varepsilon_t$$

### 9.9. Тест Гранжера (Granger) на причинно-следственную зависимость

Нередко возникает вопрос о причинно-следственной связи между факторами. Например, верно ли, что увеличение цен на нефть влечёт за собой инфляцию?

Идея теста такова: если  $x$  является причиной  $y$ , то изменения  $x$  предшествуют изменениям  $y$ , но не наоборот.

Другими словами, регрессия с лагами  $y$  на  $x$  должна быть значимой, а регрессия с лагами  $x$  на  $y$  – незначимой. Если же обе регрессии значимы, то, скорее всего, есть другой фактор, влияющий и на  $x$  и на  $y$ .

Чтобы тестировать гипотезу  $H_0$ : “ $X$  не влияет на  $Y$ ”, рассматривается ADL модель вида

$$y_t = c + \alpha_1 y_{t-1} + \dots + \alpha_p y_{t-p} + \beta_1 x_{t-1} + \dots + \beta_q x_{t-q} + \varepsilon_t,$$

Гипотеза  $H_0$  записывается в виде  $\beta_1 = \beta_2 = \dots = \beta_q = 0$ .

Проверяется она обычным F-тестом.

Чтобы тестировать гипотезу  $H_0$ : “ $Y$  не влияет на  $X$ ”, рассматривается аналогичная модель



$$x_t = c + \alpha_1 x_{t-1} + \dots + \alpha_p x_{t-p} + \beta_1 y_{t-1} + \dots + \beta_q y_{t-q} + \varepsilon_t,$$

Если первая гипотеза  $H_0$  отклоняется, а вторая принимается, то, возможно,  $X$  является *причиной* для  $Y$ .

## 10. Системы эконометрических уравнений

### 10.1. Введение

При моделировании часто приходится вводить не одно, а несколько связанных между собой эконометрических уравнений, т.е. описывать модель **системой уравнений**.

**Пример.**

$$\begin{cases} Y_1 = a_0 + a_1 Y_2 + a_2 X_1 + \varepsilon_1 \\ Y_2 = b_0 + b_1 Y_1 + b_2 X_2 + \varepsilon_2 \end{cases}$$

Классическим примером является формирование спроса  $Q^D$  и предложения  $Q^S$  товара в зависимости от его цены  $P$  и доходов  $I$

$$\begin{cases} Q_t^S = a_0 + a_1 P_t + \varepsilon_1 \\ Q_t^D = b_0 + b_1 P_t + b_2 I_t + \varepsilon_2 \end{cases}$$

В ситуации равновесия на рынке  $Q^D = Q^S$  и получается система уравнений

$$\begin{cases} Q_t^D = a_0 + a_1 P_t + \varepsilon_1 \\ Q_t^S = b_0 + b_1 P_t + b_2 I_t + \varepsilon_2 \\ Q_t^D = Q_t^S \end{cases}$$

Уравнения, не содержащие неизвестных коэффициентов (как третье в этом примере), называются *тождествами*.

## 10.2. Структурная и приведённая форма записи системы уравнений

Роль переменных в системе уравнений различна.

**Эндогенные** переменные – это **зависимые** переменные, значения которых определяются **внутри** модели. Как правило, их количество равно числу уравнений.

**Экзогенные** переменные – это **независимые** переменные, значения которых формируются вне модели.

С математической точки зрения их особенность в том, что они не коррелируют со случайным фактором  $\varepsilon$ .

Разделение переменных на эндогенные и экзогенные зависит от теоретической концепции.

Внеэкономические переменные (например, климатические условия) являются экзогенными.

Близкую к экзогенным роль играют значения эндогенных переменных в предыдущие моменты времени (так как они уже сформировались), т.е. лаговые переменные.

Экзогенные и лаговые эндогенные переменные называются *предопределёнными*

Так, потребление текущего года  $Y_t$  зависит не только от одновременных экономических факторов, но и от потребления предыдущего года  $Y_{t-1}$  (и это уже *предопределённая* переменная).

Исходная форма записи системы, когда в правой присутствуют и предопределённые и эндогенные переменные, называется **структурной**.

Как правило, каждое уравнение в левой части содержит свою эндогенную переменную, а в правой – какие-то другие переменные, эндогенные и предопределённые.

$$\begin{cases} y_1 = a_{12}y_2 + \dots + b_{11}x_1 + b_{12}x_2 + \dots, \\ y_2 = a_{21}y_1 + \dots + b_{21}x_1 + b_{22}x_2 + \dots, \\ \dots, \\ y_m = a_{m1}y_1 + \dots + b_{m1}x_1 + b_{m2}x_2 + \dots \end{cases}$$

Из-за наличия в правой части эндогенных переменных для нахождения коэффициентов нельзя применять МНК непосредственно к структурной форме.

С этой целью её преобразуют в приведенную форму.

В **приведённой** форме записи эндогенные переменные выражены через предопределённые, в каждом уравнении в левой части стоит своя эндогенная переменная.

$$\begin{cases} y_1 = \pi_{11}x_1 + \pi_{12}x_2 + \dots, \\ y_2 = \pi_{21}x_1 + \pi_{22}x_2 + \dots, \\ \dots, \\ y_m = \pi_{m1}x_1 + \pi_{m2}x_2 + \dots \end{cases}$$

**Пример.** Рассмотрим модель

$$\begin{cases} Q_t = b_0 + a_1 P_t + \varepsilon_1 \\ Q_t = b_1 + a_2 P_t + b_2 I_t + \varepsilon_2 \end{cases}$$

Здесь  $Q$  – равновесный спрос-предложение (эндогенная переменная),  $P$  – цена (эндогенная переменная),  $I$  – доход (экзогенная переменная).

Это *структурная* форма записи, числа  $a_j, b_j$  – *структурные коэффициенты*.

Чтобы получить приведённую форму записи, выразим  $P, Q$  через  $I$ .

Перенесём все эндогенные переменные в левую часть

$$\begin{cases} Q_t - a_1 P_t = b_0 + \varepsilon_1 \\ Q_t - a_2 P_t = b_1 b_2 I_t + \varepsilon_2 \end{cases}$$

и запишем в матричной форме

$$AY = BX + \varepsilon,$$

$$A = \begin{pmatrix} 1 & -a_1 \\ 1 & -a_2 \end{pmatrix}, Y = \begin{pmatrix} Q_t \\ P_t \end{pmatrix}, B = \begin{pmatrix} b_0 & 0 \\ b_1 & b_2 \end{pmatrix}, X = \begin{pmatrix} 1 \\ I_t \end{pmatrix},$$

Если существует обратная матрица

$$A^{-1} = \begin{pmatrix} 1 & -a_1 \\ 1 & -a_2 \end{pmatrix}^{-1}$$

то, помножив на неё равенство  $AX + BY = \varepsilon$ , получим

$$Y = A^{-1}BX + A^{-1}\varepsilon \Rightarrow$$

$$Y = PX + \xi, \quad \xi = B^{-1}\varepsilon,$$

где

$$\Pi = \begin{pmatrix} \pi_1 & \pi_2 \\ \pi_3 & \pi_4 \end{pmatrix} = A^{-1}B$$

т.е. получаем приведённую форму модели

$$\begin{cases} Q_t = \pi_1 + \pi_2 I_t + \xi_{1t} \\ P_t = \pi_3 + \pi_4 I_t + \xi_{2t} \end{cases}$$

В правой части только predetermined переменные, так что оценку коэффициентов здесь уже можно проводить методом наименьших квадратов.

### 10.3. Косвенный метод наименьших квадратов

#### 10.3.1. Идея метода

Попробуем применить подход из предыдущего примера в общем случае.

1) Пусть исходная система дана в структурной форме

$$AY = BX + \varepsilon.$$

2) Перепишем её (если возможно) в приведённой форме

$$Y = \Pi X + \xi, \quad \Pi = A^{-1}B.$$

3) Рассчитаем оценки коэффициентов матрицы  $\Pi$  методом наименьших квадратов, после чего

4) Выразим через них (**если это окажется возможным!**) исходные структурные коэффициенты (т.е. матрицы  $A$  и  $B$ ).

Такой способ оценивания структурных коэффициентов называется **косвенным методом наименьших квадратов**.

### 10.3.2. Проблема идентификации

При использовании косвенного метода наименьших квадратов возникает проблема идентификации: удастся ли выразить коэффициенты исходной модели через коэффициенты приведённой модели?

**Пример.** Рассмотрим модель

$$\begin{cases} Q_t = a_0 + a_1 P_t + \varepsilon_1 \\ Q_t = b_0 + b_1 P_t + b_2 I_t + \varepsilon_2 \end{cases}$$

В приведённой форме она имеет вид

$$\begin{cases} Q_t = \pi_0 + \pi_1 I_t + \xi_{1t} \\ P_t = \pi_2 + \pi_3 I_t + \xi_{2t} \end{cases}$$

Ясно, что **пять** неизвестных параметров исходной модели  $a_0, a_1, b_0, b_1, b_2$  невозможно выразить через **четыре** коэффициента приведённой модели  $\pi_0, \pi_1, \pi_2, \pi_3$ .

**Структурный параметр** называется

- ❖ **идентифицируемым**, если он может быть однозначно выражен через коэффициенты приведённой формы.
- ❖ **неидентифицируемым**, если его невозможно выразить через коэффициенты приведённой формы.
- ❖ **сверхидентифицируемым**, если он может быть выражен через коэффициенты приведённой формы несколькими способами.

**Уравнение** называется **идентифицируемым**, если идентифицируемы все входящие в него коэффициенты.

В случае неидентифицируемости косвенный метод наименьших квадратов, очевидно, неприменим.

### **Признаки идентифицируемости**

Пусть исходная система дана в структурной форме, всего  $m$  уравнений,  $m_1$  эндогенных переменных  $Y$  и  $r$  предопределённых переменных  $X$ .

Все тождества можно исключить из системы уравнений.

**Первое необходимое условие идентифицируемости.**

Если тождества исключены, то число уравнений системы должно быть равно количеству эндогенных переменных:

$$m=m_1$$

**Второе необходимое условие идентифицируемости.**

Матрица наблюдений предопределённых переменных  $X$  должна иметь ранг  $p$  (т.е. ранг, равный количеству этих переменных).

**Третье необходимое условие идентифицируемости.**

Наборы переменных в различных уравнениях системы должны быть разными.

**Четвёртое необходимое условие идентифицируемости.**

Для каждого уравнения системы количество отсутствующих в нём предопределённых переменных должно быть не меньше числа включённых в него эндогенных переменных минус 1.

Если равно, то это уравнение идентифицируемо, если больше, то сверхидентифицируемо.

Четвёртое необходимое условие называется «правило порядка». Сформулируем его ещё раз.

Общий вид  $j$ -го уравнения модели в структурной форме можно записать как ( $y$  – эндогенные переменные,  $x$  – предопределённые):

$$a_{j1} \cdot y_{1t} + \dots + a_{jm} \cdot y_{mt} + b_{j1} \cdot x_{1t} + \dots + b_{jp} \cdot x_{pt} = \varepsilon_{jt}$$

Коэффициент при  $y_{jt}$  обычно равен 1 (“нормировка”).



Часть коэффициентов может быть равна нулю.

Обозначим через  $cnt_j$  количество на самом деле присутствующих в  $j$ -м уравнении переменных.

**Теорема (правило порядка).** Пусть  $j$ -ое уравнение системы идентифицируемо или сверхидентифицируемо,  $p$  – количество predetermined переменных в системе. Тогда справедливо неравенство

$$p \geq cnt_j - 1$$

**Пример.** Рассмотрим следующую эконометрическую модель:

$$\begin{cases} C_t = a_0 + a_1 Y_t + a_2 S_t + a_3 t + \varepsilon_{1t}, \\ I_t = b_0 + b_1 Y_{t-1} + \varepsilon_{2t}, \\ S_t = c_0 + c_1 Y_t + c_2 Y_{t-1} + \varepsilon_{3t}, \\ Y_t = C_t + I_t + G_t, \end{cases}$$

где  $C_t$  – расходы на конечное потребление в период  $t$ ;  
 $Y_t, Y_{t-1}$  – совокупный доход в периоды  $t$  и  $t-1$  соответственно;  
 $I_t$  – валовые инвестиции периода  $t$ ;  
 $S_t$  – расходы на зарплату в период  $t$ ;  
 $G_t$  – государственные расходы период  $t$ ;  
 $\varepsilon_1, \varepsilon_2, \varepsilon_3$  — случайные факторы.

В модели четыре эндогенных переменных:  $C_t, I_t, S_t, Y_t$ .

Остальные три переменные модели:  $1, t$  и  $G_t$  – экзогенные.

Кроме того, модель содержит лаговую эндогенную переменную  $Y_{t-1}$ . Таким образом, общее количество эндогенных переменных  $m = 4$ , количество predetermined переменных  $p = 4$ .

Для первого уравнения общее количество включенных в него переменных  $\text{cnt}_1=5$  (в него входят переменные  $C_t, 1, Y_t, S_t, t$ ). Имеем:  $p=4, \text{cnt}_1-1=5-1=4$ .

$4 = 4$ , следовательно, первое уравнение идентифицировано.

$$\begin{cases} C_t = a_0 + a_1 Y_t + a_2 S_t + a_3 t + \varepsilon_{1t}, \\ I_t = b_0 + b_1 Y_{t-1} + \varepsilon_{2t}, \\ S_t = c_0 + c_1 Y_t + c_2 Y_{t-1} + \varepsilon_{3t}, \\ Y_t = C_t + I_t + G_t, \end{cases}$$

Для второго уравнения:  $\text{cnt}_2 = 3 (I_t, 1, Y_{t-1})$ . Имеем:  $\text{cnt}_2-1=3-1=2$ .  $p=4 > 2$ , следовательно, второе уравнение свержидентифицировано.

Для третьего уравнения:  $\text{cnt}_3 = 4 (S_t, 1, Y_t, Y_{t-1})$ .

Имеем:  $\text{cnt}_3 - 1 = 4 - 1 = 3$ .

$p=4 > 3$ , следовательно, третье уравнение свержидентифицировано.

Последнее уравнение модели представляет собой тождество, его не надо проверять на идентификацию.

Приведенная форма модели будет иметь вид:

$$\begin{cases} C_t = \pi_0 + \pi_1 Y_{t-1} + \pi_2 G_t + \pi_3 t + \xi_{1t}, \\ I_t = \pi_4 + \pi_5 Y_{t-1} + \pi_6 G_t + \pi_7 t + \xi_{2t}, \\ S_t = \pi_8 + \pi_9 Y_{t-1} + \pi_{10} G_t + \pi_{11} t + \xi_{3t}, \\ Y_t = \pi_{12} + \pi_{13} Y_{t-1} + \pi_{14} G_t + \pi_{15} t + \xi_{4t}. \end{cases}$$

Здесь 12 неизвестных коэффициентов.

Когда они будут найдены, через них потребуется выразить 9 коэффициентов исходной модели ( $a_0, a_1, a_2, a_3, b_0, b_1, c_0, c_1, c_2$ ).

Эта процедура неоднозначна, что лишний раз подтверждает свержидентифицируемость модели

Правило порядка – это необходимое условие идентифицируемости, но оно не является достаточным. Это означает, что, когда неравенство несправедливо, то j-ое уравнение заведомо неидентифицируемо. Однако при выполнении неравенства ещё нельзя сделать вывод о идентифицируемости данного уравнения. Достаточное условие выражается в терминах рангов матриц.

**Необходимое и достаточное условие идентифицируемости.**

Уравнение идентифицируемо, если матрица, составленная по остальным уравнениям из коэффициентов отсутствующих в нём переменных имеет **ранг не меньше, чем общее количество эндогенных переменных минус 1**.

**Пример.** Проверим достаточное условие идентифицируемости.

Перенесём все в левую часть и выпишем матрицу коэффициентов при переменных

$C_t$	1	$Y_t$	$S_t$	$t$	$I_t$	$Y_{t-1}$	$G_t$
1	$-a_0$	$-a_1$	$-a_2$	$-a_3$	0	0	0
0	$-b_0$	0	0	0	1	$-b_1$	0
0	$-c_0$	$-c_1$	1	0	0	$-c_2$	0
-1	0	1	0	0	-1	0	-1

$$\begin{cases} C_t = a_0 + a_1 Y_t + a_2 S_t + a_3 t + \varepsilon_{1t}, \\ I_t = b_0 + b_1 Y_{t-1} + \varepsilon_{2t}, \\ S_t = c_0 + c_1 Y_t + c_2 Y_{t-1} + \varepsilon_{3t}, \\ Y_t = C_t + I_t + G_t, \end{cases}$$

Матрица коэффициентов при переменных, не входящих в первое уравнение:

$$\begin{pmatrix} 1 & -b_1 & 0 \\ 0 & -c_2 & 0 \\ -1 & 0 & -1 \end{pmatrix}$$

Её определитель не равен 0, так что ранг равен 3. Общее количество эндогенных переменных равно  $m=4$ . Имеем:  $4-1=3$ , так что первое уравнение идентифицируемо.

Матрица коэффициентов при переменных, не входящих в второе уравнение:

$$\begin{pmatrix} -1 & a_1 & a_2 & a_3 & 0 \\ 0 & c_1 & -1 & 0 & 0 \\ 1 & -1 & 0 & 0 & 1 \end{pmatrix}$$

$C_t$	1	$Y_t$	$S_t$	$t$	$I_t$	$Y_{t-1}$	$G_t$
1	$-a_0$	$-a_1$	$-a_2$	$-a_3$	0	0	0
0	$-b_0$	0	0	0	1	$-b_1$	0
0	$-c_0$	$-c_1$	1	0	0	$-c_2$	0
-1	0	1	0	0	-1	0	-1

Ранг матрицы равен 3, так как определитель следующей подматрицы 3x3 не равен нулю. Второе уравнение идентифицируемо.

$$\begin{pmatrix} -1 & a_2 & 0 \\ 0 & -1 & 0 \\ 1 & 0 & 1 \end{pmatrix}$$

Матрица коэффициентов при переменных, не входящих в третье уравнение:

$$\begin{pmatrix} -1 & a_3 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 1 & 0 & 1 & 1 \end{pmatrix}$$

Матрица коэффициентов при переменных, не входящих в четвертое уравнение:

$$\begin{pmatrix} -a_0 & -a_2 & -a_3 & 0 \\ -b_0 & 0 & 0 & -b_1 \\ -c_0 & -1 & 0 & -c_2 \end{pmatrix} \dots$$

Ранг матрицы аналогичным образом равен 3. Третье уравнение идентифицируемо.

### 10.3.3. Алгоритм косвенного МНК

Подведём итоги.

Косвенный метод наименьших квадратов (КМНК) применяется в случае точной идентифицируемости уравнений модели

Алгоритм применения КМНК:

1. От структурной формы модели переходят к приведенной
2. Определяются МНК-оценки параметров приведенной формы модели
3. По МНК-оценкам приведенной формы вычисляются оценки параметров структурной формы модели.

Исходная система в структурной форме:

$$AY = BX + \varepsilon,$$

В приведённой форме:

$$Y = \Pi X + \xi, \quad \Pi = A^{-1}B$$

Следовательно,

$$A\Pi = A \cdot (A^{-1}B) = B$$

Это – система  $m$  х  $p$  уравнений.

Из неё, в случае идентифицируемости, по матрице  $\Pi$  можно вычислить значения структурных коэффициентов  $A$  и  $B$ .

**Задача.** Дана модель в структурной форме. Найти оценки структурных коэффициентов по исходным данным косвенным методом наименьших квадратов

$$\begin{cases} y_{1t} = c_{10} + a_{12} y_{2t} + b_{11} x_{1t} + \varepsilon_{1t}, \\ y_{2t} = c_{20} + a_{21} y_{1t} + b_{22} x_{2t} + \varepsilon_{2t}, \end{cases}$$

$t$	$y_1$	$y_2$	$x_1$	$x_2$
1	22	24	2	7
2	24	15	3	5
3	17	23	2	5
4	19	21	3	4
5	22	21	4	5
6	25	13	5	3
7	25	15	4	4
8	31	12	6	3
9	33	8	7	2
10	31	5	8	1

**Решение.** 1) Проверим на идентифицируемость.

Перенесём все в левую часть и выпишем матрицу коэффициентов при переменных

$y_1$	$y_2$	1	$x_1$	$x_2$
1	$-a_{12}$	$-c_{10}$	$-b_{11}$	0
$-a_{21}$	1	$-c_{20}$	0	$-b_{22}$

Количество эндогенных переменных  $m=2$  ( $y_1$  и  $y_2$ ).

Матрица коэффициентов при переменных, не входящих в первое уравнение состоит из одного числа  $-b_{21}$ , ранг = 1 =  $m-1$ , уравнение точно идентифицируемо. Аналогично со вторым уравнением.

2) В приведенной форме модель примет вид:

$$\begin{cases} y_{1t} = \pi_{10} + \pi_{11} x_{1t} + \pi_{12} x_{2t} + \xi_{1t}, \\ y_{2t} = \pi_{20} + \pi_{21} x_{1t} + \pi_{22} x_{2t} + \xi_{2t}, \end{cases}$$

$t$	$y_1$	$y_2$	$x_1$	$x_2$
1	22	24	2	7
2	24	15	3	5
3	17	23	2	5
4	19	21	3	4
5	22	21	4	5
6	25	13	5	3
7	25	15	4	4
8	31	12	6	3
9	33	8	7	2
10	31	5	8	1

Оценивая уравнения по отдельности обычным МНК, получаем

$$\begin{cases} y_{1t} = 0,71 + 3,8 x_{1t} + 1,91 x_{2t}, & R^2 = 0,88 \\ y_{2t} = 20,44 - 2,06 x_{1t} + 1,11 x_{2t}, & R^2 = 0,88 \end{cases}$$

3) Записываем равенство для нахождения структурных коэффициентов

$$A\Pi = B$$

$$\begin{pmatrix} 1 & -a_{12} \\ -a_{21} & 1 \end{pmatrix} \begin{pmatrix} 0,71 & 3,8 & 1,91 \\ 20,44 & -2,06 & 1,11 \end{pmatrix} = \begin{pmatrix} c_{10} & b_{11} & 0 \\ c_{20} & 0 & b_{22} \end{pmatrix}$$

Из шести уравнений есть два с нулевой правой частью, откуда находятся коэффициенты матрицы А.

$$\begin{cases} 1,91 - 1,11a_{12} = 0 \\ 3,8a_{21} + 2,06 = 0 \end{cases} \Rightarrow \begin{cases} a_{12} \approx 1,72 \\ a_{21} \approx -0,54 \end{cases} \Rightarrow A = \begin{pmatrix} 1 & -1,72 \\ 0,54 & 1 \end{pmatrix}$$

Далее находим матрицу В

$$\begin{aligned} \begin{pmatrix} c_{10} & b_{11} & 0 \\ c_{20} & 0 & b_{22} \end{pmatrix} &= \begin{pmatrix} 1 & -1,72 \\ 0,54 & 1 \end{pmatrix} \cdot \begin{pmatrix} 0,71 & 3,8 & 1,91 \\ 20,44 & -2,06 & 1,11 \end{pmatrix} = \\ &= \begin{pmatrix} -34,44 & 7,34 & 0 \\ 20,82 & 0 & 2,14 \end{pmatrix} \end{aligned}$$

Теперь выписываем модель в структурной форме.

$$\begin{cases} y_{1t} = c_{10} + a_{12} y_{2t} + b_{11} x_{1t} + \varepsilon_{1t}, \\ y_{2t} = c_{20} + a_{21} y_{1t} + b_{22} x_{2t} + \varepsilon_{2t}, \end{cases}$$

$$\begin{cases} y_{1t} = -34,44 + 1,72 y_{2t} + 7,34 x_{1t} + \varepsilon_{1t}, \\ y_{2t} = 20,82 - 0,54 y_{1t} + 2,14 x_{2t} + \varepsilon_{2t}, \end{cases}$$

**Задача.** Дана модель в структурной форме

$$\begin{cases} Y_1 = a_0 + a_1 X_1 + a_2 X_2 + a_3 Y_3 + \varepsilon_{1t}, \\ Y_2 = b_0 + b_1 X_2 + b_2 X_3 + b_3 Y_1 + \varepsilon_{2t}, \\ Y_3 = c_0 + c_1 X_1 + c_2 X_3 + c_3 Y_1 + \varepsilon_{3t} \end{cases}$$

В приведённой форме  
получены коэффициенты

$$\begin{cases} Y_1 = 4 + 6X_1 + X_2 + 4X_3, \\ Y_2 = 6 - 12X_1 - 7X_2 + 8X_3, \\ Y_3 = 10 - 5X_1 - 2X_2 + 5X_3 \end{cases}$$

Определить значения всех структурных коэффициентов, которые можно найти.

**Решение.**

Общий вид структурной формы

$$AY = BX + \varepsilon,$$

Общий вид приведённой формы

$$Y = \Pi X + \xi,$$

Выпишем матрицы А, В

$$A = \begin{pmatrix} 1 & 0 & -a_3 \\ -b_3 & 1 & 0 \\ -c_3 & 0 & 1 \end{pmatrix},$$

$$\begin{cases} Y_1 = a_0 + a_1 X_1 + a_2 X_2 + a_3 Y_3 + \varepsilon_{1t}, \\ Y_2 = b_0 + b_1 X_2 + b_2 X_3 + b_3 Y_1 + \varepsilon_{2t}, \\ Y_3 = c_0 + c_1 X_1 + c_2 X_3 + c_3 Y_1 + \varepsilon_{3t} \end{cases}$$

$$B = \begin{pmatrix} a_0 & a_1 & a_2 & 0 \\ b_0 & 0 & b_1 & b_2 \\ c_0 & c_1 & 0 & c_2 \end{pmatrix}$$

Выпишем матрицу П

$$\Pi = \begin{pmatrix} 4 & 6 & 1 & 4 \\ 6 & -12 & -7 & 8 \\ 10 & -5 & -2 & 5 \end{pmatrix}$$

$$\begin{cases} Y_1 = 4 + 6X_1 + X_2 + 4X_3, \\ Y_2 = 6 - 12X_1 - 7X_2 + 8X_3, \\ Y_3 = 10 - 5X_1 - 2X_2 + 5X_3 \end{cases}$$

3) Запишем равенство для  
нахождения структурных  
коэффициентов

$$A\Pi = B$$

$$\begin{pmatrix} 1 & 0 & -a_3 \\ -b_3 & 1 & 0 \\ -c_3 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} 4 & 6 & 1 & 4 \\ 6 & -12 & -7 & 8 \\ 10 & -5 & -2 & 5 \end{pmatrix} = \begin{pmatrix} a_0 & a_1 & a_2 & 0 \\ b_0 & 0 & b_1 & b_2 \\ c_0 & c_1 & 0 & c_2 \end{pmatrix}$$

Из 12 уравнений три имеют правую часть равную нулю.

$$\begin{cases} 4 - 5a_3 = 0, \\ -6b_3 - 12 = 0, \\ -c_3 - 2 = 0 \end{cases}$$

Отсюда находим все коэффициенты матрицы А.

$$\begin{cases} a_3 = 0,8, \\ b_3 = -2, \\ c_3 = -2 \end{cases}$$

Перемножая матрицы, находим остальные коэффициенты:

$$\begin{pmatrix} a_0 & a_1 & a_2 & 0 \\ b_0 & 0 & b_1 & b_2 \\ c_0 & c_1 & 0 & c_2 \end{pmatrix} = \begin{pmatrix} 1 & 0 & -0,8 \\ 2 & 1 & 0 \\ 2 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} 4 & 6 & 1 & 4 \\ 6 & -12 & -7 & 8 \\ 10 & -5 & -2 & 5 \end{pmatrix}$$

### 10.3.4. Двухшаговый метод наименьших квадратов

Если система сверхидентифицируема, то косвенный МНК не применим, структурные коэффициенты не выражаются однозначно через коэффициенты приведённой системы.

**Пример.** Рассмотрим следующую модель в структурной форме.

$$\begin{cases} y_{1t} = a_1 y_{2t} + b_1 x_{1t} + b_2 x_{2t} + \varepsilon_{1t}, \\ y_{2t} = a_2 y_{1t} + b_3 x_{3t} + \varepsilon_{2t}, \end{cases}$$

В приведённой форме она будет иметь вид

$$\begin{cases} y_{1t} = \pi_1 x_{1t} + \pi_2 x_{2t} + \pi_3 x_{3t} + \xi_{1t}, \\ y_{2t} = \pi_4 x_{1t} + \pi_5 x_{2t} + \pi_6 x_{3t} + \xi_{2t}, \end{cases}$$

Пусть получены оценки параметров приведённой формы

$$\begin{cases} y_{1t} = 2 x_{1t} + 3 x_{2t} + 3 x_{3t}, \\ y_{2t} = 3 x_{1t} + 4 x_{2t} + 2 x_{3t}, \end{cases}$$

Попробуем найти оценки структурной формы:

$$AP = B$$



$$\begin{pmatrix} 1 & -a_1 \\ -a_2 & 1 \end{pmatrix} \begin{pmatrix} 2 & 3 & 3 \\ 3 & 4 & 2 \end{pmatrix} = \begin{pmatrix} b_1 & b_2 & 0 \\ 0 & 0 & b_{23} \end{pmatrix}$$

Из шести уравнений есть три с нулевой правой частью, откуда получается три уравнения для нахождения  $a_1, a_2$ .

Видно, что  $a_2$  выражается двумя способами, противоречащими друг другу – это и есть сверхидентифицируемость.

$$\begin{cases} -2a_2 + 3 = 0 \\ -3a_2 + 4 = 0 \\ 3 - 2a_1 = 0 \end{cases}$$

### Алгоритм оценки коэффициентов структурной формы двухшаговым МНК

1. Оцениваются параметры приведенной формы модели при помощи МНК
2. Оцениваются параметры структурной формы модели, в правую часть которой вместо значений эндогенных переменных подставляются их оценки, рассчитанные по приведенной форме

**Пример (продолжение).** Пусть в предыдущем примере исходные данные имели вид:

Подставляя в найденные формулы, мы можем найти расчётные значения:

$$\begin{cases} \hat{y}_{1t} = 2x_{1t} + 3x_{2t} + 3x_{3t}, \\ \hat{y}_{2t} = 3x_{1t} + 4x_{2t} + 2x_{3t}, \end{cases}$$

$t$	$y_1$	$y_2$	$x_1$	$x_2$	$x_3$	$\hat{y}_1$	$\hat{y}_2$
1	33,4	40,0	2	7	3	34	40
2	24,0	30,3	3	5	1	24	31
3	21,3	28,6	2	5	1	22	28
4	23,5	29,8	3	4	2	24	29
5	24,9	34,7	4	5	1	26	34
6	21,8	28,9	5	3	1	22	29
7	26,1	31,9	4	4	2	26	32
8	24,5	32	6	3	1	24	32
9	32,7	36,8	7	2	4	32	37
10	33,9	37,9	8	1	5	34	38

$t$	$y_1$	$y_2$	$x_1$	$x_2$	$x_3$
1	33,4	40,0	2	7	3
2	24,0	30,3	3	5	1
3	21,3	28,6	2	5	1
4	23,5	29,8	3	4	2
5	24,9	34,7	4	5	1
6	21,8	28,9	5	3	1
7	26,1	31,9	4	4	2
8	24,5	32	6	3	1
9	32,7	36,8	7	2	4
10	33,9	37,9	8	1	5

Эти рассчитанные значения используем в качестве новых переменных в правой части структурной формы

$$\begin{cases} y_{1t} = a_1 \hat{y}_{2t} + b_1 x_{1t} + b_2 x_{2t} + \varepsilon_{1t}, \\ y_{2t} = a_2 \hat{y}_{1t} + b_3 x_{3t} + \varepsilon_{2t}, \end{cases}$$

Опять обычным методом получаем коэффициенты структурной формы

$$\begin{cases} y_{1t} = 1,49 \hat{y}_{2t} - 2,4 x_{1t} - 3,09 x_{2t} + \varepsilon_{1t}, & R^2 = 0,999 \\ y_{2t} = 1,38 \hat{y}_{1t} - 1,9 x_{3t} + \varepsilon_{2t}, & R^2 = 0,999 \end{cases}$$

Таким образом, двухшаговый МНК даёт следующий ответ:

$$\begin{cases} y_{1t} = 1,49 y_{2t} - 2,4 x_{1t} - 3,09 x_{2t} + \varepsilon_{1t}, & R^2 = 0,999 \\ y_{2t} = 1,38 y_{1t} - 1,9 x_{3t} + \varepsilon_{2t}, & R^2 = 0,999 \end{cases}$$

#### **Замечания.**

- 1) В случае точной идентифицируемости двухшаговый МНК даёт те же результаты, что и косвенный МНК.
- 2) Двухшаговый МНК является частным случаем метода инструментальных переменных, когда “неудобные” эндогенные переменные в правой части заменяются на новые переменные.

## Литература

1. Магнус Я.Р., Катышев П.К., Пересецкий А.А. Эконометрика. Начальный курс: Учебник. – М.: Дело, 2004.
2. Эконометрика: Учебник./Под ред. И.И. Елисеевой. – М.: Финансы и статистика, 2005.
3. Практикум по эконометрике: Учеб. пособие. /И.И. Елисеева, С.В. Курьшева, Н.М. Гордеенко и др. Под ред. И.И. Елисеевой. – М.: Финансы и статистика, 2005.
4. Кремер Н.Ш., Путко Б.А. Эконометрика: учебник для вузов. /Под ред. проф. Н.Ш. Кремера. – М.: ЮНИТИ–ДАНА, 2002.
5. Кобзарь А.И. Прикладная математическая статистика. Для инженеров и научных работников. – М.: ФИЗМАТЛИТ, 2006.
6. Берндт Э. Р. Практика эконометрики: классика и современность: Учебник. – М.: ЮНИТИ–ДАНА, 2005.
7. Вербик М. Путеводитель по современной эконометрике. Пер. с англ. В.А.Банникова. – М.: Научная книга, 2008.

## Оглавление

1. Введение	3
2. Корреляционный анализ количественных зависимостей	7
3. Метод наименьших квадратов	16
4. Классическая модель множественной линейной регрессии	19
5. Спецификация модели	41
6. Некоторые дополнительные вопросы	54
7. Нарушения допущений классической линейной модели	63
8. Анализ временных рядов	83
9. Взаимосвязь временных рядов	125
10. Системы эконометрических уравнений	135
Литература	153