

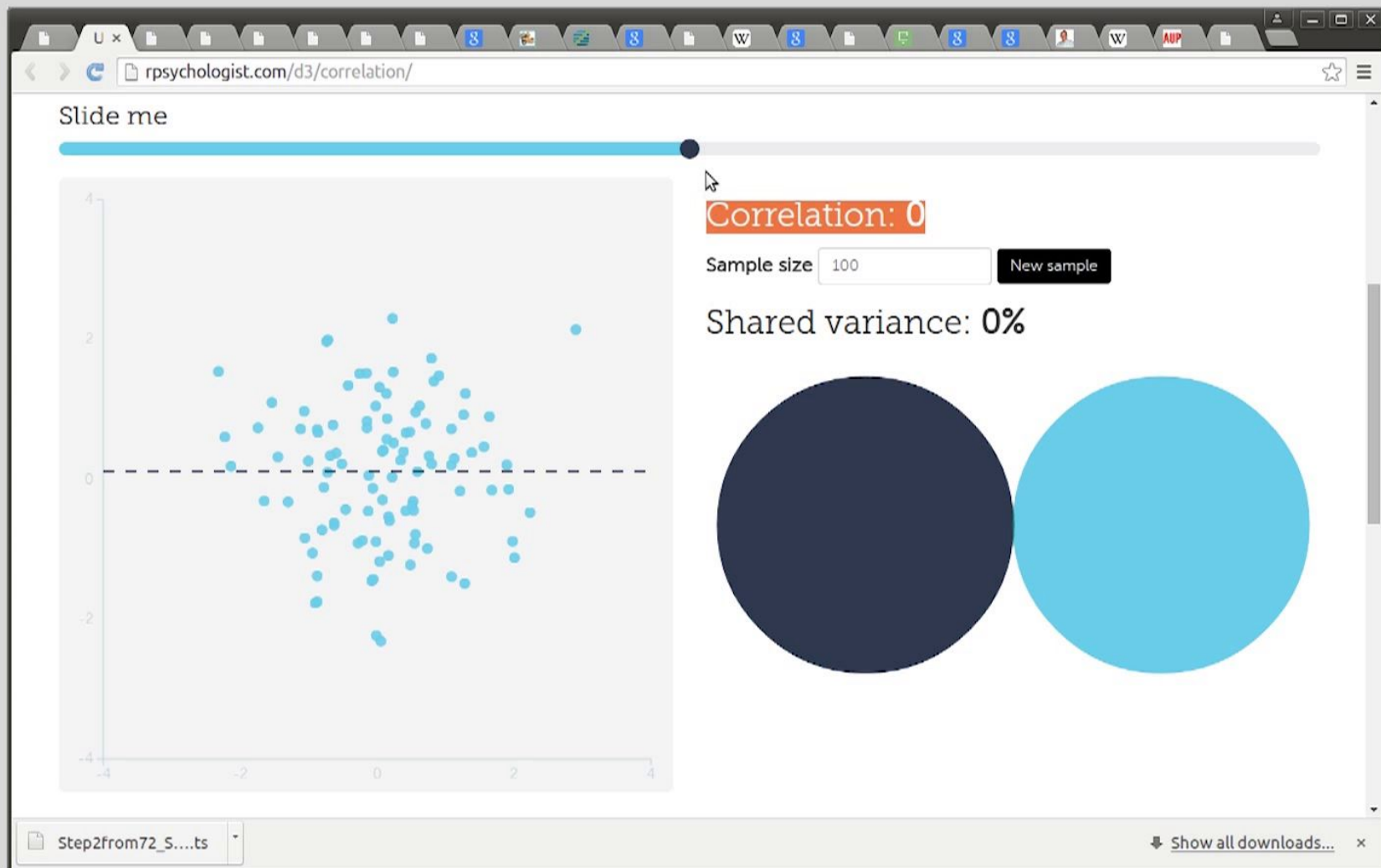
$$+ \begin{cases} (x_1 - \bar{x}) \cdot (y_1 - \bar{y}) > 0 \\ (x_2 - \bar{x}) \cdot (y_1 - \bar{y}) > 0 \\ \vdots \\ (x_i - \bar{x}) \cdot (y_i - \bar{y}) \end{cases}$$

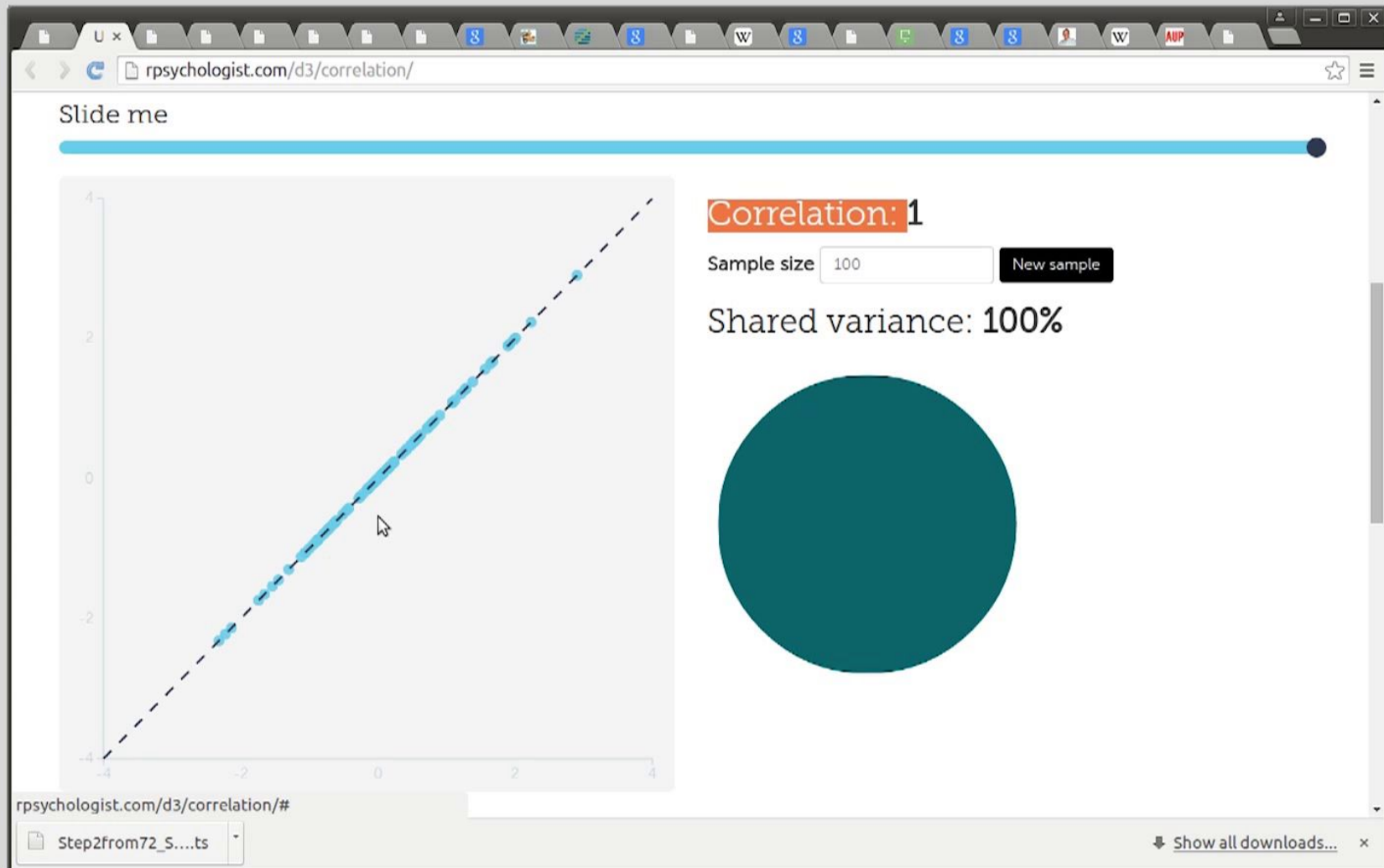
$$\underline{cov} = \frac{\sum (x_i - \bar{x}) \cdot (y_i - \bar{y})}{N-1} + \begin{matrix} \nearrow 7200 \\ \nwarrow \end{matrix}$$

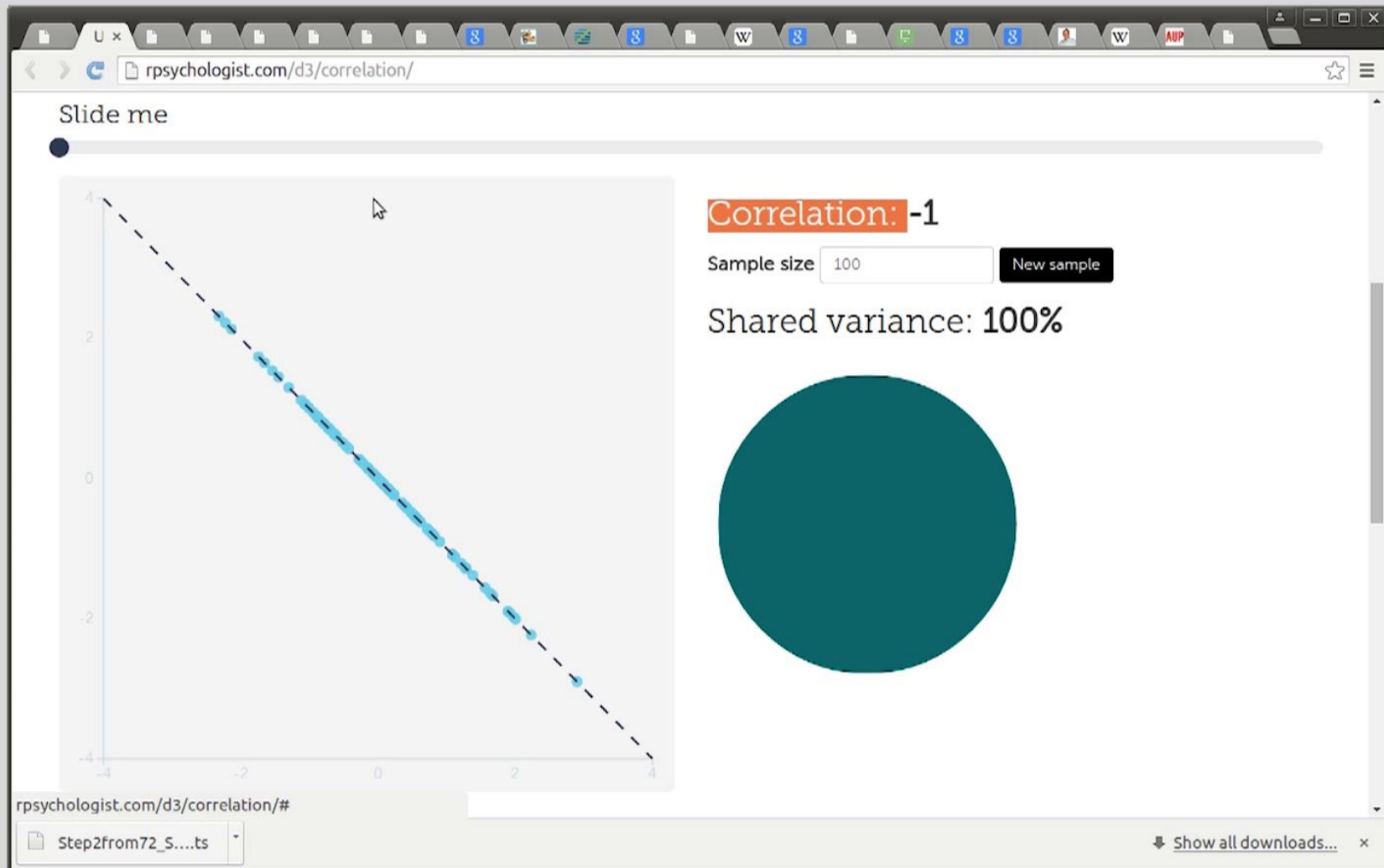
$$r_{xy} = \frac{cov}{\sigma_x \cdot \sigma_y} \quad [-1, 1]$$

"-"  $\nwarrow$   $\nearrow$  "+"









# Коэффициент корреляции

$r_{xy}$  — показатель силы и направления взаимосвязи двух количественных переменных



Принимает значения  $[-1, 1]$

Знак коэффициента корреляции показывает направление взаимосвязи



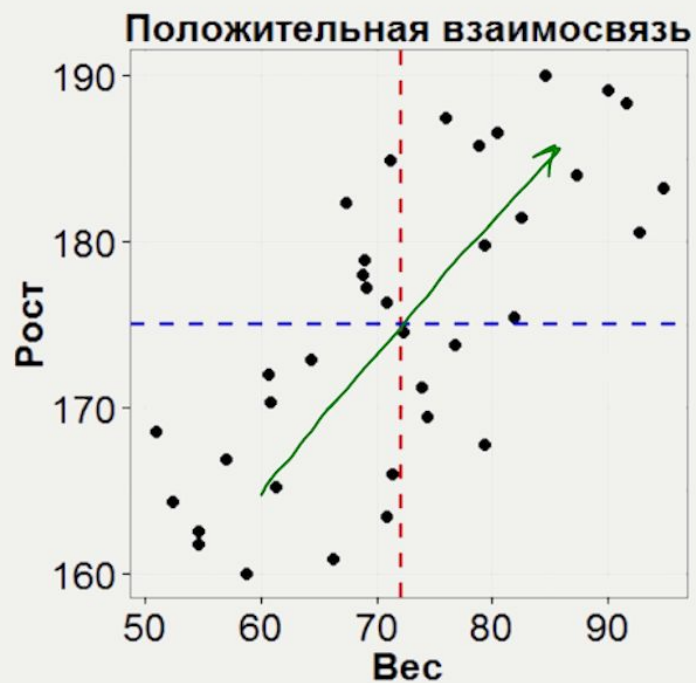
# Коэффициент детерминации

$R^2$  — показывает, в какой степени дисперсия одной переменной обусловлена влиянием другой переменной

Равен квадрату коэффициента корреляции

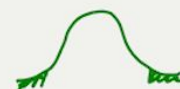
Принимает значения  $[0, 1]$





$$H_0 \quad r_{xy} = 0$$

$$H_1 \quad r_{xy} \neq 0$$

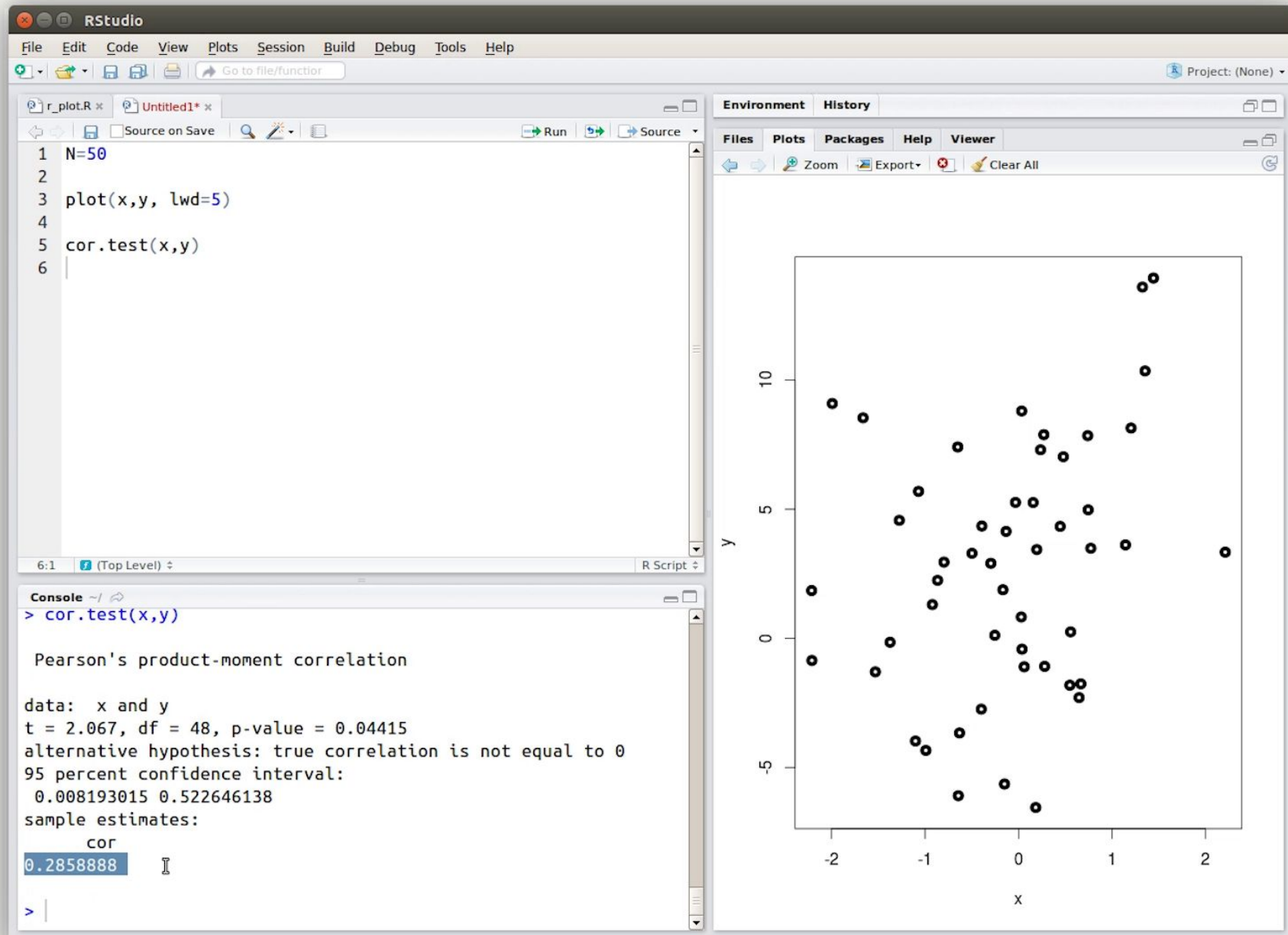


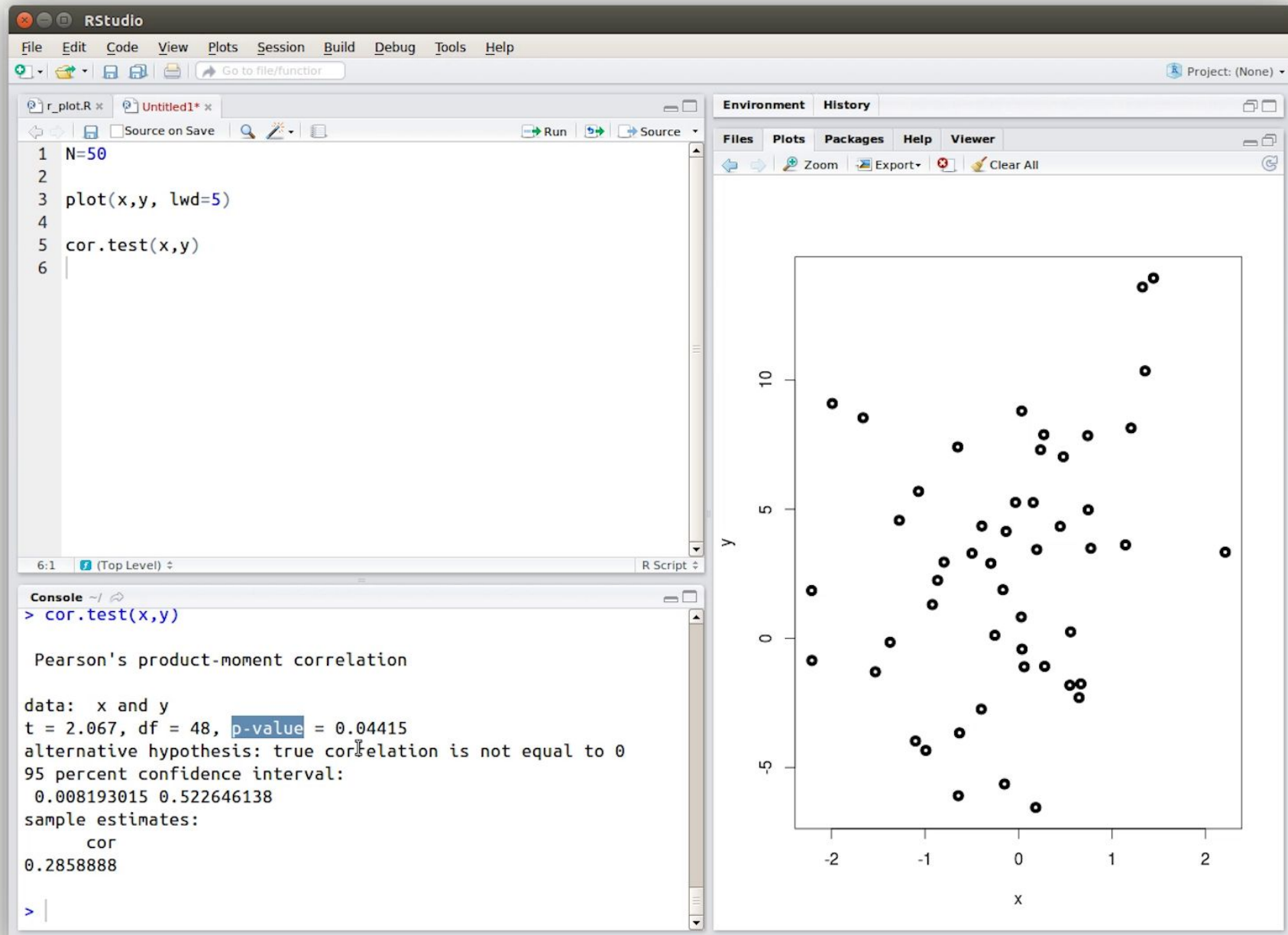
$$t \quad df = N - 2$$

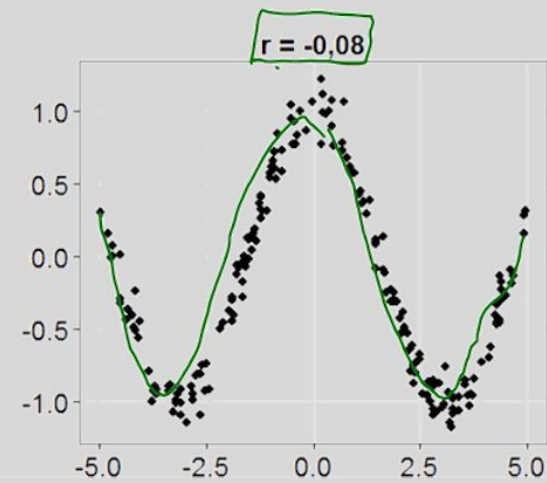
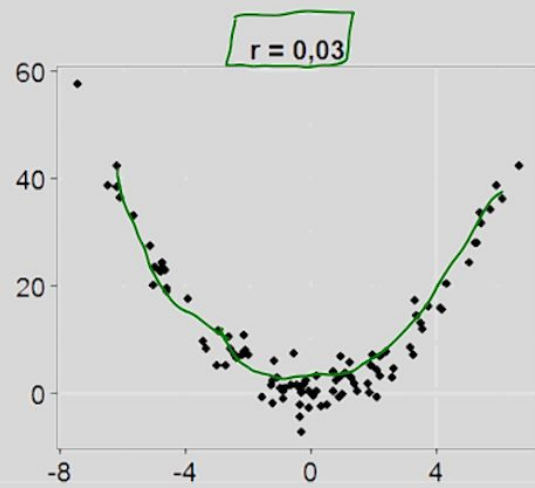
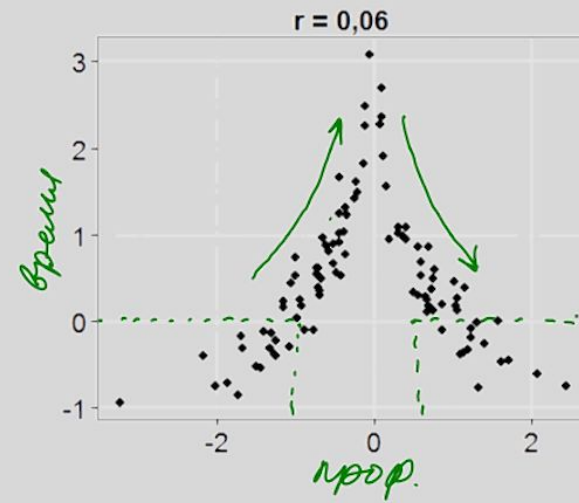
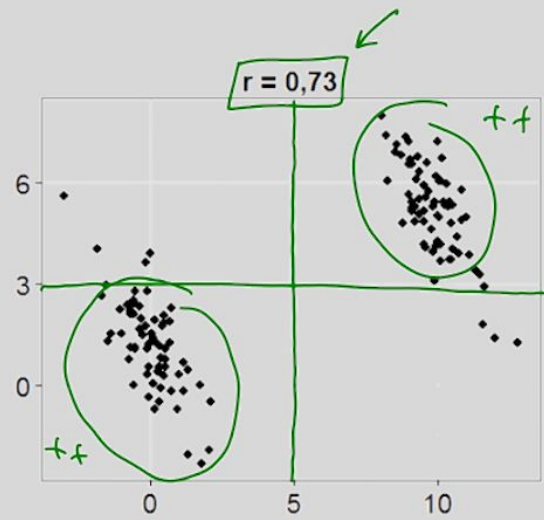
P



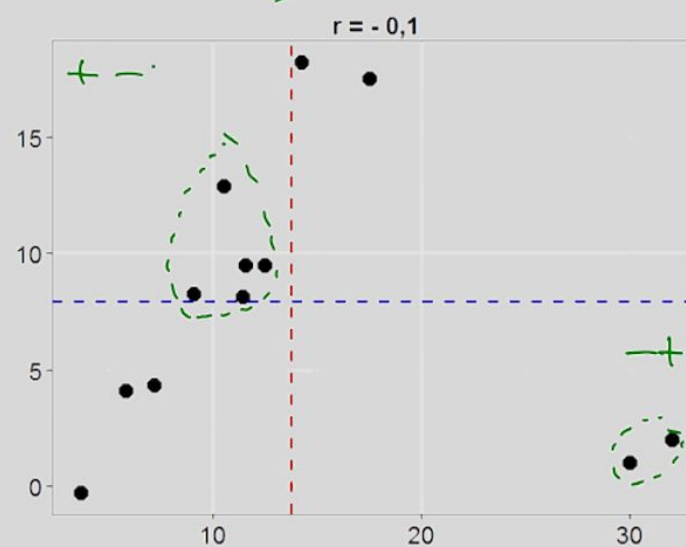
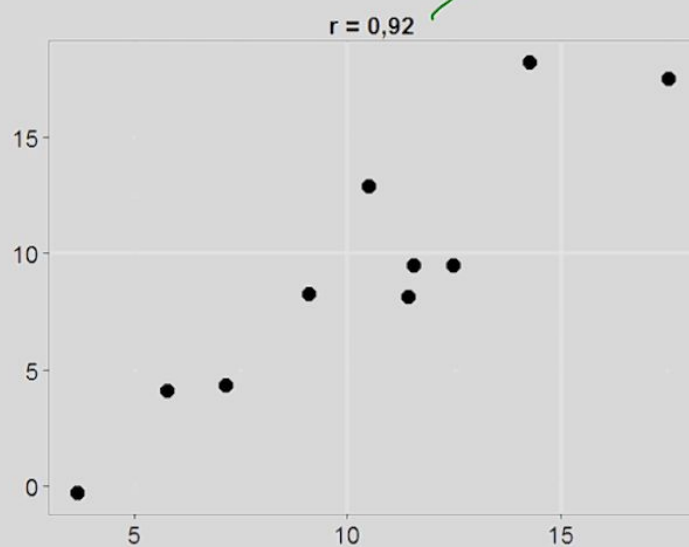




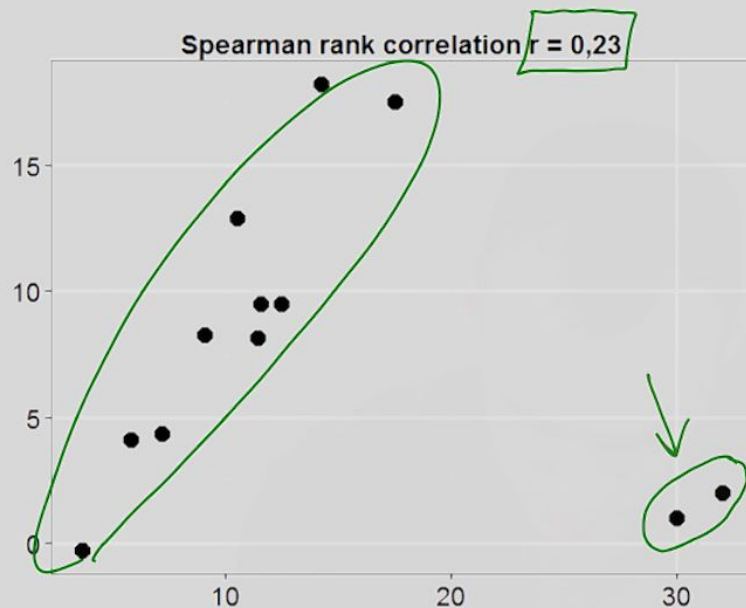




# Опасность выбросов!!!



## Коэффициент корреляции Спирмена



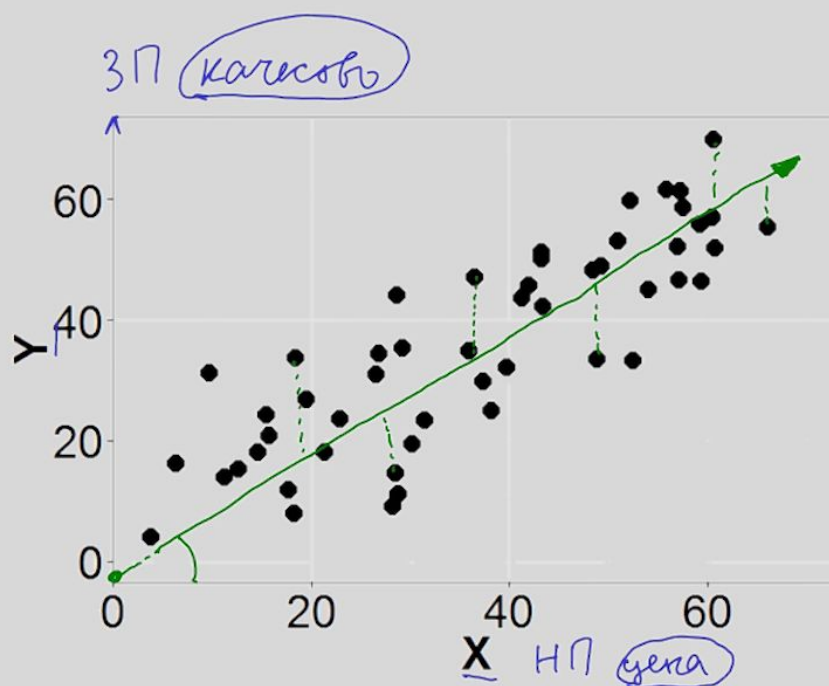
$$r_s = 1 - \frac{6 \sum_i d_i^2}{N(N^2 - 1)}$$

<u>X</u>	<u>Y</u>		<u>X</u>	<u>Y</u>	<u>d<sup>2</sup></u>
3,7	-0,3		1	1	0
5,8	4,1		2	4	<del>(-2)<sup>2</sup></del> 4
7,1	4,3		3	5	4
9,1	8,3		4	7	9
10,5	12,9		5	10	25
11,4	8,1		6	6	0
11,6	9,5		7	9	4
12,5	9,5		8	8	0
14,3	18,2		9	12	9
17,5	17,5		10	11	1
30,0	1,0		11	2	81
32,0	2,0		12	3	81

Σ



# Линия регрессии



Statistic	N	Mean	St. Dev.
x	51	36.7	17.3
y	51	35.9	16.9

$$y = \underline{\underline{b_0}} + \underline{\underline{b_1}} x$$

intercept slope

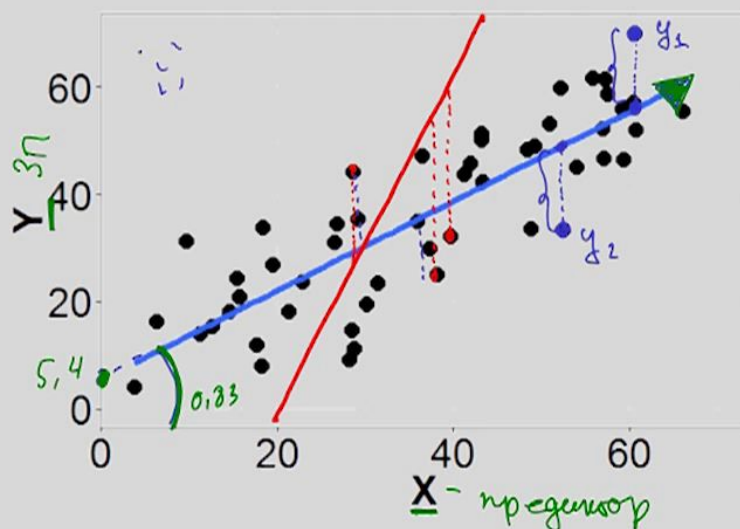


# Метод наименьших квадратов

$b_0$   $b_1$

МНК — метод нахождения оптимальных параметров линейной регрессии, таких, что сумма квадратов ошибок (остатков) была минимальна

$$\hat{y} = 5,4 + 0,83 \cdot x$$



$$e_1 = y_1 - \hat{y}_1 \quad \sum e_i^2$$

$$e_2 = y_2 - \hat{y}_2$$

$$b_1 = \frac{sd_y}{sd_x} \cdot r_{xy} \quad b_1 = 0,83$$

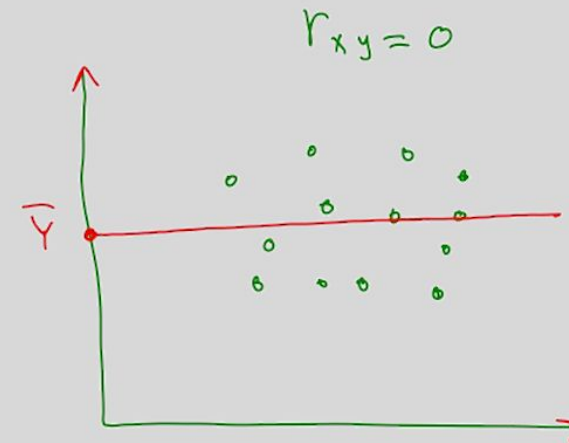
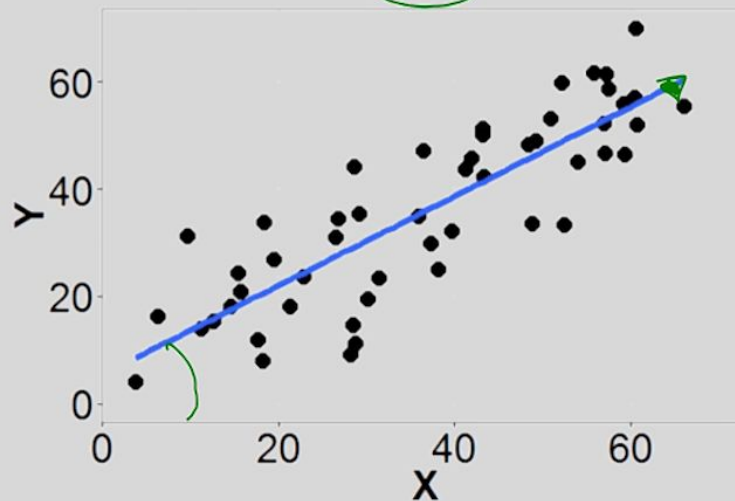
$$b_0 = \bar{y} - b_1 \cdot \bar{x} \quad b_0 = 5,4$$





# Коэффициенты линейной регрессии

$$\hat{y} = 5.4 + \overset{\beta_1}{\underset{\text{0.83}}{\circledast}} x$$

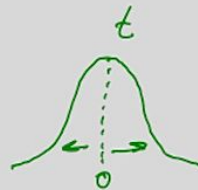


$$\beta_1 = \frac{sd_y}{sd_x} \cdot r_{xy} = 0$$

$$\beta_0 = \bar{Y} - \beta_1 \cdot \bar{X} = \bar{Y}$$

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$



$$t = \frac{\beta_1 - 0}{se}$$

$$t = \frac{\beta_1}{se} \quad df = N - 2$$

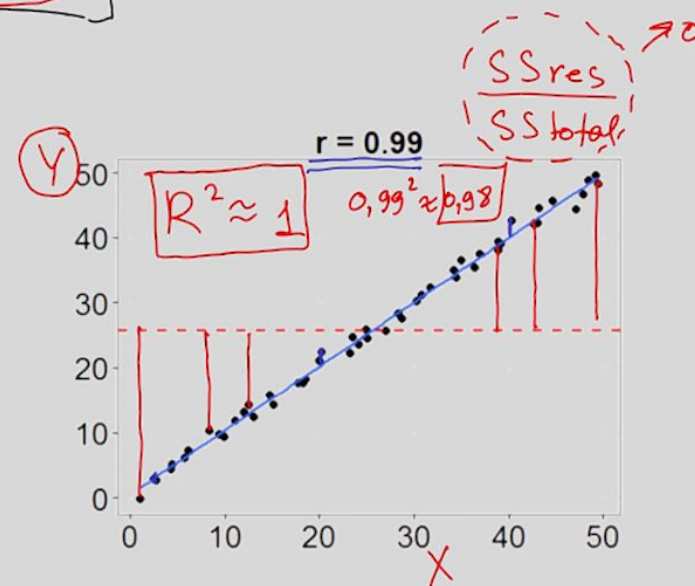
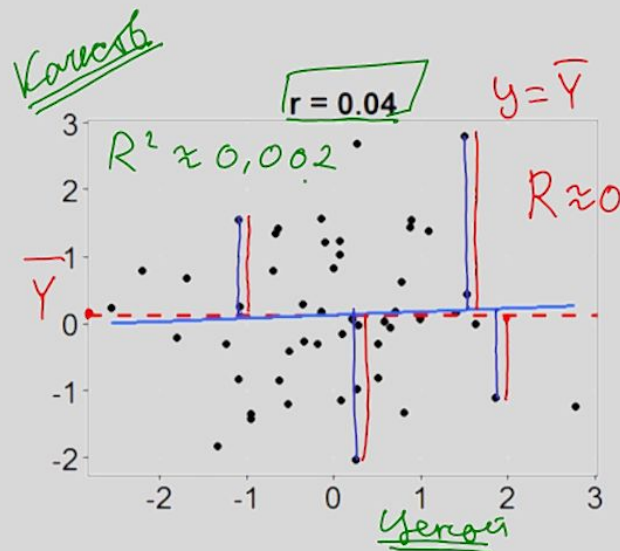




# Коэффициент детерминации

$R^2$  – доля дисперсии зависимой переменной (Y), объясняемая регрессионной моделью.

$$R^2 = 1 - \frac{SS_{res}}{SS_{total}}$$



# Условия применения

- Линейная взаимосвязь  $X$  и  $Y$
- Нормальное распределение остатков
- Гомоскедастичность - постоянная изменчивость остатков на всех уровнях независимой переменной

[http://bitly.com/slr\\_diag](http://bitly.com/slr_diag)



# Diagnostics for simple linear regression

Select a trend:

- ☒ Linear up
- ☐ Linear down
- ☐ Curved up
- ☐ Curved down
- ☐ Fan-shaped

☒ Show residuals

This applet uses ordinary least squares (OLS) to fit a regression line to the data with the selected trend. The applet is designed to help you practice evaluating whether or not the linear model is an appropriate fit to the data. The three diagnostic plots on the lower half of the page are provided to help you identify undesirable patterns in the residuals that may arise from non-linear trends in the data.

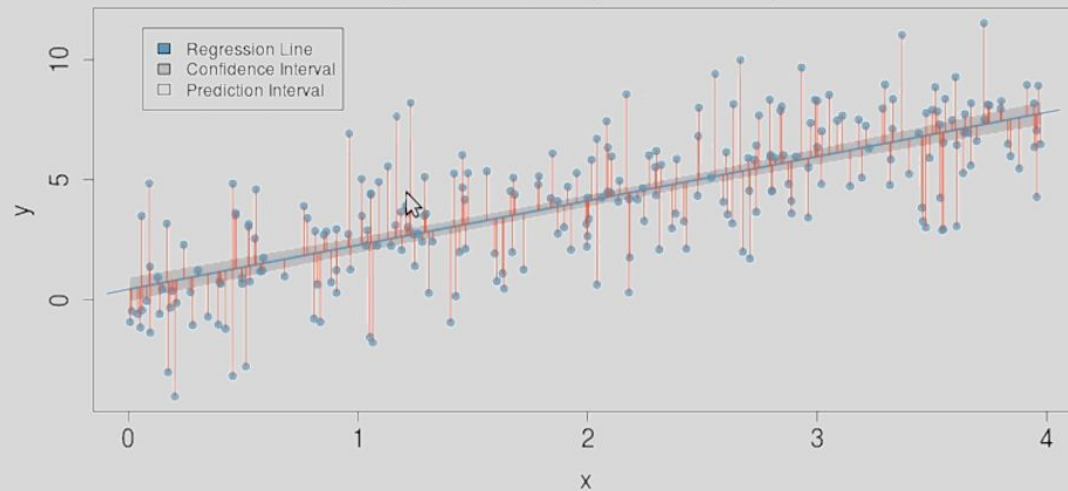
[Rate this app!](#)

[View code](#)

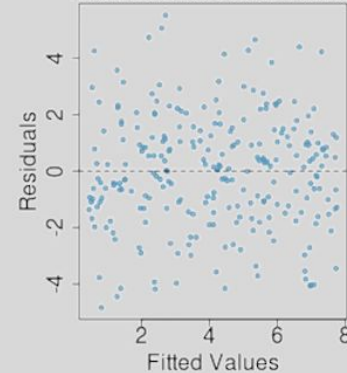
[Check out other apps](#)

[Want to learn more for free?](#)

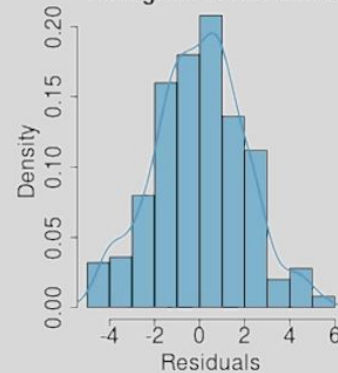
Regression Model  
( $R = 0.7346$ ,  $R\text{-squared} = 0.5397$ )



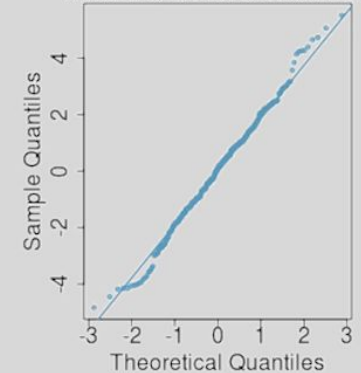
Residuals vs. Fitted Values



Histogram of Residuals



Normal Q-Q Plot of Residuals



# Diagnostics for simple linear regression

Select a trend:

- ☐ Linear up
- ☒ Linear down
- ☐ Curved up
- ☐ Curved down
- ☐ Fan-shaped

☒ Show residuals

This applet uses ordinary least squares (OLS) to fit a regression line to the data with the selected trend. The applet is designed to help you practice evaluating whether or not the linear model is an appropriate fit to the data. The three diagnostic plots on the lower half of the page are provided to help you identify undesirable patterns in the residuals that may arise from non-linear trends in the data.

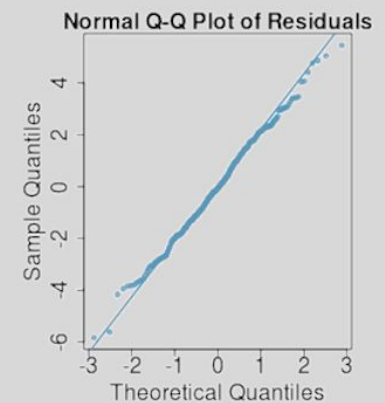
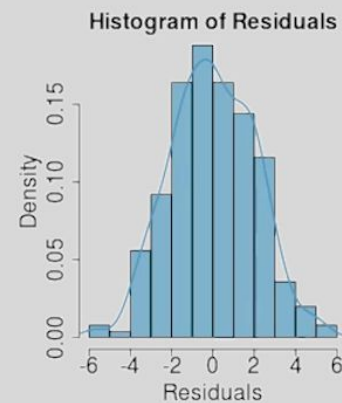
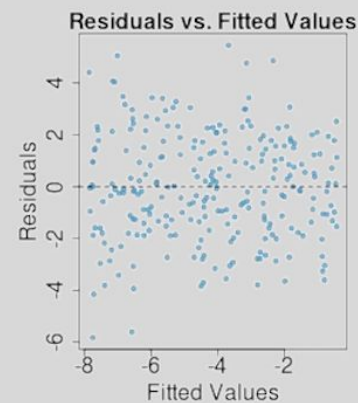
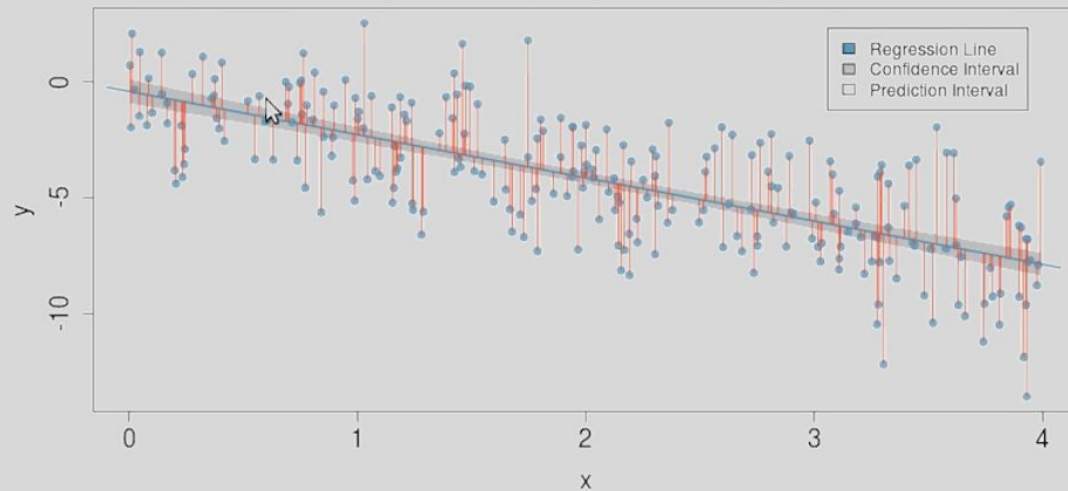
[Rate this app!](#)

[View code](#)

[Check out other apps](#)

[Want to learn more for free?](#)

Regression Model  
( $R = -0.7284$ ,  $R\text{-squared} = 0.5306$ )



# Diagnostics for simple linear regression

Select a trend:

- ☐ Linear up
- ☐ Linear down
- ☒ Curved up
- ☐ Curved down
- ☐ Fan-shaped

☒ Show residuals

This applet uses ordinary least squares (OLS) to fit a regression line to the data with the selected trend. The applet is designed to help you practice evaluating whether or not the linear model is an appropriate fit to the data. The three diagnostic plots on the lower half of the page are provided to help you identify undesirable patterns in the residuals that may arise from non-linear trends in the data.

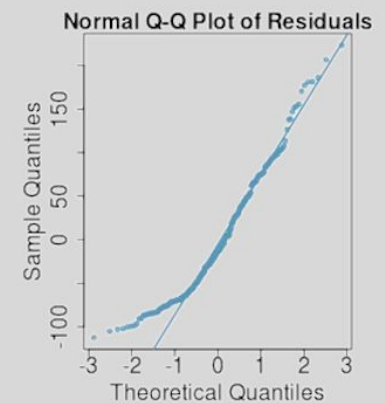
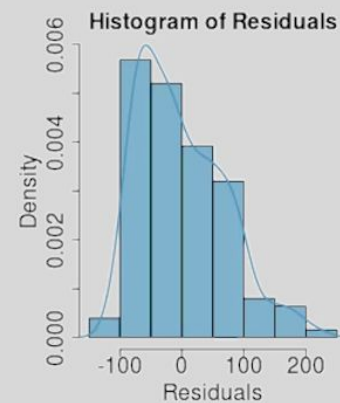
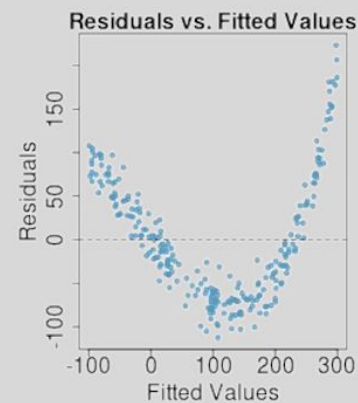
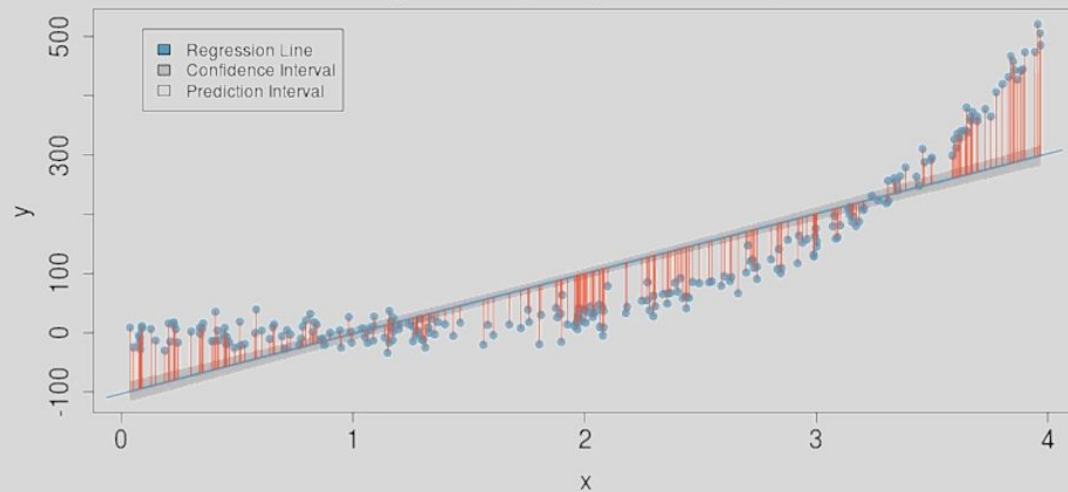
[Rate this app!](#)

[View code](#)

[Check out other apps](#)

[Want to learn more for free?](#)

Regression Model  
( $R = 0.8581$ ,  $R\text{-squared} = 0.7363$ )





# Diagnostics for simple linear regression

Select a trend:

- ☐ Linear up
- ☐ Linear down
- ☐ Curved up
- ☒ Curved down
- ☐ Fan-shaped

☒ Show residuals

This applet uses ordinary least squares (OLS) to fit a regression line to the data with the selected trend. The applet is designed to help you practice evaluating whether or not the linear model is an appropriate fit to the data. The three diagnostic plots on the lower half of the page are provided to help you identify undesirable patterns in the residuals that may arise from non-linear trends in the data.

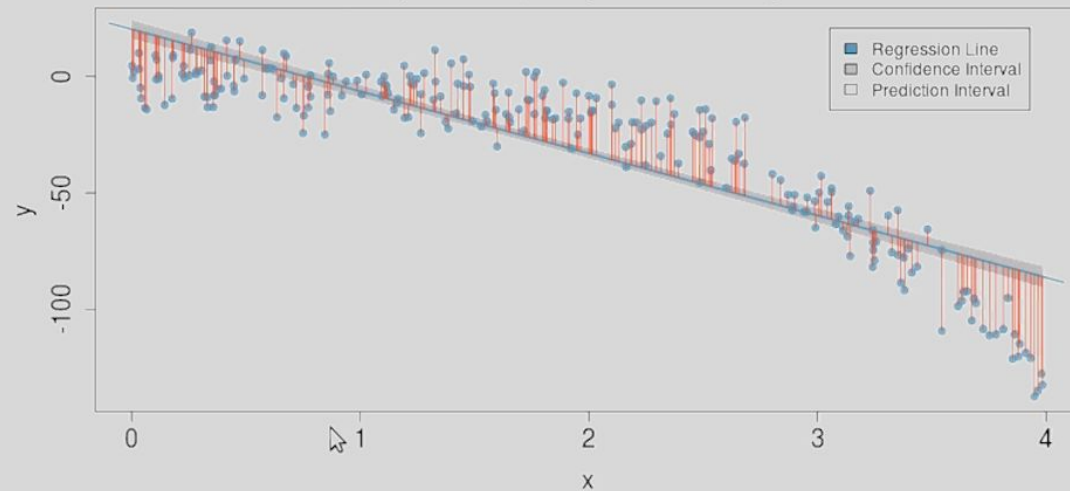
[Rate this app!](#)

[View code](#)

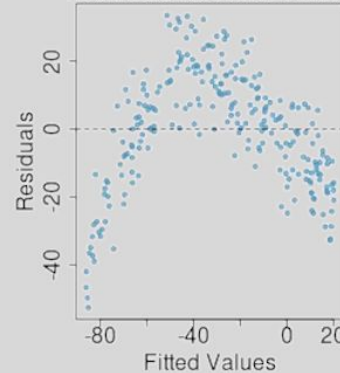
[Check out other apps](#)

[Want to learn more for free?](#)

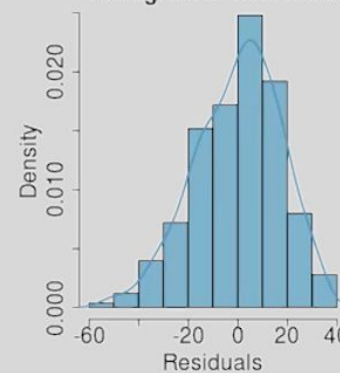
Regression Model  
( $R = -0.874$ ,  $R\text{-squared} = 0.7638$ )



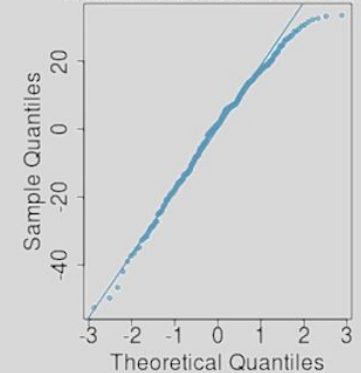
Residuals vs. Fitted Values



Histogram of Residuals



Normal Q-Q Plot of Residuals



# Diagnostics for simple linear regression

Select a trend:

- ☐ Linear up
- ☐ Linear down
- ☐ Curved up
- ☐ Curved down
- ☒ Fan-shaped

☒ Show residuals

This applet uses ordinary least squares (OLS) to fit a regression line to the data with the selected trend. The applet is designed to help you practice evaluating whether or not the linear model is an appropriate fit to the data. The three diagnostic plots on the lower half of the page are provided to help you identify undesirable patterns in the residuals that may arise from non-linear trends in the data.

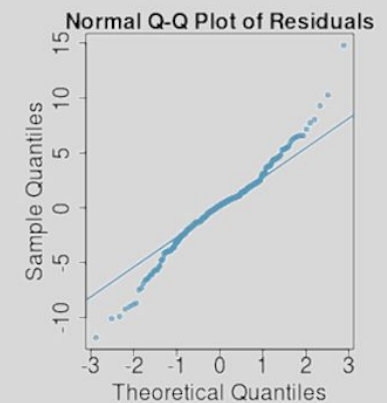
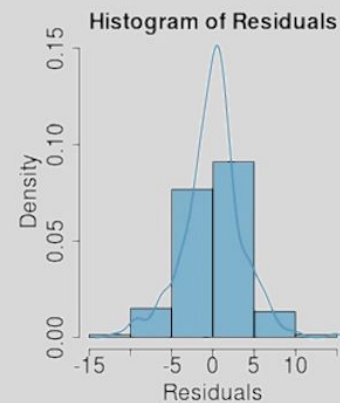
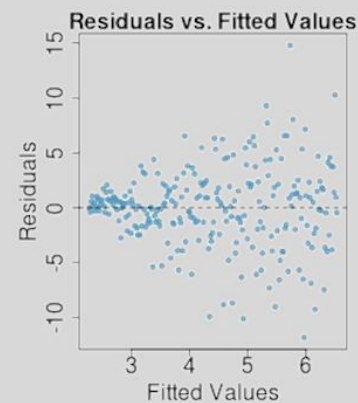
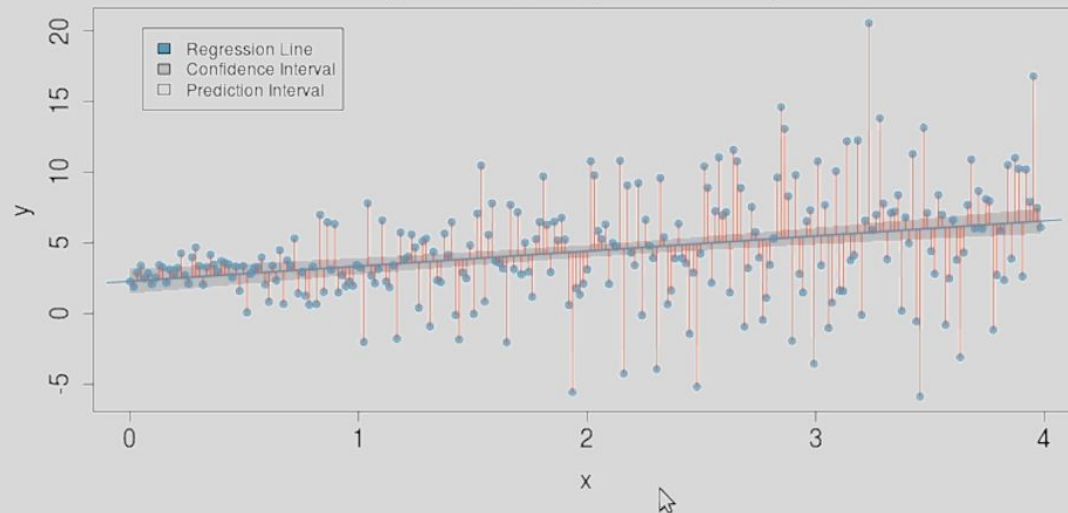
[Rate this app!](#)

[View code](#)

[Check out other apps](#)

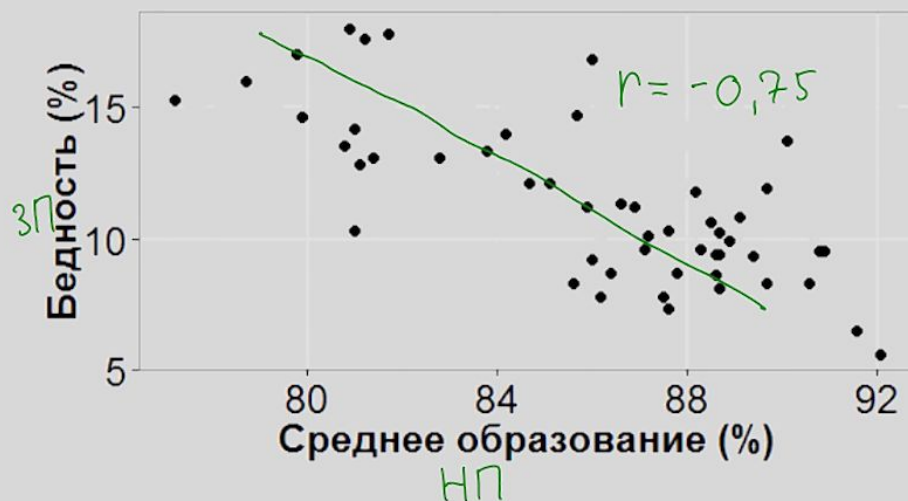
[Want to learn more for free?](#)

Regression Model  
( $R = 0.3231$ ,  $R\text{-squared} = 0.1044$ )



# Регрессионный анализ с одной независимой переменной

Связь бедности и уровня образования



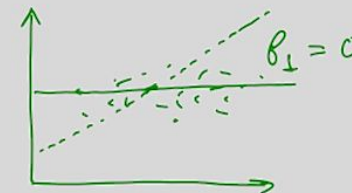
## Descriptive statistics

Statistic	N	Mean	St. Dev.	Min	Max
poverty	51	11.3	3.1	5.6	18.0
hs_grad	51	86.0	3.7	77.2	92.1

1.  $\hat{y} = \underline{\underline{\beta_0}} + \underline{\underline{\beta_1}}x$

2.  $R^2$

$\beta_1 = 0 : H_0$

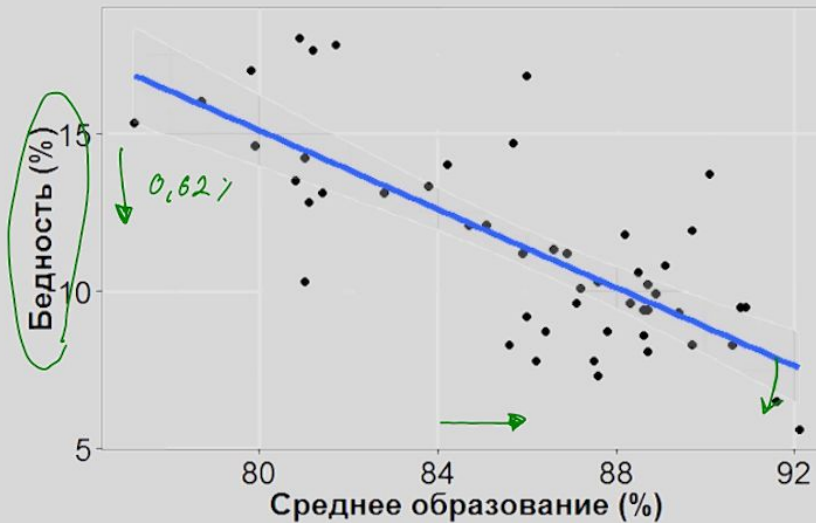


3.  $H_0 \rightarrow Z_{\alpha/2}$





### Связь бедности и уровня образования



$$\hat{y} = 64,78 - 0,62 \cdot \text{hs\_grad}$$

$$H_0 \quad \beta_1 = 0$$

	<u>Estimate</u>	<u>Std. Error</u>	<u>t value</u>	<u>Pr(&gt; t )</u>
(Intercept) $\beta_0$	<u>64.7810</u>	<u>6.8026</u>	<u>9.52</u>	<u>0.0000</u>
hs_grad $\beta_1$	<u>-0.6212</u>	<u>0.0790</u>	<u>-7.86</u>	<u>0.0000</u>

$p < 0,05$

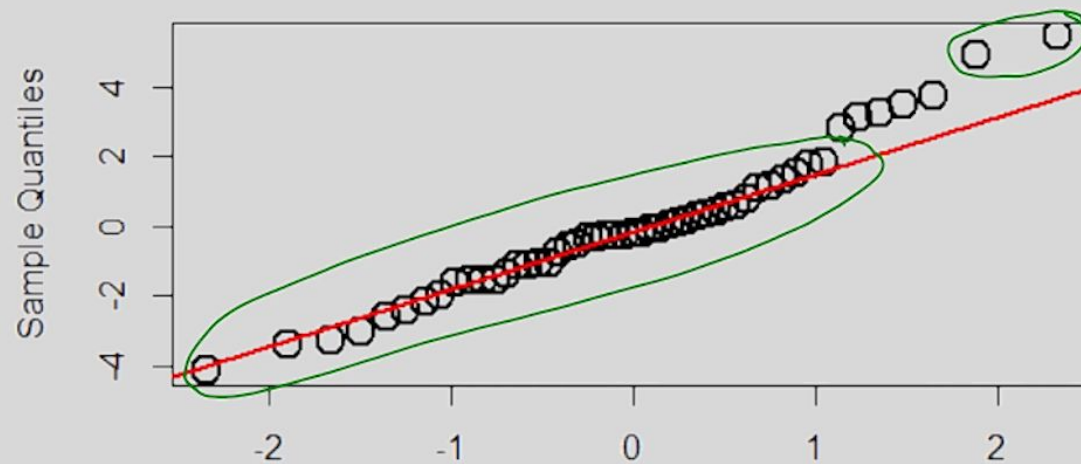
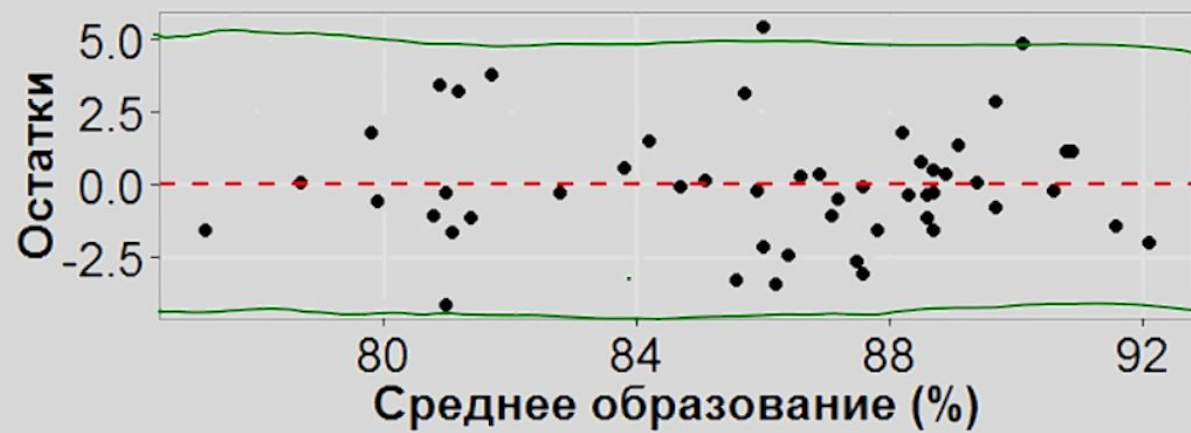
$p < 0,05$

Multiple R-squared: 0.5578

F-statistic (1, 49) = 61.81, p-value < 0.01



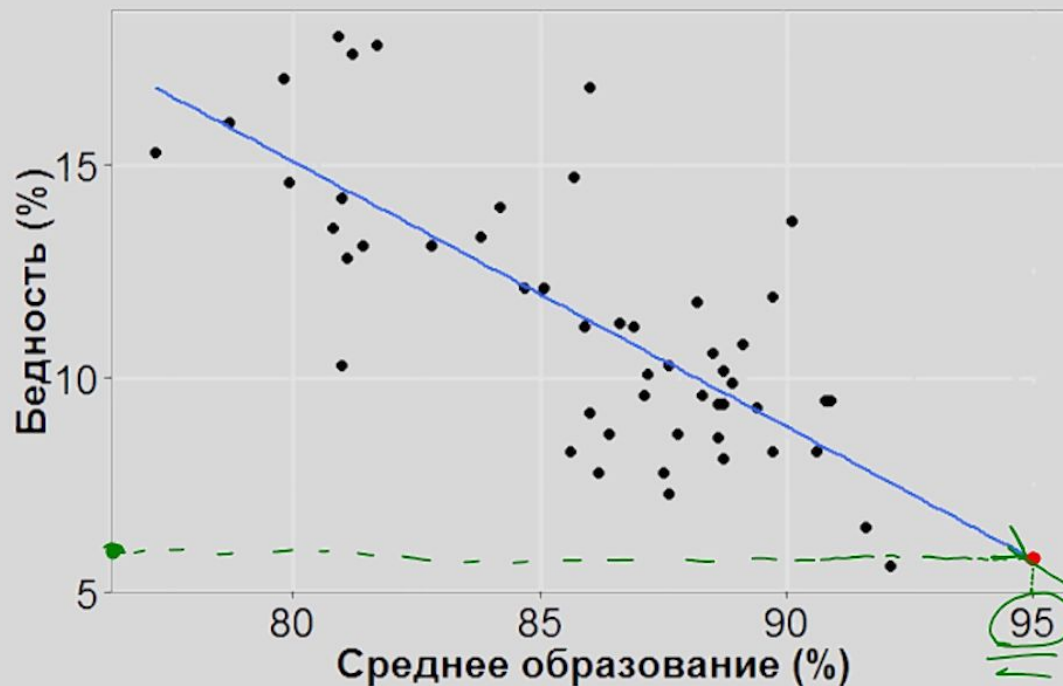
## Анализ остатков



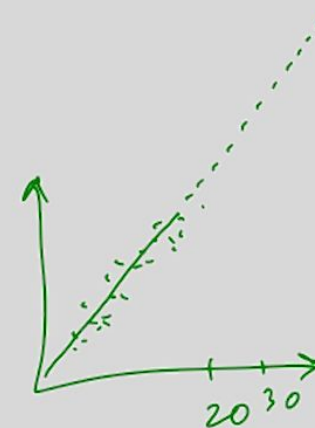
# Предсказание значений

$$\hat{\text{Бедность}} = 64,78 - 0,62 \cdot \text{образование}^{95}$$

Связь бедности и уровня образования



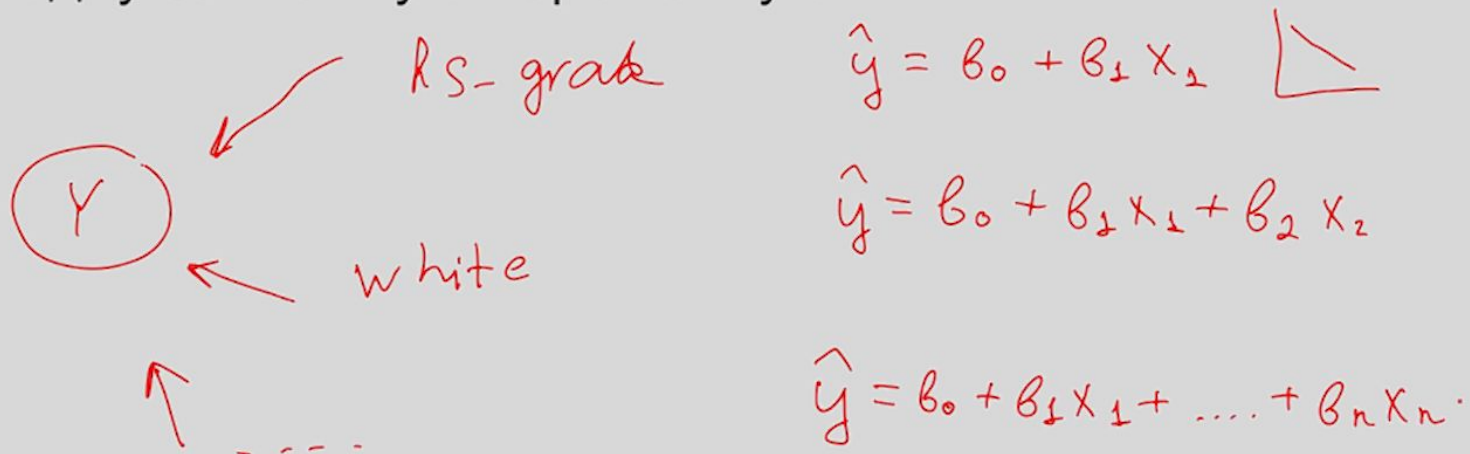
$$\hat{y} = 64,78 - 0,62 \cdot 95 = 5,77$$



# Множественная регрессия

## Multiple regression

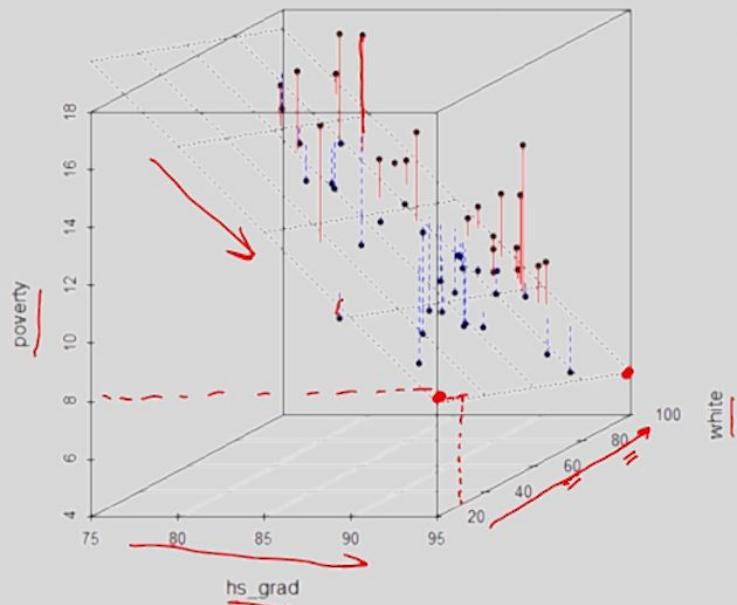
Множественная регрессия позволяет исследовать влияние сразу нескольких независимых переменных на одну зависимую переменную.



# Множественная регрессия

## Multiple regression

Множественная регрессия позволяет исследовать влияние сразу нескольких независимых переменных на одну зависимую переменную.



$$\hat{y} = \underline{b_0} + \underline{b_1} x_1 + \underline{b_2} x_2$$



# Требования к данным

Линейная зависимость переменных

Нормальное распределение остатков

Гетероскедастичность

Проверка на мультиколлинеарность

Нормальное распределение переменных  
(желательно)

3П





# Множественная регрессия

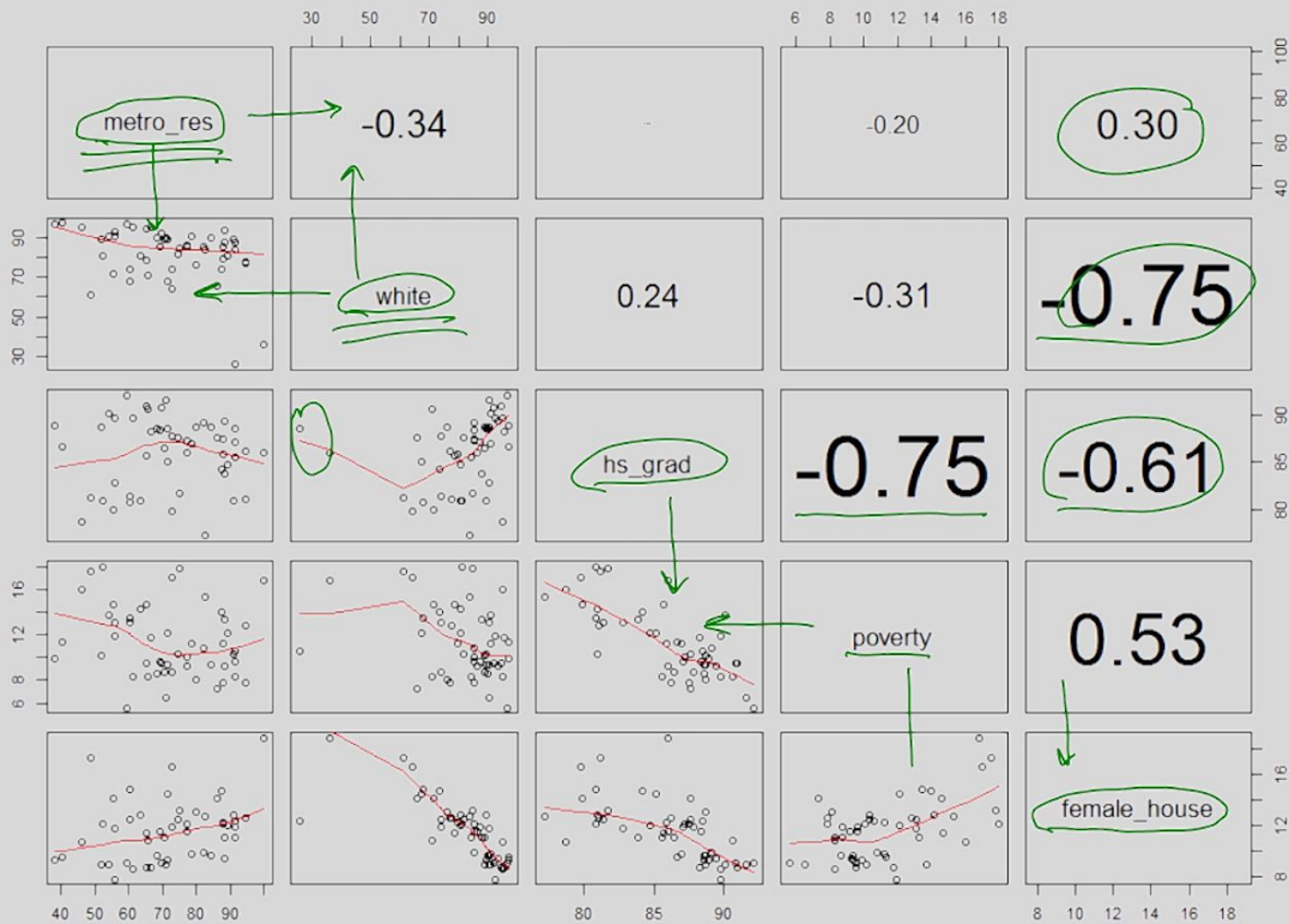
	<u>Estimate</u>	<u>Std. Error</u>	<u>t value</u>	<u>Pr(&gt; t )</u>
(Intercept)	<u>66.47</u>	12.5	5.28	<u>0.0000</u>
<u>metro_res</u>	<u>-0.06</u>	0.02	-2.88	<u>0.0060</u>
white	<u>-0.05</u>	0.03	-1.46	<u>0.1522</u>
hs_grad	<u>-0.55</u>	0.1	-5.29	<u>0.0000</u>
female_house	<u>0.05</u>	0.24	0.21	0.8363

Multiple R-squared: 0.6416, Adjusted R-squared: 0.6104

F-statistic (4, 46) = 20.58, p-value < 0,01

**Исправленный R - квадрат (adjusted R-squared) - скорректированный коэффициент детерминации.**  
Рассчитывается при включении в модель дополнительных независимых переменных.







# Выбор модели

Удаляем	Модель	Adj R-squared
	Poverty ~ hs_grad + white + metro_res + female_house	0,61
female_house	Poverty ~ hs_grad + white + metro_res	0,62
metro_res	Poverty ~ hs_grad + white + female_house	0,54
white	Poverty ~ hs_grad + metro_res + female_house	0,60
hs_grad	Poverty ~ white + metro_res + female_house	0,38

Удаляем	Модель	Adj R-squared
	Poverty ~ hs_grad + white + metro_res	0,62
metro_res	Poverty ~ hs_grad + white	0,55
white	Poverty ~ hs_grad + metro_res	0,57
<u>hs_grad</u>	Poverty ~ white + metro_res	0,17



## Итоговая модель

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	68.7	6.39	10.76	0.0000
white	<u>-0.05</u>	0.02	-2.48	<u>0.0167</u>
metro_res	<u>-0.05</u>	0.02	-2.93	0.0053
hs_grad	<u>-0.57</u>	0.08	-7.57	0.0000

Multiple R-squared: 0.6412, Adjusted R-squared: 0.62

F-statistic (3,47) = 28, p-value < 0,01

$$\hat{Y} = 68,7 - 0,05 \cdot \text{white} - 0,05 \cdot \text{metro\_res} - 0,57 \cdot \text{hs\_grad}$$

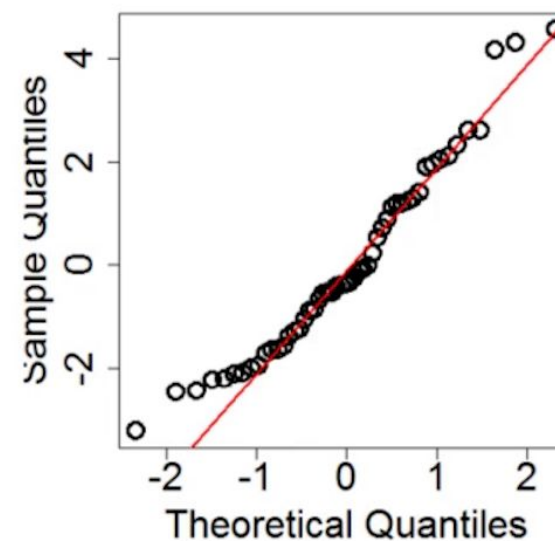
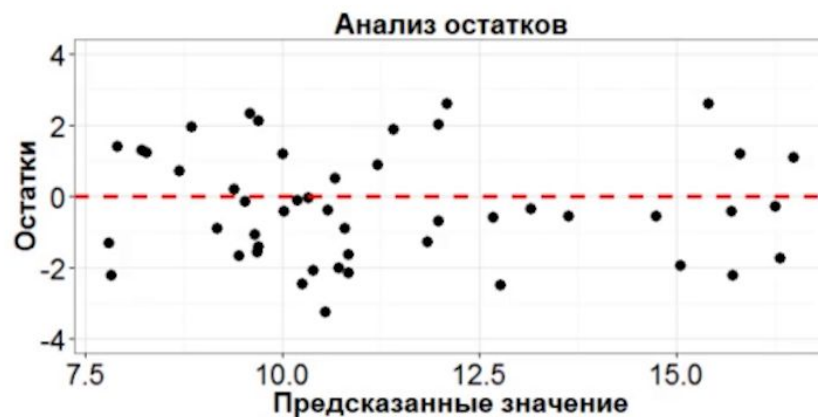


## Итоговая модель

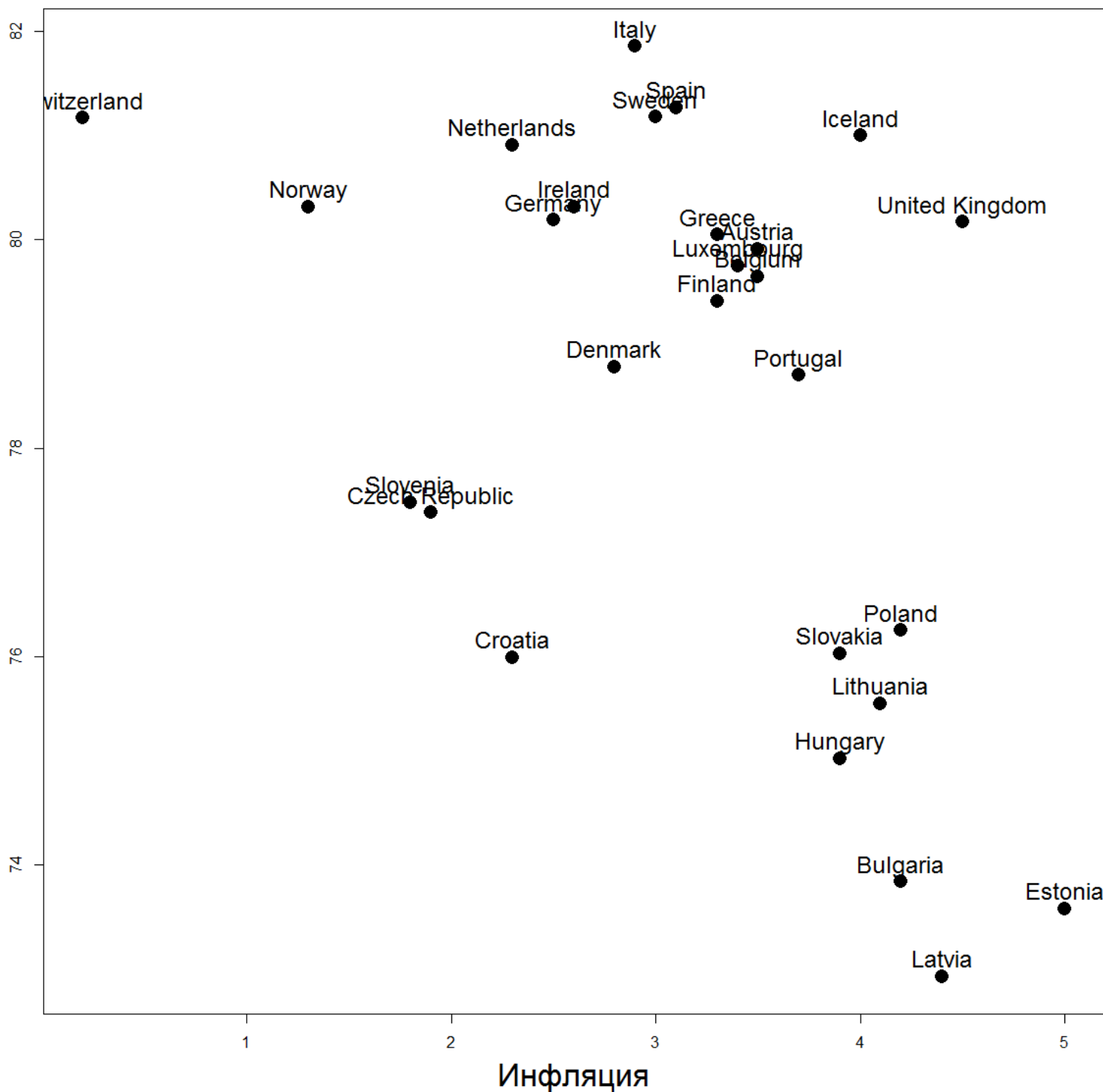
	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	68.7	6.39	10.76	0.0000
white	-0.06	0.02	-2.48	0.0167
metro_res	-0.05	0.02	-2.93	0.0053
hs_grad	-0.57	0.08	-7.57	0.0000

Multiple R-squared: 0.6412, Adjusted R-squared: 0.62

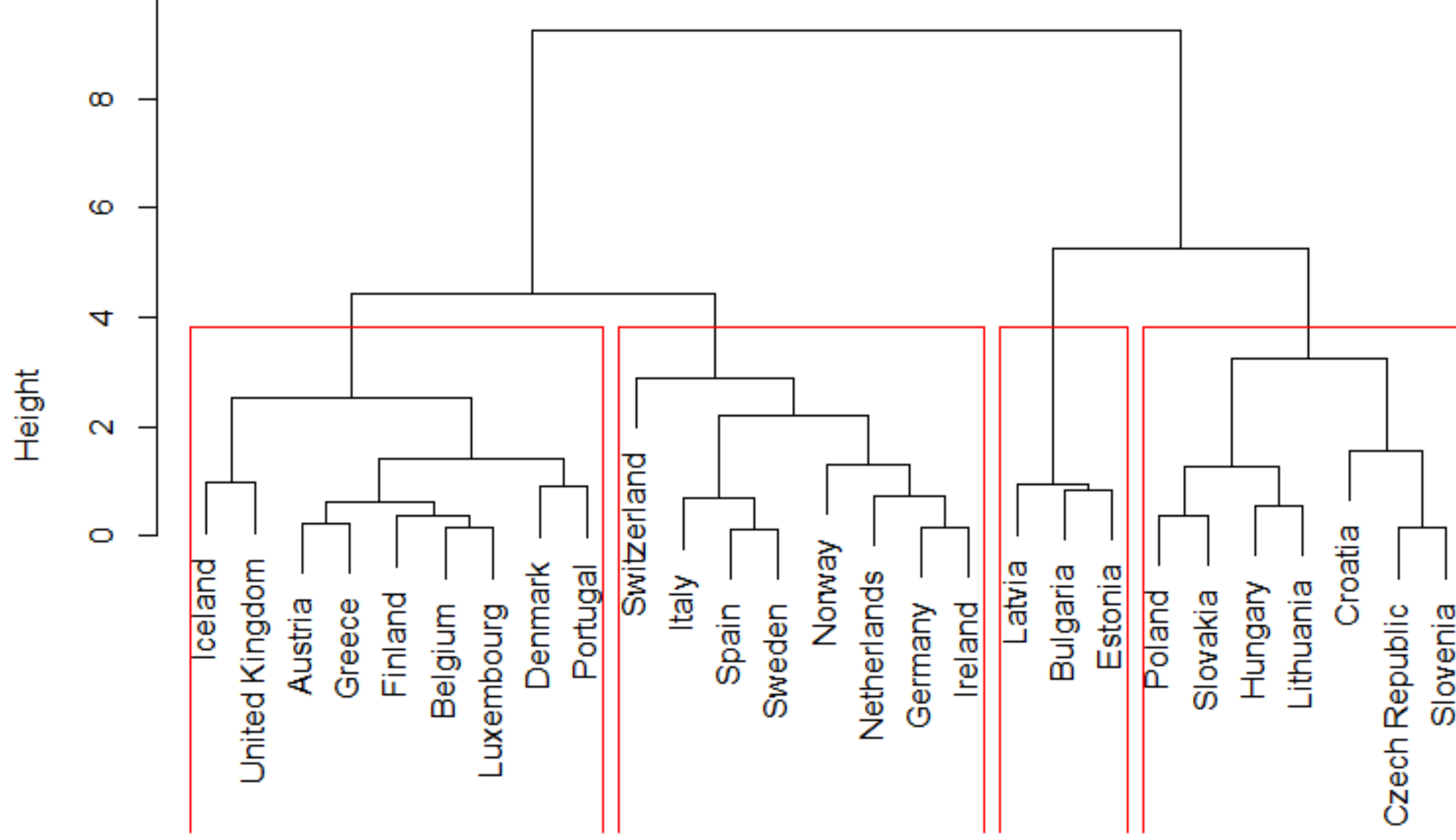
F-statistic (3,47) = 28, p-value < 0,01



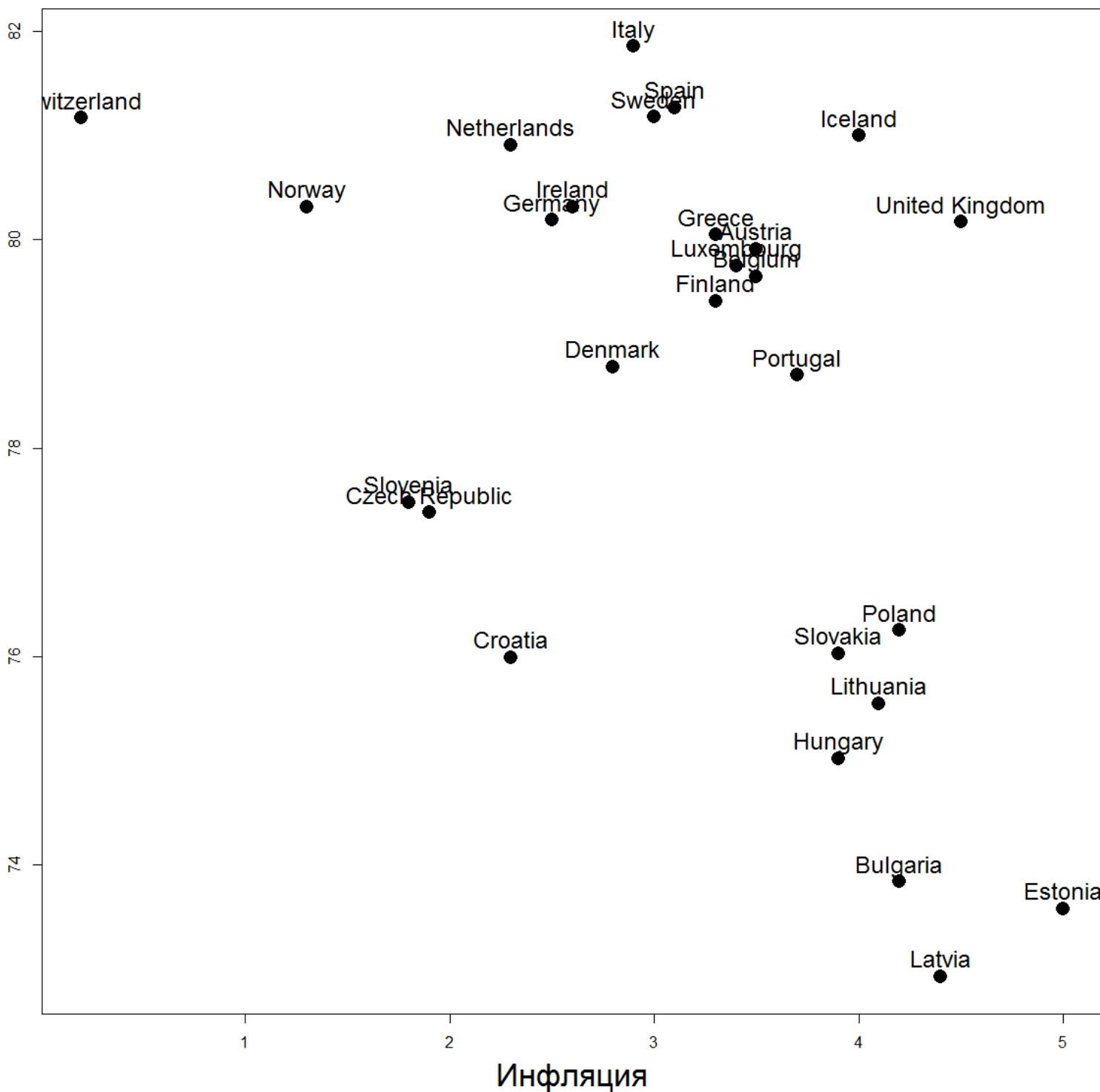
І тродолжителность жизни



# Cluster Dendrogram

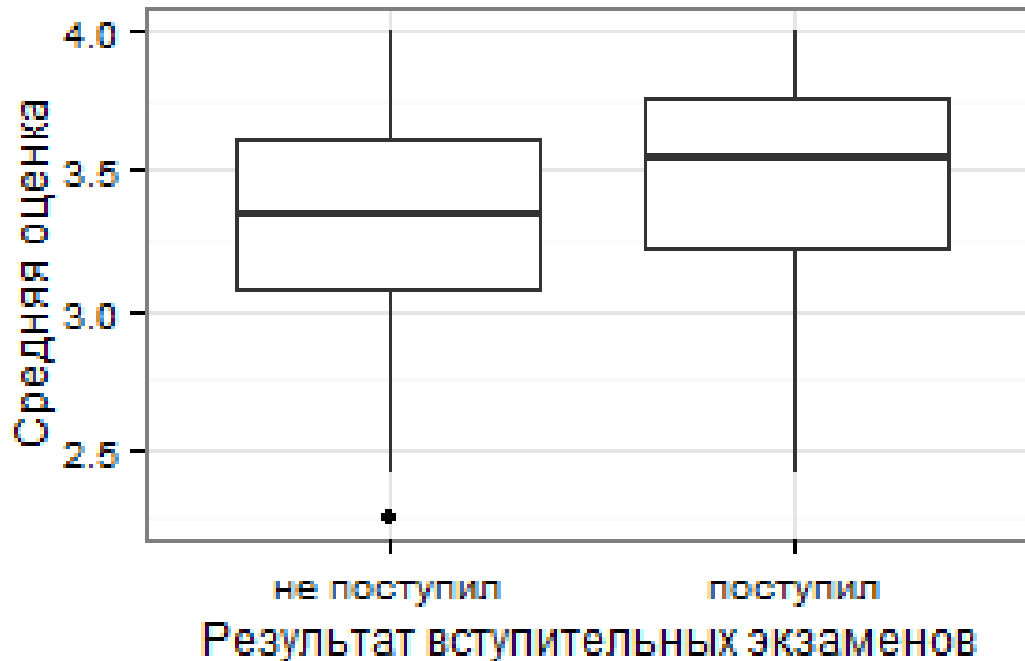


І тродолжителность жизни



# Логистическая регрессия

	admit	mark
1	0	3.61
2	1	3.67
3	1	4.00
4	1	3.19
5	0	2.93
6	1	3.00
7	1	2.98
8	0	3.08
9	1	3.39
10	0	3.92
11	0	4.00
12	0	3.22
13	1	4.00



	Estimate	Std. Error	z value	Pr(> z )
Intercept	-4.3576	1.0353	-4.21	0.0000
mark	1.0511	0.2989	3.52	0.0004