



# КРАТКИЙ ОБЗОР МЕТОДОВ СТАТИСТИЧЕСКОГО АНАЛИЗА КОЛИЧЕСТВЕННЫХ ПЕРЕМЕННЫХ

**МОРДОВСКИЙ ЭДГАР АРТУРОВИЧ**  
**К.М.Н., ДОЦЕНТ**

# ПЛАН

- ОБЩИЕ ТРЕБОВАНИЯ К ВЫПОЛНЕНИЮ СТАТИСТИЧЕСКИХ ТЕСТОВ
- СРАВНЕНИЕ 2-Х СРЕДНИХ ВЕЛИЧИН
- СРАВНЕНИЕ 3-Х И БОЛЕЕ СРЕДНИХ ВЕЛИЧИН
- КОРЕЛЯЦИОННЫЙ АНАЛИЗ
- ОСНОВЫ РЕГРЕССИОННОГО АНАЛИЗА

# ОБЩИЕ ТРЕБОВАНИЯ К ВЫПОЛНЕНИЮ СТАТИСТИЧЕСКИХ ТЕСТОВ

# ПЕРЕМЕННЫЕ

```
graph TD; A[ПЕРЕМЕННЫЕ] --> B[КОЛИЧЕСТВЕННЫЕ]; A --> C[КАТЕГОРИАЛЬНЫЕ]; B --> D[НЕПРЕРЫВНЫЕ  
(CONTINUOUS)]; B --> E[ДИСКРЕТНЫЕ  
(DISCRETE)]; C --> F[ПОРЯДКОВЫЕ  
(ORDINAL)]; C --> G[НОМИНАЛЬНЫЕ  
(NOMINAL)];
```

The diagram is a hierarchical flowchart. At the top is a light orange box labeled 'ПЕРЕМЕННЫЕ'. Two red arrows point down from this box to two intermediate boxes: a dark grey box labeled 'КОЛИЧЕСТВЕННЫЕ' on the left and an orange box labeled 'КАТЕГОРИАЛЬНЫЕ' on the right. From 'КОЛИЧЕСТВЕННЫЕ', two red arrows point down to a blue box labeled 'НЕПРЕРЫВНЫЕ (CONTINUOUS)' and an orange box labeled 'ДИСКРЕТНЫЕ (DISCRETE)'. From 'КАТЕГОРИАЛЬНЫЕ', two red arrows point down to a blue box labeled 'ПОРЯДКОВЫЕ (ORDINAL)' and an orange box labeled 'НОМИНАЛЬНЫЕ (NOMINAL)'. A red dashed rounded rectangle encloses the 'КОЛИЧЕСТВЕННЫЕ' box and its two sub-categories.

## КОЛИЧЕСТВЕННЫЕ

## КАТЕГОРИАЛЬНЫЕ

НЕПРЕРЫВНЫЕ  
(CONTINUOUS)

ДИСКРЕТНЫЕ  
(DISCRETE)

ПОРЯДКОВЫЕ  
(ORDINAL)

НОМИНАЛЬНЫЕ  
(NOMINAL)

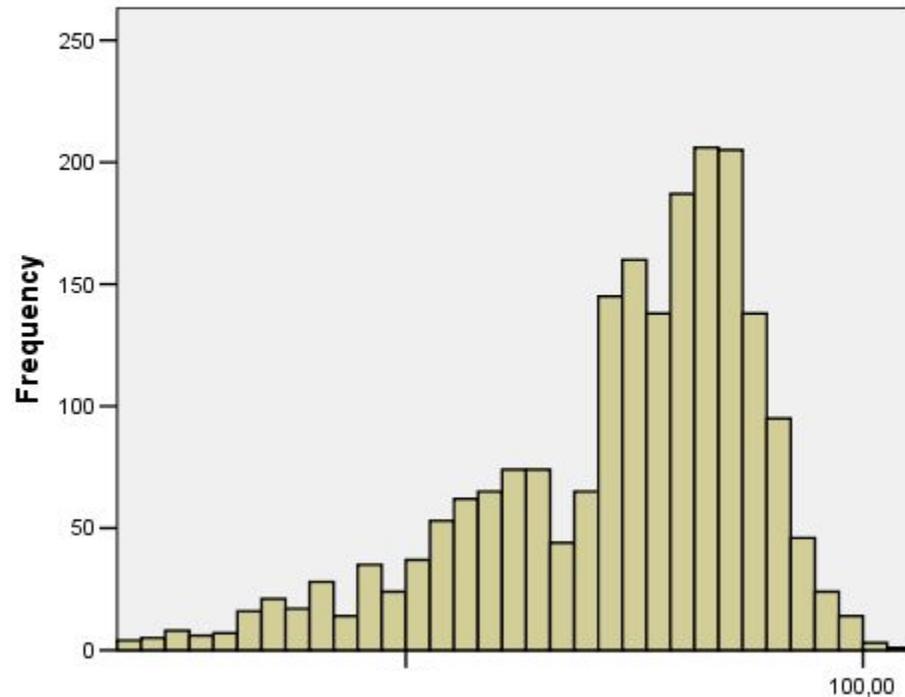
**ТИП ПЕРЕМЕННОЙ ОПРЕДЕЛЯЕТ  
НАБОР МЕТОДОВ СТАТИСТИЧЕСКОГО АНАЛИЗА**

**ПРИМЕР: ФАКТИЧЕСКАЯ СРЕДНЯЯ ПРОДОЛЖИТЕЛЬНОСТЬ ЖИЗНИ В ВЫБОРКЕ МУЖЧИН И ЖЕНЩИН, - ЖИТЕЛЕЙ АРХАНГЕЛЬСКОЙ ОБЛАСТИ, УМЕРШИХ В 2012 Г.**

**ЖЕНЩИНЫ**

**Ы**

for Gender= woman



Mean = 73,2945  
Std. Dev. = 15,37406  
N = 2 021

**X = 73,3 лет** Age

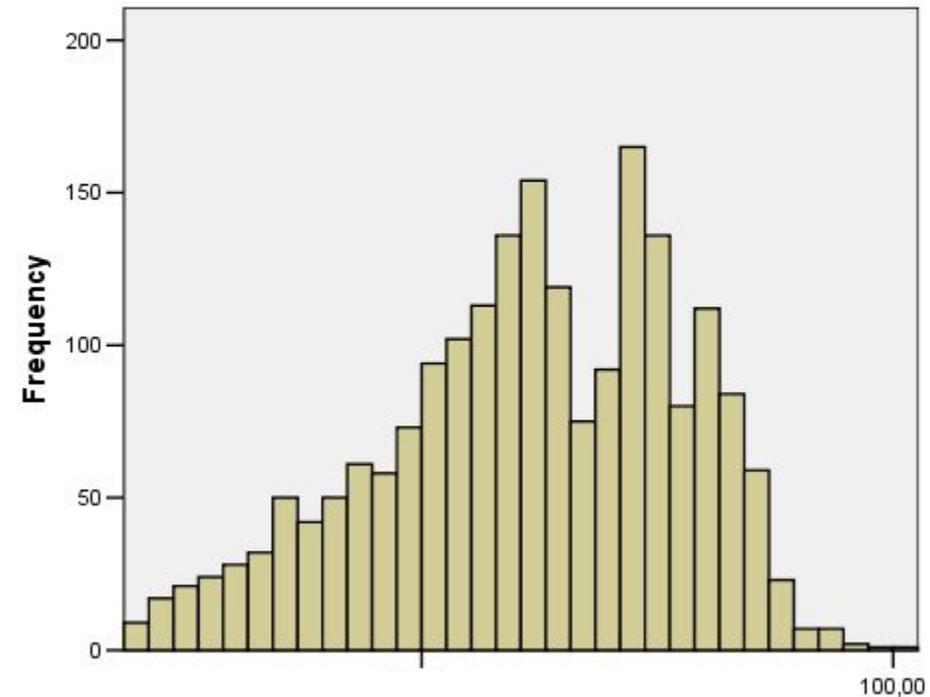
**SD = 15,4**

**N = 2021**

**МУЖЧИНЫ**

**Ы**

for Gender= man



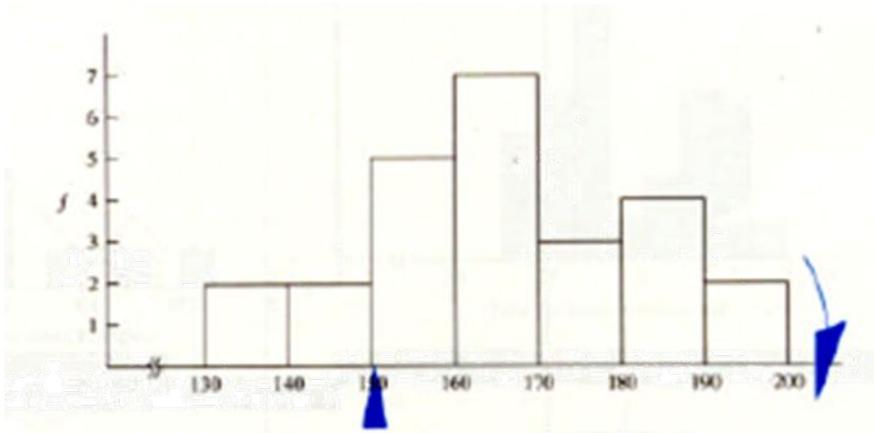
Mean = 61,4049  
Std. Dev. = 15,92564  
N = 2 027

**X = 61,4 лет** Age

**SD = 15,9**

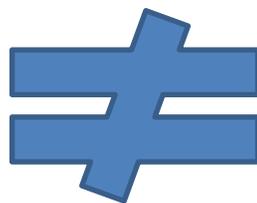
**N = 2027**

**СРЕДНЕЕ  
АРИФМЕТИЧЕСКИЕ  
ДЛЯ ВЫБОРКИ (X / m)**

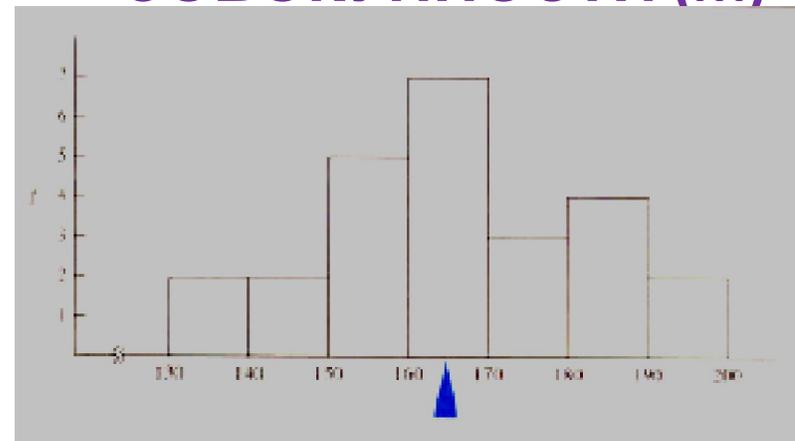


$$\bar{X}_{\text{арифм}} = \frac{\sum x_i}{n},$$

**X (женщины) =  
73,3  
SD = 15,4  
N = 2021**



**СРЕДНЕЕ  
АРИФМЕТИЧЕСКИЕ  
ДЛЯ ГЕНЕРАЛЬНОЙ  
СОВОКУПНОСТИ (M)**



$$\mu = \frac{\sum X}{N}$$

**X (мужчины) =  
61,4  
SD = 15,9  
N = 2027**



# НУЛЕВАЯ И АЛЬТЕРНАТИВНАЯ ГИПОТЕЗА

**ГИПОТЕЗА (HYPOTHESIS)** – *предположение* о свойстве популяции (параметре...)

**ФОРМУЛИРУЕМ ДВЕ ВЗАИМОИСКЛЮЧАЮЩИЕ ГИПОТЕЗЫ:**

ГИПОТЕЗЫ	ФОРМУЛИРОВКА
----------	--------------

**Н<sub>0</sub>** (нулевая гипотеза)

Распределение признака **СЛУЧАЙНОЕ**  
(категориальные переменные)  
**НЕТ** отличий в сравниваемых величинах  
(количественные непрерывные переменные)

**Н<sub>а</sub>** (альтернативная гипотеза)

Распределение признака **НЕСЛУЧАЙНОЕ**  
(категориальные переменные)  
**ЕСТЬ** отличия в сравниваемых величинах  
(количественные непрерывные переменные)

# НУЛЕВАЯ И АЛЬТЕРНАТИВНАЯ ГИПОТЕЗА

ГИПОТЕЗЫ	ФОРМУЛИРОВКА
<b>H<sub>0</sub></b> (нулевая гипотеза)	Распределение признака <b>СЛУЧАЙНОЕ</b> <b>НЕТ</b> отличий в сравниваемых величинах
<b>H<sub>a</sub></b> (альтернативная гипотеза)	Распределение признака <b>НЕСЛУЧАЙНОЕ</b> <b>ЕСТЬ</b> отличия в сравниваемых величинах

**X (женщины) = 73,3 года**  
**SD = 15,4**  
**N = 2021**

**X (мужчины) = 61,4**  
**года**  
**SD = 15,9**

**N = 2027** ФОРМУЛИРОВКА

ГИПОТЕЗЫ	ФОРМУЛИРОВКА
<b>H<sub>0</sub></b> (нулевая гипотеза)	<b>X (женщины) = X (мужчины)</b> средняя продолжительность жизни женщин <b>НЕ</b> отличается от средней продолжительности жизни мужчин <b>(т.е. 73,3 = 61,4 в популяции)</b>
<b>H<sub>a</sub></b> (альтернативная гипотеза)	<b>X (женщины) ≠ X (мужчины)</b> средняя продолжительность жизни женщин <b>ОТЛИЧАЕТСЯ</b> от средней продолжительности жизни мужчин <b>(т.е. 73,3 ≠ 61,4 в популяции)</b>

## 2 ВИДА АЛЬТЕРНАТИВНЫХ ГИПОТЕЗ

ГИПОТЕЗЫ	ФОРМУЛИРОВКА
Двусторонняя альтернатива (two-tailed hypothesis)	$H_0: X \text{ (женщины)} = X \text{ (мужчины)}$ $H_a: X \text{ (женщины)} \neq X \text{ (мужчины)}$
Односторонняя альтернатива (one-tailed hypothesis)	$H_0: X \text{ (женщины)} \geq X \text{ (мужчины)}$ $H_a: X \text{ (женщины)} < X \text{ (мужчины)}$

ЕСТЬ КАКИЕ-ТО ДОПОЛНИТЕЛЬНЫЕ СВЕДЕНИЯ / ГИПОТЕЗЫ  
(ИСПОЛЬЗУЕТСЯ РЕДКО)

# ТЕСТИРОВАНИЕ ГИПОТЕЗ

## ИСТИНА

***H<sub>0</sub> - ВЕРНА***

***H<sub>a</sub> - ВЕРНА***

**МЫ ПРИНИМАЕМ**

***H<sub>0</sub>***

**ПРАВИЛЬНЫЙ  
РЕЗУЛЬТАТ**

это чувствительность  
теста  
(1- $\alpha$ )

**ОШИБКА 2 ТИПА ( $\beta$ )**

(вероятность НЕ найти  
то, чего ЕСТЬ)

**МЫ ОТВЕРГАЕМ *H<sub>0</sub>*  
(ПРИНИМАЕМ *H<sub>a</sub>*)**

**ОШИБКА 1 ТИПА ( $\alpha$ )**  
(уровень значимости –  
significance (Sig.)  
“p”

(вероятность найти то,  
чего НЕТ)

**ПРАВИЛЬНЫЙ  
РЕЗУЛЬТАТ**

это «мощность теста»  
(1- $\beta$ )

# СТАТИСТИЧЕСКАЯ ОБРАБОТКА ДАННЫХ

**СТАТИСТИЧЕСКАЯ ОБРАБОТКА ДАННЫХ (методы статистического анализа)** – математические расчеты, позволяющие оценить **ВЕРОЯТНОСТЬ ОШИБКИ 1 ТИПА (p / significance (Sig.))**

**СТАТИСТИЧЕСКАЯ ОБРАБОТКА ДАННЫХ (методы статистического анализа)** – математические расчеты, результаты которых позволяют с определенной долей вероятности принять нулевую гипотезу (accept) или ее отвергнуть (reject)

**«Приемлемая» вероятность ошибки 1 типа ( $\alpha$ -ошибки) = 0.05 (5%)**

**«КОНСЕНСУС ФИШЕРА»**

**ЭТО ОТНОСИТЕЛЬНАЯ ВЕЛИЧИНА !!!!!!!!!**

# СТАТИСТИЧЕСКАЯ ОБРАБОТКА ДАННЫХ

## ПОЧЕМУ ВАЖНО ???

А) ПРАВИЛЬНО **РАССЧИТАТЬ ОБЪЕМ  
ВЫБОРКИ** ДО НАЧАЛА ИССЛЕДОВАНИЯ  
???



ЧТОБЫ МИНИМИЗИРОВАТЬ  
ВЕРОЯТНОСТЬ **ОШИБКИ 1 ТИПА**

Б) ПРАВИЛЬНО **СФОРМИРОВАТЬ  
ВЫБОРКУ**  
И ПРАВИЛЬНО **ВЫБРАТЬ  
СТАТИСТИЧЕСКИЙ МЕТОД** АНАЛИЗА  
(СТАТИСТИЧЕСКИЙ КРИТЕРИЙ)



ЧТОБЫ МИНИМИЗИРОВАТЬ  
ВЕРОЯТНОСТЬ **ОШИБКИ 2 ТИПА**

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**1 ЭТАП:  
ФОРМУЛИРУЕМ  $H_0$  и  
 $H_a$**

## ГИПОТЕЗЫ

## ФОРМУЛИРОВКА

**$H_0$**  (нулевая гипотеза)

**$X$  (женщины) =  $X$  (мужчины)**

средняя продолжительность жизни женщин **НЕ** отличается от средней продолжительности жизни мужчин

**$H_a$**  (альтернативная гипотеза)

**$X$  (женщины)  $\neq$   $X$  (мужчины)**

средняя продолжительность жизни женщин **ОТЛИЧАЕТСЯ** от средней продолжительности жизни мужчин

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

## 2 ЭТАП:

ОПРЕДЕЛЯЕМ УСЛОВИЯ,  
ПРИ КОТОРЫХ ПРИМЕМ  
**Ha** (ОТВЕРГНЕМ **Ho**)

**БУДЕМ** считать результаты теста «статистически значимыми» (т.е. примем **Ha**) при вероятности ошибки 1 типа ( $\alpha$ -ошибки) менее 0.05 (5%)  
**«КОНСЕНСУС ФИШЕРА»**

$p < 0.05$  «достаточно», если имеем дело с социологическими исследованиями, «ориентировочными» исследованиями, «пилотными» исследованиями

В клинических испытаниях “ $p$ ” устанавливается индивидуально (в зависимости от клинической значимости искомого результата) – в т.ч. устанавливается в «SD»

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**3 ЭТАП:  
ВЫБОР  
СТАТИСТИЧЕСКОГО  
КРИТЕРИЯ / МЕТОДА**

**ВЫБОР СТАТИСТИЧЕСКОГО КРИТЕРИЯ ОПРЕДЕЛЯЕТСЯ НАБОРОМ ПАРАМЕТРОВ !!!**

**4 ЭТАП:  
МАТЕМАТИЧЕСКИЕ  
РАССЧЕТЫ**

**СТРОГО ИНДИВИДУАЛЬНО**

**СТАТИСТИЧЕСКИЕ ПРОГРАММЫ  
(IBM SPSS, STATA, STATISTICA, PASW, R)**

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**5 ЭТАП:  
ИНТЕРПРЕТАЦИЯ  
РЕЗУЛЬТАТОВ**

**ПРИНИМАЕМ  $H_0$  / ОТВЕРГАЕМ  $H_a$  (если " $p$ " < 0.05)**

**ПРИНИМАЕМ  $H_a$  / ОТВЕРГАЕМ  $H_0$  (если " $p$ "  $\geq$  0.05)**

**+ ОЦЕНИВАЕМ ВОЗМОЖНОСТЬ ЭКСТРАПОЛЯЦИИ ПОЛУЧЕННЫХ РЕЗУЛЬТАТОВ**

**НА ГЕНЕРАЛЬНУЮ СОВОКУПНОСТЬ**

**+ ОЦЕНИВАЕМ СТАТИСТИЧЕСКУЮ МОЩНОСТЬ РЕЗУЛЬТАТА**

**+ ОЦЕНИВАЕМ ПРАКТИЧЕСКУЮ ЗНАЧИМОСТЬ РЕЗУЛЬТАТОВ**

# СРАВНЕНИЕ 2-Х СРЕДНИХ ВЕЛИЧИН

# ПЕРЕМЕННЫЕ

## КОЛИЧЕСТВЕННЫЕ

НЕПРЕРЫВН  
ЫЕ  
(CONTINUOUS)

ДИСКРЕТНЫ  
Е  
(DISCRETE)

## КАТЕГОРИАЛЬНЫЕ

ПОРЯДКОВЫ  
Е  
(ORDINAL)

НОМИНАЛЬНЫ  
Е  
(NOMINAL)

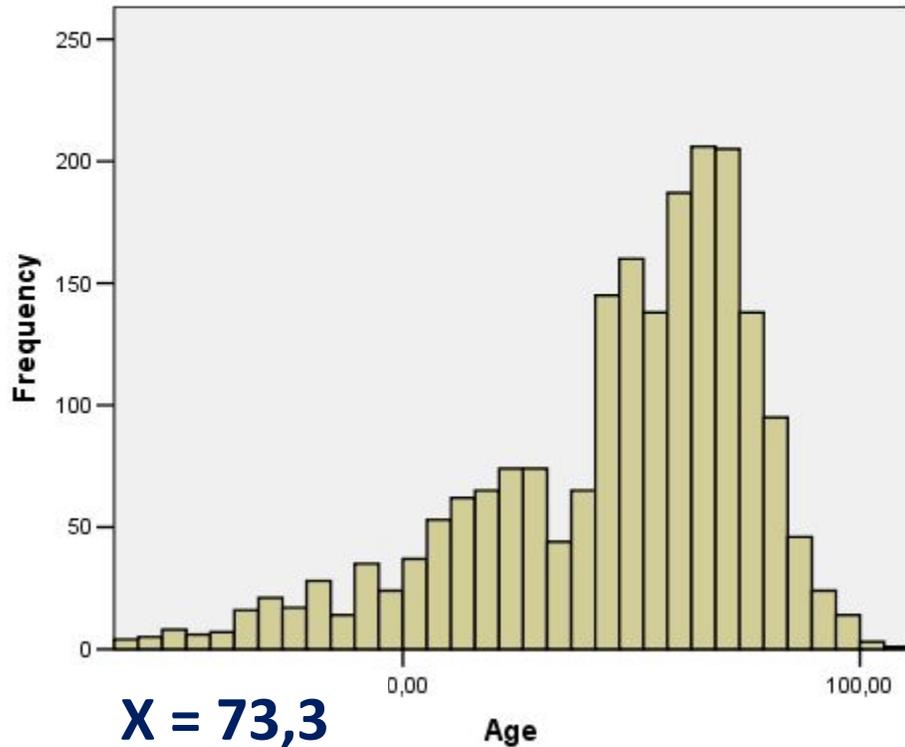
**СРЕДНИЕ ВЕЛИЧИНЫ МОЖНО ВЫЧИСЛИТЬ ТОЛЬКО ДЛЯ  
КОЛИЧЕСТВЕННЫХ НЕПРЕРЫВНЫХ ВЕЛИЧИН**

# ПРИМЕР: СРЕДНЯЯ ПРОДОЛЖИТЕЛЬНОСТЬ ЖИЗНИ В ВЫБОРКЕ МУЖЧИН И ЖЕНЩИН, - ЖИТЕЛЕЙ АРХАНГЕЛЬСКОЙ ОБЛАСТИ, УМЕРШИХ В

## ЖЕНЩИНЫ

Ы

for Gender= woman



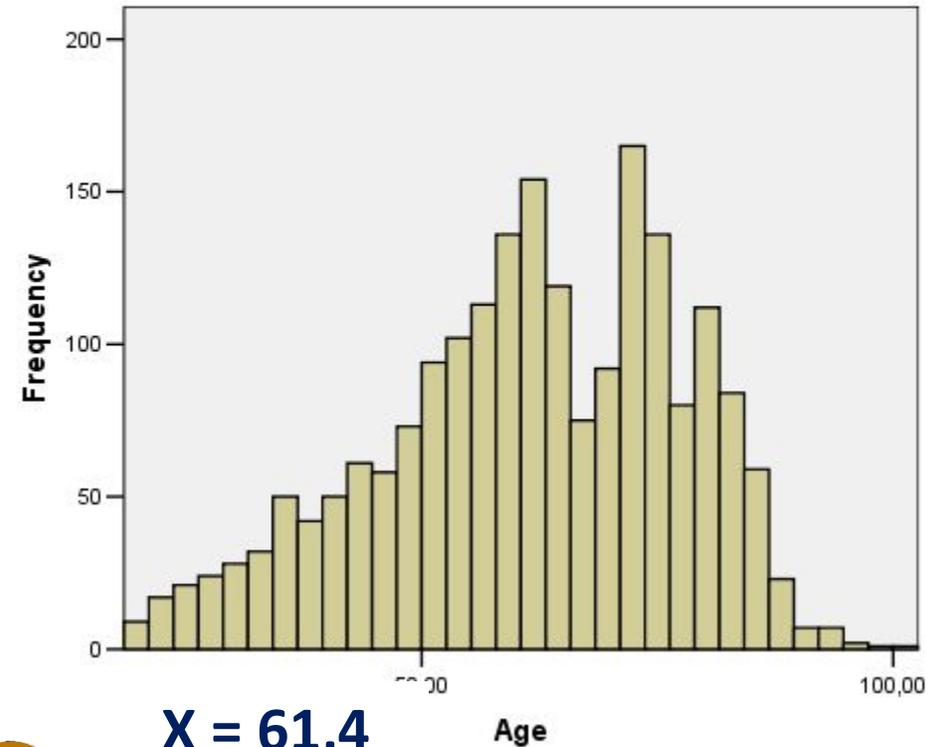
Mean = 73,2945  
Std. Dev. = 15,37406  
N = 2 021

$X = 73,3$   
 $SD = 15,4$   
 $N = 2021$

## МУЖЧИНЫ

Ы

for Gender= man



Mean = 61,4049  
Std. Dev. = 15,92564  
N = 2 027

$X = 61,4$   
 $SD = 15,9$   
 $N = 2027$



# ВЫБОР КОНКРЕТНОГО СТАТИСТИЧЕСКОГО МЕТОДА ПРИ СРАВНЕНИИ СРЕДНИХ ВЕЛИЧИН ОПРЕДЕЛЯЕТСЯ:

## УСЛОВИЕ

1	КОЛИЧЕСТВО СРАВНИВАЕМЫХ ГРУПП	2 / 3+
2	РАСПРЕДЕЛЕНИЕ ПРИЗНАКА <b>В КАЖДОЙ</b> ИЗ СРАВНИВАЕМЫХ ГРУПП	нормальное или скошенное
3	ТИП ВЫБОРКИ	зависимые выборки («до и после») / независимые выборки (простое сравнение)

**ПОПРАВКА БОНФЕРРОНИ: 2 / 3+ групп**

**ГОМОГЕННОСТЬ / ГОМОСКЕДАСТИЧНОСТЬ ДИСПЕРСИИ: НЕ КРИТИЧНОЕ ТРЕБОВАНИЕ; ПРИ РАВЕНСТВЕ ОБЪЕМОВ ВЫБОРОК «ПОЧТИ НЕКРИТИЧНОЕ»**

# СРАВНЕНИЕ 2-Х СРЕДНИХ ВЕЛИЧИН

	<b>НОРМАЛЬНОЕ РАСПРЕДЕЛЕНИЕ В КАЖДОЙ ИЗ СРАВНИВАЕМЫХ ВЫБОРОК (РАЗНИЦЫ ПРИЗНАКА В ПАРАХ ДО-ПОСЛЕ)</b>	<b>СКОШЕННОЕ РАСПРЕДЕЛЕНИЕ В 1 ИЛИ 2 СРАВНИВАЕМЫХ ВЫБОРКАХ (РАЗНИЦЫ ПРИЗНАКА В ПАРАХ ДО- ПОСЛЕ)</b>
<b>НЕЗАВИСИМЫЕ ВЫБОРКИ</b>	<b>Independent Samples T-test (Student T-test) тест Стьюдента</b> <i><u>для независимых выборок</u></i>	<b>2-Independent Samples test (Mann-Whitney U test) тест Манна-Уитни</b> <i><u>для независимых выборок</u></i>
<b>ЗАВИСИМЫЕ ВЫБОРКИ (ПОВТОРНЫЕ ИЗМЕРЕНИЯ)</b>	<b>Dependent (Paired Samples) T-test тест Стьюдента</b> <i><u>для парных выборок</u></i>	<b>2-Related Samples test (Wilcoxon signed-rank test) тест Вилкоксона</b> <i><u>для парных выборок</u></i>

# Independent Samples T-test (Student test)

## T-тест Стьюдента

### ASSUMPTIONS / УСЛОВИЯ ПРИМЕНЕНИЯ

### КАК ПРОВЕРИТЬ?

1. Сравниваем 2 выборки

см. характеристики собранных данных

2. Выборки д.б. независимыми

см. характеристики собранных данных

3. Количественный непрерывный тип данных в каждой из сравниваемых выборок

см. тип данных

4. Нормальное распределение изучаемого признака в каждой из выборок

Test Shapiro-Wilk / Kolmogorov-Smirnov

5. Равенство дисперсий

**Levene's test for Equality of Variances**  
(sig. (p)  $\geq 0,05$ )

**Ho:  $v_1 = v_2$**

**Ha:  $v_1 \neq v_2$**

**Если дисперсии не равны ( $p < 0,05$ )**

**= проблема БЕРЕНСА-ФИШЕРА**

# 2-Independent Samples test (Mann-Whitney U test)

## U-тест Манна-Уитни

### ASSUMPTIONS / УСЛОВИЯ ПРИМЕНЕНИЯ

1. Сравниваем 2 выборки
2. Выборки д.б. независимыми
3. Количественный непрерывный тип данных в каждой из сравниваемых выборок
4. Скошенное распределение данных в одной или обеих сравниваемых выборок

### КАК ПРОВЕРИТЬ?

- см. характеристики собранных данных
- см. характеристики собранных данных
- см. тип данных

***ВНИМАНИЕ: несмотря на то, что распределение скошенное, тест Манна-Уитни оценивает именно СРЕДНИЕ АРИФМЕТИЧЕСКИЕ, А НЕ МЕДИАНЫ !!!***

Test Shapiro-Wilk / Kolmogorov-Smirnov

**ДИСПЕРСИЯ НЕ ПРОВЕРЯЕТСЯ**

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**1 ЭТАП:**  
ФОРМУЛИРУЕМ  $H_0$  и  $H_a$

ГИПОТЕЗЫ	ФОРМУЛИРОВКА
<b><math>H_0</math></b> (нулевая гипотеза)	<b><math>X</math> (женщины) = <math>X</math> (мужчины)</b> средняя продолжительность жизни женщин <b>НЕ</b> отличается от средней продолжительности жизни мужчин
<b><math>H_a</math></b> (альтернативная гипотеза)	<b><math>X</math> (женщины) <math>\neq</math> <math>X</math> (мужчины)</b> средняя продолжительность жизни женщин <b>ОТЛИЧАЕТСЯ</b> от средней продолжительности жизни мужчин

**2 ЭТАП:**  
ОПРЕДЕЛЯЕМ УСЛОВИЯ,  
ПРИ КОТОРЫХ ПРИМЕМ  
 **$H_a$**  (ОТВЕРГНЕМ  **$H_0$** )

**БУДЕМ считать результаты теста «статистически значимыми» (т.е. примем  $H_a$ ) при вероятности ошибки 1 типа ( $\alpha$ -ошибки) менее 0.05 (5%)**

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**3 ЭТАП:  
ВЫБОР  
СТАТИСТИЧЕСКОГО  
КРИТЕРИЯ / МЕТОДА**

Tests of Normality

	Gender	Kolmogorov-Smirnov <sup>a</sup>			Shapiro-Wilk		
		Statistic	df	Sig.	Statistic	df	Sig.
Age	woman	,118	2021	,000	,923	2021	,000
	man	,066	2027	,000	,978	2027	,000

a. Lilliefors Significance Correction

**Н<sub>0</sub>: РАСПРЕДЕЛЕНИЕ В ВЫБОРКЕ НЕ ОТЛИЧАЕТСЯ ОТ НОРМАЛЬНОГО**

**Н<sub>а</sub>: РАСПРЕДЕЛЕНИЕ В ВЫБОРКЕ ОТЛИЧАЕТСЯ ОТ НОРМАЛЬНОГО**

## ИСТИНА

	ИСТИНА	
	<b>Н<sub>0</sub> - ВЕРНА</b>	<b>Н<sub>а</sub> - ВЕРНА</b>
<b>МЫ ПРИНИМАЕМ Н<sub>0</sub></b>	<b>ПРАВИЛЬНЫЙ РЕЗУЛЬТАТ</b> (= чувствительность теста)	<b>ОШИБКА 2 ТИПА</b> (вероятность НЕ найти то, чего ЕСТЬ)
<b>МЫ ОТВЕРГАЕМ Н<sub>0</sub></b> (ПРИНИМАЕМ Н <sub>а</sub> )	<b>ОШИБКА 1 ТИПА</b> (уровень значимости – significance (Sig.) “p” (вероятность найти то, чего НЕТ)	<b>ПРАВИЛЬНЫЙ РЕЗУЛЬТАТ</b> (= мощность теста)

**$p$  (женщины) < 0,0001**

**$p$  (мужчины) < 0,0001**

т.е. **МОЖЕМ** принять **Н<sub>а</sub>**  
вероятность ошибки 1 типа  
(ошибочно принять Н<sub>а</sub> - найти то, чего нет) < 0,1%

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**3 ЭТАП:  
ВЫБОР  
СТАТИСТИЧЕСКОГО  
КРИТЕРИЯ / МЕТОДА**

Tests of Normality

	Gender	Kolmogorov-Smirnov <sup>a</sup>			Shapiro-Wilk		
		Statistic	df	Sig.	Statistic	df	Sig.
Age	woman	,118	2021	,000	,923	2021	,000
	man	,066	2027	,000	,978	2027	,000

a. Lilliefors Significance Correction

ASSUMPTIONS / УСЛОВИЯ ПРИМЕНЕНИЯ	КАК ПРОВЕРИТЬ?
1. Сравниваем 2 выборки	см. характеристики собранных данных
2. Выборки д.б. независимыми	см. характеристики собранных данных
3. Количественный непрерывный тип данных в каждой из сравниваемых выборок	см. тип данных
4. Скошенное распределение данных в одной или обеих сравниваемых выборок	Test Shapiro-Wilk / Kolmogorov-Smirnov
<p><b>ВНИМАНИЕ: несмотря на то, что распределение скошенное, тест Манна-Уитни сравнивает именно СРЕДНИЕ АРИФМЕТИЧЕСКИЕ, А НЕ МЕДИАНЫ !!!</b></p> <p><b>ДИСПЕРСИЯ НЕ ПРОВЕРЯЕТСЯ</b></p>	

**2-Independent Samples test (Mann-Whitney U test)  
U-тест Манна-Уитни**

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

## 4 ЭТАП: МАТЕМАТИЧЕСКИЕ РАССЧЕТЫ

формулируем  $H_0$  и  $H_a$  для теста Манна-Уитни

**$H_0$ :**  $m_1 = m_2$  (средняя продолжительность жизни мужчин не отличается от средней продолжительности жизни женщин)

**$H_a$ :**  $m_1 \neq m_2$  (средняя продолжительность жизни мужчин отличается от средней продолжительности жизни женщин)

Test Statistics<sup>a</sup>

	Age
Mann-Whitney U	1144664
Wilcoxon W	3200042
Z	-24,305
Asymp. Sig. (2-tailed)	,000

a. Grouping Variable: Gender

$$p < 0,0001$$

т.е. **МОЖЕМ** принять  $H_a$   
вероятность ошибки 1 типа (ошибочно  
принять  $H_a$  - найти то, чего нет)  $< 0,1\%$

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**5 ЭТАП:  
ИНТЕРПРЕТАЦИЯ  
РЕЗУЛЬТАТОВ**

**+ ОЦЕНИВАЕМ ПРАКТИЧЕСКУЮ ЗНАЧИМОСТЬ РЕЗУЛЬТАТОВ**

Средняя продолжительность жизни мужчин меньше, чем средняя продолжительность жизни женщин на 11,9 лет

$X = 73,3$   
 $SD = 15,4$   
 $N = 2021$

$X = 61,4$   
 $SD = 15,9$   
 $N = 2027$

## 2-Independent Samples test (*Mann-Whitney U test*) *тест Манна-Уитни*

### КАК ПРЕДСТАВИТЬ РЕЗУЛЬТАТ («АКАДЕМИЧЕСКАЯ ВЕРСИЯ»)

X (мужчины) = 61,4 лет (95% ДИ: 60,7 – 62,1)

X (женщины) = 73,3 лет (95% ДИ: 72,6 – 74,0)

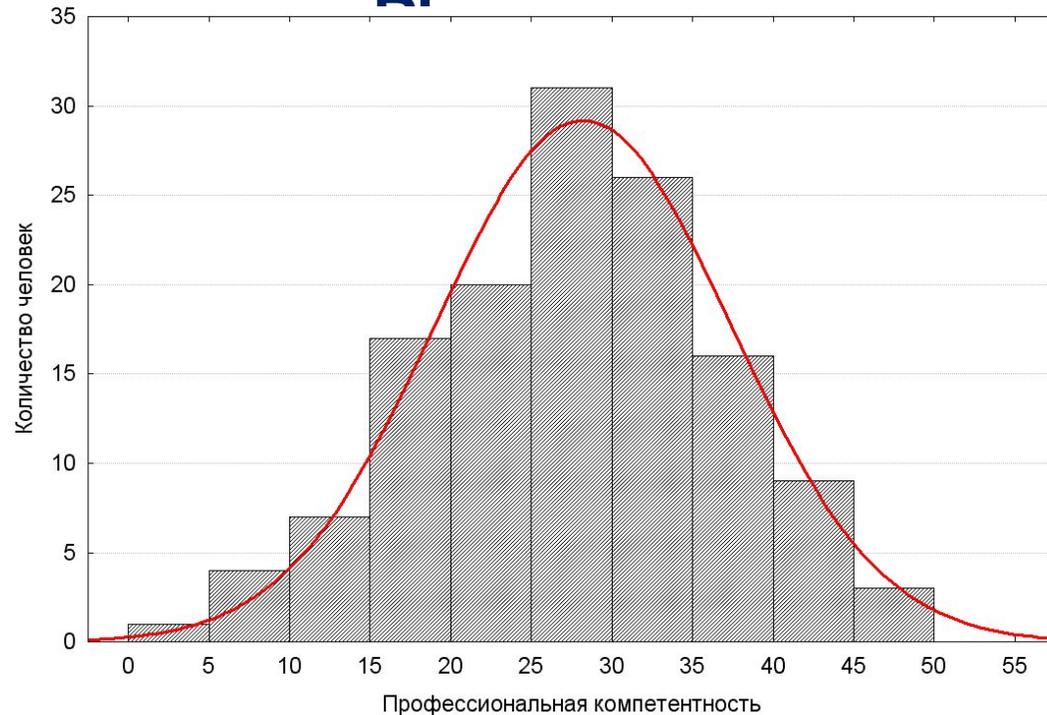
Различия являются статистически значимыми ( $p < 0,0001$ )

**РЕКОМЕНДУЕТСЯ УКАЗЫВАТЬ ТОЧНОЕ ЗНАЧЕНИЕ «p»**

(необходимо продемонстрировать вероятность ошибки)

# ПРИМЕР: СРЕДНЯЯ ПРОДОЛЖИТЕЛЬНОСТЬ ЖИЗНИ В ВЫБОРКЕ МУЖЧИН И ЖЕНЩИН, - ЖИТЕЛЕЙ АРХАНГЕЛЬСКОЙ ОБЛАСТИ, УМЕРШИХ В 2012 Г.

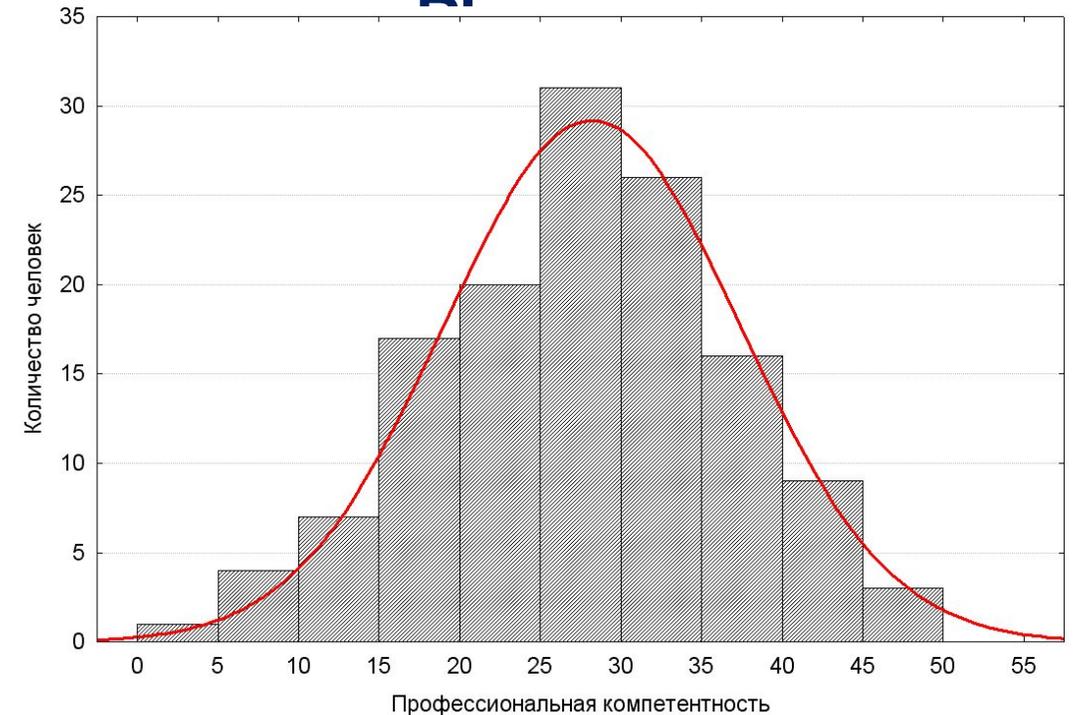
## ЖЕНЩИНЫ



**$X = 73,3$**   
 **$SD = 15,4$**   
 **$N = 2021$**



## МУЖЧИНЫ



**$X = 61,4$**   
 **$SD = 15,9$**   
 **$N = 2027$**

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**3 ЭТАП:  
ВЫБОР  
СТАТИСТИЧЕСКОГО  
КРИТЕРИЯ / МЕТОДА**

Tests of Normality

	Gender	Kolmogorov-Smirnov <sup>a</sup>			Shapiro-Wilk		
		Statistic	df	Sig.	Statistic	df	Sig.
Age	woman	,118	2021	<b>,298</b>	,923	2021	,000
	man	,066	2027	<b>,345</b>	,978	2027	,000

a. Lilliefors Significance Correction

**Н<sub>0</sub>: РАСПРЕДЕЛЕНИЕ В ВЫБОРКЕ НЕ ОТЛИЧАЕТСЯ ОТ НОРМАЛЬНОГО**  
**Н<sub>а</sub>: РАСПРЕДЕЛЕНИЕ В ВЫБОРКЕ ОТЛИЧАЕТСЯ ОТ НОРМАЛЬНОГО**

	ИСТИНА	
	<b>Н<sub>0</sub> - ВЕРНА</b>	<b>Н<sub>а</sub> - ВЕРНА</b>
<b>МЫ ПРИНИМАЕМ Н<sub>0</sub></b>	<b>ПРАВИЛЬНЫЙ РЕЗУЛЬТАТ</b> (= чувствительность теста)	<b>ОШИБКА 2 ТИПА</b> (вероятность НЕ найти то, чего ЕСТЬ)
<b>МЫ ОТВЕРГАЕМ Н<sub>0</sub></b> (ПРИНИМАЕМ Н <sub>а</sub> )	<b>ОШИБКА 1 ТИПА</b> (уровень значимости – significance (Sig.) “p” (вероятность найти то, чего НЕТ)	<b>ПРАВИЛЬНЫЙ РЕЗУЛЬТАТ</b> (= мощность теста)

**$p$  (женщины) = 0,298**  
 **$p$  (мужчины) = 0,345**

т.е. **НЕ МОЖЕМ** принять **Н<sub>а</sub>**  
 вероятность ошибки 1 типа (ошибочно принять Н<sub>а</sub> - найти то, чего нет) = 29,8% и 34,5%

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**3 ЭТАП:  
ВЫБОР  
СТАТИСТИЧЕСКОГО  
КРИТЕРИЯ / МЕТОДА**

Tests of Normality

	Gender	Kolmogorov-Smirnov <sup>a</sup>			Shapiro-Wilk		
		Statistic	df	Sig.	Statistic	df	Sig.
Age	woman	,118	2021	<b>,298</b>	,923	2021	,000
	man	,066	2027	<b>,345</b>	,978	2027	,000

a. Lilliefors Significance Correction

ASSUMPTIONS / УСЛОВИЯ ПРИМЕНЕНИЯ	КАК ПРОВЕРИТЬ?
1. Сравниваем 2 выборки	см. характеристики собранных данных
2. Выборки д.б. независимыми	см. характеристики собранных данных
3. Количественный непрерывный тип данных в каждой из сравниваемых выборок	см. тип данных
4. Нормальное распределение изучаемого признака в каждой из выборок	Test Shapiro-Wilk / Kolmogorov-Smirnov
5. Равенство дисперсий	<p><b>Levene's test for Equality of Variances</b> (sig. (p) <math>\geq</math> 0,05)</p> <p>Если дисперсии не равны (p &lt; 0,05)</p> <p><b>= проблема БЕРЕНСА-ФИШЕРА</b></p>

**Independent Samples T-test (Student test)  
Т-тест Стьюдента**

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

## 4 ЭТАП: МАТЕМАТИЧЕСКИЕ РАСЧЕТЫ

формулируем  $H_0$  и  $H_a$  для **теста Стьюдента**

**$H_0$ :  $m_1 = m_2$**  (средняя продолжительность жизни мужчин не отличается от средней продолжительности жизни женщин)

**$H_a$ :  $m_1 \neq m_2$**  (средняя продолжительность жизни мужчин отличается от средней продолжительности жизни женщин)

ASSUMPTIONS / УСЛОВИЯ ПРИМЕНЕНИЯ	КАК ПРОВЕРИТЬ?
1. Сравниваем 2 выборки	см. характеристики собранных данных
2. Выборки д.б. независимыми	см. характеристики собранных данных
3. Количественный непрерывный тип данных в каждой из сравниваемых выборок	см. тип данных
4. Нормальное распределение изучаемого признака в каждой из выборок	Test Shapiro-Wilk / Kolmogorov-Smirnov
5. Равенство дисперсий	<b>Levene's test for Equality of Variances</b> (sig. (p) $\geq 0,05$ )  Если дисперсии не равны (p < 0,05)  <b>= проблема БЕРЕНСА-ФИШЕРА</b>

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

## 4 ЭТАП: МАТЕМАТИЧЕСКИЕ РАССЧЕТЫ

формулируем  $H_0$  и  $H_a$  для **теста ЛЕВЕНЕ**  
(тест равенства дисперсий)

$H_0: \sigma_1 = \sigma_2$  (дисперсия средней продолжительности жизни мужчин не отличается от дисперсии средней продолжительности жизни женщин)

$H_a: \sigma_1 \neq \sigma_2$  (дисперсия средней продолжительности жизни мужчин отличается от дисперсии средней продолжительности жизни женщин)

Independent Samples Test

		Levene's Test for Equality of Variances		t-test for Equality of Means						
		F	Sig.	t	df	Sig. (2-tailed)	Mean Difference	Std. Error Difference	95% Confidence Interval of the Difference	
									Lower	Upper
Age	Equal variances assumed	8,286	,004	24,164	4046	,000	11,88957	,49204	10,92491	12,85424
	Equal variances not assumed			24,165	4041,791	,000	11,88957	,49201	10,92496	12,85419

# Independent Samples T-test (Student test) тест Стьюдента

## КАК ПРЕДСТАВИТЬ РЕЗУЛЬТАТ

X (мужчины) = 61,4 лет (95% ДИ: 60,7 – 62,1)

X (женщины) = 73,3 лет (95% ДИ: 72,6 – 74,0)

Средняя продолжительность жизни мужчин на 11,9 лет меньше (95% ДИ: 11,9 – 12,9), чем женщин ( $p < 0,0001$ )

## РЕКОМЕНДУЕТСЯ УКАЗЫВАТЬ ТОЧНОЕ ЗНАЧЕНИЕ «p»

(необходимо продемонстрировать вероятность ошибки)

# СРАВНЕНИЕ 2-Х СРЕДНИХ ВЕЛИЧИН

	НОРМАЛЬНОЕ РАСПРЕДЕЛЕНИЕ В КАЖДОЙ ИЗ СРАВНИВАЕМЫХ ВЫБОРОК	СКОШЕННОЕ РАСПРЕДЕЛЕНИЕ В 1 ИЛИ 2 СРАВНИВАЕМЫХ ВЫБОРКАХ
НЕЗАВИСИМЫЕ ВЫБОРКИ	Independent Samples T-test (Student T-test) <b>тест Стьюдента</b>	2-Independent Samples test (Mann-Whitney U test) <b>тест Манна-Уитни</b>
ЗАВИСИМЫЕ ВЫБОРКИ (ПОВТОРНЫЕ ИЗМЕРЕНИЯ)	Dependent (Paired Samples) T-test <b>тест Стьюдента</b> <i><u>для парных выборок</u></i>	2-Related Samples test (Wilcoxon signed-rank test) <b>тест Вилкоксона</b> <i><u>для парных выборок</u></i>

# Paired Samples T-test

## тест Стьюдента для парных выборок

### ASSUMPTIONS / УСЛОВИЯ ПРИМЕНЕНИЯ

1. Сравниваем 2 выборки
2. Выборки д.б. **зависимыми**  
(одни и те же участники, но в разное время)
3. Количественный непрерывный тип данных в каждой из сравниваемых выборок
4. Нормальное распределение **разности** между значениями изучаемого признака в парах

### КАК ПРОВЕРИТЬ?

см. характеристики собранных данных

см. характеристики собранных данных

см. тип данных

Test Shapiro-Wilk / Kolmogorov-Smirnov

ДО	<i>(до-после)</i>	ПОСЛЕ	РАЗНОСТЬ
167		134	-33
156		160	4
177		129	-48
...		...	...

## 2-Related Samples test (Wilcoxon) тест Вилкоксона

### ASSUMPTIONS / УСЛОВИЯ ПРИМЕНЕНИЯ

1. Сравниваем 2 выборки
2. Выборки д.б. **зависимыми**  
(одни и те же участники в разное время)
3. Количественный непрерывный тип данных  
в каждой из сравниваемых выборок
4. Скошенное распределение **разности**  
между значениями изучаемого признака

### КАК ПРОВЕРИТЬ?

см. характеристики собранных данных

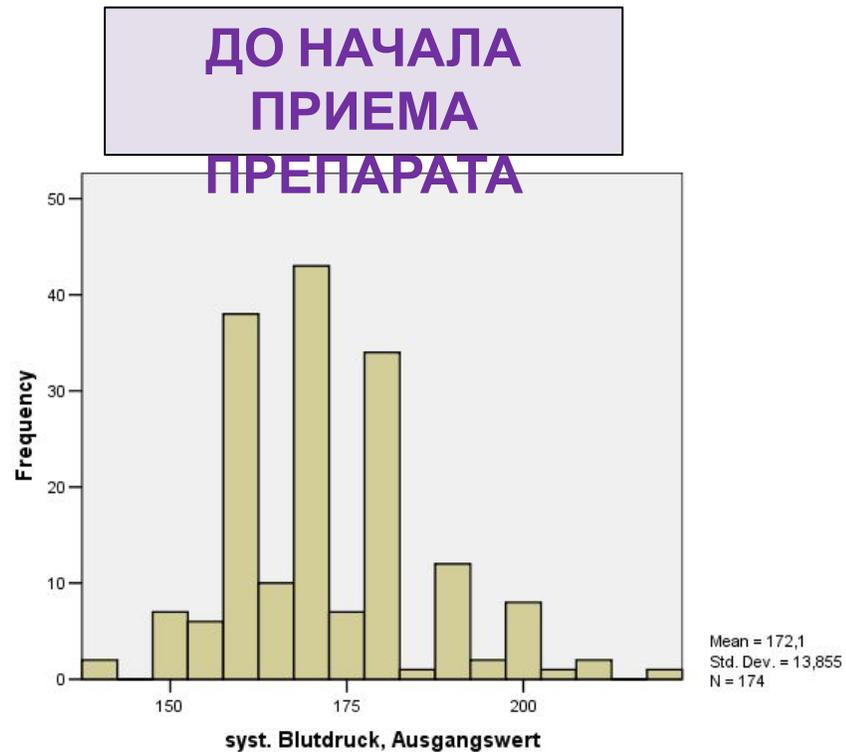
см. характеристики собранных данных

см. тип данных

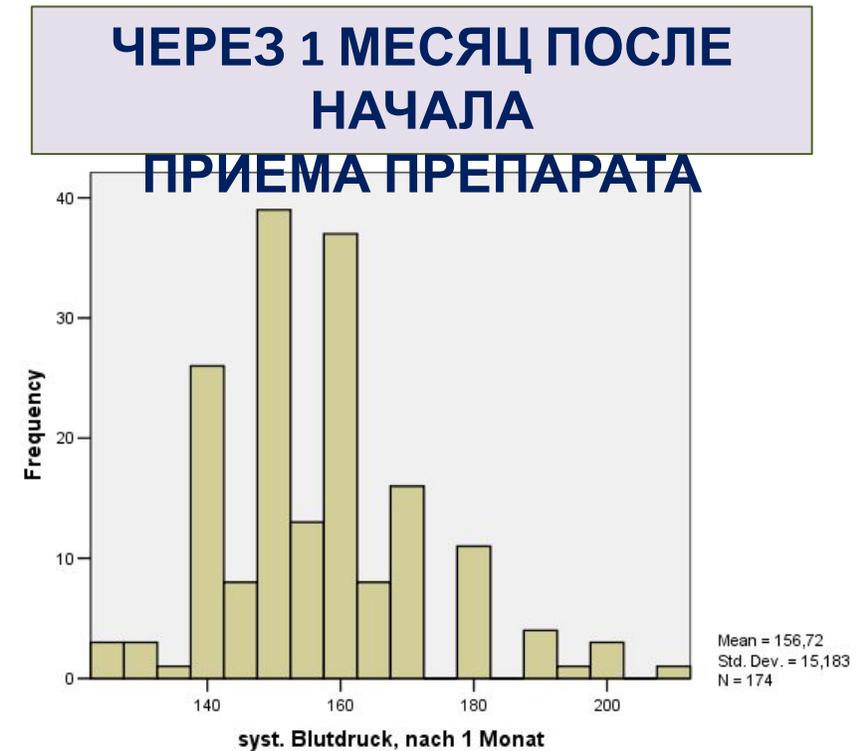
Test Shapiro-Wilk / Kolmogorov-Smirnov

ДО	ПОСЛЕ	РАЗНОСТЬ
167	134	-33
156	160	4
177	129	-48
...	...	...

# ПРИМЕР: УРОВЕНЬ АРТЕРИАЛЬНОГО ДАВЛЕНИЯ В ГРУППЕ ПАЦИЕНТОВ, ПРИНИМАЮЩИХ АНТИГИПЕРТЕНЗИВНЫЙ ПРЕПАРАТ



**$X = 172,1$**   
 **$SD = 13,9$**   
 **$N = 174$**



**$X = 156,7$**   
 **$SD = 15,2$**   
 **$N = 174$**

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**1 ЭТАП:**  
ФОРМУЛИРУЕМ  $H_0$  и  $H_a$

ГИПОТЕЗЫ	ФОРМУЛИРОВКА
$H_0$ (нулевая гипотеза)	$X(\text{ДО}) = X(\text{ПОСЛЕ})$ средний уровень артериального давления в группе пациентов до начала приема препарата $H_1$ отличается от среднего уровня артериального давления в группе пациентов после начала приема препарата
$H_a$ (альтернативная гипотеза)	$X(\text{ДО}) \neq X(\text{ПОСЛЕ})$ средний уровень артериального давления в группе пациентов до начала приема препарата <b>ОТЛИЧАЕТСЯ</b> от среднего уровня артериального давления в группе пациентов после начала приема препарата

**2 ЭТАП:**  
ОПРЕДЕЛЯЕМ УСЛОВИЯ, ПРИ КОТОРЫХ ПРИМЕМ  $H_a$  (И ОТВЕРГНЕМ  $H_0$ )

**БУДЕМ считать результаты теста «статистически значимыми» (т.е. примем  $H_a$ ) при вероятности ошибки 1 типа ( $\alpha$ -ошибки) менее 0.01 / 0.05 (1% / 5%)**

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**3 ЭТАП:  
ВЫБОР  
СТАТИСТИЧЕСКОГО  
КРИТЕРИЯ / МЕТОДА**

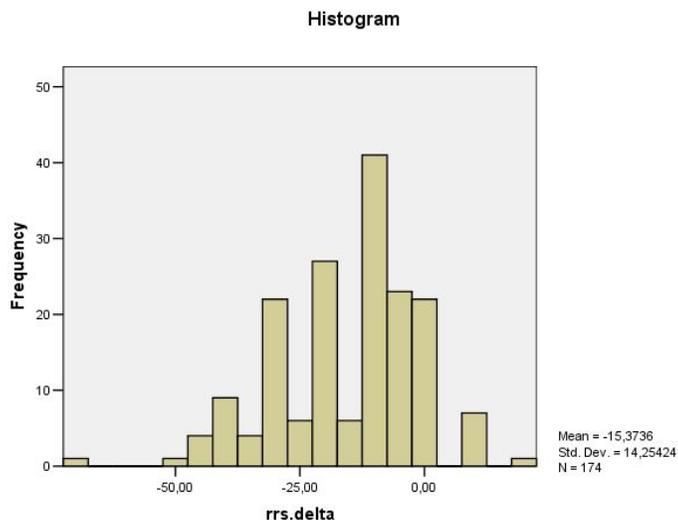
Tests of Normality

	Kolmogorov-Smirnov <sup>a</sup>			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
rrs.delta	,187	174	,000	,954	174	,000

a. Lilliefors Significance Correction

**Но:** РАСПРЕДЕЛЕНИЕ РАЗНИЦЫ СРЕДНИХ ВЕЛИЧИН (ДО-ПОСЛЕ) НЕ ОТЛИЧАЕТСЯ ОТ НОРМАЛЬНОГО

**На:** РАСПРЕДЕЛЕНИЕ РАЗНИЦЫ СРЕДНИХ ВЕЛИЧИН (ДО-ПОСЛЕ) В ВЫБОРКЕ ОТЛИЧАЕТСЯ ОТ НОРМАЛЬНОГО



$p$  (женщины) < 0,0001

$p$  (мужчины) < 0,0001

т.е. **МОЖЕМ** принять На  
вероятность ошибки 1 типа  
(ошибочно принять На - найти то,  
чего нет) < 0,1%

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**3 ЭТАП:  
ВЫБОР  
СТАТИСТИЧЕСКОГО  
КРИТЕРИЯ / МЕТОДА**

Tests of Normality

	Kolmogorov-Smirnov <sup>a</sup>			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
rrs_delta	,187	174	,000	,954	174	,000

a. Lilliefors Significance Correction

## ASSUMPTIONS / УСЛОВИЯ ПРИМЕНЕНИЯ

1. Сравниваем 2 выборки
2. Выборки д.б. **зависимыми**  
(одни и те же участники в разное время)
3. Количественный непрерывный тип данных в каждой из сравниваемых выборок
4. Скошенное распределение **разности** между значениями изучаемого признака

## КАК ПРОВЕРИТЬ?

1. см. характеристики собранных данных
2. см. характеристики собранных данных
3. см. тип данных

Test Shapiro-Wilk / Kolmogorov-Smirnov

**2-Related Samples test (Wilcoxon)  
тест Вилкоксона**

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

## 4 ЭТАП: МАТЕМАТИЧЕСКИЕ РАСЧЕТЫ

формулируем  $H_0$  и  $H_a$  для **теста Вилкоксона**

**$H_0$ :**  $m_1 = m_2$  (среднее АД до начала приема препарата не отличается от среднего АД через 1 месяц после начала приема препарата)

**$H_a$ :**  $m_1 \neq m_2$  (среднее АД до начала приема препарата отличается от среднего АД через 1 месяц после начала приема препарата)

Test Statistics<sup>b</sup>

	syst. Blutdruck, nach 1 Monat - syst. Blutdruck, Ausgangswert
Z	-9,970 <sup>a</sup>
Asymp. Sig. (2-tailed)	,000

a. Based on positive ranks.

b. Wilcoxon Signed Ranks Test

$$p < 0,0001$$

т.е. **МОЖЕМ** принять  $H_a$   
вероятность ошибки 1 типа (ошибочно  
принять  $H_a$  - найти то, чего нет) < 0,1%

## 2-Related Samples test (Wilcoxon) тест Вилкоксона

### **КАК ПРЕДСТАВИТЬ РЕЗУЛЬТАТ («АКАДЕМИЧЕСКАЯ ВЕРСИЯ»)**

М (до) = 172,1 мм рт.ст.

М (после) = 156,7 мм рт.ст.

Различия являются статистически значимыми  
( $p < 0,0001$ )

**РЕКОМЕНДУЕТСЯ УКАЗЫВАТЬ ТОЧНОЕ ЗНАЧЕНИЕ «p»**

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**3 ЭТАП:  
ВЫБОР  
СТАТИСТИЧЕСКОГО  
КРИТЕРИЯ / МЕТОДА**

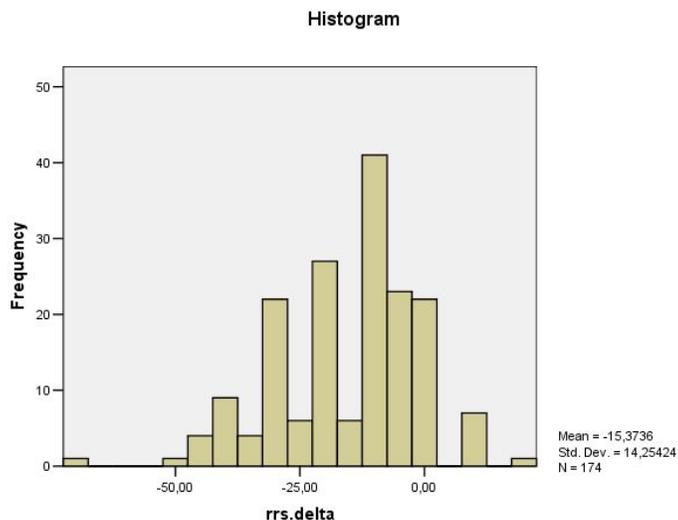
Tests of Normality

	Kolmogorov-Smirnov <sup>a</sup>			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
rrs.delta	,187	174	<b>,298</b>	,954	174	,000

a. Lilliefors Significance Correction

**Но:** РАСПРЕДЕЛЕНИЕ РАЗНИЦЫ СРЕДНИХ ВЕЛИЧИН (ДО-ПОСЛЕ) НЕ ОТЛИЧАЕТСЯ ОТ НОРМАЛЬНОГО

**На:** РАСПРЕДЕЛЕНИЕ РАЗНИЦЫ СРЕДНИХ ВЕЛИЧИН (ДО-ПОСЛЕ) В ВЫБОРКЕ ОТЛИЧАЕТСЯ ОТ НОРМАЛЬНОГО



**$p$  (мужчины) = 0,298**

т.е. **НЕ МОЖЕМ** принять **На**  
вероятность ошибки 1 типа (ошибочно  
принять **На** - найти то, чего нет)  $< 0,1\%$

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**3 ЭТАП:  
ВЫБОР  
СТАТИСТИЧЕСКОГО  
КРИТЕРИЯ / МЕТОДА**

Tests of Normality

	Kolmogorov-Smirnov <sup>a</sup>			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
rrs. delta	,187	174	,298	,954	174	,000

a. Lilliefors Significance Correction

## ASSUMPTIONS / УСЛОВИЯ ПРИМЕНЕНИЯ

## КАК ПРОВЕРИТЬ?

- |   |  |
|---|--|
| 1. Сравниваем 2 выборки   | см. характеристики собранных данных    |
| 2. Выборки д.б. <b>зависимыми</b><br>(одни и те же участники в разное время)                        | см. характеристики собранных данных    |
| 3. Количественный непрерывный тип данных в каждой из сравниваемых выборок                           | см. тип данных                         |
| 4. Нормальное распределение <b>разности</b> между значениями изучаемого признака в парах (до-после) | Test Shapiro-Wilk / Kolmogorov-Smirnov |



**Paired Samples T-test  
тест Стьюдента для парных  
выборок**

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

## 4 ЭТАП: МАТЕМАТИЧЕСКИЕ РАСЧЕТЫ

формулируем  $H_0$  и  $H_a$  для **парного теста Стьюдента**

**$H_0$ :**  $m_1 = m_2$  (среднее АД до начала приема препарата не отличается от среднего АД через 1 месяц после начала приема препарата)

**$H_a$ :**  $m_1 \neq m_2$  (среднее АД до начала приема препарата отличается от среднего АД через 1 месяц после начала приема препарата)

Paired Samples Test

		Paired Differences					t	df	Sig. (2-tailed)
		Mean	Std. Deviation	Std. Error Mean	95% Confidence Interval of the Difference				
					Lower	Upper			
Pair 1	syst. Blutdruck, Ausgangswert - syst. Blutdruck, nach 1 Monat	15,374	14,254	1,081	13,241	17,506	14,227	173	,000

# Paired Samples T-test

## тест Стьюдента для парных выборок

### **КАК ПРЕДСТАВИТЬ РЕЗУЛЬТАТ («АКАДЕМИЧЕСКАЯ ВЕРСИЯ»)**

М (до) = 172,1 мм рт.ст.

М (после) = 156,7 мм рт.ст.

Различия являются статистически значимыми  
( $p < 0,0001$ )

**РЕКОМЕНДУЕТСЯ УКАЗЫВАТЬ ТОЧНОЕ ЗНАЧЕНИЕ «p»**

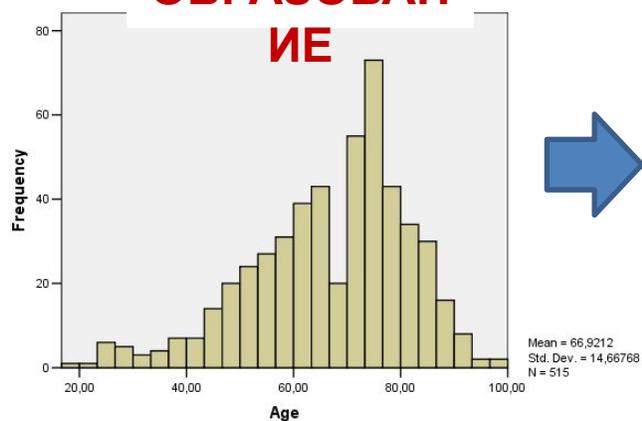
# СРАВНЕНИЕ 3-Х И БОЛЕЕ СРЕДНИХ ВЕЛИЧИН

# СРАВНЕНИЕ 2-х СРЕДНИХ ВЕЛИЧИН СРАВНЕНИЕ 3-х И БОЛЕЕ СРЕДНИХ ВЕЛИЧИН

**ИССЛЕДОВАТЕЛЬСКИЙ ВОПРОС:  
УРОВЕНЬ ОБРАЗОВАНИЯ ВЛИЯЕТ  
НА ПРОДОЛЖИТЕЛЬНОСТЬ ЖИЗНИ?**

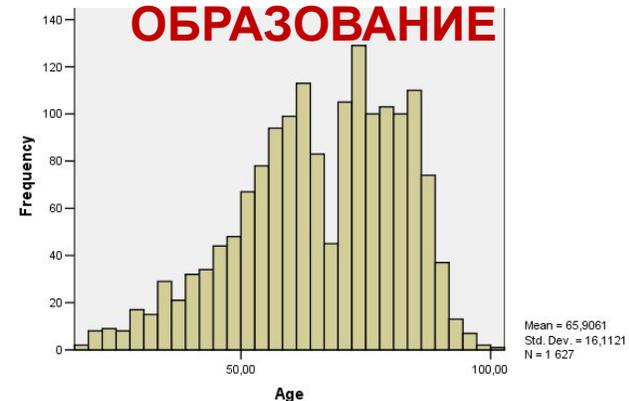


## ВЫСШЕЕ ОБРАЗОВАНИЕ



$X = 66,9$   
 $SD = 14,7$   
 $N = 515$

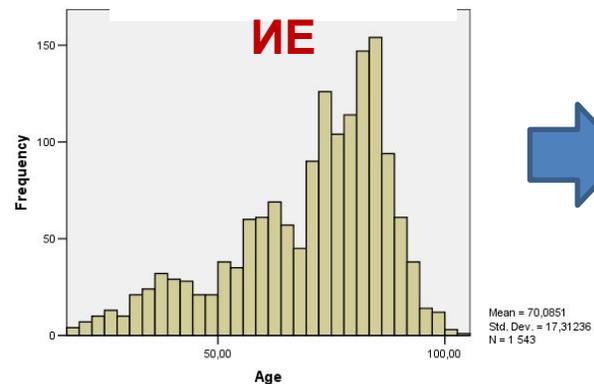
## СРЕДНЕЕ СПЕЦИАЛЬНОЕ ОБРАЗОВАНИЕ



$X = 65,9$   
 $SD = 16,1$   
 $N = 1627$



## СРЕДНЕЕ ОБРАЗОВАНИЕ



$X = 70,1$   
 $SD = 17,3$   
 $N = 1543$



# СРАВНЕНИЕ 2-х СРЕДНИХ ВЕЛИЧИН СРАВНЕНИЕ 3-х И БОЛЕЕ СРЕДНИХ ВЕЛИЧИН

**ИССЛЕДОВАТЕЛЬСКИЙ ВОПРОС:  
УРОВЕНЬ ОБРАЗОВАНИЯ ВЛИЯЕТ  
НА ПРОДОЛЖИТЕЛЬНОСТЬ ЖИЗНИ?**

**ВЫСШЕЕ  
ОБРАЗОВАНИЕ**

**ИЕ**

**X = 66,9**

**SD = 14,7**

**N = 515**



**СРЕДНЕЕ  
СПЕЦИАЛЬНОЕ  
ОБРАЗОВАНИЕ**

**X = 65,9**

**SD = 16,1**

**N = 1627**



**СРЕДНЕЕ  
ОБРАЗОВАНИЕ**

**ИЕ**

**X = 70,1**

**SD = 17,3**

**N = 1543**

**Почему нельзя сравнить группы попарно с помощью *t*-критерия  
Стьюдента?**

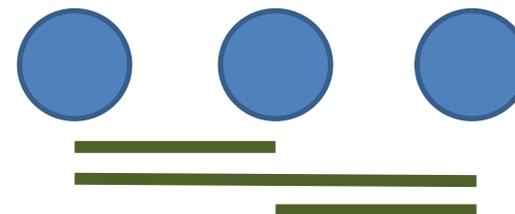


## **ЭФФЕКТ МНОЖЕСТВЕННЫХ СРАВНЕНИЙ**

При уровне значимости  $\alpha = 0,05$   
**вероятность ошибиться хотя бы в  
одном** из  $k$  сравнений

$$P_{\text{ошибки}} = 1 - (1 - 0,05)^k$$

$$P_{\text{ошибки}} = 1 - (1 - 0,05)^3 = 1 - (1 - 0,05)^3 = 14.3\%$$



**ВЫПОЛНЯЯ СЕРИЮ ПОПАРНЫХ  
СРАВНЕНИЙ, В КАЖДОМ СЛУЧАЕ  
МЫ УМЕНЬШАЕМ ОБЪЕМ**

# СРАВНЕНИЕ 3-Х И БОЛЕЕ СРЕДНИХ ВЕЛИЧИН

		НОРМАЛЬНОЕ РАСПРЕДЕЛЕНИЕ В КАЖДОЙ ИЗ СРАВНИВАЕМЫХ ВЫБОРОК	СКОШЕННОЕ РАСПРЕДЕЛЕНИЕ В 1 ИЛИ БОЛЕЕ СРАВНИВАЕМЫХ ВЫБОРОК
НЕЗАВИСИМЫЕ ВЫБОРКИ	РУС. ВЕРСИЯ	ONE-WAY ANOVA ДИСПЕРСИОННЫЙ АНАЛИЗ	K-Independent Samples test (Kruskall-Wallis H test) Тест Крускалла-Уоллиса
ЗАВИСИМЫЕ ВЫБОРКИ (ПОВТОРНЫЕ ИЗМЕРЕНИЯ)	РУС. ВЕРСИЯ	REPEATED MEASURES ANOVA (GLM-4) Дисперсионный анализ для повторных измерений	Friedman's test (Friedman's ANOVA) Дисперсионный анализ Фридмана

НО ! Считается, что нарушение нормальности распределения не оказывает существенного влияния на результаты)

НО ! Считается, что нарушение равенства дисперсии выборок оказывает значимое влияние в том случае, если сравниваемые выборки отличаются по численности)

# ONE-WAY ANOVA

## ДИСПЕРСИОННЫЙ АНАЛИЗ

### ASSUMPTIONS / УСЛОВИЯ ПРИМЕНЕНИЯ

1. Сравниваем 3 и более выборки
2. Выборки д.б. независимыми
3. **Количественный непрерывный тип данных** в каждой из сравниваемых выборок
4. **Нормальное распределение** изучаемого признака в сравниваемых группах
5. **Равенство дисперсий** изучаемого признака в сравниваемых группах (гомоскедастичность)

### КАК ПРОВЕРИТЬ?

см. характеристики собранных данных

см. характеристики собранных данных

см. тип данных

Test Shapiro-Wilk / Kolmogorov-Smirnov

Levene's test for Equality of Variances  
(Sig. (p)  $\geq 0,05$ )

Если дисперсии не равны (p < 0,05)

**= поправки Brown-Forsythe / Welch**

# K-Independent Samples test (Kruskall-Wallis H test)

## Тест Краскелла-Уоллиса

### ASSUMPTIONS / УСЛОВИЯ ПРИМЕНЕНИЯ

1. Сравниваем 3 и более выборок
2. Выборки д.б. независимыми
3. **Количественный непрерывный тип данных** в каждой из сравниваемых выборок
4. **Скошенное распределение** данных хотя бы в одной из сравниваемых выборок

### КАК ПРОВЕРИТЬ?

см. характеристики собранных данных

см. характеристики собранных данных

см. тип данных

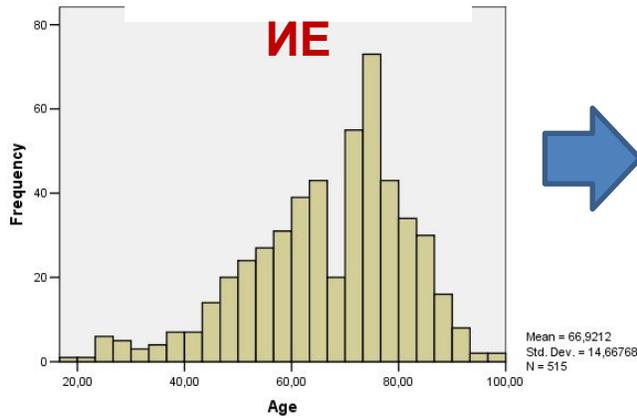
Test Shapiro-Wilk / Kolmogorov-Smirnov

**ДИСПЕРСИЯ НЕ  
ПРОВЕРЯЕТСЯ**

# ИССЛЕДОВАТЕЛЬСКИЙ ВОПРОС: УРОВЕНЬ ОБРАЗОВАНИЯ ВЛИЯЕТ НА ПРОДОЛЖИТЕЛЬНОСТЬ ЖИЗНИ?

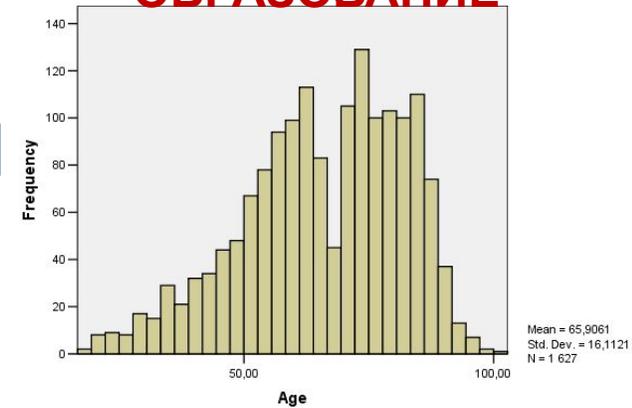


## ВЫСШЕЕ ОБРАЗОВАНИЕ



$X = 66,9$   
 $SD = 14,7$   
 $N = 515$

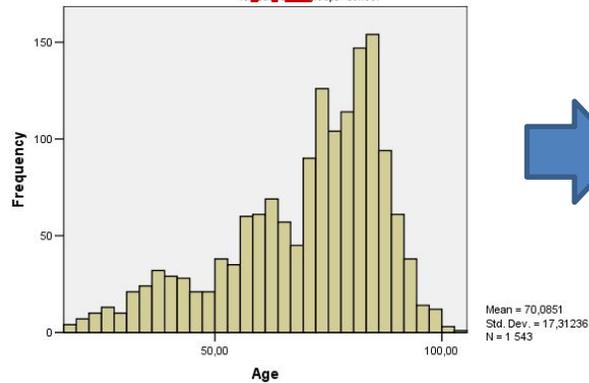
## СРЕДНЕЕ СПЕЦИАЛЬНОЕ ОБРАЗОВАНИЕ



$X = 65,9$   
 $SD = 16,1$   
 $N = 1627$



## СРЕДНЕЕ ОБРАЗОВАНИЕ



$X = 70,1$   
 $SD = 17,3$   
 $N = 1543$



# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**1 ЭТАП:**  
ФОРМУЛИРУЕМ  $H_0$  и  $H_a$

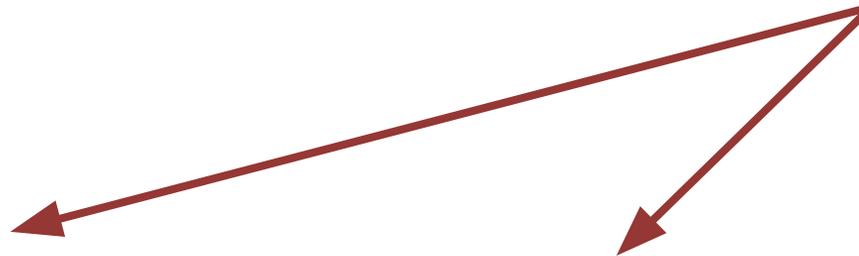
ГИПОТЕЗЫ	ФОРМУЛИРОВКА
<b><math>H_0</math></b> (нулевая гипотеза)	<b><math>X</math> (высшее) = <math>X</math> (ср.спец.) = <math>X</math> (среднее)</b> средняя продолжительность жизни не зависит от уровня образования
<b><math>H_a</math></b> (альтернативная гипотеза)	<b><math>X</math> (высшее) <math>\neq</math> <math>X</math> (ср.спец.)</b> <b><math>X</math> (высшее) <math>\neq</math> <math>X</math> (среднее)</b> <b><math>X</math> (ср.спец.) <math>\neq</math> <math>X</math> (среднее)</b>  мы отвергаем $H_0$ гипотезу если верна <b>хотя бы одна</b> из частных $H_a$

**2 ЭТАП:**  
ОПРЕДЕЛЯЕМ УСЛОВИЯ,  
ПРИ КОТОРЫХ ПРИМЕМ  
 **$H_a$**  (ОТВЕРГНЕМ  **$H_0$** )

**БУДЕМ** считать результаты теста  
«статистически значимыми» (т.е. примем  $H_a$ )  
при вероятности ошибки 1 типа ( $\alpha$ -ошибки)  
менее 0.05 (5%)

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**3 ЭТАП:  
ВЫБОР  
СТАТИСТИЧЕСКОГО  
КРИТЕРИЯ / МЕТОДА**



		НОРМАЛЬНОЕ РАСПРЕДЕЛЕНИЕ В КАЖДОЙ ИЗ СРАВНИВАЕМЫХ ВЫБОРОК	СКОШЕННОЕ РАСПРЕДЕЛЕНИЕ В 1 ИЛИ БОЛЕЕ СРАВНИВАЕМЫХ ВЫБОРОК
<b>НЕЗАВИСИМЫЕ ВЫБОРКИ</b>	РУС.ВЕРСИЯ	<b>ONE-WAY ANOVA ДИСПЕРСИОННЫЙ АНАЛИЗ</b>	<b>K-Independent Samples test (Kruskall-Wallis H test) Тест Краскелла-Уоллиса</b>
<b>ЗАВИСИМЫЕ ВЫБОРКИ (ПОВТОРНЫЕ ИЗМЕРЕНИЯ)</b>	РУС.ВЕРСИЯ	<b>REPEATED MEASURES ANOVA (GLM-4) Дисперсионный анализ для повторных измерений</b>	<b>Friedman's test (Friedman's ANOVA) тест Фридмана</b>

**НО !** Считается, что нарушение нормальности распределения не оказывает существенного влияния на результаты)

**НО !** Считается, что нарушение равенства дисперсии выборок оказывает значимое влияние в том случае, если сравниваемые выборки отличаются по численности)

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**3 ЭТАП:  
ВЫБОР  
СТАТИСТИЧЕСКОГО  
КРИТЕРИЯ / МЕТОДА**

Tests of Normality

Education groups	Kolmogorov-Smirnov <sup>a</sup>			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
Age university	,098	515	,000	,965	515	,000
college	,074	1627	,000	,972	1627	,000
school	,118	1543	,000	,931	1543	,000

a. Lilliefors Significance Correction

Но: РАСПРЕДЕЛЕНИЕ ПРИЗНАКА В ГРУППАХ **НЕ ОТЛИЧАЕТСЯ** ОТ НОРМАЛЬНОГО

На: РАСПРЕДЕЛЕНИЕ В ГРУППАХ **ОТЛИЧАЕТСЯ** ОТ НОРМАЛЬНОГО

K-Independent Samples test  
(Kruskal-Wallis H test)

**Тест Краскелла-Уоллиса**

$p$  (высшее) < 0,0001

$p$  (сред. спец.) < 0,0001

$p$  (среднее) < 0,0001

т.е. **МОЖЕМ** принять **На**  
вероятность ошибки 1 типа < 0,1%  
(ошибочно принять **На** - найти то, чего нет)

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

## 4 ЭТАП: МАТЕМАТИЧЕСКИЕ РАСЧЕТЫ

формулируем  $H_0$  и  $H_a$  для **теста Краскелла-Уоллиса**

$H_0: m_1 = m_2 = m_3$

$H_a: m_1 \neq m_2 / m_1 \neq m_3 / m_2 \neq m_3$

Test Statistics<sup>a,b</sup>

	Age
Chi-Square	79,561
df	2
Asymp. Sig.	,000

a. Kruskal Wallis Test

b. Grouping Variable: Education.groups

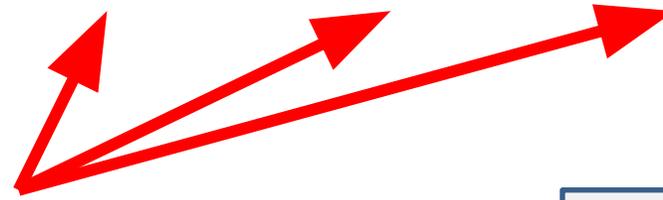
$p < 0,0001$

т.е. **МОЖЕМ** принять  $H_a$   
вероятность ошибки 1 типа (ошибочно  
принять  $H_a$  - найти то, чего нет)  $< 0,1\%$

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**4 ЭТАП:  
МАТЕМАТИЧЕСКИЕ  
РАСЧЕТЫ**

**$H_a: m_1 \neq m_2 / m_1 \neq m_3 / m_2 \neq m_3$**



**ДАЛЕЕ НЕОБХОДИМА СЕРИЯ  
ПРОЦЕДУР  
*ТЕСТА МАННА-УИТНИ***

**1**

**$H_0: m_1 = m_2$   
 $H_a: m_1 \neq m_2$**

**2**

**$H_0: m_1 = m_3$   
 $H_a: m_1 \neq m_3$**

**3**

**$H_0: m_2 = m_3$   
 $H_a: m_2 \neq m_3$**

# ПОПРАВКА БОНФЕРРОНИ: критический уровень "p" < 0.05/3 = < 0.017

Test Statistics<sup>a</sup>

	Age
Mann-Whitney U	407457,0
Wilcoxon W	1731835
Z	-,940
Asymp. Sig. (2-tailed)	,347

a. Grouping Variable: Education.groups



**H0: m1 = m2**

Средняя продолжительность жизни лиц с высшим образованием не отличается от средней продолжительности жизни лиц со средним специальным образованием

Test Statistics<sup>a</sup>

	Age
Mann-Whitney U	331066,0
Wilcoxon W	463936,0
Z	-5,674
Asymp. Sig. (2-tailed)	,000

a. Grouping Variable: Education.groups



**Ha: m1 ≠ m3**

Средняя продолжительность жизни лиц с высшим образованием отличается от средней продолжительности жизни лиц со средним образованием

Test Statistics<sup>a</sup>

	Age
Mann-Whitney U	1038009
Wilcoxon W	2362387
Z	-8,434
Asymp. Sig. (2-tailed)	,000

a. Grouping Variable: Education.groups



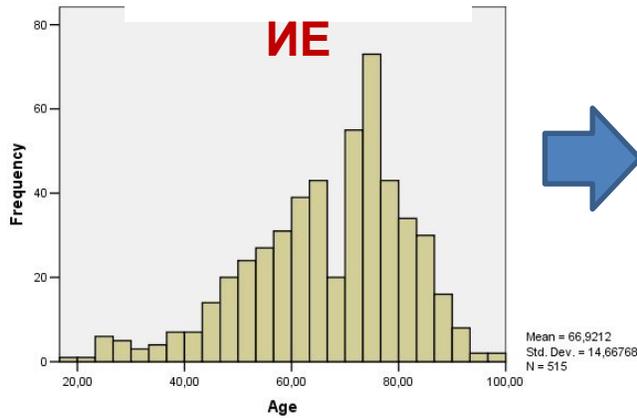
**Ha: m2 ≠ m3**

Средняя продолжительность жизни лиц со средним специальным образованием отличается от средней продолжительности жизни лиц со средним образованием

# ИССЛЕДОВАТЕЛЬСКИЙ ВОПРОС: УРОВЕНЬ ОБРАЗОВАНИЯ ВЛИЯЕТ НА ПРОДОЛЖИТЕЛЬНОСТЬ ЖИЗНИ?

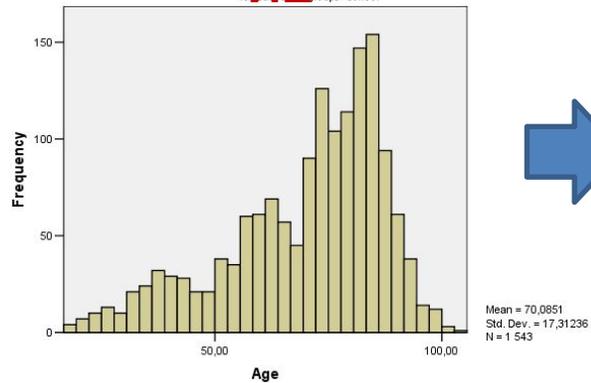


## ВЫСШЕЕ ОБРАЗОВАНИЕ



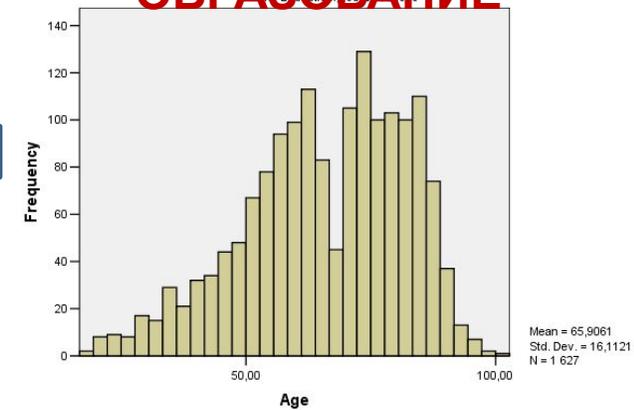
$X = 66,9$   
 $SD = 14,7$   
 $N = 515$

## СРЕДНЕЕ ОБРАЗОВАНИЕ



$X = 70,1$   
 $SD = 17,3$   
 $N = 1543$

## СРЕДНЕЕ СПЕЦИАЛЬНОЕ ОБРАЗОВАНИЕ



$X = 65,9$   
 $SD = 16,1$   
 $N = 1627$



# K-Independent Samples test (Kruskall-Wallis H test)

## Тест Краскелла-Уоллиса

### КАК ПРЕДСТАВИТЬ РЕЗУЛЬТАТ («АКАДЕМИЧЕСКАЯ ВЕРСИЯ»)

$m_1 = 66,9$  (95% ДИ: 65,7 – 68,2)

$m_2 = 65,9$  (95% ДИ: 65,1 – 66,7)

$m_3 = 70,1$  (95% ДИ: 69,2 – 70,9)

«...средняя продолжительность жизни зависит от уровня образования человека ( $H = 79,6$ ;  $p < 0,0001$ ). Продолжительность жизни лиц, имевших среднее образование, была статистически значимо выше, чем у лиц, имевших высшее и среднее специальное образование; средняя продолжительность жизни лиц, имевших высшее и среднее специальной

**ПОПРАВКА БОНФЕРРОНИ:**

**ОШИБКА 1 ТИПА:  $\alpha / n = 0.05/3 = 0,017$**

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**3 ЭТАП:  
ВЫБОР  
СТАТИСТИЧЕСКОГО  
КРИТЕРИЯ / МЕТОДА**

Tests of Normality

Education groups	Kolmogorov-Smirnov <sup>a</sup>			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
Age university	,098	515	<b>,298</b>	,965	515	,000
college	,074	1627	<b>,345</b>	,972	1627	,000
school	,118	1543	<b>,455</b>	,931	1543	,000

a. Lilliefors Significance Correction

Но: РАСПРЕДЕЛЕНИЕ ПРИЗНАКА В ГРУППАХ **НЕ ОТЛИЧАЕТСЯ** ОТ НОРМАЛЬНОГО

На: РАСПРЕДЕЛЕНИЕ ПРИЗНАКА В ГРУППАХ **ОТЛИЧАЕТСЯ** ОТ НОРМАЛЬНОГО

**ONE-WAY ANOVA  
ДИСПЕРСИОННЫЙ  
АНАЛИЗ**

$p$  (высшее) = 0,298  
 $p$  (сред. спец.) < 0,345  
 $p$  (среднее) < 0,455

т.е. **ОТКЛОНЯЕМ**  $H_0$   
вероятность ошибки 1 типа > 5%

# ONE-WAY ANOVA

## ДИСПЕРСИОННЫЙ АНАЛИЗ

- **ЦЕЛЬ:** с помощью **ДА** исследуют влияние *одной (одномерный анализ)* или нескольких (*многомерный анализ*) независимых переменных на *одну зависимую* переменную или на *несколько зависимых* переменных
- Независимые переменные **КАК ПРАВИЛО** принимают только дискретные значения (относятся к номинальной или порядковой шкале) - это **ФАКТОРНЫЙ АНАЛИЗ**
- Если независимые переменные принадлежат к интервальной шкале или к шкале отношений, то их называют **ковариациями** - это **КОВАРИАЦИОННЫЙ АНАЛИЗ**

# ДИСПЕРСИОННЫЙ АНАЛИЗ: ОСНОВНАЯ ИДЕЯ

- $SD = \sigma =$  СТАНДАРТНОЕ ОТКЛОНЕНИЕ
- $SD^2 =$  ДИСПЕРСИЯ (VARIANCE)

Оценка общей дисперсии по разбросу МЕЖДУ группами

средние в каждой группе

общее среднее

$$MS_B = s_{\bar{X}}^2 n = \frac{\sum (\bar{X}_j - \bar{X}_G)^2}{k-1} n$$

размер группы

число групп

$MS_B$  – mean square between groups  
оценка расстояния между средними в группах

ВЫСШЕЕ      СРЕД.СПЕЦ.      СРЕДНЕЕ

34	32	43
56	44	56
76	57	43
46	87	35
89	91	53
51	43	47
60	74	48
67	73	40
76	68	44
43	35	46
54	63	56
71	49	80
80	21	16
24	67	37
59	78	50

66,9

65,9

70,1

# ДИСПЕРСИОННЫЙ АНАЛИЗ: ОСНОВНАЯ ИДЕЯ

- $SD = \sigma =$  СТАНДАРТНОЕ ОТКЛОНЕНИЕ
- $SD^2 =$  ДИСПЕРСИЯ (VARIANCE)

Оценка общей дисперсии по разбросу ВНУТРИ групп

сумма квадратов  
стандартных отклонений  
внутри групп

$$MS_W = \frac{s_1^2 + s_2^2 + \dots + s_k^2}{k}$$

число групп

$$df_W = n_G - k$$

ВЫСШЕЕ	СРЕД.СПЕЦ.	СРЕДНЕЕ
34	32	43
56	44	56
76	57	43
46	87	35
89	91	53
51	43	47
60	74	48
67	73	40
76	68	44
43	35	46
54	63	56
71	49	80
80	21	16
24	67	37
59	78	50
...	...	...
<b>66,9</b>	<b>65,9</b>	<b>70,1</b>

**ВЫСШЕЕ      СРЕД.СПЕЦ.      СРЕДНЕЕ**

34	32	43
56	44	56
76	57	43
46	87	35
89	91	53
51	43	47
60	74	48
67	73	40
76	68	44
43	35	46
54	63	56
71	49	80
80	21	16
24	67	37
59	78	50
...	...	...
<b>66,9</b>	<b>65,9</b>	<b>70,1</b>

## ДИСПЕРСИОННЫЙ АНАЛИЗ: ОСНОВНАЯ ИДЕЯ

- $SD = \sigma =$  СТАНДАРТНОЕ ОТКЛОНЕНИЕ
- $SD^2 =$  ДИСПЕРСИЯ (VARIANCE)

**Расчет F-статистики  
ANOVA**

$F = \frac{\text{оценка дисперсии между группами}}{\text{оценка дисперсии внутри групп}}$

$$F = \frac{MS_B}{MS_W}$$

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**4 ЭТАП:  
МАТЕМАТИЧЕСКИЕ  
РАССЧЕТЫ**

формулируем  $H_0$  и  $H_a$  для **теста ЛЕВЕНЕ**  
(тест равенства дисперсий)

**$H_0$ :  $\sigma_1 = \sigma_2 = \sigma_3$**  (дисперсии средней продолжительности жизни в группах лиц в зависимости от уровня образования равны между собой)

**$H_a$ :  $\sigma_1 \neq \sigma_2 \neq \sigma_3$**  (дисперсии средней продолжительности жизни в группах лиц в зависимости от уровня образования **НЕ** равны между собой)

Test of Homogeneity of Variances

Age			
Levene Statistic	df1	df2	Sig.
9,963	2	3682	,000

***NB:***

**НЕОБХОДИМА  
ПОПРАВКА БРОУНА-  
ФОРСИТА / УЭЛЧА**

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**4 ЭТАП:  
МАТЕМАТИЧЕСКИЕ  
РАСЧЕТЫ**

формулируем  $H_0$  и  $H_a$  для **ANOVA**

**$H_0$ :  $m_1 = m_2 = m_3$**

**$H_a$ :  $m_1 \neq m_2 / m_1 \neq m_3 / m_2 \neq m_3$**

Robust Tests of Equality of Means

Age	Statistic <sup>a</sup>	df1	df2	Sig.
Welch	25,418	2	1508,262	,000
Brown-Forsythe	28,339	2	2562,151	,000

a. Asymptotically F distributed.

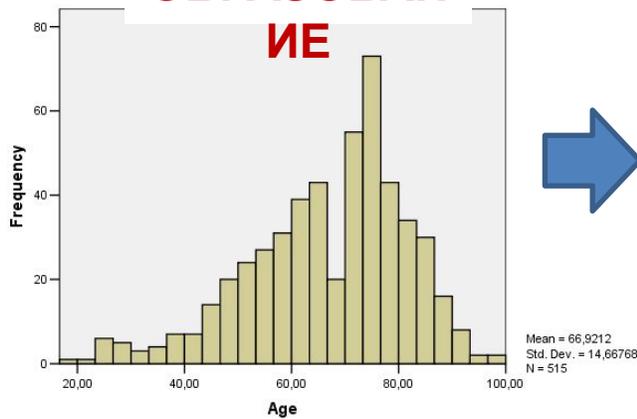
**$p < 0,0001$**

т.е. **МОЖЕМ** принять  **$H_a$**   
вероятность ошибки 1 типа  
(ошибочно принять  $H_a$  - найти  
то, чего нет)  $< 0,1\%$

# ИССЛЕДОВАТЕЛЬСКИЙ ВОПРОС: УРОВЕНЬ ОБРАЗОВАНИЯ ВЛИЯЕТ НА ПРОДОЛЖИТЕЛЬНОСТЬ ЖИЗНИ?



## ВЫСШЕЕ ОБРАЗОВАНИЕ

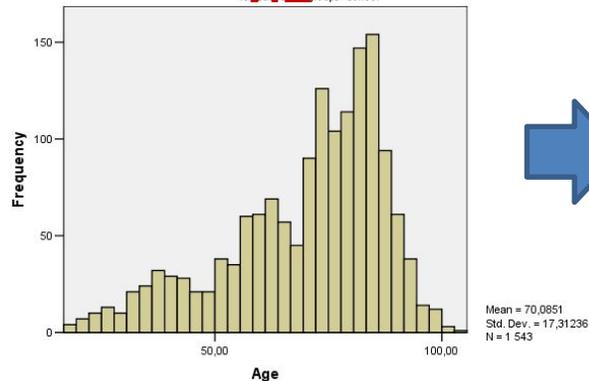


$X = 66,9$   
 $SD = 14,7$   
 $N = 515$

Но:  $m_1 = m_2 = m_3$

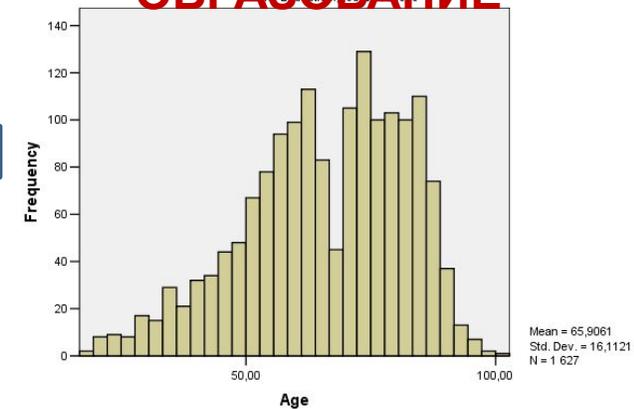
На:  $m_1 \neq m_2 / m_1 \neq m_3 / m_2 \neq m_3$

## СРЕДНЕЕ ОБРАЗОВАНИЕ



$X = 70,1$   
 $SD = 17,3$   
 $N = 1543$

## СРЕДНЕЕ СПЕЦИАЛЬНОЕ ОБРАЗОВАНИЕ



$X = 65,9$   
 $SD = 16,1$   
 $N = 1627$

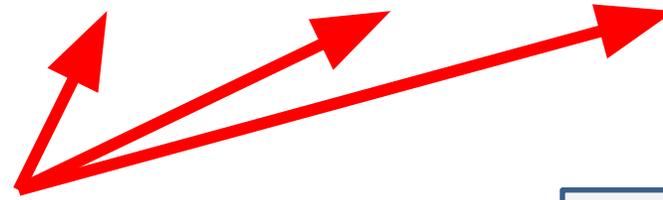


**В КАКОЙ ИМЕННО ПАРЕ СРЕДНЯЯ  
ПРОДОЛЖИТЕЛЬНОСТЬ ЖИЗНИ  
ОТЛИЧАЕТСЯ ???**

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**4 ЭТАП:  
МАТЕМАТИЧЕСКИЕ  
РАССЧЕТЫ**

**$H_a: m_1 \neq m_2 / m_1 \neq m_3 / m_2 \neq m_3$**



**1**

**$H_0: m_1 = m_2$   
 $H_a: m_1 \neq m_2$**

**2**

**$H_0: m_1 = m_3$   
 $H_a: m_1 \neq m_3$**

**3**

**$H_0: m_2 = m_3$   
 $H_a: m_2 \neq m_3$**

***ДАЛЕЕ НЕОБХОДИМА СЕРИЯ  
POST HOC тестов***

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**4 ЭТАП:  
МАТЕМАТИЧЕСКИЕ  
РАСЧЕТЫ**

**ДАЛЕЕ НЕОБХОДИМА СЕРИЯ  
*POST HOC* тестов**

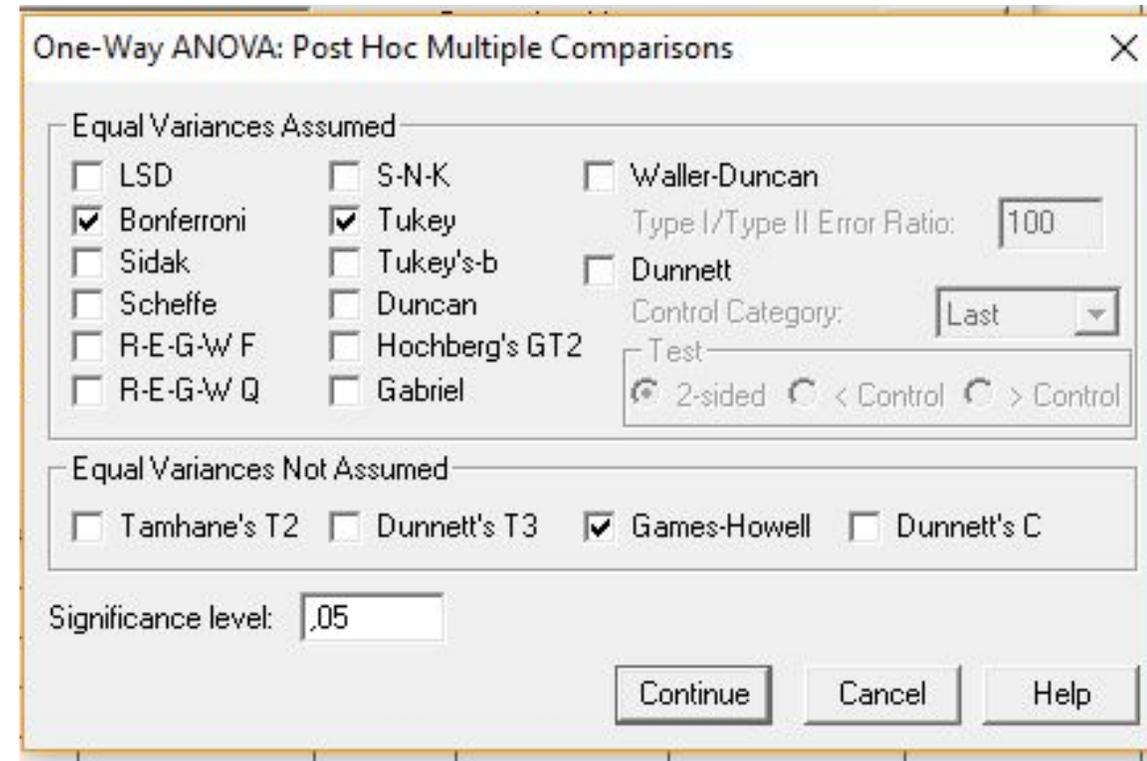
**УСЛОВИЕ О РАВЕНСТВЕ ДИСПЕРСИЙ  
СОБЛЮДЕНО**

**Bonferroni** – если число групп не более 5

**Tukey** – если число групп более 5

**УСЛОВИЕ О РАВЕНСТВЕ ДИСПЕРСИЙ НЕ  
СОБЛЮДЕНО**

**Games-Howell** – если группы равны, большие группы



# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

## ПРОБЛЕМА БОНФЕРРОНИ:

**НЕОБХОДИМО ВНЕСТИ ПОПРАВКУ НА КОЛИЧЕСТВО ГРУПП**

**ОШИБКА 1 ТИПА:  $\alpha / n = 0.05 / 3 = 0,017$**

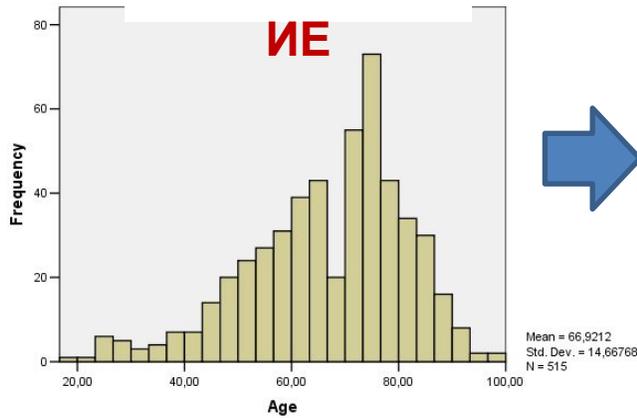
Dependent Variable				Difference		95% Confidence Interval		
Bonferroni	school	university		3,16388*	,83652	,000	1,2025	5,1252
		college		4,17901*	,58411	,000	2,8095	5,5485
	university	college		1,01513	,83109	,666	-,9754	3,0057
		school		-3,16388*	,83652	,000	-5,1674	-1,1604
	college	university		-1,01513	,83109	,666	-3,0057	,9754
		school		-4,17901*	,58411	,000	-5,5780	-2,7800
Games-Howell	school	university		3,16388*	,83652	,000	1,1604	5,1674
		college		4,17901*	,58411	,000	2,7800	5,5780
	university	college		1,01513	,75981	,376	-,7685	2,7987
		school		-3,16388*	,78230	,000	-5,0000	-1,3277
	college	university		-1,01513	,75981	,376	-2,7987	,7685
		school		-4,17901*	,59481	,000	-5,5757	-2,7845
school	university		3,16388*	,78230	,000	1,3277	5,0000	
	college		4,17901*	,59481	,000	2,7843	5,5737	

\*. The mean difference is significant at the .05 level.

# ИССЛЕДОВАТЕЛЬСКИЙ ВОПРОС: УРОВЕНЬ ОБРАЗОВАНИЯ ВЛИЯЕТ НА ПРОДОЛЖИТЕЛЬНОСТЬ ЖИЗНИ?

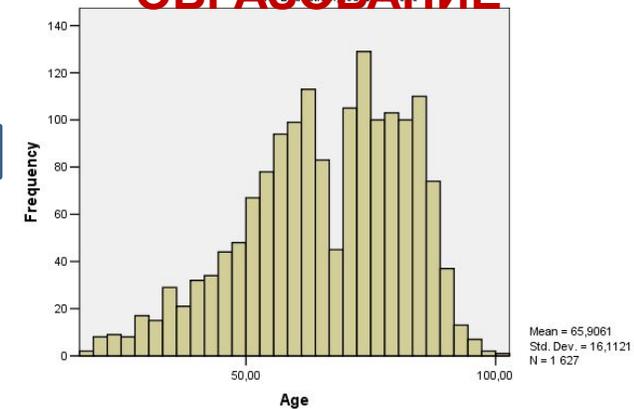


## ВЫСШЕЕ ОБРАЗОВАНИЕ



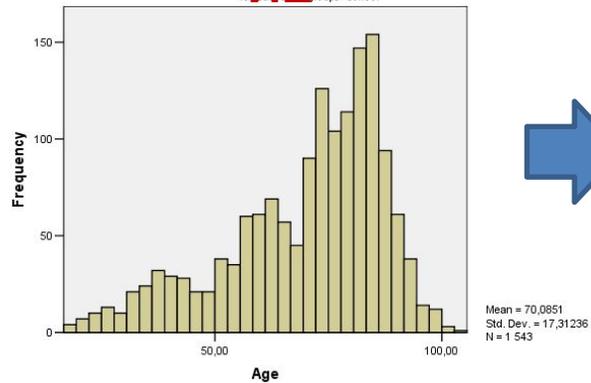
$X = 66,9$   
 $SD = 14,7$   
 $N = 515$

## СРЕДНЕЕ СПЕЦИАЛЬНОЕ ОБРАЗОВАНИЕ



$X = 65,9$   
 $SD = 16,1$   
 $N = 1627$

## СРЕДНЕЕ ОБРАЗОВАНИЕ



$X = 70,1$   
 $SD = 17,3$   
 $N = 1543$



# ONE-WAY ANOVA

## ДИСПЕРСИОННЫЙ АНАЛИЗ

### КАК ПРЕДСТАВИТЬ РЕЗУЛЬТАТ («АКАДЕМИЧЕСКАЯ ВЕРСИЯ»)

$m_1 = 66,9$  (95% ДИ: 65,7 – 68,2)

$m_2 = 65,9$  (95% ДИ: 65,1 – 66,7)

$m_3 = 70,1$  (95% ДИ: 69,2 – 70,9)

«...средняя продолжительность жизни зависит от уровня образования человека ( $F = 25,4$  (Welch);  $p < 0,0001$ ). Продолжительность жизни лиц, имевших среднее образование, была статистически значимо выше, чем у лиц, имевших высшее и среднее специальное образование»; средняя продолжительность жизни лиц, имевших высшее и среднее специальное образование, была равной

# ONE-WAY ANOVA

## ДИСПЕРСИОННЫЙ АНАЛИЗ

$$R^2 = \frac{SS_B}{SS_T}$$

**SS** - суммы квадратов отклонений (sum of squares):

**SS<sub>B</sub>** - средних в группах от общего среднего = **Effect**

**SS<sub>W</sub>** - измерений от средних в группах = **Error**

**«доля объяснённой  
вариабельности»**

$R^2 = 0.01$  – «незначительный» эффект

$R^2 = 0.06$  – «средний» эффект

$R^2 = 0.14$  – «значительный» эффект

**5 ЭТАП:  
ИНТЕРПРЕТАЦИЯ  
РЕЗУЛЬТАТОВ / оценка  
практической значимости**

# ONE-WAY ANOVA

## ДИСПЕРСИОННЫЙ АНАЛИЗ

Общая дисперсия по разбросу ВНУТРИ  
групп

$$f = \frac{s_{\bar{X}}}{\sqrt{MS_W}}$$

$$MS_W = \frac{s_1^2 + s_2^2 + \dots + s_k^2}{k}$$

**«практическая значимость»**

**результата:**

$f = 0,1$  – «незначительный» эффект

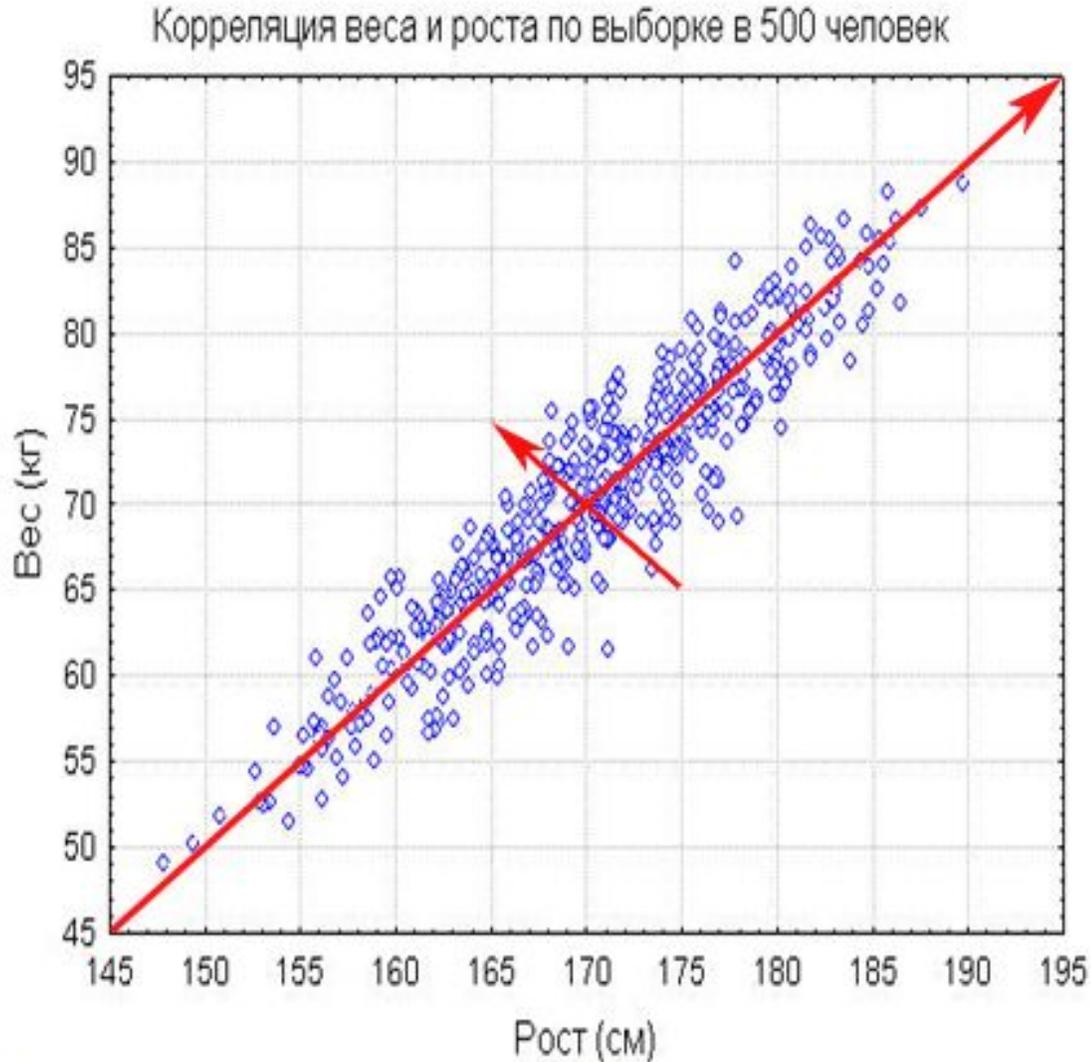
$f = 0.25$  – «средний» эффект

$f = 0.4$  – «значительный» эффект

**5 ЭТАП:**  
**ИНТЕРПРЕТАЦИЯ**  
**РЕЗУЛЬТАТОВ / оценка**  
**практической значимости**

# КОРРЕЛЯЦИОННЫЙ АНАЛИЗ

# Математическая зависимость величин



## НАПРАВЛЕНИЕ ЗАВИСИМОСТИ:

- Положительная

- Отрицательная

## СИЛА ЗАВИСИМОСТИ:

- Отсутствует

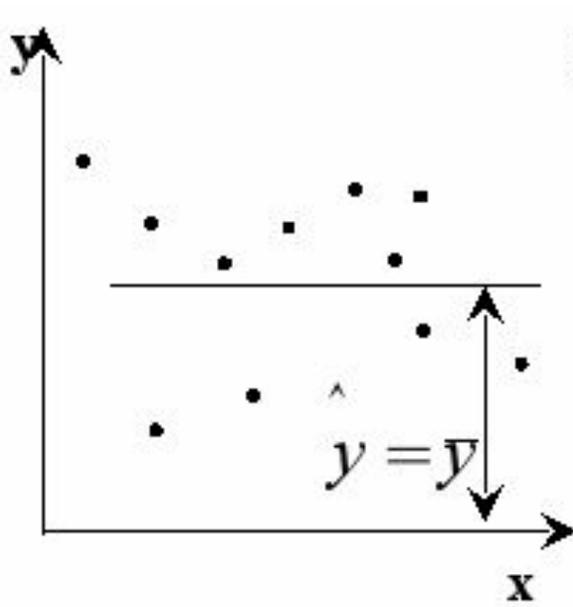
- Слабая

- Средняя

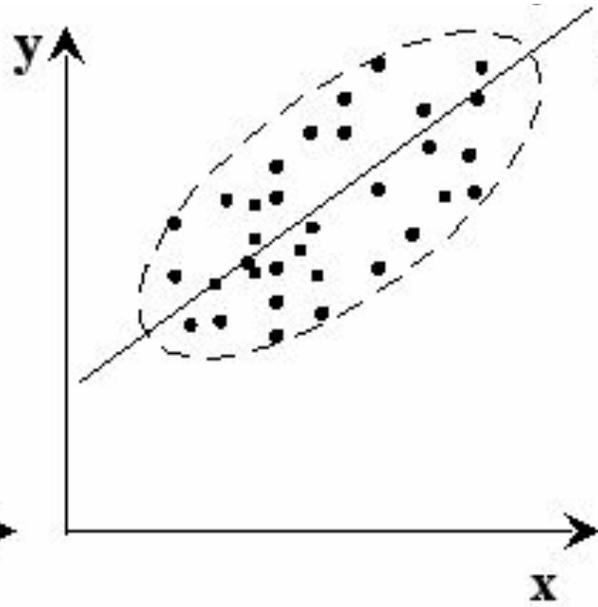
- Сильная

- Абсолютная

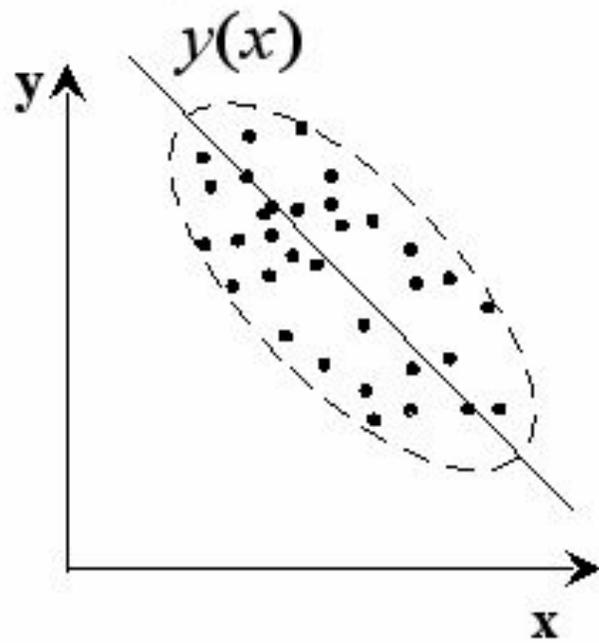
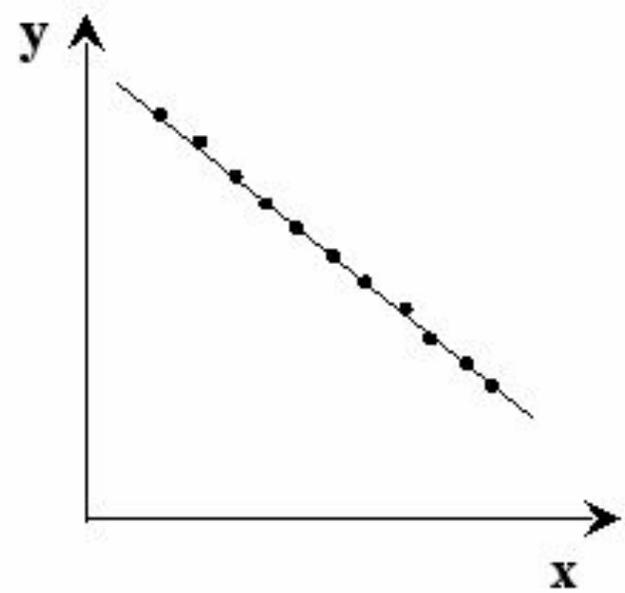
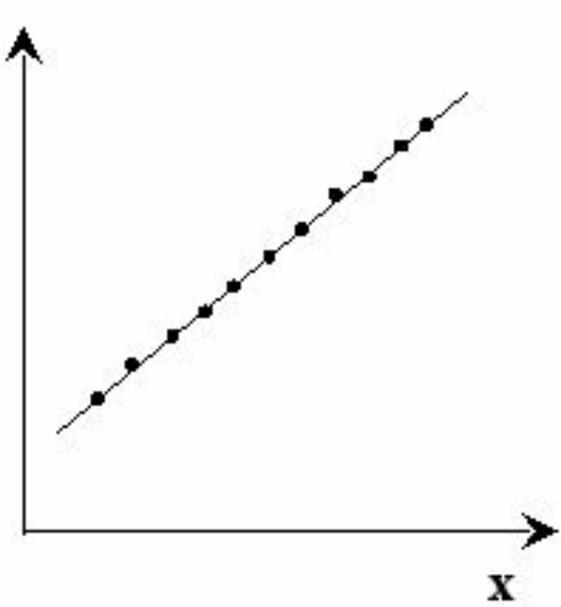
Наличие математической зависимости / корреляции **НЕ ОЗНАЧАЕТ** наличия **ПРИЧИННО-СЛЕДСТВЕННОЙ** взаимосвязи между переменными



$$\tau_{y,x} = -1$$



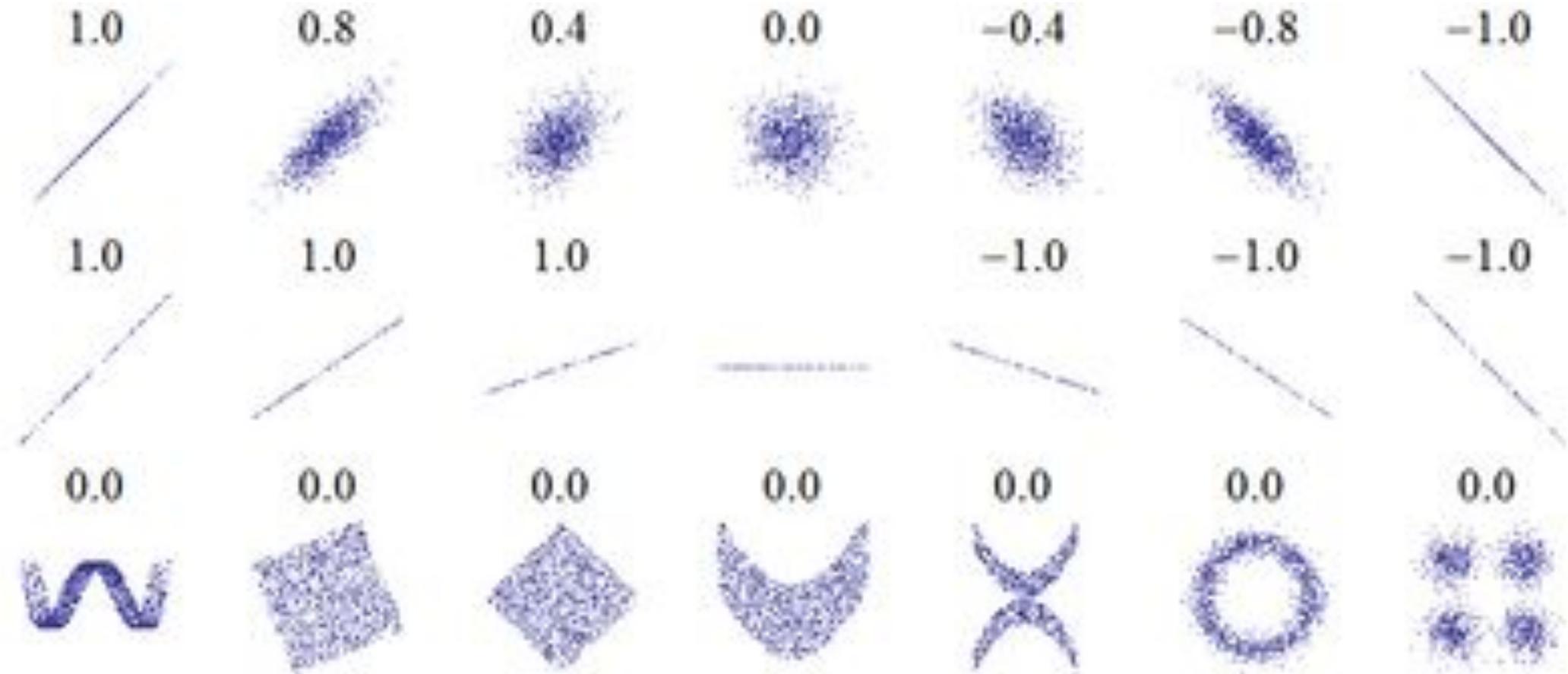
$$-1 < \tau_{y,x} < 0$$



**ЗАДАНИЕ:**

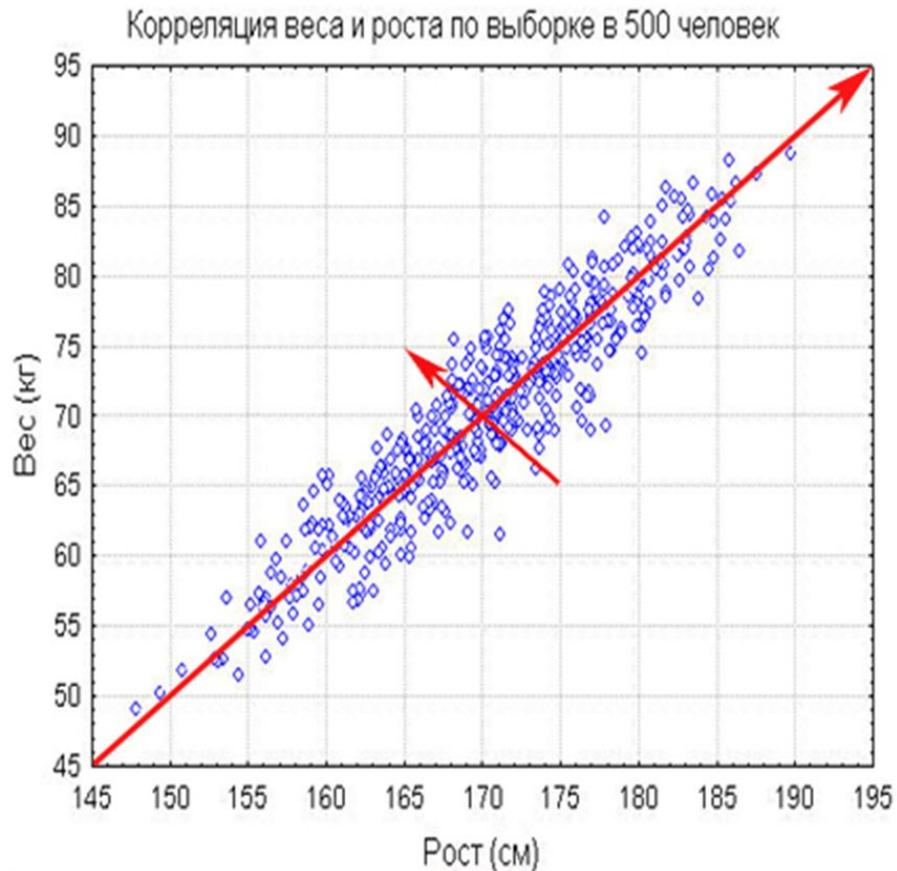
**ОПРЕДЕЛИТЬ  
НАПРАВЛЕНИЕ  
И СИЛУ  
ЗАВИСИМОСТИ  
ПЕРЕМЕННЫХ**

# МНОЖЕСТВО КОРРЕЛЯЦИОННЫХ ПОЛЕЙ



*Множество корреляционных полей.*  
<https://ru.wikipedia.org/wiki/Корреляция>

# Как можно количественно выразить математическую зависимость 2-х величин ?



**КОВАРИАЦИЯ**

**КОРРЕЛЯЦИЯ**

**КОВАРИАЦИЯ – это степень согласованности отклонений двух переменных**

$$\text{cov}(x,y) = \Sigma[(\mathbf{X} - \text{среднее}\mathbf{X})(\mathbf{Y} - \text{среднее}\mathbf{Y})]$$

Смысл: если 1 варианта отклоняется от средней, можно ожидать, что 2-я отклонится в ту же сторону

**КОРРЕЛЯЦИЯ – это ковариация стандартизованных переменных**

$$r = \text{cov}(x,y) / \text{SD}_{xy}$$

Смысл: отношение наблюдаемой ковариации двух стандартизованных переменных к максимально возможной ковариации

# КОРРЕЛЯЦИЯ

**КОРРЕЛЯЦИЯ** – это двумерное измерение силы и направления математической взаимосвязи между двумя переменными

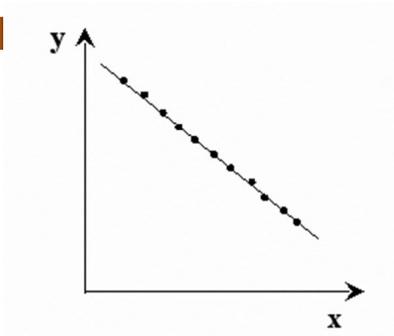
**-1**

**0**

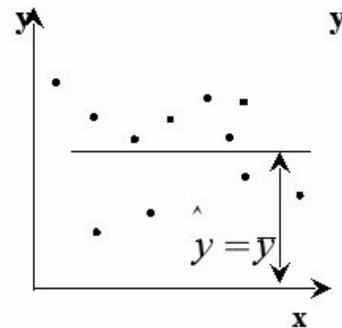
**+1**

абсолютная  
негативная  
линейная

СВЯЗЬ

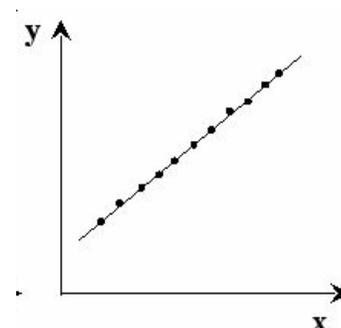


случайная  
связь



абсолютная  
положительная  
линейная

СВЯЗЬ



# КОЭФФИЦИЕНТЫ КОРРЕЛЯЦИИ

*Единственный  
параметрический критерий*

	Непрерывные	Порядковые	Дихотомические
Непрерывные	<p>Pearson's r</p> <p>Spearman's rho</p> <p>Kendall's tau</p>	<p>Spearman's rho</p> <p>Kendall's tau</p> <p>Polyserial correlation</p>	<p>Polyserial correlation</p> <p>Point-biserial correlation (истинная дихотомия)</p> <p>Biserial correlation (ложная дихотомия)</p>
Порядковые	<p>Spearman's rho</p> <p>Kendall's tau</p> <p>Polyserial correlation</p>	<p>Spearman's rho</p> <p>Kendall's tau</p> <p>Polychoric correlation</p>	<p>Rank biserial correlation</p>
Дихотомические	<p>Polyserial correlation</p> <p>Point-biserial correlation (истинная дихотомия)</p> <p>Biserial correlation (ложная дихотомия)</p>	<p>Rank biserial correlation</p>	<p>Polychoric correlation (tetrachoric correlation)</p> <p>phi</p>

# Пример расчета коэффициента корреляции Пирсона

<b>N</b>	<b>Содержание тестостерона в крови, нг/дл (X)</b>	<b>Процент мышечной массы, % (Y)</b>
<b>1.</b>	<b>951</b>	<b>83</b>
<b>2.</b>	<b>874</b>	<b>76</b>
<b>3.</b>	<b>957</b>	<b>84</b>
<b>4.</b>	<b>1084</b>	<b>89</b>
<b>5.</b>	<b>903</b>	<b>79</b>

## 1 ЭТАП. Расчет суммы значений переменных X и Y:

$$\Sigma(X) = 951 + 874 + 957 + 1084 + 903 = 4769$$

$$\Sigma(Y) = 83 + 76 + 84 + 89 + 79 = 441$$

# Пример расчета коэффициента корреляции Пирсона

<b>N</b>	<b>Содержание тестостерона в крови, нг/дл (X)</b>	<b>Процент мышечной массы, % (Y)</b>
<b>1.</b>	<b>951</b>	<b>83</b>
<b>2.</b>	<b>874</b>	<b>76</b>
<b>3.</b>	<b>957</b>	<b>84</b>
<b>4.</b>	<b>1084</b>	<b>89</b>
<b>5.</b>	<b>903</b>	<b>79</b>

## 2 ЭТАП. Расчет средних арифметических для X и Y:

$$M_x = \Sigma(X) / n = 4769 / 5 = 953.8$$

$$M_y = \Sigma(Y) / n = 441 / 5 = 82.2$$

# Пример расчета коэффициента корреляции Пирсона

<b>N</b>	<b>Содержание тестостерона в крови, нг/дл (X)</b>	<b>Процент мышечной массы, % (Y)</b>	<b>Отклонение содержания тестостерона от среднего значения (<math>d_x</math>)</b>	<b>Отклонение % мышечной массы от среднего значения (<math>d_y</math>)</b>
1.	951	83	-2.8	0.8
2.	874	76	-79.8	-6.2
3.	957	84	3.2	1.8
4.	1084	89	130.2	6.8
5.	903	79	-50.8	-3.2

**3 ЭТАП. Расчет для каждого значения сопоставляемых показателей величину отклонения от среднего арифметического**

$$dx = X - M_x$$

$$dy = Y - M_y$$

# Пример расчета коэффициента корреляции Пирсона

N	Содержание тестостерона в крови, нг/дл (X)	Процент мышечной массы, % (Y)	Отклонение содержания тестостерона от среднего значения ( $d_x$ )	Отклонение % мышечной массы от среднего значения ( $d_y$ )	$d_x^2$	$d_y^2$
1.	951	83	-2.8	0.8	7.84	0.64
2.	874	76	-79.8	-6.2	6368.04	38.44
3.	957	84	3.2	1.8	10.24	3.24
4.	1084	89	130.2	6.8	16952,04	46.24
5.	903	79	-50.8	-3.2	2580,64	10.24

**4 ЭТАП. Возвести в квадрат каждое значение отклонения  $d_x$  и  $d_y$**

## Пример расчета коэффициента корреляции Пирсона

N	Содержание тестостерона в крови, нг/дл (X)	Процент мышечной массы, % (Y)	Отклонение содержания тестостерона от среднего значения ( $d_x$ )	Отклонение % мышечной массы от среднего значения ( $d_y$ )	$d_x^2$	$d_y^2$	$d_x \times d_y$
1.	951	83	-2.8	0.8	7.84	0.64	-2.24
2.	874	76	-79.8	-6.2	6368.04	38.44	494.76
3.	957	84	3.2	1.8	10.24	3.24	5.76
4.	1084	89	130.2	6.8	16952,04	46.24	885.36
5.	903	79	-50.8	-3.2	2580,64	10.24	162.56

**5 ЭТАП. Расчет для каждой пары анализируемых значений произведение отклонений  $d_x \times d_y$ :**

# Пример расчета коэффициента корреляции Пирсона

**6 ЭТАП. Расчет значения суммы квадратов отклонений  $\Sigma$**

$$(d_x^2) \text{ и } \Sigma(d_y^2)$$

$$\Sigma(d_x^2) = 25918.8$$

$$\Sigma(d_y^2) = 98.8$$

**7 ЭТАП. Расчет значения суммы произведений отклонений  $\Sigma$**

$$(d_x \times d_y)$$

$$\Sigma(d_x \times d_y) = 1546.2$$

**8 ЭТАП. Расчет значения коэффициента корреляции Пирсона**

$$r_{xy} = \frac{\Sigma(d_x \times d_y)}{\sqrt{(\Sigma d_x^2 \times \Sigma d_y^2)}} = \frac{1546.2}{\sqrt{(25918.8 \times 98.8)}} = 0.966$$

# Пример расчета коэффициента корреляции Пирсона

9 ЭТАП. Оценка достоверности результата – расчет t-критерия

$$t_r = \frac{r_{xy} \sqrt{n - 2}}{\sqrt{1 - r_{xy}^2}} = \frac{0.97 \sqrt{5 - 2}}{\sqrt{1 - 0.97^2}} = 7.0$$

***Критическое значение*** t-критерия можно найти по специальной статистической таблице

# УСЛОВИЯ ПРИМЕНЕНИЯ КОРРЕЛЯЦИИ ПИРСОНА

ASSUMPTIONS / УСЛОВИЯ ПРИМЕНЕНИЯ	КАК ПРОВЕРИТЬ?
1. Сравниваем 2 выборки	см. характеристики собранных данных
2. Выборки д.б. <u>независимыми</u>	см. характеристики собранных данных
3. <u>Количественный непрерывный</u> тип данных в каждой из сравниваемых выборок	см. тип данных
4. <u>Нормальное распределение</u> изучаемого признака в каждой из выборок	Test Shapiro-Wilk / Kolmogorov-Smirnov
5. <u>Гомоскедастичность</u> - предполагается, что дисперсия ошибки остается той же самой в любой точке на протяжении всей линейной связи (иначе коэффициент корреляции будет завышаться или, наоборот, занижаться)	<b><u>обычно не проверяется</u></b>
6. <u>Линейная связь</u>	Graphs – Scatter/Dot (точечный график)
7. <u>Отсутствие «выбросов»</u>	

# КОЭФФИЦИЕНТ КОРРЕЛЯЦИИ ПИРСОНА

Значение коэффициента корреляции $r$	Интерпретация
$0 < r \leq 0,2$	Очень слабая корреляция
$0,2 < r \leq 0,5$	Слабая корреляция
$0,5 < r \leq 0,7$	Средняя корреляция
$0,7 < r \leq 0,9$	Сильная корреляция
$0,9 < r \leq 1$	Очень сильная корреляция

**Корреляция является симметричной, поэтому она не может говорить о направлении каузальной СВЯЗИ**

## Коэффициент детерминации $R^2$

$R^2$  - коэффициент детерминации - доля дисперсии переменной X, объясняемая вариабельностью переменной Y

$$r_{xy} = 0,5$$

$$R^2 = 0,25$$

Таким образом, вариабельность переменной X объясняет 25% вариабельности переменной Y

# УСЛОВИЯ ПРИМЕНЕНИЯ КОЭФФИЦИЕНТА РАНГОВОЙ КОРРЕЛЯЦИИ СПИРМЕНА, КОЭФФИЦИЕНТА КОРРЕЛЯЦИИ КЕНДАЛЛА ( $\tau$ )

## ASSUMPTIONS / УСЛОВИЯ ПРИМЕНЕНИЯ

1. Сравниваем 2 выборки
2. Выборки д.б. **независимыми**
3. **Количественный непрерывный / порядковый** тип данных в каждой из сравниваемых выборок
4. **Нормальное / скошенное** распределение изучаемого признака

## КАК ПРОВЕРИТЬ?

- см. характеристики собранных данных
- см. характеристики собранных данных
- см. тип данных

**можно не проверять**

# ОСНОВНОЙ НЕДОСТАТОК КОЭФФИЦИЕНТА КОРРЕЛЯЦИИ

Коэффициент корреляции демонстрирует

- А) направление взаимосвязи переменных
- Б) силу взаимосвязи переменных

**НО** коэффициент корреляции бесполезен, если мы хотим ПРЕДСКАЗАТЬ значение переменной X по значению переменной Y

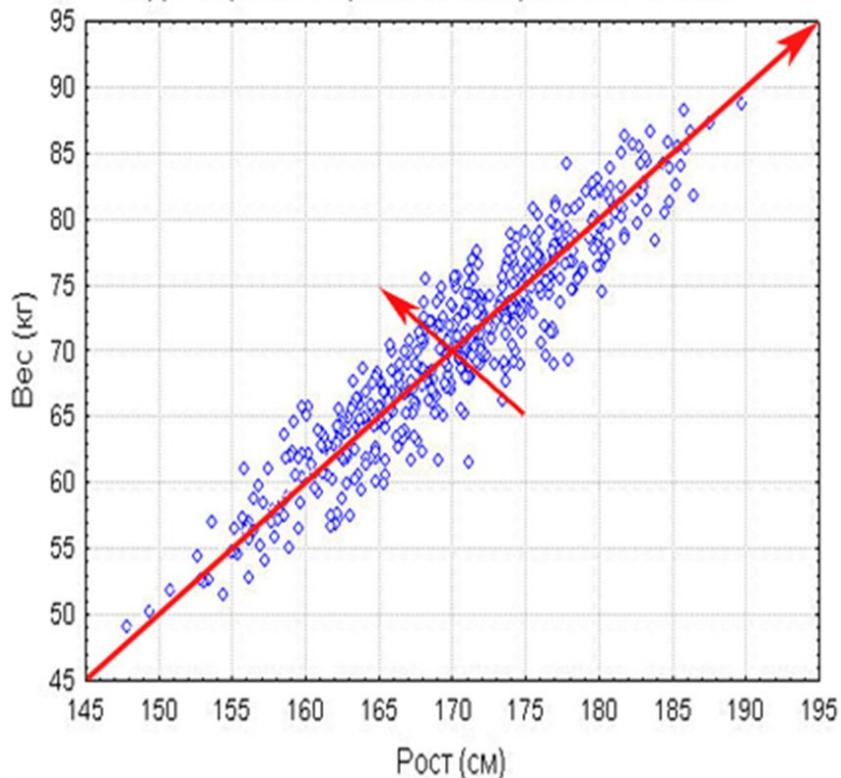


**РЕГРЕССИОННЫЙ  
АНАЛИЗ**

# ОСНОВЫ РЕГРЕССИОННОГО АНАЛИЗА

# КОРРЕЛЯЦИЯ vs. РЕГРЕССИЯ

Корреляция веса и роста по выборке в 500 человек



**МЕЖДУ ПЕРЕМЕННЫМИ ЕСТЬ ЗАВИСИМОСТЬ?**

**КОРРЕЛЯЦИОННЫЙ АНАЛИЗ** – демонстрирует лишь направление взаимосвязи переменных и силу взаимосвязи переменных

**ИССЛЕДОВАТЕЛЯ МОГУТ ДОПОЛНИТЕЛЬНО ИНТЕРЕСОВАТЬ ВОПРОСЫ:**

- 1) как сильно влияет на зависимую (1) переменную  
А) другая (1) независимая переменная?  
Б) одновременно 2 и > независимых переменных?
- 2) какие именно переменные влияют на зависимую переменную (отсеять из набора переменных «лишние»)?
- 3) какие именно переменные влияют одновременно на 2 и более зависимых переменных из набора?
- 4) можно ли по значениям одной (нескольких) переменных ПРЕДСКАЗАТЬ значение другой (других)

переменных

# РЕГРЕССИЯ: ОСНОВНАЯ ИДЕЯ

$$Y = f(X)$$



Зависимость между переменными может быть выражена **УРАВНЕНИЕМ**

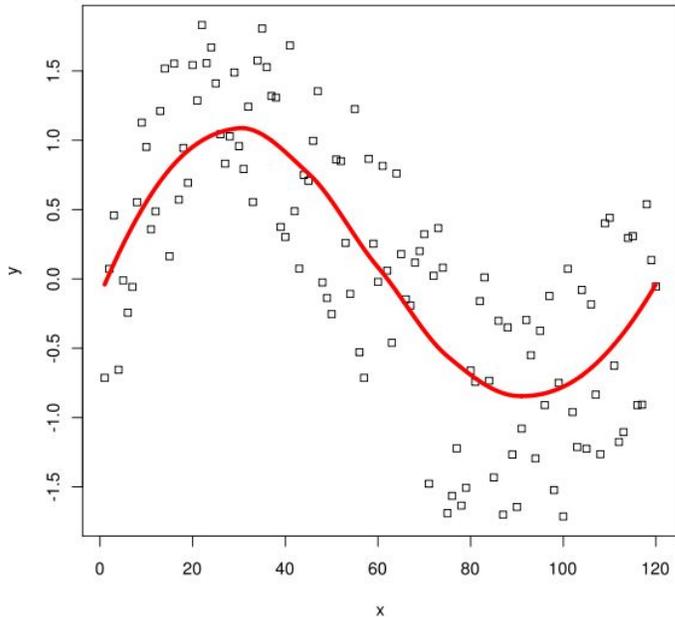


**ОСНОВНАЯ ИДЕЯ РЕГРЕССИОННОГО АНАЛИЗА:**

*математически рассчитать параметры  
УРАВНЕНИЯ РЕГРЕССИИ*

*(с какой силой / в каком направлении переменные влияют на зависимую переменную)*

# РЕГРЕССИЯ: ОСНОВНАЯ ПРОБЛЕМА

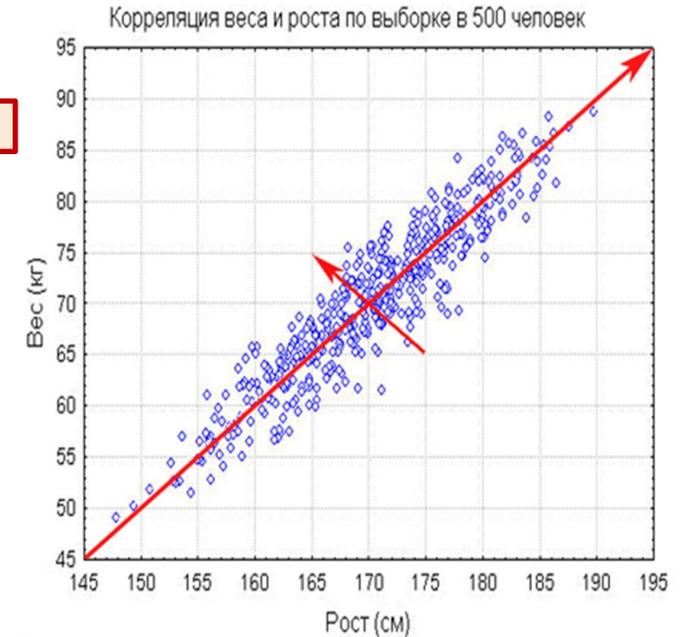
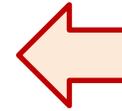


**нелинейная  
зависимость**



**ЛИНЕЙНЫЙ  
РЕГРЕССИОННЫЙ  
АНАЛИЗ**

**НЕЛИНЕЙНЫЙ  
РЕГРЕССИОННЫЙ  
АНАЛИЗ**

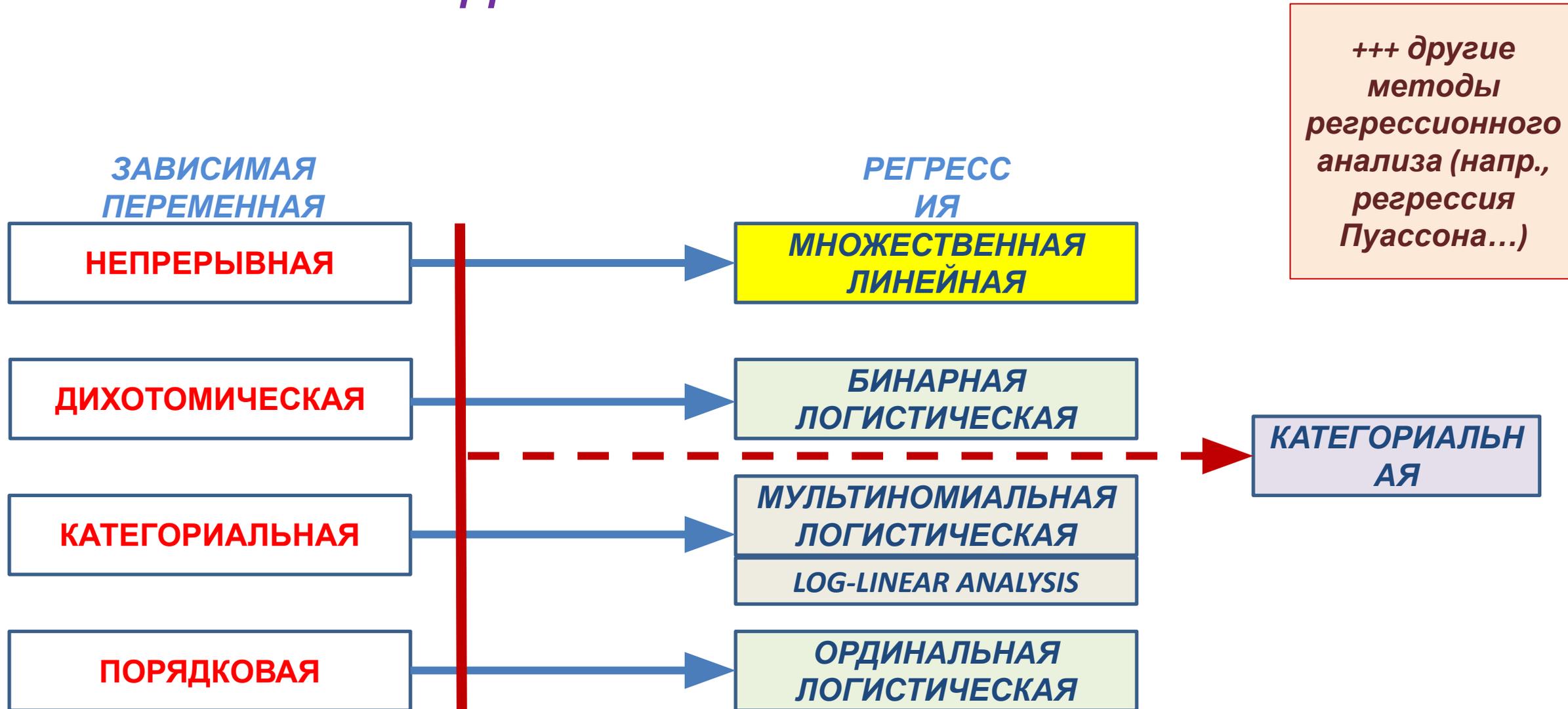


**линейная  
зависимость**

**КАКАЯ ФОРМА ЗАВИСИМОСТИ ОДНОЙ ПЕРЕМЕННОЙ ОТ ДРУГОЙ ПЕРЕМЕННОЙ?**

**КАКАЯ ФОРМА ЗАВИСИМОСТИ ОДНОЙ ПЕРЕМЕННЫХ ОТ НЕСКОЛЬКИХ ПЕРЕМЕННЫХ?**

# ВЫБОР МОДЕЛИ РЕГРЕССИОННОГО АНАЛИЗА



# ЛИНЕЙНЫЙ РЕГРЕССИОННЫЙ АНАЛИЗ

$$Y = b_0 + b_1 X_1 + b_2 X_2 + \dots + b_N X_N + \epsilon$$

ПРОСТАЯ ЛИНЕЙНАЯ  
РЕГРЕССИЯ

МНОЖЕСТВЕННАЯ ЛИНЕЙНАЯ  
РЕГРЕССИЯ

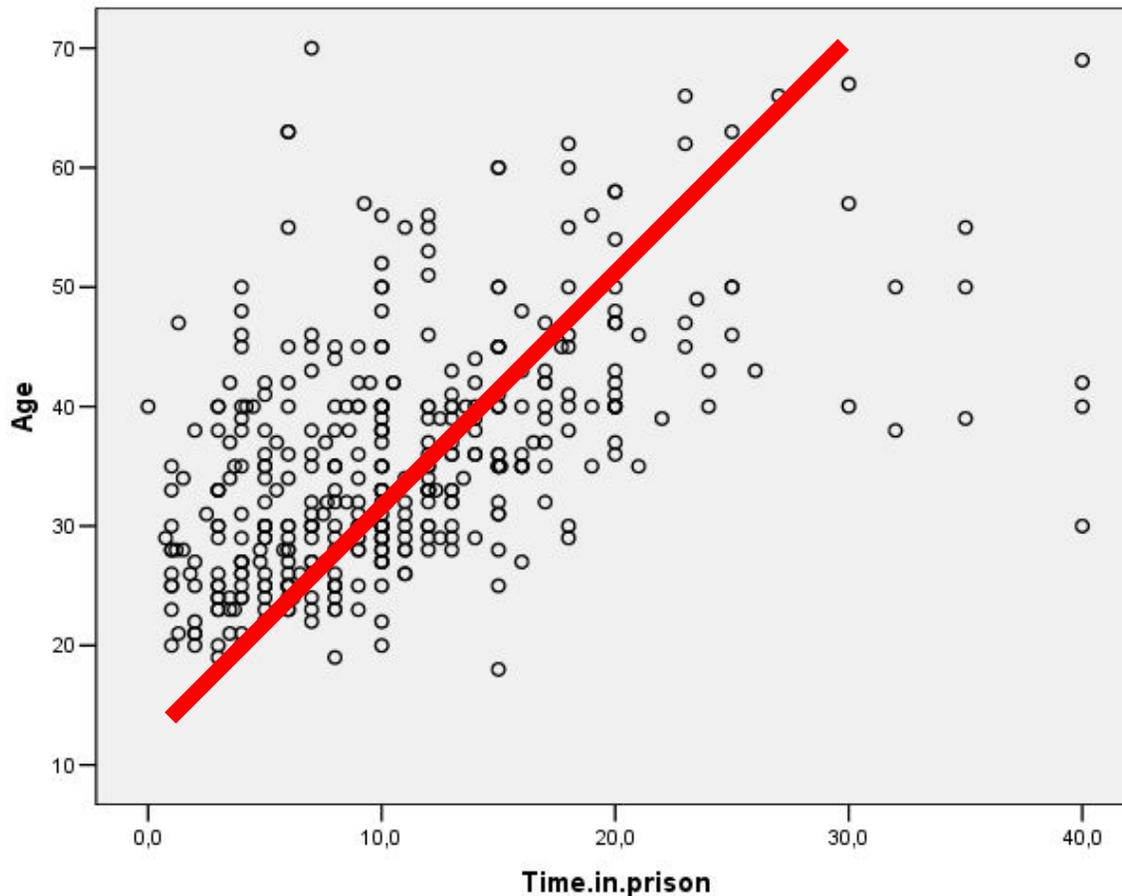
**Y** – зависимая переменная / переменная отклика

**b<sub>0</sub>** – константа

**b<sub>n</sub>** – коэффициент регрессии / градиент

**ε** - ошибка

# ЛИНЕЙНЫЙ РЕГРЕССИОННЫЙ АНАЛИЗ

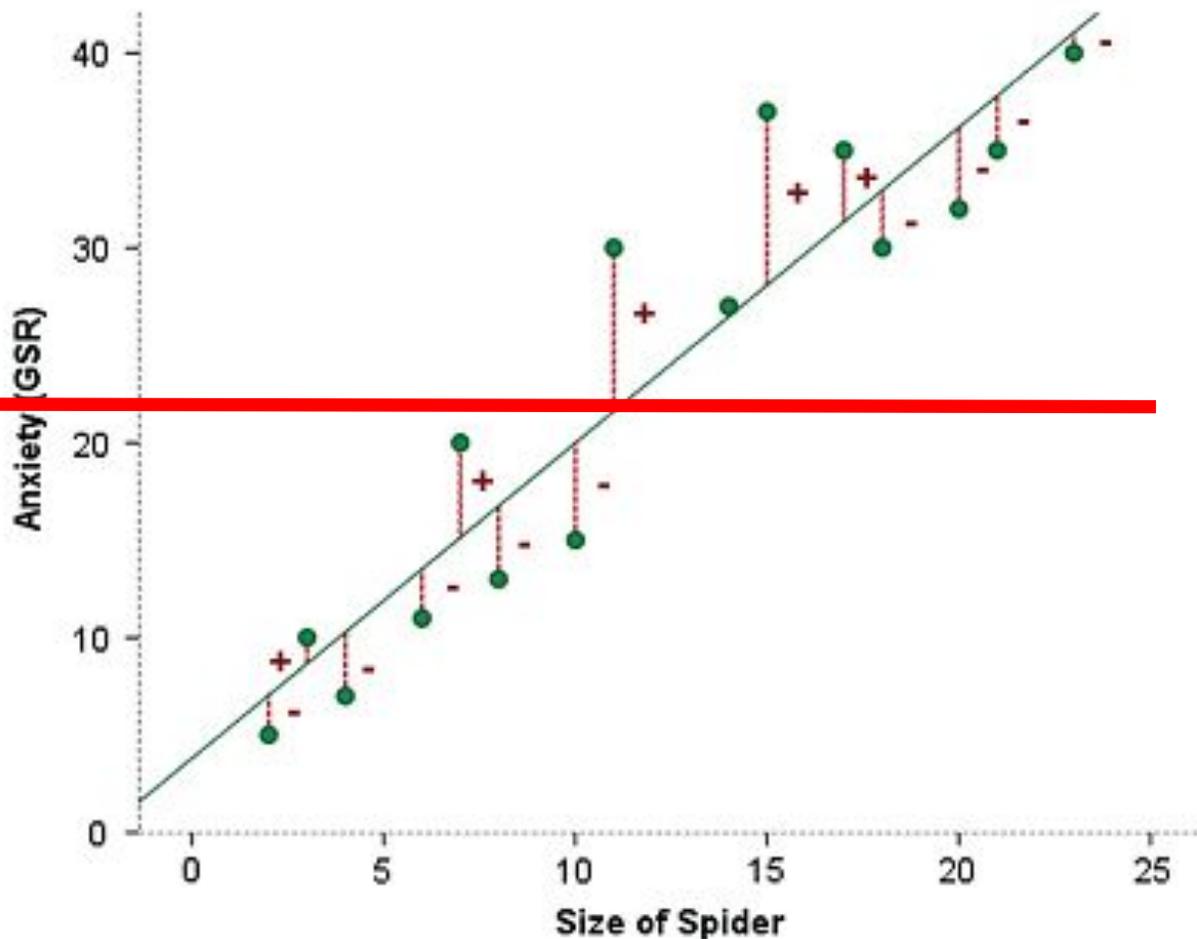


**Идея LRA:** построить прямую, наиболее точно предсказывающую значение зависимой переменной от предиктора (-ов) (и рассчитать ее параметры, т.е. ФОРМУЛУ)

– «линейный» анализ

**В ЭТОМ «МИНУС» ЛРА** – в природе нет линейной зависимости (тем более 1 зависимой переменной от нескольких)

# ЛИНЕЙНЫЙ РЕГРЕССИОННЫЙ АНАЛИЗ



## $H_0$ (LRA):

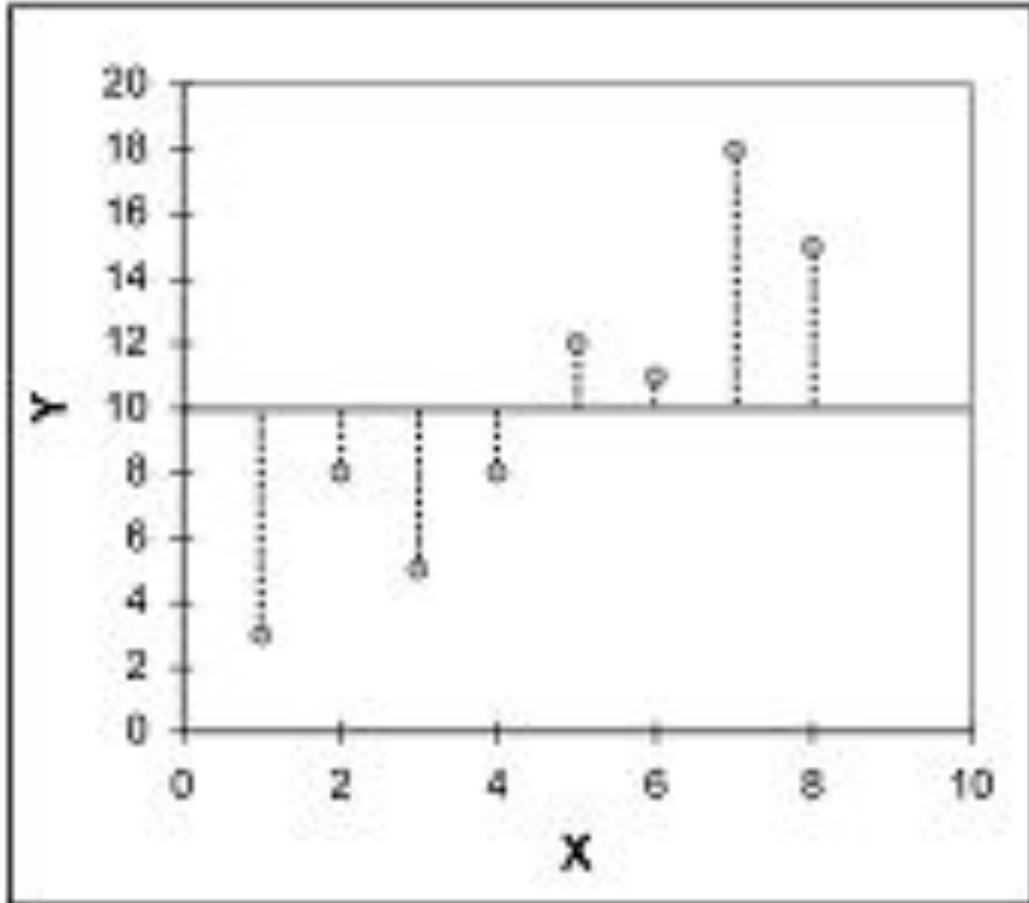
- Зависимая переменная лучше всего описывается средней арифметической

## $H_a$ (LRA):

- Зависимая переменная лучше всего описывается некоторой линейной моделью

Далее программа (по методу «наименьших квадратов») «подбирает» линию (модель), которая наилучшим образом «предсказывает» зависимую переменную по значению независимого предиктора

# ЛИНЕЙНЫЙ РЕГРЕССИОННЫЙ АНАЛИЗ



СНАЧАЛА ПРОГРАММА АНАЛИЗИРУЕТ,

НАСКОЛЬКО ХОРОШО СРЕДНЯЯ АРИФМЕТИЧЕСКАЯ ( $\bar{y}$ ) ПРЕДСКАЗЫВАЕТ ЗАВИСИМУЮ ПЕРЕМЕННУЮ:

**SST**

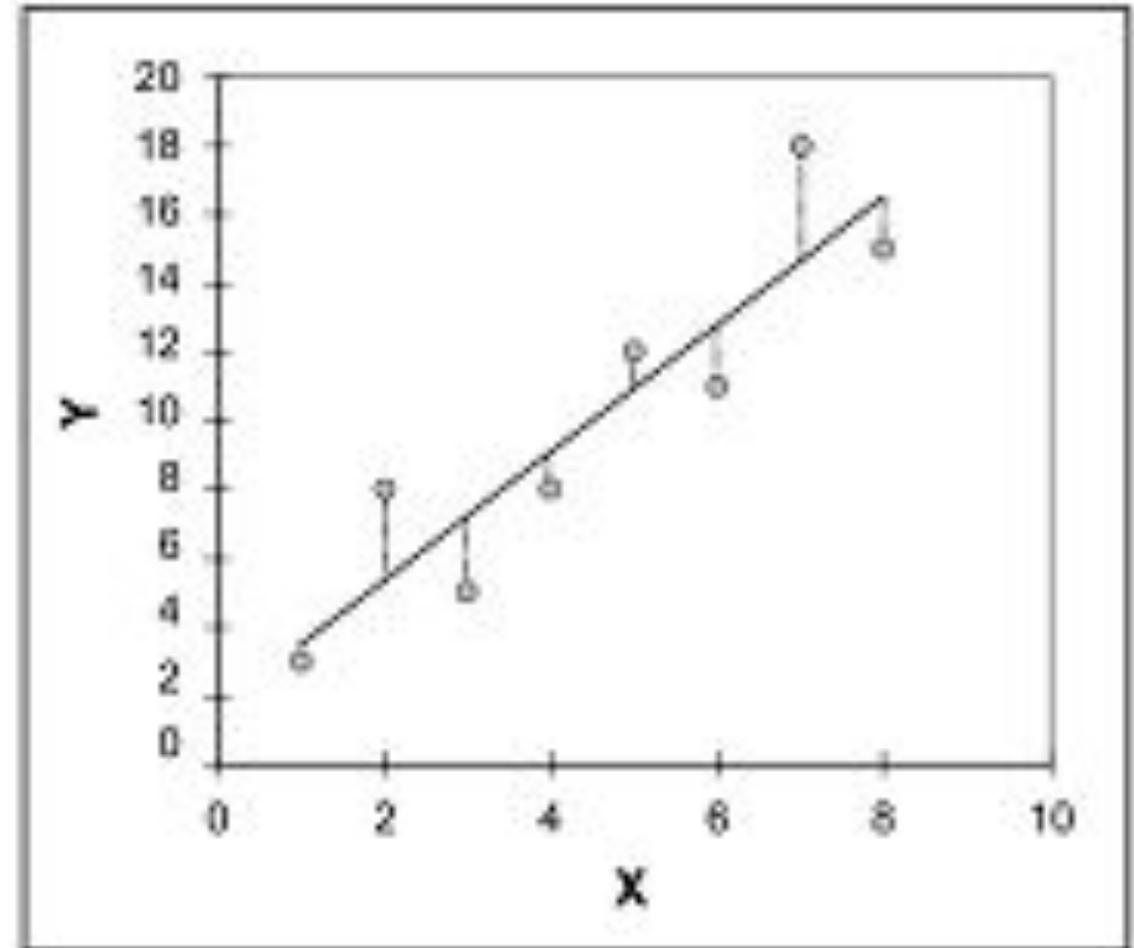
общая сумма различий между фактическими данными и средней арифметической

# ЛИНЕЙНЫЙ РЕГРЕССИОННЫЙ АНАЛИЗ

ДАЛЕЕ ПРОГРАММА АНАЛИЗИРУЕТ,  
НАСКОЛЬКО ХОРОШО МОДЕЛЬ (на)  
ПРЕДСКАЗЫВАЕТ ЗАВИСИМУЮ  
ПЕРЕМЕННУЮ

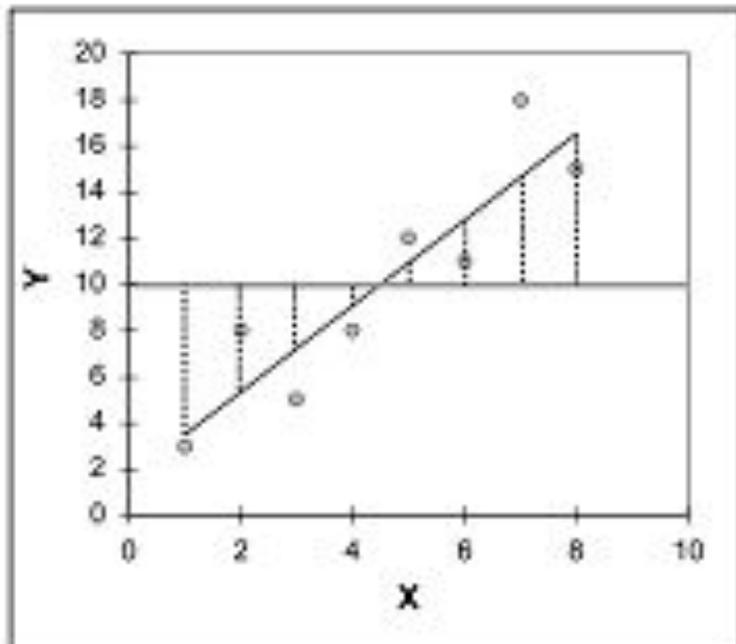
**SSR**

общая сумма различий между  
фактическими данными и моделью



# ЛИНЕЙНЫЙ РЕГРЕССИОННЫЙ АНАЛИЗ

ДАЛЕЕ ПРОГРАММА АНАЛИЗИРУЕТ,  
НАСКОЛЬКО ХОРОШО МОДЕЛЬ (на)  
ПРЕДСКАЗЫВАЕТ ЗАВИСИМУЮ  
ПЕРЕМЕННУЮ  
В СРАВНЕНИИ С ПРОСТОЙ СРЕДНЕЙ  
АРИФМЕТИЧЕСКОЙ (H0)



$SS_M = SST - SSR$   
ПОКАЗЫВАЕТ УЛУЧШЕНИЕ В  
ПРЕДСКАЗАТЕЛЬНОЙ СИЛЕ МОДЕЛИ В  
СРАВНЕНИИ С ПРОСТОЙ СРЕДНЕЙ  
АРИФМЕТИЧЕСКОЙ

$$RR^2 = \frac{SS(M)}{SS(T)}$$

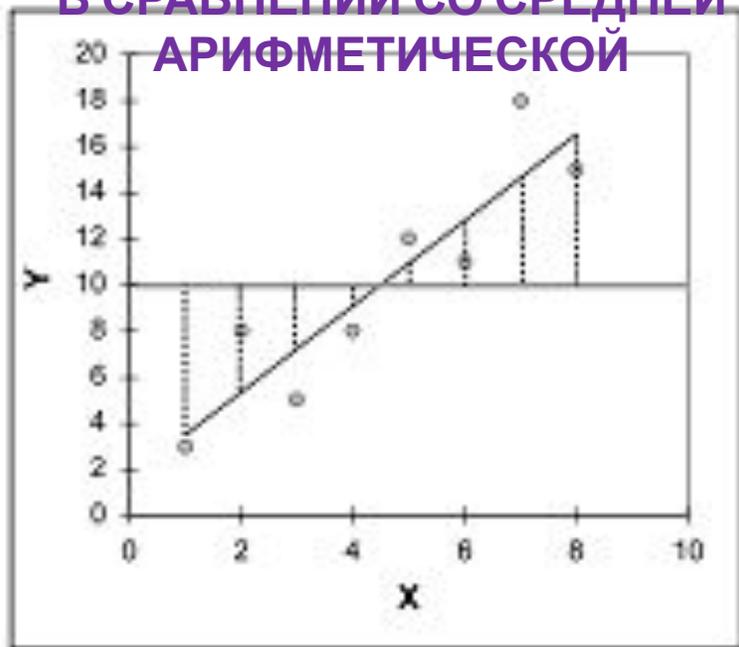
$RR^2$  - показывает количество дисперсии,  
которая объясняется моделью

$$1 - RR^2 = \text{ERROR}$$

# ЛИНЕЙНЫЙ РЕГРЕССИОННЫЙ АНАЛИЗ

**ПРОГРАММА РАССЧИТЫВАЕТ  
СТАТИСТИКУ РЕГРЕССИОННОЙ МОДЕЛИ  
(F – TEST)**

**СПОСОБНОСТЬ МОДЕЛИ УЛУЧШАТЬ  
ПРЕДСКАЗАНИЕ ЗАВИСИМОЙ ПЕРЕМЕННОЙ  
В СРАВНЕНИИ СО СРЕДНЕЙ  
АРИФМЕТИЧЕСКОЙ**



$$F_{\text{test}} = \frac{MS(M)}{MS(R)}$$

$$MS(M) = \frac{SS(M)}{df(M)}$$

$$MS(R) = \frac{SS(R)}{df(R)}$$

**$p(F\text{-test}) < 0,05$**   
**МОДЕЛЬ «РАБОТАЕТ», т.е.**  
предсказывает зависимую переменную  
лучше, чем средняя арифметическая

( $H_0$ )

# ЛИНЕЙНЫЙ РЕГРЕССИОННЫЙ АНАЛИЗ

$$Y = b_0 + b_1 X_1 + b_2 X_2 + \dots + b_N X_N + E$$

СТАТИСТИКА РЕГРЕССИОННОЙ МОДЕЛИ (F – TEST)

демонстрирует статистическую значимость  
всего уравнения регрессии

**$b_n$  – коэффициент регрессии / градиент** - демонстрирует изменение значения зависимой переменной ( $Y$ ) при изменении предиктора ( $X_n$ ) на “1” (единицу)

**Статистическую значимость каждого коэффициента регрессии необходимо оценить**

**$H_0: b_1 = 0$**

**$H_a: b_1 \neq 0$**

# УСЛОВИЯ ПРИМЕНЕНИЯ (ASSUMPTIONS) ЛИНЕЙНОГО РЕГРЕССИОННОГО АНАЛИЗА

- А) ЗАВИСИМАЯ ПЕРЕМЕННАЯ: **количественная непрерывная (неограниченная)**
- Б) НЕЗАВИСИМЫЕ ПЕРЕМЕННЫЕ (ПРЕДИКТОРЫ): **количественные непрерывные и дихотомические (0;1)**
- В) ЛИНЕЙНАЯ СВЯЗЬ: Graphs – Scatter/Dot (можно проверить для простой регрессии)
- Г) ГОМОСКЕДАСТИЧНОСТЬ - предполагается, что дисперсия ошибки остается той же самой в любой точке на протяжении всей линейной связи
- Д) НЕЗАВИСИМЫЕ НАБЛЮДЕНИЯ (DURBIN-WATSON  $\approx 2$  (DW  $\in [1;3]$ ))
- Е) НОРМАЛЬНО РАСПРЕДЕЛЕННЫЕ ОСТАТКИ (residuals)
- Ж) НЕ Д.Б. МУЛЬТИКОЛЛИНЕАРНОСТИ ( $R > 0,8$  – проблема;  $VIF > 10$  – проблема)

# УСЛОВИЯ ПРИМЕНЕНИЯ (ASSUMPTIONS) ЛИНЕЙНОГО РЕГРЕССИОННОГО АНАЛИЗА

OUTLIER: случаи, значительно влияющие на тренд ( $>2,58$  – проблема)

INFLUENTIAL CASES: случаи, заметно влияющие на модель (ее значимость)

COOK'S DISTANCE – д.б.  $< 1$  – мера влияния случая на модель

MAHALANOBIS – разница м/д случаем и средней арифметической

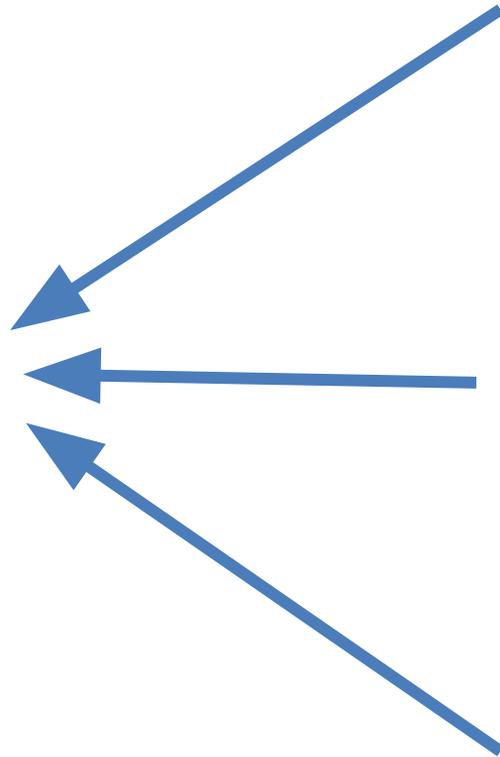
$N = 500$  – д.б.  $< 25$

$N = 100$  – д.б.  $< 15$

$N = 30$  – д.б.  $< 11$

# ПРИМЕР ЛИНЕЙНОГО РЕГРЕССИОННОГО АНАЛИЗА

**ВЕЛИЧИНА  
РАСХОДОВ  
ПАЦИЕНТОВ НА  
МЕДИКАМЕНТЫ**



**ПОЛ ПАЦИЕНТА**



**ВОЗРАСТ ПАЦИЕНТА**



**ДОХОД ПАЦИЕНТА**

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**1 ЭТАП:**  
ФОРМУЛИРУЕМ  $H_0$  и  $H_a$

ГИПОТЕЗЫ	ФОРМУЛИРОВКА
$H_0$ (нулевая гипотеза)	простая средняя арифметическая предсказывает исход лучше, чем модель регрессии
$H_a$ (альтернативная гипотеза)	модель регрессии предсказывает исход лучше, чем простая средняя арифметическая

**2 ЭТАП:**  
ОПРЕДЕЛЯЕМ УСЛОВИЯ, ПРИ КОТОРЫХ ПРИМЕМ  $H_a$  (ОТВЕРГНЕМ  $H_0$ )

**БУДЕМ** считать результаты теста «статистически значимыми» (т.е. примем  $H_a$ ) при вероятности ошибки 1 типа ( $\alpha$ -ошибки) менее 0.05 (5%)

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**3 ЭТАП:  
ВЫБОР  
СТАТИСТИЧЕСКОГО  
КРИТЕРИЯ / МЕТОДА**

зависимая переменная:  
количественная  
непрерывная  
**ВЕЛИЧИНА РАСХОДОВ  
ПАЦИЕНТОВ НА  
МЕДИКАМЕНТЫ**

**ПОДХОДИТ  
МНОЖЕСТВЕННАЯ  
ЛИНЕЙНАЯ РЕГРЕССИЯ**

предикторы:  
количественная  
непрерывная /  
дихотомическая  
**ПОЛ ПАЦИЕНТА:  
дихотомическая**  
**ВОЗРАСТ ПАЦИЕНТА:  
количественная  
непрерывная**  
**ВМІ ПАЦИЕНТА:  
количественная  
непрерывная**

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**4 ЭТАП:  
МАТЕМАТИЧЕСКИЕ  
РАССЧЕТЫ**

формулируем  $H_0$  и  $H_a$

**$H_0$ : F-статистика модели стат.незначима**

**$H_a$ : F-статистика модели стат.значима**

**МОДЕЛЬ РЕГРЕССИИ «РАБОТАЕТ»  
(описывает данные лучше, чем средняя  
арифметическая)**

ANOVA<sup>c</sup>

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	4E+008	3	116900381,9	30,579	,000 <sup>a</sup>
	Residual	4E+009	1024	3822872,805		
	Total	4E+009	1027			
2	Regression	3E+008	2	173363139,4	45,347	,000 <sup>b</sup>
	Residual	4E+009	1025	3823021,092		
	Total	4E+009	1027			

a. Predictors: (Constant), BMI, Gender\_new, Age

b. Predictors: (Constant), Gender\_new, Age

c. Dependent Variable: Drug\_spendings

**$p < 0,0001$**

т.е. **МОЖЕМ** принять  **$H_a$**   
вероятность ошибки 1 типа  
(ошибочно принять  $H_a$  - найти  
то, чего нет)  $< 0,1\%$

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

**4 ЭТАП:  
МАТЕМАТИЧЕСКИЕ  
РАССЧЕТЫ**

**МОДЕЛЬ ОБЪЯСНЯЕТ  
8,1% ДИСПЕРСИИ ЗАВИСИМОЙ ПЕРЕМЕННОЙ**

Model Summary

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	Change Statistics					Durbin-Watson
					R Square Change	F Change	df1	df2	Sig. F Change	
1	,287 <sup>a</sup>	,082	,080	1955,217	,082	30,579	3	1024	,000	
2	,285 <sup>b</sup>	,081	,079	1955,255	-,001	1,040	1	1024	,308	1,997

- a. Predictors: (Constant), BMI, Gender\_new, Age
- b. Predictors: (Constant), Gender\_new, Age
- c. Dependent Variable: Drug\_spendings

# ПОРЯДОК ТЕСТИРОВАНИЯ СТАТИСТИЧЕСКОЙ ГИПОТЕЗЫ

## 4 ЭТАП: МАТЕМАТИЧЕСКИЕ РАСЧЕТЫ

формулируем  $H_0$  и  $H_a$  для t-статистики коэффициентов  $b$

$H_0$ : t-статистика  $b$  стат.незначима  
 $H_a$ : t-статистика  $b$  стат.значима

Coefficients<sup>a</sup>

Model		Unstandardized Coefficients		Standardized Coefficients	t	Sig.	95% Confidence Interval for B		Collinearity Statistics	
		B	Std. Error	Beta			Lower Bound	Upper Bound	Tolerance	VIF
1	(Constant)	288,542	344,777		,837	,403	-388,009	965,093		
	Age	29,726	3,986	,235	7,457	,000	21,903	37,548	,901	1,109
	Gender_new	-470,182	123,483	-,115	-3,808	,000	-712,490	-227,874	,977	1,023
	BMI	12,029	11,797	,032	1,020	,308	-11,120	35,179	,909	1,100
2	(Constant)	551,149	229,234		2,404	,016	101,328	1000,970		
	Age	30,885	3,821	,244	8,083	,000	23,387	38,383	,981	1,019
	Gender_new	-478,291	123,229	-,117	-3,881	,000	-720,100	-236,481	,981	1,019

a. Dependent Variable: Drug\_spending

# ИНТЕРПРЕТАЦИЯ РЕЗУЛЬТАТОВ

Coefficients<sup>a</sup>

Model		Unstandardized Coefficients		Standardized Coefficients	t	Sig.	95% Confidence Interval for B		Collinearity Statistics	
		B	Std. Error	Beta			Lower Bound	Upper Bound	Tolerance	VIF
1	(Constant)	288,542	344,777		,837	,403	-388,009	965,093		
	Age	29,726	3,986	,235	7,457	,000	21,903	37,548	,901	1,109
	Gender_new	-470,182	123,483	-,115	-3,808	,000	-712,490	-227,874	,977	1,023
	BMI	12,029	11,797	,032	1,020	,308	-11,120	35,179	,909	1,100
2	(Constant)	551,149	229,234		2,404	,016	101,328	1000,970		
	Age	30,885	3,821	,244	8,083	,000	23,387	38,383	,981	1,019
	Gender_new	-478,291	123,229	-,117	-3,881	,000	-720,100	-236,481	,981	1,019

a. Dependent Variable: Drug\_spending

$$Y = b_0 + b_1 X_1 + b_2 X_2 + \dots + b_N X_N$$

Y = расходы на медикаменты

B0 = CONSTANT = 551,1

B1 = ВОЗРАСТ = 30,9

B2 = ПОЛ = -478,3 (для мужчин)

**ДЛЯ 50-ЛЕТНЕГО МУЖЧИНА ВЕЛИЧИНА РАСХОДОВ  
НА МЕДИКАМЕНТЫ**

**РАСХОДЫ = 551,1 + 30,9 × 50 – 478,3 = 1617,8 руб. + ERROR**

The background features a sunset over the ocean with a blue-to-orange gradient. Overlaid on this are several technical graphics: a large circular gauge with numerical markings (100, 110, 120, 130, 140, 150, 160, 170, 180, 190, 200, 210) and arrows in the upper right; a smaller circular gauge in the upper left; and a circular arrow graphic in the lower left.

**КРАТКИЙ ОБЗОР МЕТОДОВ  
СТАТИСТИЧЕСКОГО АНАЛИЗА  
КОЛИЧЕСТВЕННЫХ ПЕРЕМЕННЫХ**