

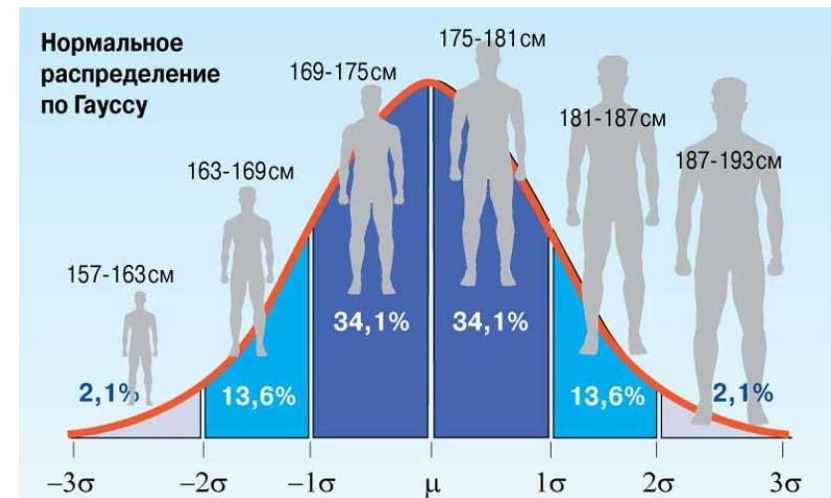
## Занятие 3

# Непараметрические критерии

**С точки зрения анализа данных их удобнее разбить на 3 группы.**

## **I. Количественные признаки с нормальным распределением**

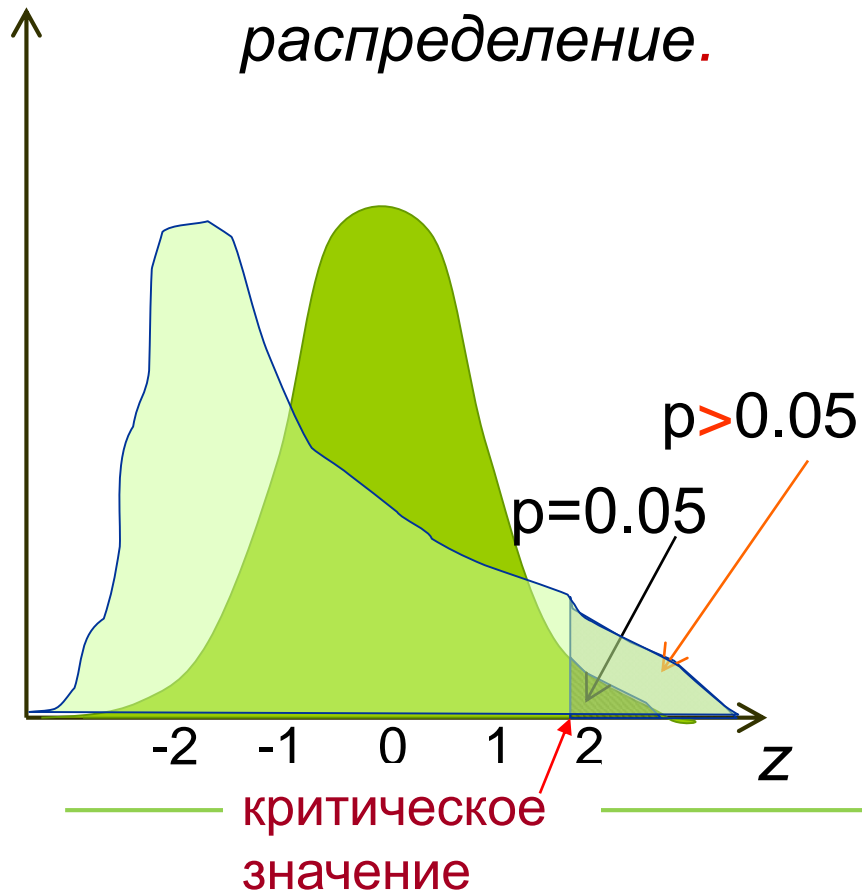
Используются для анализа **параметрические** тесты



Они задействуют в расчётах **параметры** известных распределений, главным образом — параметры нормального распределения (математическое ожидаемое  $\mu$  «мю» и стандартное отклонение  $\sigma$  «сигма»).

Почему при проведении параметрических тестов важно соблюдать **условие нормальности распределения**?

Распределение **статистики критерия не будет нормальным**, если в выборке не нормальное распределение.



Вероятность, что среднее в выборке попадёт в критическую область (рассчитанную для нормального распределения), будет выше, чем 0.05 – **увеличится ошибка 1-го рода**

## Основной вывод:

пренебрежения условиями использования параметрических тестов может **увеличивать ошибку 1-го рода** (отвергнута верная нулевая гипотеза об отсутствии связи между явлениями или искомого эффекта)

(Неизвестно, насколько)

**Примечание:** слабые отклонения от нормального распределения не очень страшны (в силу Центральной предельной теоремы), а **для больших выборок ими можно пренебречь** (кроме регрессионного анализа).

ANOVA устойчива к отклонениям от нормального распределения, особенно если выборки одинаковы по размеру.

## II. Количественные признаки с ненормальным распределением и порядковые признаки

Если нет уверенности в нормальности распределения признака или распределение неизвестно анализ данных можно провести 3-мя способами:

- 1) **нормализовать данные** с помощью преобразований шкалы (логарифмирование, преобразование арксинуса, Бокса — Кокса и др.) и использовать далее параметрические методы;
- 2) использовать **непараметрические методы**. Способ традиционен и популярен (медианы и квартили, корреляция Спирмена, критерии Уилкоксона — Манна — Уитни, Краскела — Уоллиса, Фридмана и др.).

3) работать с исходными непреобразованными данными методами, устойчивыми к отклонениям от нормальности.

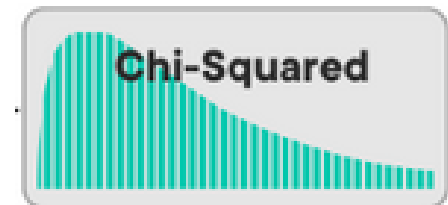
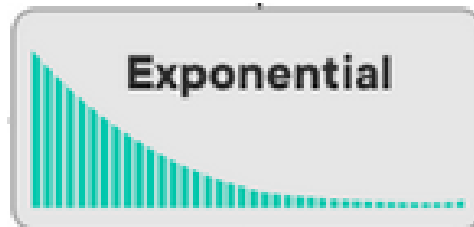
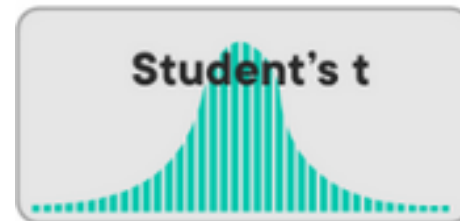
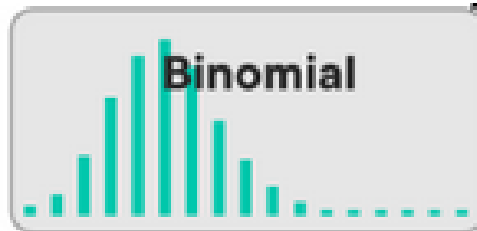
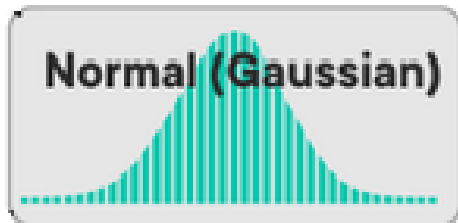
Это **методы робастной статистики** (усечённые средние или отличные от среднего М-оценки, средние абсолютные отклонения вместо среднеквадратичных и т. д.), или современные **ресэмплинг-техники**, основанные на вычислительных возможностях компьютеров (складной нож, бутстреп, рандомизационные методы Монте-Карло).

**III. Качественные признаки** (см. занятие «Анализ качественных»)

# Распределения бывают

Природные (нормальное, биномиальное, Пуассона, экспоненциальное, лог-нормальное и пр.)

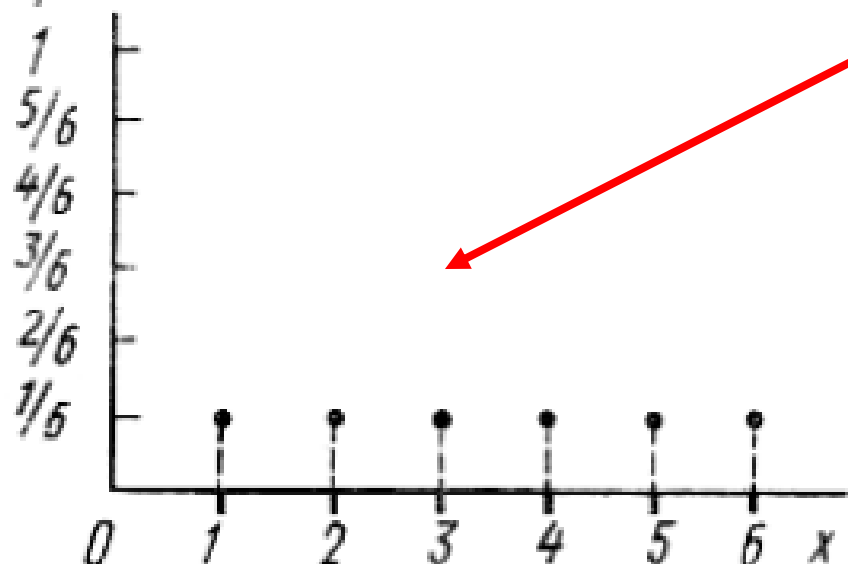
Распределения статистик критериев (t, F, U ...)



## Равномерное (uniform)

**Случайная величина** имеет **дискретное равномерное распределение**, если она принимает конечное **число значений с равными вероятностями**.

Вероятность



Распределение вероятностей дискретного равномерного распределения ( $n=6$ ).

Может быть и дискретным, и непрерывным



## **Биномиальное распределение**

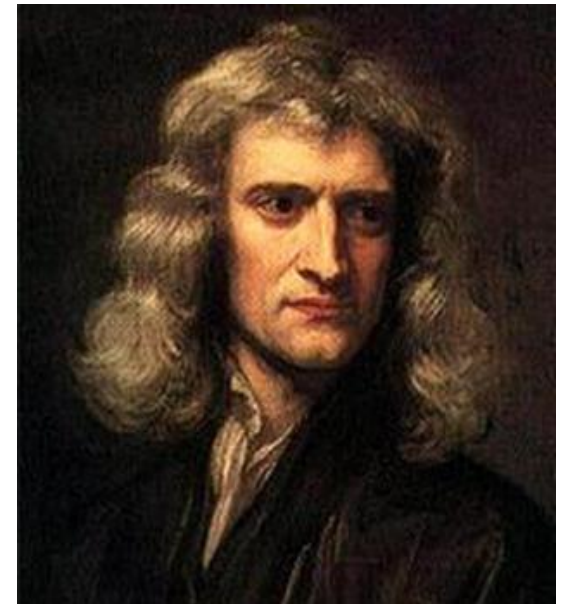
Одно из важнейших распределений вероятностей дискретно изменяющейся случайной величины. Введено в науку швейцарским математиком и одним из основателей теории вероятностей - **Якобом Бернулли** и опубликовано в **1713 году**.



**Якоб Бернулли**  
(*Jakob Bernoulli* 1655–1705)

**Вероятности** представляют собой **члены бинома Ньютона**, благодаря чему распределение и получило своё название.

Биномиальному распределению обычно соответствуют **доли, частоты**, пропорции



**Исаак Ньютон**  
(*Isaac Newton* 1643-1727)

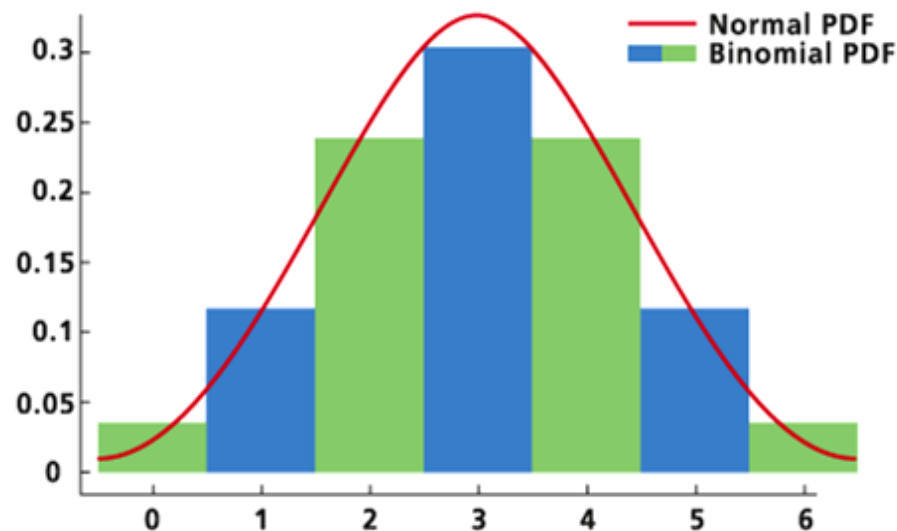
**Пример:** рассмотрим выводки из 6 детёнышей каждый.  
Возможное соотношение самцов и самок в выводке:

6:0; 5:1; 4:2; 3:3;  
2:4; 1:5; 0:6



Распределение количества самцов в  $N$  выводков (независимых случайных экспериментов) из  $n = 6$  зверьков, таких что вероятность рождения самца постоянна и равна  $p$ , а вероятность рождения самки  $q = 1 - p$ .

Вероятность такого выводка



Количество самцов в выводке из 6 зверьков

Если  $p$  мало, ситуация лучше описывается распределением Пуассона

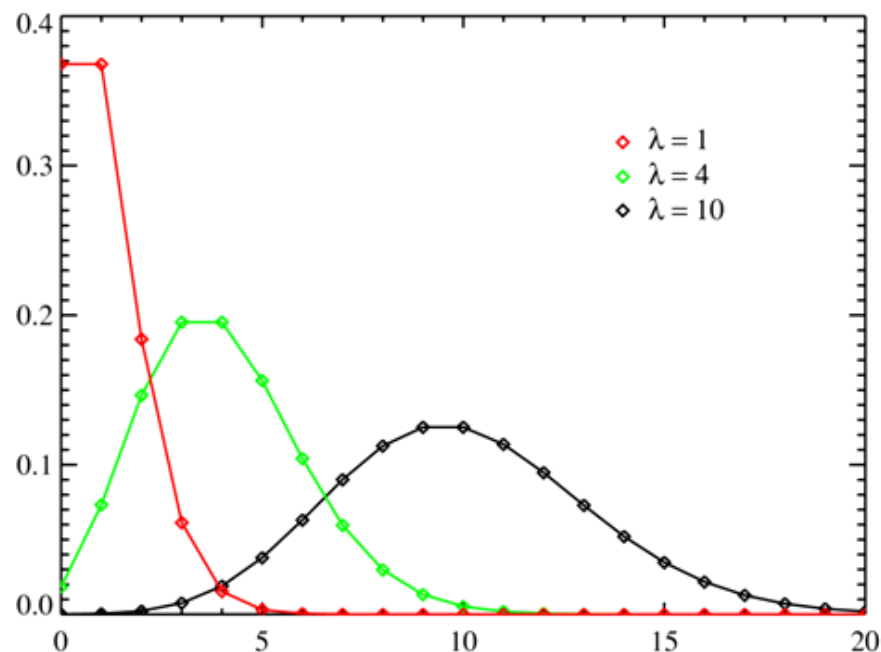
## Распределение Пуассона

Показывает вероятность того или иного количества независимых друг от друга **редких и случайных** событий (особей и пр.) на заданном интервале времени (участке пространства, объёме...).



$$\mu = \sigma^2$$

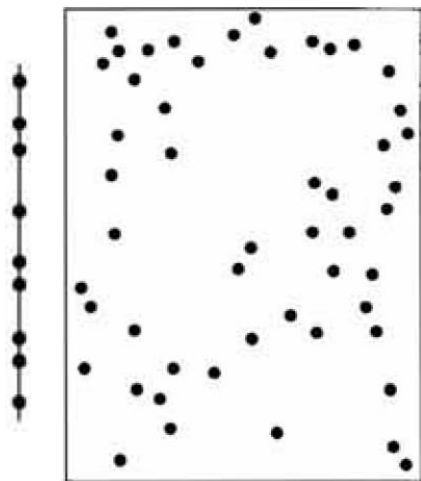
**Симеон Дени Пуассон**  
(*Siméon Denis Poisson 1781-1840*)



Распределение имеет один **параметр  $\lambda$**  (греческая буква «лямбда») – **среднее количество успешных испытаний в заданной области возможных исходов.**

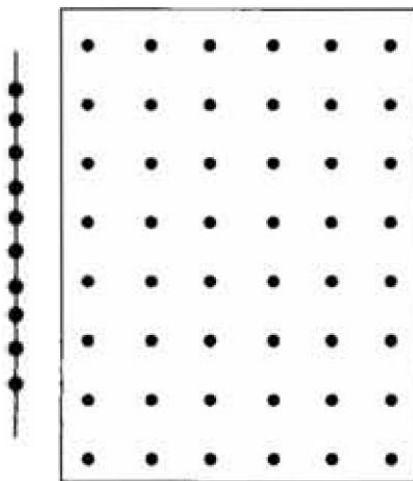
Распределению Пуассона соответствуют **частоты**, количества случайно распределённых объектов

# Сравнение распределения объектов во времени и пространстве со случайным распределением (testing for randomness)



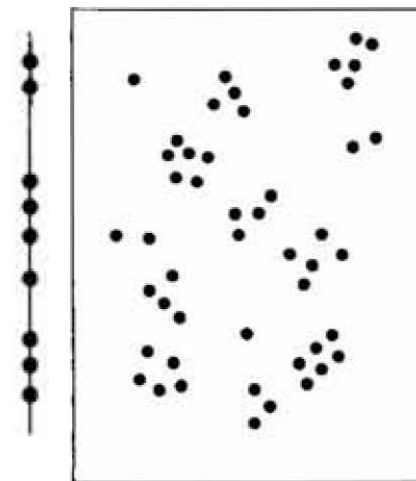
**случайное**

$$\sigma^2 = \mu$$



**равномерное**

$$\sigma^2 < \mu$$



**групповое**

$$\sigma^2 > \mu$$

Важно: следует задавать размер элементарной единицы пространства (времени и пр.), напр., квадрата, так, чтобы  $\mu \approx 1$

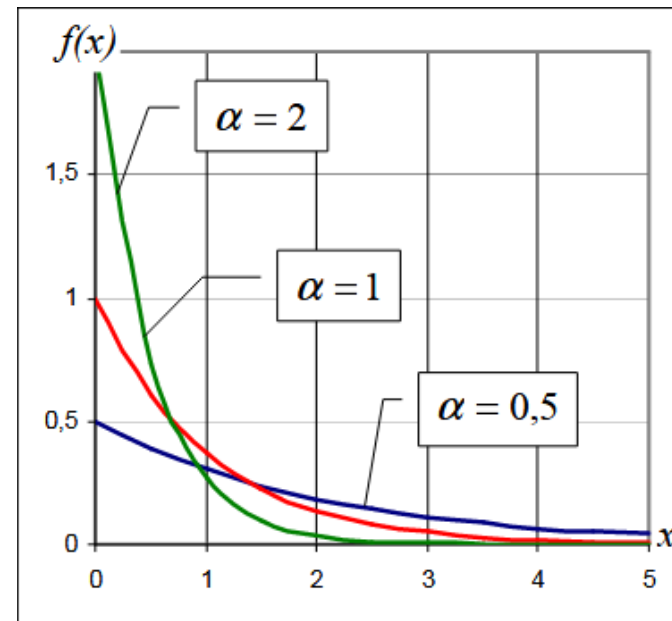
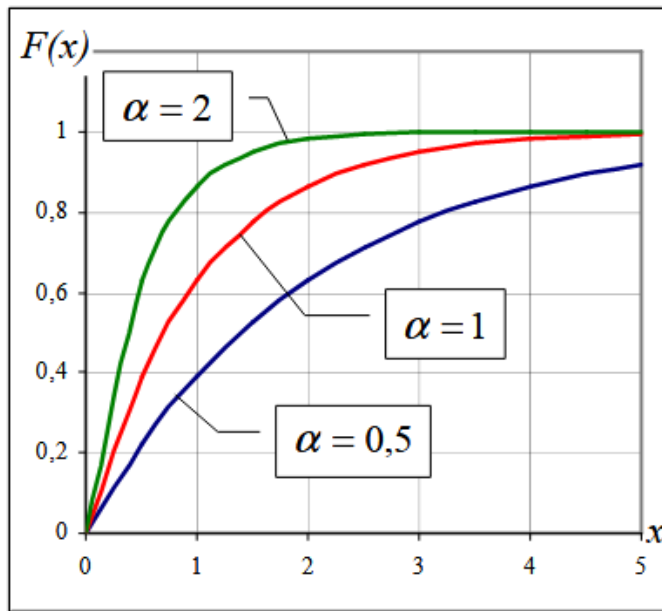
## Экспоненциальное распределение

Хорошо описывает распределение **промежутков времени** (расстояний) между случайными событиями с заданной средней частотой событий.

$$F(x) = 1 - e^{-\alpha x},$$

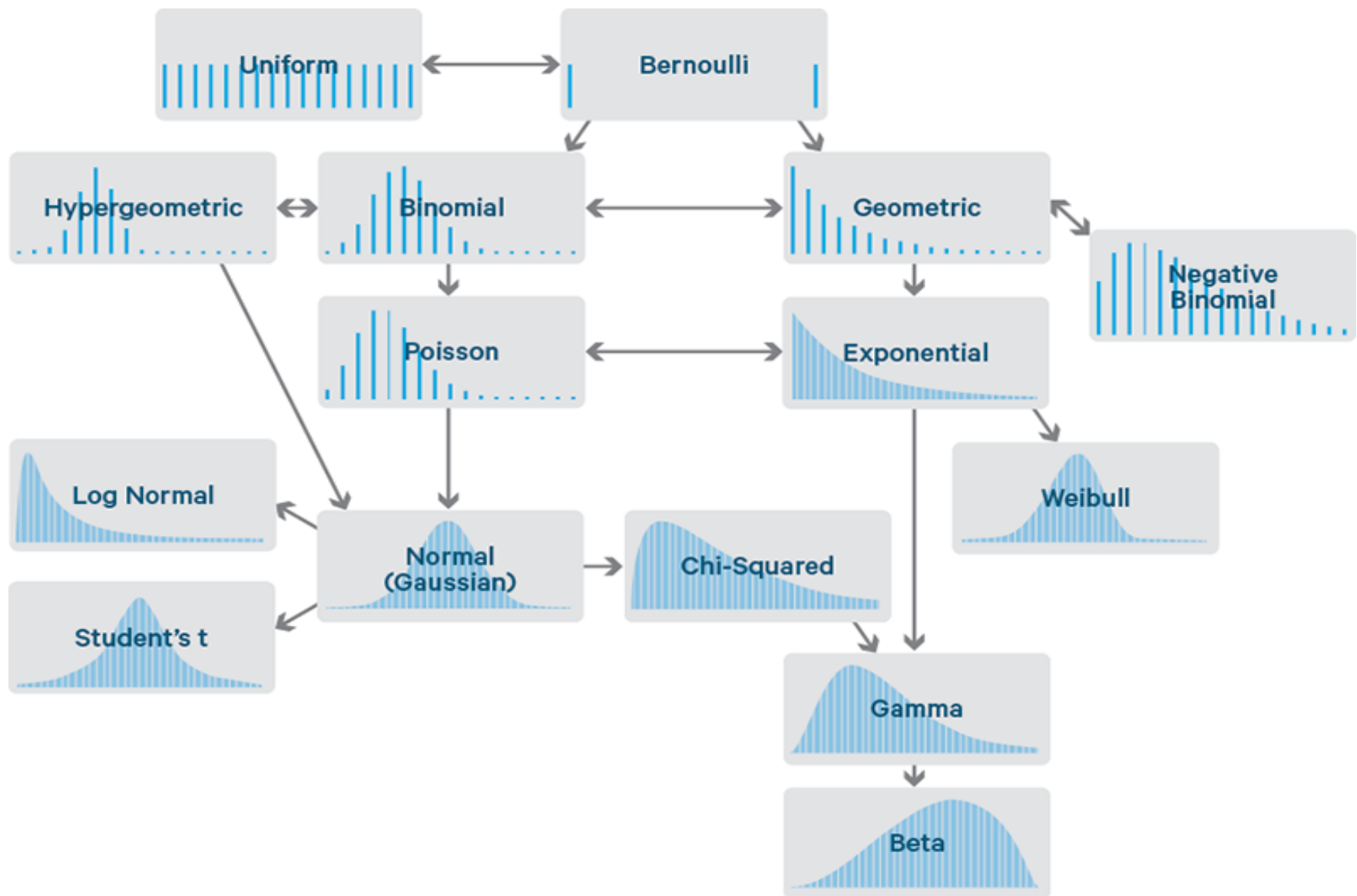
$$f(x) = \alpha e^{-\alpha x},$$

где  $\alpha > 0$  – параметр распределения;  $x \geq 0$  – непрерывная случайная величина.



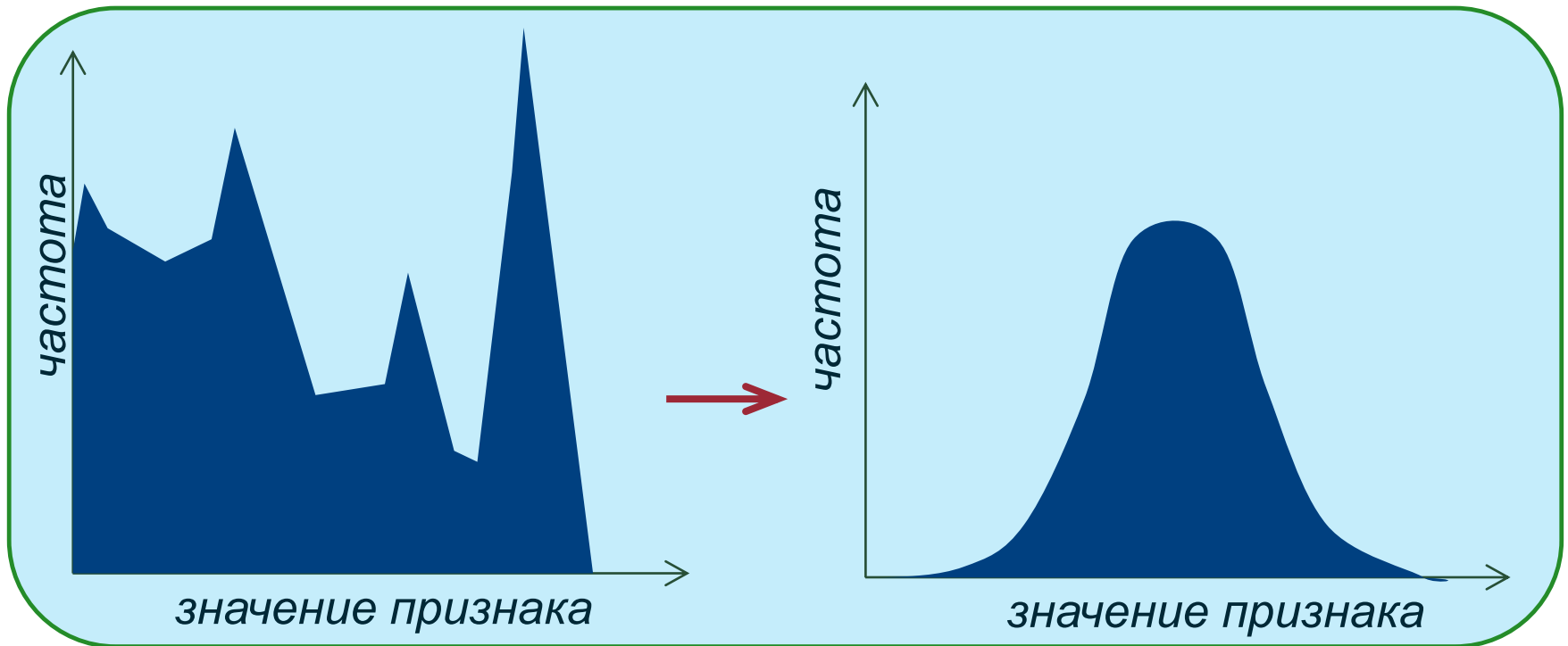
Функция и плотность экспоненциального распределения для трех значений параметра:  $\alpha=5,0$  ;  $\alpha=1$  ;  $\alpha=2$

# Типы распределений



# Трансформация данных

Если распределение отлично от нормального, выборки не гомогенны, факторы мультипликативны, можно **ТРАНСФОРМИРОВАТЬ** данные



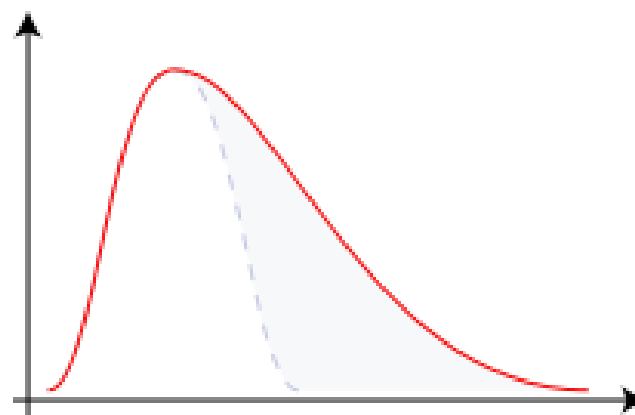
**Прекрасное свойство:** часто трансформация данных приводит одновременно к нормальному распределению, гомогенности и аддитивности

## 1. Логарифмическая трансформация (*logarithmic transformation*):

- Делает симметричным скошенное вправо (positively skewed) распределение.
- Используется в случае, когда среднее значение в группе прямо пропорционально стандартному отклонению.

$$X'_i = \lg X_i$$

$$X'_i = \lg(X_i + 1)$$



Positive Skew

Если в результате логарифмирования получилось нормальное распределение, исходное распределение было **логнормальным**.

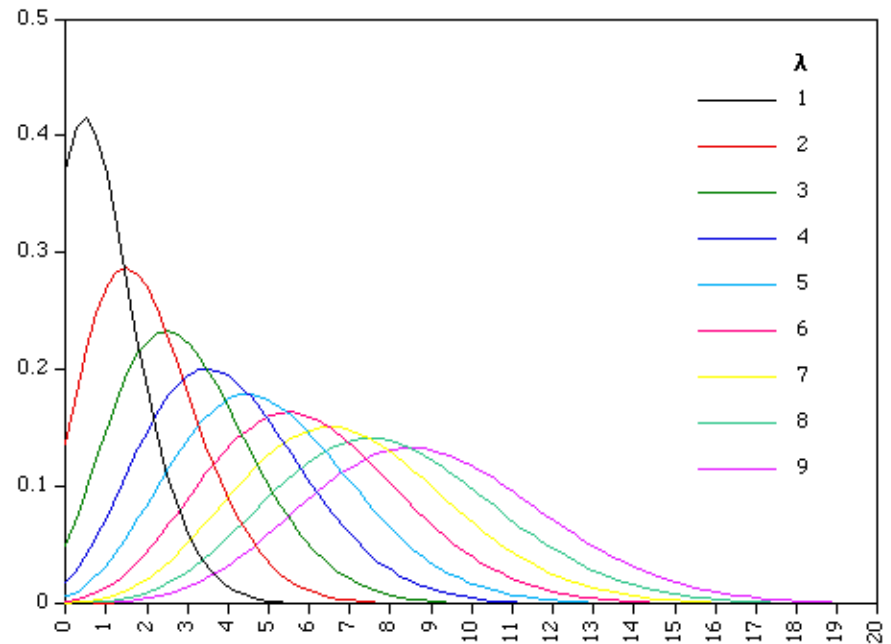


## 2. Извлечение квадратного корня (*square root transformation*)

- Используется, когда среднее значение в группе прямо пропорционально дисперсии.
- обычно такое явление свойственно выборкам из **распределения Пуассона** (т.е., данные представляют собой количества случайных событий, объектов...)

$$X'_i = \sqrt{X_i}$$

$$X'_i = \sqrt{X_i + 0,5}$$

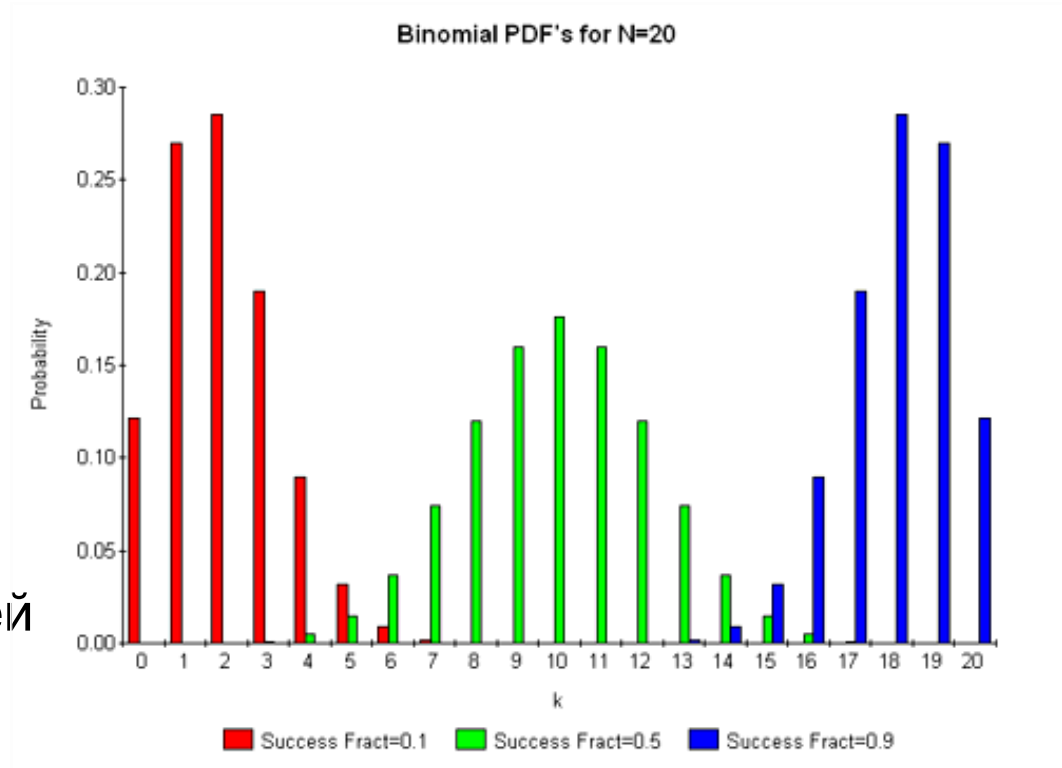


### 3. Арксинусная трансформация (*arcsine transformation*)

применяется для процентов и долей ( $X_i \leq 1$ ), которые обычно формируют биномиальное распределение.

$$X'_i = \arcsin \sqrt{X_i}$$

Например, мы исследуем долю самцов или долю переживших зиму детёнышей в выводках сурков.



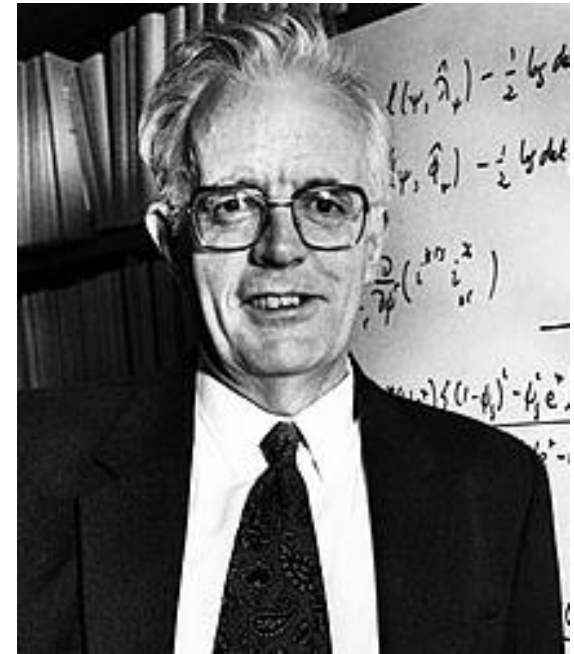
#### 4. Бокс-Кокс преобразование (*Box-Cox transformation*)

Среди множества методов преобразований одним из лучших (при неизвестном типе распределения) считается Бокс-Кокс преобразование.



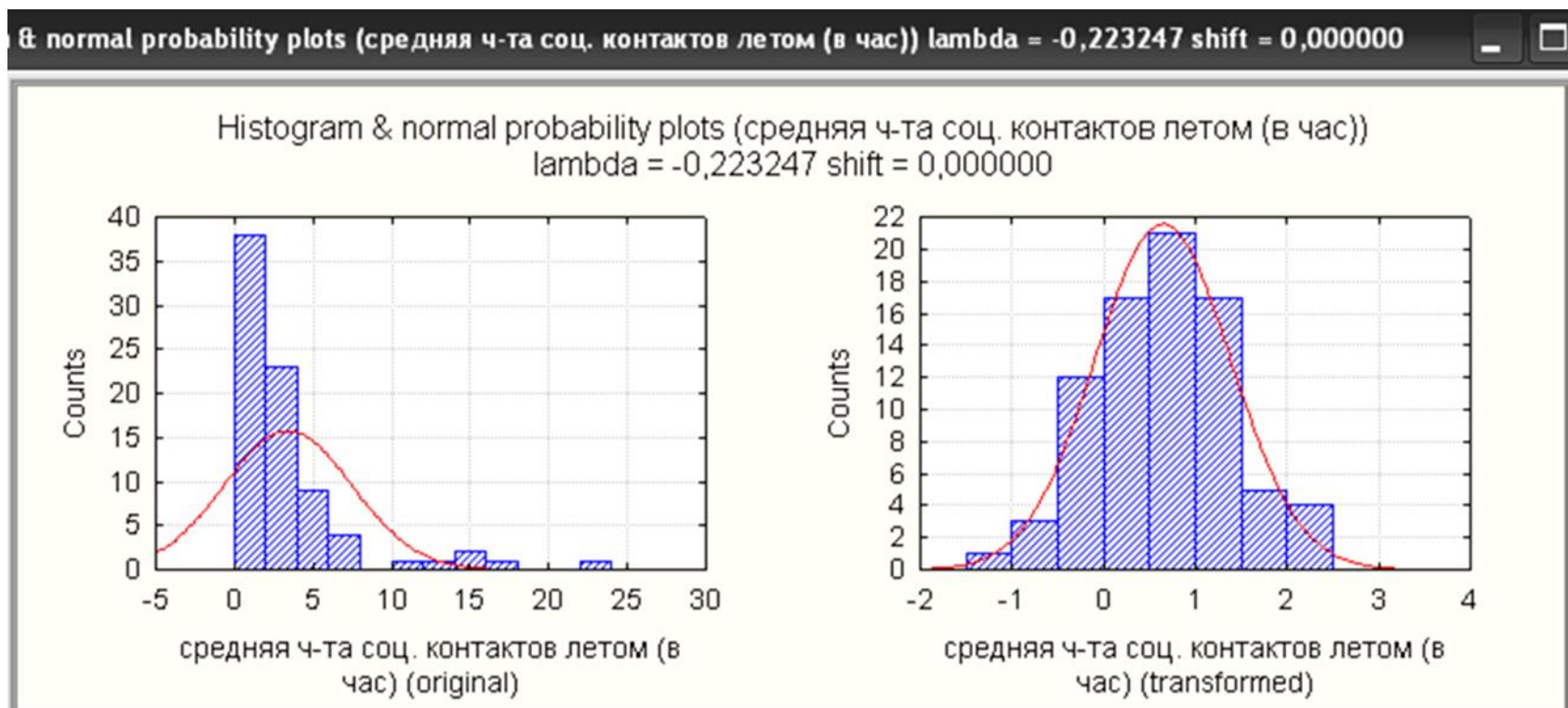
**Джордж Бокс**  
(*George E. P. Box 1919-2013*)

Сущность метода впервые была изложена в **1964 году**, в Журнале Королевского статистического общества, известными статистиками — Джорджем Боксом и сэром Дэвидом Коксом



**Дэвид Кокс**  
(*David Roxbee Cox 1927*)

Универсальная трансформация данных, в которой программа методом проб подбирает наилучшие параметры и способ трансформации для конкретных данных (ищется особый параметр  $\lambda$ )



В зависимости от значения лямбда, преобразование Бокса-Кокса включает в себя следующие частные случаи:

$\lambda = -1.0,$	$x_i(\lambda) = \frac{1}{x_i}$
$\lambda = -0.5,$	$x_i(\lambda) = \frac{1}{\sqrt{x_i}}$
$\lambda = 0.0,$	$x_i(\lambda) = \ln(x_i)$
$\lambda = 0.5,$	$x_i(\lambda) = \sqrt{x_i}$
$\lambda = 2.0,$	$x_i(\lambda) = x_i^2$

При использовании Бокс-Кокс преобразования необходимо, чтобы все значения входной последовательности были положительными и отличными от нуля.

# Box-Cox transformation

Трансформация Box-Cox.dat

File Edit Transform Plot Univariate Multivariate Model Diversity Timeseries Geometry Stratigraphy

Show

☐ Row at a time

☒ Column

Type

Name

1

2

3

4

Log

Subtract mean

Remove trend

Convert to ranks

Row percentage

Row normalize length

Box-Cox

Compositional data transforms

Remove size from distances

Edit

Cut

Copy

Paste

Select all

View

☐ Bands

☐ Binary

Исходные данные	Трансформированные
-	-
Исходные данные	Трансформированные
151	5,11761602183518
145	5,07544422330792
99	4,67918786483493
123	4,90443520177034
134	4,99342555438293
147	5,08969082103865
135	5,00115283513417
139	5,03150431271104
140	5,03895676730543
143	5,06100131726685
78	4,42822814264722

Box-Cox transformation

Lambda: 0,0433171

Log likelihood: -385

Close Help Transform

---

*Если распределение не удовлетворяет условиям параметрических тестов и трансформация не помогает или невозможна, спользуем*

## Непараметрические методы (nonparametric methods)

= “distribution-free” tests

- ✓ Свойства распределения неизвестны, и **параметры** **распределения** (среднее, дисперсию и т. п.) мы использовать не можем
  - ✓ Основной подход – ранжирование (*ranking*) наблюдений (выстраиваем их по порядку от самого маленького значения к наибольшему).
  - ✓ подразумевается, что сравниваемые распределения имеют одинаковую форму и дисперсию.
-

## Сравнение 2-х независимых групп

### Манн-Уитни тест (*Mann-Whitney U-test*)

В 1947 году двумя американскими математиками – **Манном** и **Уитни** для сравнения 2-х независимых выборок был предложен не параметрический тест.



**Генри Манн**  
(*Henry Berthold Mann*  
1900-2000)

Непараметрический аналог теста Стьюдента.

Является развитием идей Франка Уилкоксона изложенных в 1945 году. Поэтому в ряде случаев называется – тест **Уилкоксона-Манна-Уитни**



**Дональд Рэнсом Уитни**  
(*Donald Ransom Whitney*  
1915—2001)



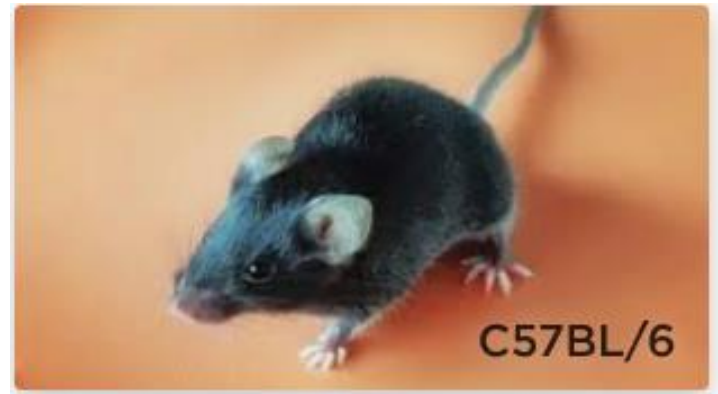
**Пример.** Мы исследуем две линии лабораторных мышей. Хотим сравнить размеры выводков у этих зверей.

**Фактор** – линия: 1. белые (BaLB/C); 2. черные (C57BL/6)

**Зависимая переменная** – размер выводка



*белые*



*черные*

$H_0$ : размер выводка у белые мышей такой же, как и у черных.

$H_1$ : размер выводка не одинаков у этих линий.

**Мы ничего не говорим про параметры распределений!**

*Тест Манна-Уитни можно использовать и для ранговых, и для непрерывных переменных.*

белые		черные	
размер	ранг	размер	ранг
8	15.5	4	5
7	13	7	13
4	5	5	8.5
7	13	8	15.5
9	17.5	3	2
3	2	3	2
5	8.5	5	8.5
6	11	4	5
9	17.5		
5	8.5		
111.5		59.5	

Ранжируем данные от меньшего к большему (**игнорируя** деление на группы).

Число 3 встретилось трижды (это называется **связанные ранги**, *tied ranks*): ранги у них будут одинаковы  $(1+2+3)/3=2$

Статистика критерия:

$$U_1 = n_1 n_2 + \frac{n_1(n_1 + 1)}{2} - R_1$$

$$U_2 = n_1 n_2 + \frac{n_2(n_2 + 1)}{2} - R_2$$

$n_1$  и  $n_2$  – размер выборок,  
 $R_1$  и  $R_2$  – суммы рангов в выборках.

Статистикой критерия  $U_{obs}$  будет **меньшее** из этих двух значений. Причём  $H_0$  мы отвергнем в случае, если оно будет МЕНЬШЕ критического значения  $U_{cv}$ . (т.е., это исключение среди прочих критериев).

Подставим наши данные в формулы:

$$U_1 = 10 \times 8 + \frac{10(8+1)}{2} - 115,5 = 135 - 115,5 = 23,5$$

$$U_2 = 10 \times 8 + \frac{8(8+1)}{2} - 59,5 = 116 - 59,5 = 56,5$$

$$U_{cv} = \mathbf{20}, \text{ при } p=0,05$$

$$U_{obs} = 23,5$$

$$U_{obs} > U_{cv}$$



$n_1$	2	3	4	5	6	7	8	9	10	11	12
$n_2$	$p=0,05$										
3	-	0									
4	-	0	1								
5	0	1	2	4							
6	0	2	3	5	7						
7	0	2	4	6	8	11					
8	1	3	5	8	10	13	15				
9	1	4	6	9	12	15	18	21			
10	1	4	7	11	14	17	20	24	27		
11	1	5	8	12	16	19	23	27	31	34	

Следовательно статистически значимых различий в величине выводка у разных линий мышей не наблюдается ( $U=23,5$ ,  $p<0,05$ )

# Тест Колмогорова-Смирнова

*(Kolmogorov-Smirnov two-sample test).*



**Колмогоров  
Андрей Николаевич**  
(1903-1987)

Отличается от М-У теста тем, что М-У более чувствителен к различиям средних значений, медианы и т.п., а К-С тест более чувствителен к различиям распределений по форме.



**Смирнов  
Николай Васильевич**  
(1900-1966)

---

Манн-Уитни тест более мощный, чем этот тест.

# Mann-Whitney U-test

## Kolmogorov-Smirnov two-sample test



В пакете PAST

Плодовитость лабораторных мышей.dat

File Edit Transform Plot **Univariate** Multivariate Model Diversity Timeseries Geometry Stratigraphy Script Help

Show

☐ Row attributes

☐ Column attributes

View

☐ Bands Recover windows

Белые мыши

1	•	8
2	•	7
3	•	4
4	•	7
5	•	9
6	•	3
7	•	5
8	•	6
9	•	9
10	•	5
11	•	
12	•	

- Summary statistics
- One-sample tests (t, Wilcoxon, single-case)
- Two-sample tests**
  - Two-sample tests (F, t, Mann-Wh, Kolm-Sm etc.)
  - Two-sample paired tests
  - F and t tests from parameters
- ANOVA etc. (several samples)
- Correlation
- Intraclass correlation
- Normality tests
- Contingency table (chi<sup>2</sup> etc.)
- Mantel-Cochran-Haenszel test
- Risk/odds
- Single proportion test
- Multiple proportion CIs
- Ratios of counts CI
- Survival analysis
- Combine errors

Не отвергаем  $H_0$ : М-У тест показал, что размеры выводков у разных линий **одинаковые**

Two-sample tests

t test   F test   Mann-Whitney   Mood median   Kolm-Smirnov

Tests for equal medians

Белые мыши	Черные мыши
N: 10	N: 8
Mean rank: 6,1944	Mean rank: 3,3056
Mann-Whitn U: 23,5	
z: 1,4396	$p$ (same med.): 0,14999
Monte Carlo permutation:	$p$ (same med.): 0,1476
Exact permutation:	$p$ (same med.): 0,14713

Two-sample tests

t test   F test   Mann-Whitnev   Mood median   Kolm-Smirnov

Kolmogorov-Smirnov test for equal distributions

Белые мыши	Черные мыши
N: 10	N: 8
D: 0,35	$p$ (same dist.): 0,54719
Monte Carlo permutation:	$p$ (same dist.): 0,408

## Сравнение 2-х связанных групп

### Критерий Уилкоксона (*Wilcoxon matched pairs test*)

W критерий Уилкоксона - это непараметрический аналог парного критерия Стьюдента (t-критерия).



**Фрэнк Уилкоксон**  
(*Frank Wilcoxon*  
1892-1965)

Предложен **в 1945 году** американским химиком и статистиком **Френком Уилкоксоном** (создатель первого университетского курса по непараметрической статистике и первой научной школы непараметрической статистики).

Мощность – около 95% мощности t-теста. При числе пар  $>100$  T аппроксимируется нормальным распределением.

**Пример.** Сравниваем 2 метода определения тестостерона в пробах, и хотим знать – различается ли его содержание в зависимости от метода определения

$H_0$ : количество тестостерона при определении первым методом, **такое же**, как и вторым.

$H_1$ : количество тестостерона не одинаково.

**Фактор** – метод определения. (Метод 1; Метод 2)

**Зависимая переменная** – содержание тестостерона в пробе.





№ пробы	Метод 1	Метод 2	$D_i = X_{i1} - X_{i2}$	Ранг
1	0,49	0,40	0,09	5
2	0,71	0,61	0,10	6
3	0,96	0,84	0,12	8
4	0,41	0,35	0,06	3,5
5	0,48	0,51	-0,03	-2
6	0,71	0,60	0,11	7
7	0,41	0,42	-0,01	-1
8	0,52	0,52	0,00	
9	0,63	0,57	0,06	3,5

$$T_{\text{эмп}} = 1 + 2 = 3$$

По таблице критических значений критерия Вилкоксона определяем, что при  $n=9$

$$T_{\text{кр}} = 8, \text{ для } p \leq 0,05$$

$$T_{\text{эмп}} < T_{\text{кр}(0,05)}$$

*Различия статистически значимы ( $p < 0,05$ )*

1. Считают разности между значениями в парах;

2. исключают нулевые разности;

3. присуждают абсолютным значениям (по модулю) разностей ранги;

4. суммируют отдельно ранги положительных и отрицательных разностей;

5. Наименьшая из этих сумм - статистика  $T_{\text{эмп}}$ .

6. Отвергаем  $H_0$ , если  $T_{\text{эмп}}$  меньше  $T_{\text{кр}}$ .

# Wilcoxon matched pair test

The screenshot shows the SPSS software interface. The 'Univariate' menu is open, and the path 'Two-sample tests' > 'Two-sample paired tests' is highlighted with red circles. The 'Two-sample paired tests' dialog box is open, showing the following data:

Метод 1 Тестостерон (мг)		Метод 2 Тестостерон (мг)	
N:	9	Mean:	0,53556
Mean:	0,59111	Median:	0,52
Median:	0,52		

**t test**

Mean difference:	0,055556	95% conf.:	(0,012567 0,098544)
t:	2,9801	p (same mean):	0,017598
Exact:		p (same mean):	0,027344

**Sign test**

r:	6	p (same median):	0,28906
----	---	------------------	---------

**Wilcoxon test :**

W:	33	p (same median):	0,035692
Normal appr. z:	2,1004	p (same median):	0,03852
Monte Carlo (n=99999):		p (same median):	0,039063
Exact:			

Содержание тестостерона при определении разными методами неодинаково ( $p=0,04$ )

## Сравнение $\geq 3$ -х независимых групп

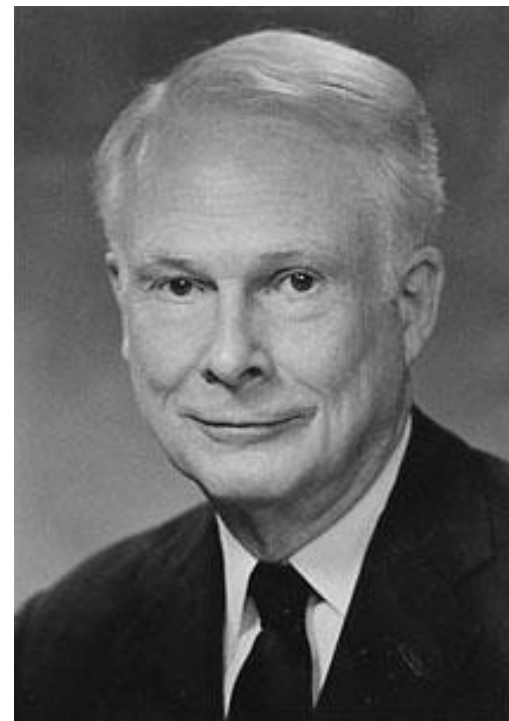
### Тест Крускала-Уоллиса (*Kruskal-Wallis test*)



**Уильям Крускал**  
(*William Kruskal*  
1919-2005)

Непараметрический аналог однофакторного дисперсионного анализа и предназначен для проверки равенства медиан нескольких выборок.

Предложенный в **1952 году** американскими учеными — математиком **Уильямом Крускалом** и экономистом **Алленом Уоллесом**



**Уилсон Аллен Уоллис**  
(*Wilson Allen Wallis*  
1912-1998)

- ✓ Непараметрический аналог One-way ANOVA
- ✓ на 95% настолько же мощный, как и ANOVA;
- ✓ для 2-х групп идентичен Манн-Уитни тесту;
- ✓ подразумевает сходство форм распределений и равенство дисперсий в группах (хотя бы на глаз)

**Пример.** Нас интересует, различается ли масса тела студентов, из разных групп разбитых по росту.

**Фактор** – рост. Группы: 1. 161-165 см.; 2. 166-170 см; 3. 171-175 см.; 4. 176-180 см.

**Зависимая переменная** – масса тела, кг.



$H_0$ : **распределение** в разных группах, из которых мы получили выборки, **одинаковое**.

$H_1$ : распределения не одинаковые.

1. все значения ранжируются от меньшего к большему (игнорируя деление на группы);
2. Считается сумма рангов в каждой группе;

Рост 161-165 см		Рост 166-170 см		Рост 171-175 см		Рост 176-180 см	
масса, кг	ранг	масса, кг	ранг	масса, кг	ранг	масса, кг	ранг
59	4,5	63	9	67	11	73	17
53	1	61	7	68	12,5	79	20
60	6	68	12,5	74	18	71	15
54	2	62	8	72	16	75	19
57	3	64	10	69	14		
59	4,5						
$\Sigma$	21		46,5		71,5		71

3. считается статистика  $H(df, N)$ .

сумма рангов в  
каждой группе

$$H = \frac{12}{N(N+1)} \sum \frac{R_j^2}{n_j} - 3(N+1)$$

общий размер  
выборки

размер группы

$$H = \frac{12}{20 \times (20 + 1)} \left[ \frac{21^2}{6} + \frac{46,5^2}{5} + \frac{71,5^2}{5} + \frac{71^2}{4} \right] - 3(20 + 1) = 16,676$$

При уровне значимости  $\alpha = 0,01$  и числе степеней свободы  $df = k - 1 = 4 - 1 = 3$ , где  $k$  – число групп;  $\chi^2_{кр} = 11,34$

$$H > \chi^2_{кр}$$

Следовательно **масса тела не одинакова** в разных группах.

Рост и масса тела студентов.dat

		Рост 161-165 см	Рост 166-170 см	Рост 171-175 см	Рост 176-180 см
1	•	59	63	67	73
2	•	53	61	68	79
3	•	60	68	74	71
4	•	54	62	72	75

# Kruskal-Wallis test

File Edit Transform Plot **Univariate** Multivariate Model Diversity Timeseries Geometry Stratigraphy Script Help

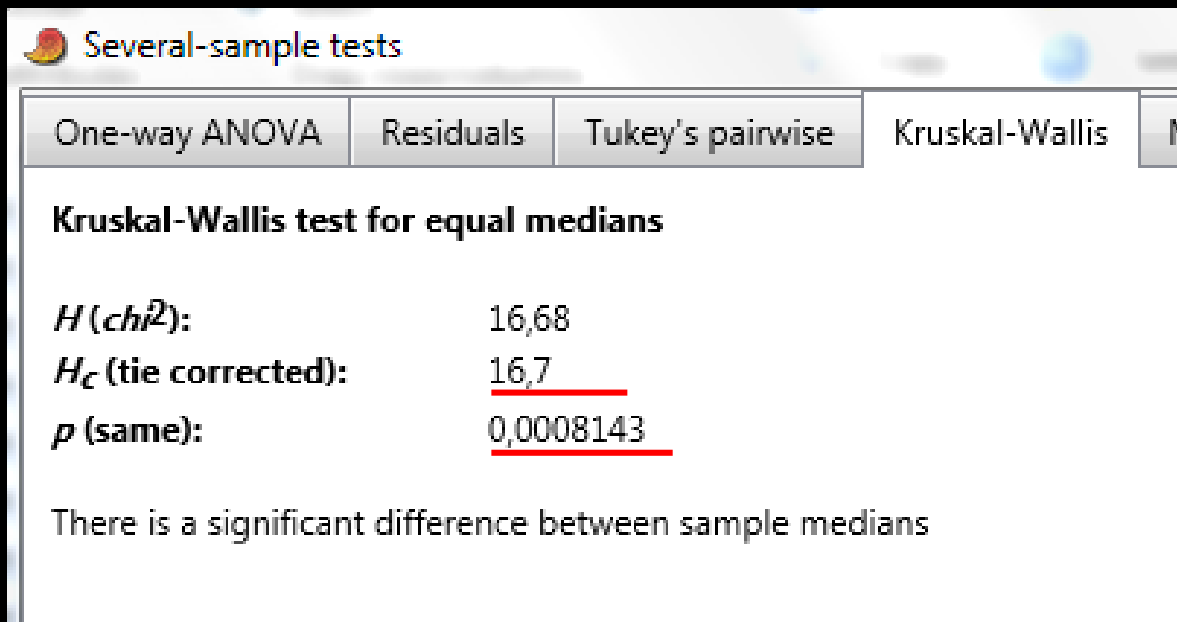
Show ☐ Row attributes ☒ Click mode ☐ Drag rows/columns Edit Cut Paste Copy Select all

**ANOVA etc. (several samples)**

- Summary statistics
- One-sample tests (t, Wilcoxon, single-case)
- Two-sample tests
- ANOVA etc. (several samples)**
  - Several-sample tests (ANOVA, Kruskal-Wallis)**
  - Several-sample repeated measures tests
  - Two-way ANOVA
  - Two-way ANOVA without replication
  - Two-way repeated measures ANOVA
  - One-way ANCOVA
- Correlation
- Intraclass correlation
- Normality tests
- Contingency table (chi<sup>2</sup> etc.)
- Mantel-Cochran-Haenszel test
- Risk/odds
- Single proportion test
- Multiple proportion CIs
- Ratios of counts CI
- Survival analysis
- Combine errors

View ☐ Bands ☐ Binary Recover windows Decimals: -

	Рост 161-165 см
1	• 59
2	• 53
3	• 60
4	• 54
5	• 57
6	• 59
7	•
8	•
9	•
10	•
11	•
12	•
13	•



$H_c$  (tie corrected) -  $H$ -критерия с поправкой на связанные значения (одинаковые значения в разных группах)

Масса тела студентов статистически значимо отличается в разных по росту группах ( $H_{(2)} = 6,70$ ;  $P < 0,001$ ).



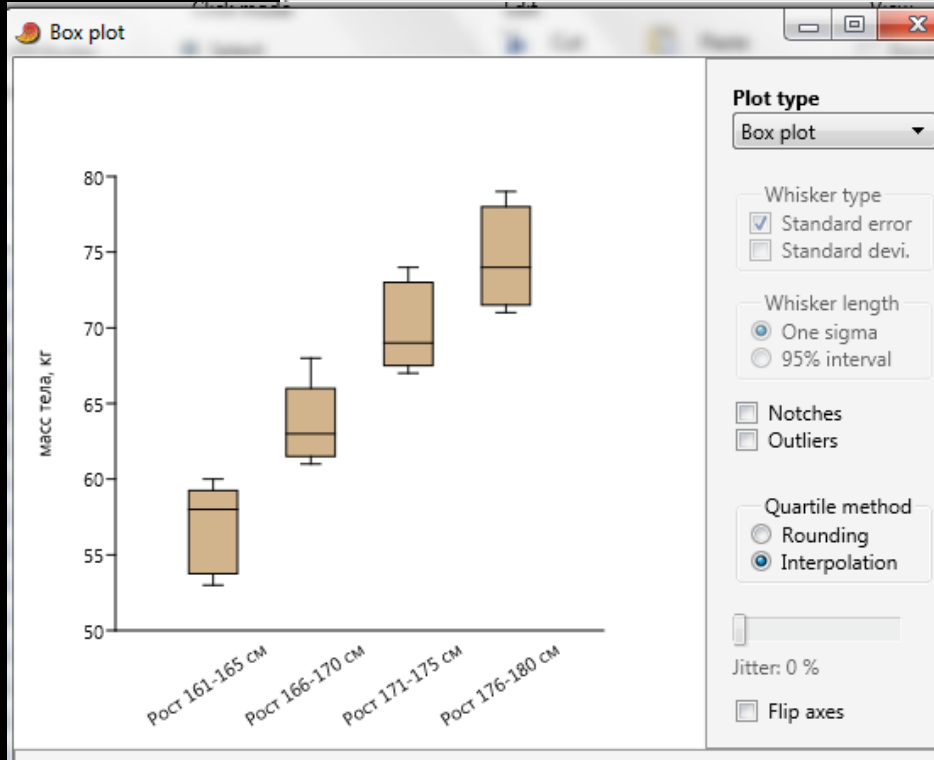
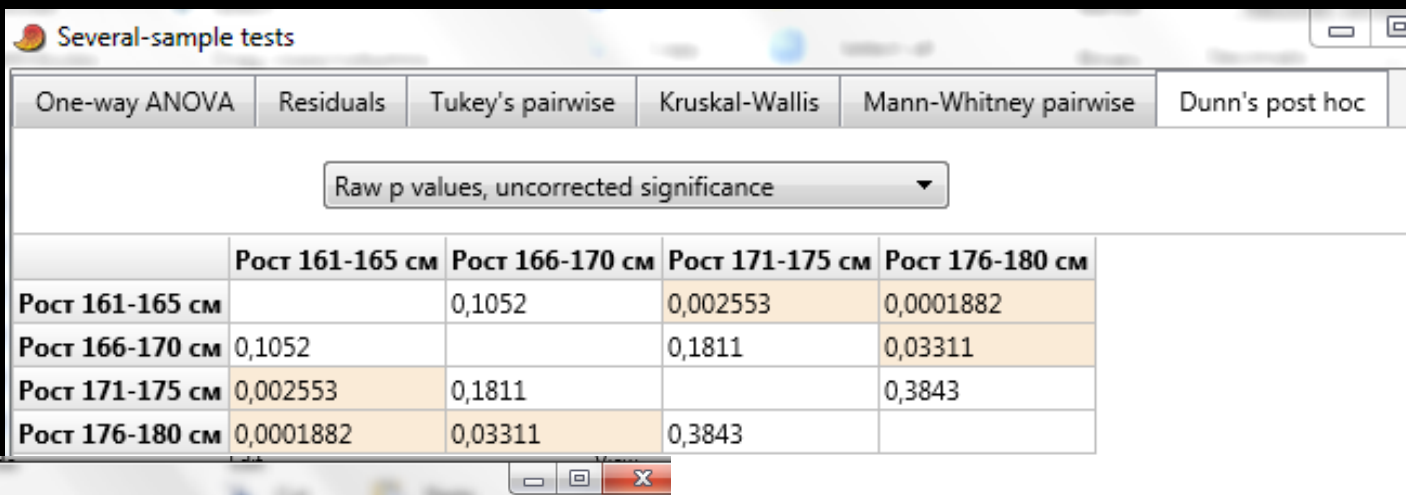
## Множественные апостериорные парные сравнения (*post-hoc comparisons*)

Как и в ANOVA, после сравнения нескольких групп имеет смысл провести **множественные апостериорные сравнения (*post-hoc comparisons*)**, по аналогии с тестом Тьюки, чтобы выяснить какие же группы различаются.

Такие тесты существуют – **Данна** (*Dunn's test*), **Манн-Уитни** (*Mann-Whitney pairwise*), **Неменьи** (*Nemenyi test*).

В ходе решения нашего примера далее для парного сравнения используем непараметрический **критерий Данна** (*Bonferroni–Dunn post hoc test, Dunn's multiple comparison post – test*).

Критерий применим для независимых групп как равной, так и различной численности.



Пост-хок тест для  
непараметрической ANOVA

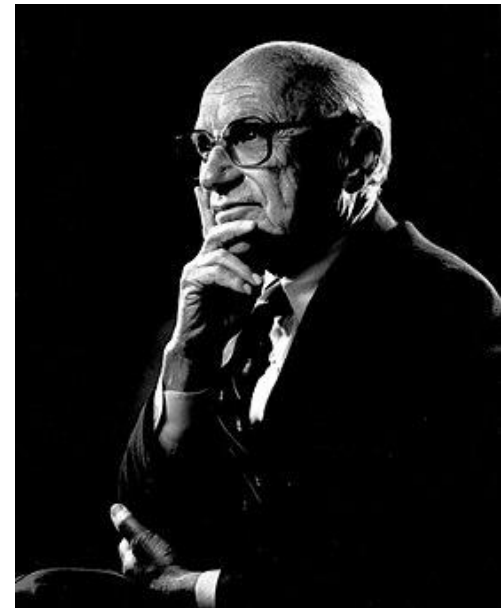
## Сравнение $\geq 3$ связанных групп

### Критерий Фридмана (*Friedman ANOVA*)

Непараметрический статистический тест, аналог дисперсионного анализа с повторными измерениями ANOVA.

Разработанный американским экономистом, нобелевским лауреатом по экономике Милтоном Фридманом.

Является обобщением критерия Уилкоксона на большее, чем 2, количество условий измерения.



**Милтон Фридман**  
(*Milton Friedman 1912-2006*)

По сравнению с аналогичными параметрическими тестами, для 2-х групп имеет всего 64% мощности, для 3-х – 72%, для 100 стремится к 95%.

**Пример.** Группа из шести человек, желающих отказаться от курения. Проводилось измерение жизненной емкости легких (ЖЕЛ) в динамике:

1. На момент включения в группу;
2. Через 1 месяц после отказа от курения;
3. Через 2 месяца после отказа от курения



**1.** Значения ранжируются меньшего к большему внутри каждой **строки**.

**2.** Суммируют ранги для каждого столбца и считают статистику  $\chi^2_r$ , которая имеет распределение  $\chi^2$ .

Группа наблюдения	На момент обследования		Через 1 месяц		Через 2 месяц	
	ЖЕЛ	ранг	ЖЕЛ	ранг	ЖЕЛ	ранг
1	2	1,5	2	1,5	3	3
2	2,1	1	3,1	2	4	3
3	2,4	2	2,1	1	3,5	3
4	2,5	1	3,5	2	4	3
5	2,3	1	3	2,5	3	2,5
6	2,6	1,5	2,6	1,5	4	3
Сумма рангов		8		10,5		17,5

$$\chi^2 = \frac{12}{nk(k+1)} \sum_{j=1}^k T_j^2 - 3n(k+1),$$

где  $n$  – число наблюдений;  $k$  – количество повторных измерений;  $T_j$  – сумма рангов для повторных измерений  $j$

$$\begin{aligned} \chi^2 &= \left[ \frac{12}{6 \times 3 \times (3 + 1)} \times (8^2 + 10,5^2 + 17,5^2) \right] - 3 \times 6 \times (3 + 1) \\ &= 80,8 - 72 = 8,08 \end{aligned}$$

**3.** Полученное значение сравнивается с критическим. При уровне значимости  $\alpha = 0,05$  и числе степеней свободы  $df = k - 1 = 3 - 1 = 2$ , где  $k$  – число групп;  $\chi^2_{кр} = 5,991$  Наша статистика больше, поэтому нулевая гипотеза отвергается.

$H_0$  и  $H_1$  - по аналогии с предыдущими тестами, о сходстве выборок.

# Friedman ANOVA

ЖЕЛ и отказ от курения.dat

File Edit Transform Plot Univariate Multivariate Model Diversity			
Show		Click mode	
<input type="checkbox"/> Row attributes		<input checked="" type="radio"/> Select	
<input type="checkbox"/> Column attributes		<input type="radio"/> Drag rows/columns	
		Edit	
		Cut	
		Copy	
	В момент обследования	через 1 месяц	через 2 месяца
1	• 2	2	3
2	• 2,1	3,1	4

ЖЕЛ и отказ от курения.dat

File Edit Transform Plot **Univariate** Multivariate Model Diversity Timeseries Geometry Stratigraphy Script Help

Show

☐ Row attributes ☒ Column attributes

В момент обследования

1	• 2
2	• 2,1
3	• 2,4
4	• 2,5
5	• 2,3
6	• 2,6
7	•
8	•
9	•
10	•
11	•
12	•
13	•

Summary statistics

One-sample tests (t, Wilcoxon, single-case)

Two-sample tests

**ANOVA etc. (several samples)**

Correlation

Intraclass correlation

Normality tests

Contingency table (chi^2 etc.)

Mantel-Cochran-Haenszel test

Risk/odds

Single proportion test

Multiple proportion CIs

Ratios of counts CI

Survival analysis

Combine errors

View

☐ Bands ☐ Binary Recover windows

Decimals: -

Several-sample tests (ANOVA, Kruskal-Wallis)

**Several-sample repeated measures tests**

Two-way ANOVA

Two-way ANOVA without replication

Two-way repeated measures ANOVA

One-way ANCOVA

Several-sample repeated measures tests

Repeated-measures ANOVA	Tukey's pairwise	Friedman test	Wilcoxon pairwise
-------------------------	------------------	---------------	-------------------

Test for equal medians

chi2:	8.0833	Degrees of freedom:	2
chi2, tie corrected:	9.2381		
chi2, continuity corrected:	8.8031		
p (same), asymptotic:	0.012258		
p (same), exact:	0.0062241		

Отвергаем  $H_0$  –  
ЖЭЛ изменилась

$$\chi^2_{(2)} = 9,24; P < 0,012$$

Several-sample repeated measures tests

Repeated-measures ANOVA	Tukey's pairwise	Friedman test	Wilcoxon pairwise
-------------------------	------------------	---------------	-------------------

Raw p values, uncorrected significance

	В момент обседа	через 1 месяц	через 2 месяца
В момент обседа		0.2941	0.04615
через 1 месяц	0.2941		0.09091
через 2 месяца	0.04615	0.09091	

*Далее попарные сравнения групп методом Вилкоксона*

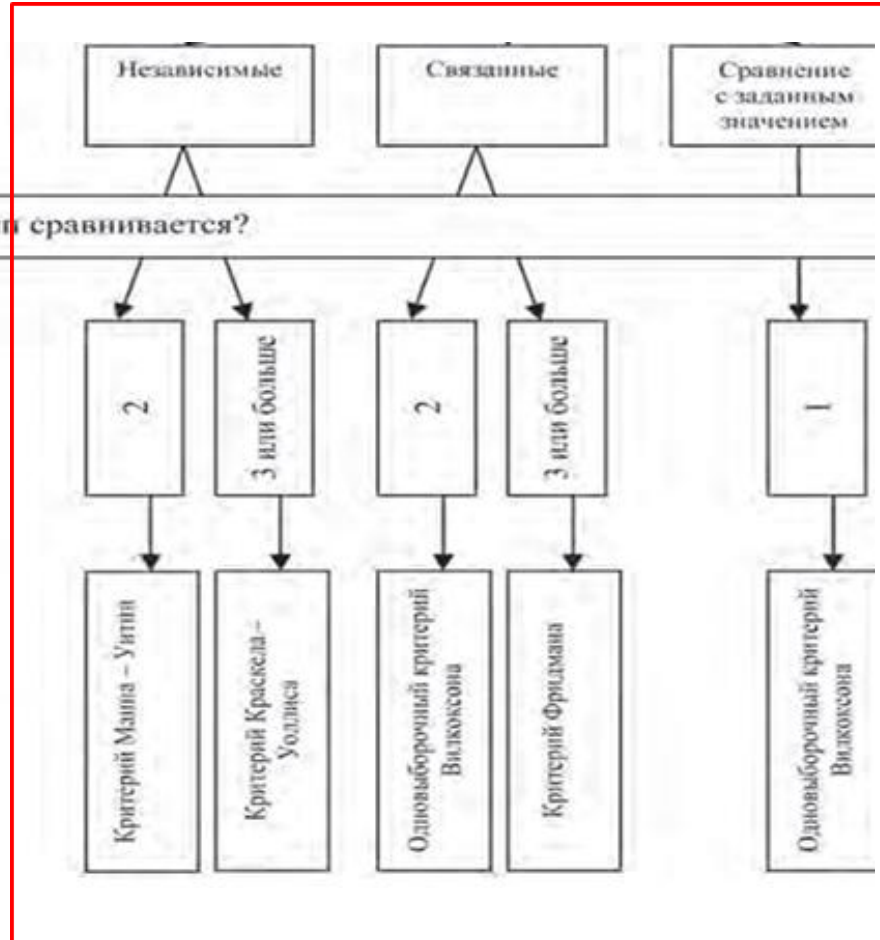
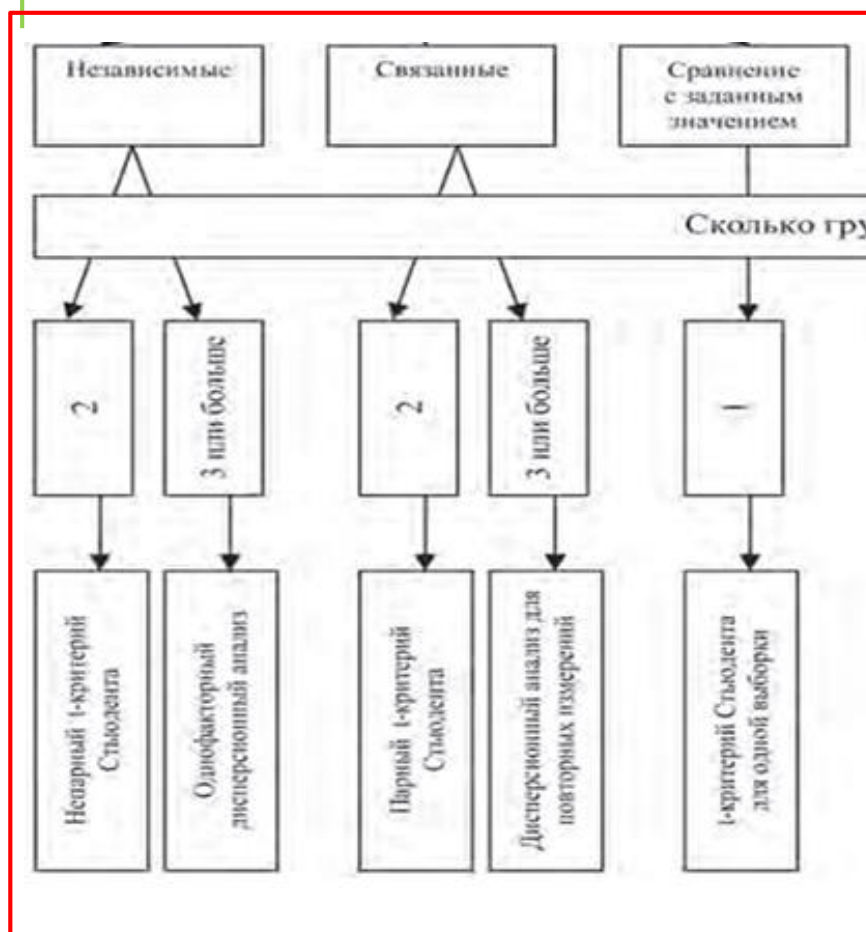


## В итоге, при выборе теста важно, что:

1. Параметрические тесты более мощные, чем непараметрические;
2. Непараметрические безопаснее в плане ошибки 1-го рода;
3. Чем больше размер выборки, тем менее критичны требования к распределению (по Центральной предельной теореме); для выборок  $N \geq 100$  используют параметрические тесты даже при больших отклонениях от нормального распределения (кроме регрессий).
4. АНОВА не очень чувствительна к отклонениям от нормального распределения (для одинаковых по размеру групп).

# Алгоритм выбора статистического критерия для сравнения *количественных данных*





# Точечная и интервальная оценка

**Термины** точечная оценка (point estimation) и интервальная оценка (interval estimation) впервые использовал в **1943** году американский математик **Генри Шеффе** в статье «*Statistical Inference in the NonParametric Case*»

**Точечная оценка** определяется одним числом (среднее значение, стандартное отклонение и т.д.)

**Интервальная оценка** определяется двумя числами – концами интервала (доверительный интервал)



**Генри Шеффе**  
(Henry Scheffé 1907-1977)

**Доверительным интервалом (ДИ)** называется интервал, в который попадают измеренные в эксперименте значения, соответствующие **доверительной вероятности**

**Метод доверительных интервалов** разработал американский статистик **Ежи Нейман** (термин доверительный введен им в **1934** году), исходя из идей Рональда Фишера



**Ежи Нейман**  
(Jerzy Neyman 1894-1981)

**ДИ**, используемый статистике при **интервальной оценке** статистических параметров, более предпочтительной при небольшом объёме выборки, чем **точечная оценка**

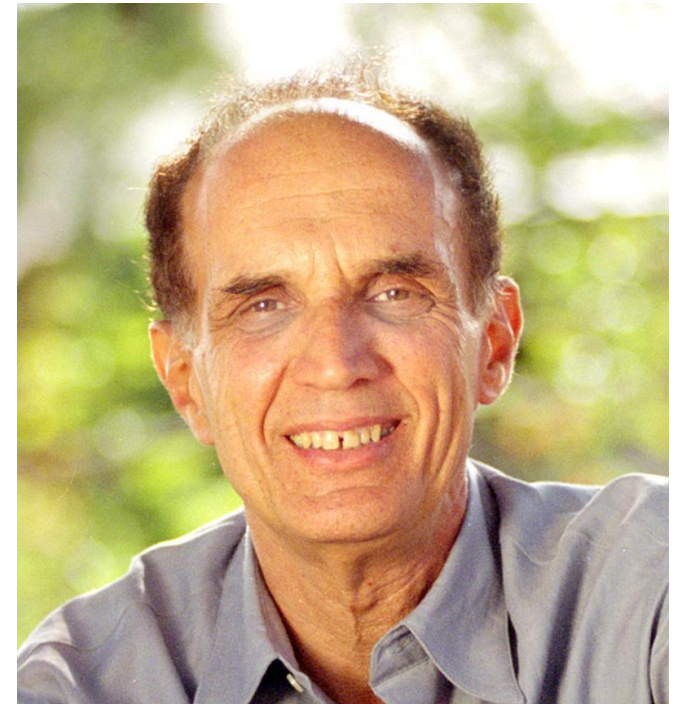
## ***Бутстреп (bootstrap, bootstrapping)***

*Доверительный интервал* можно построить с использованием *ресэмплинг-техник (resampling)*: методом складного ножа или более современным методом бутстрепа.

***Бутстреп*** — это современная ресэмплинг-техника, то есть техника, основанная на взятии повторных (re...) выборок (...sample).

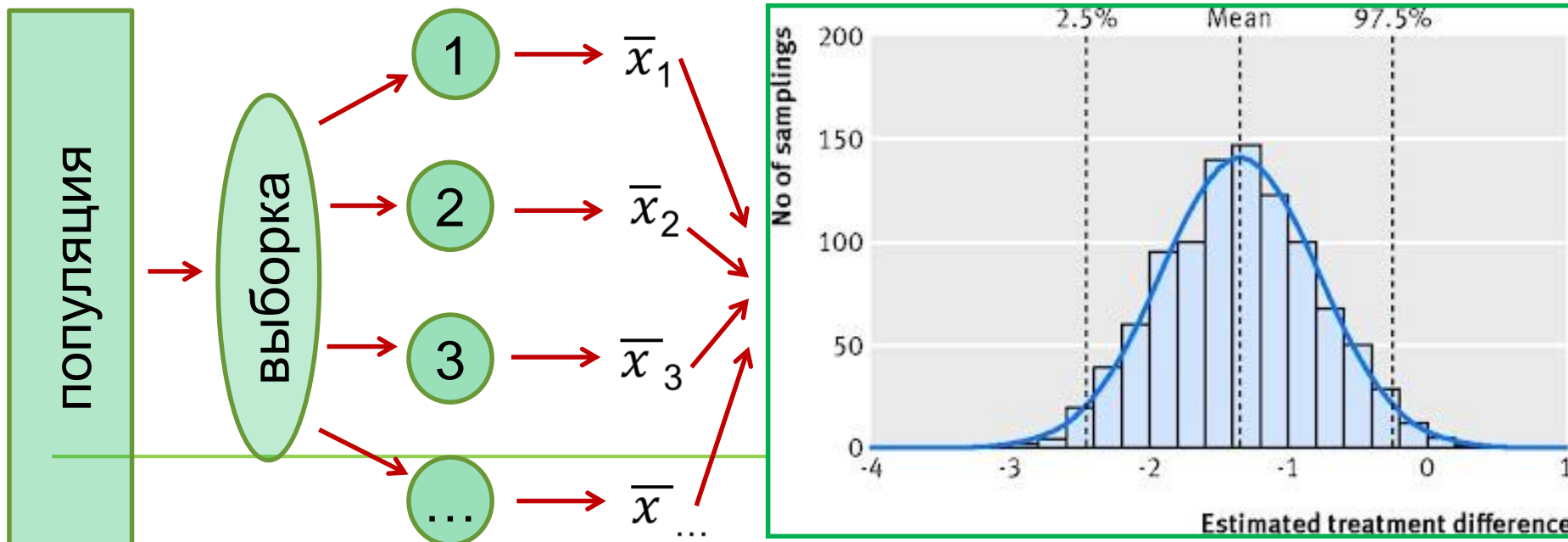
Понятие введено в **1977** году **Брэдли Эфроном** (первая публикация относится к 1979 году).

Данный метод является методом ***непараметрической статистики.***



**Брэдли Эфрон**  
(Bradley Efron p.1938)

**Например - выборка** 25 студентов, измерили рост. Выписываем значения на карточки, перемешиваем и случайным образом вытягиваем одну. Записываем результат и возвращаем карточку обратно. Повторяем процедуру ещё 24 раза. Получаем *сгенерированную выборку*. Рассчитываем интересующую нас статистику. В нашем случае – среднее значение. Повторяем множество раз (1 тыс. и более). Получаем *множество средних значений*. Строим распределение средних и с концов отрезаем по 2,5% площади. Остаётся 95% ДИ для среднего, вычисленный с помощью процедуры бутстрепа **процентильным методом**.





Один из лучших бутстрепа — **метод ВСа** (**Bias Corrected accelerated** — *ускоренный бутстреп с поправкой на смещение*).

***Наиболее известные алгоритмы ресэмплинг-техник :***

- *Перестановочный тест (permutation)*
- *Бутстреп (bootstrap)*
- *Метод «складного ножа» (jackknife)*
- *Кросс-проверка (cross-validation)*

Слово происходит от выражения: «To pull oneself over a fence by one's bootstraps.» (дословно — «перебраться через ограду, потянув за ремешки на ботинках»)



Название «складной нож», потому что его действие напоминает складной нож — простой инструмент, которым можно решить множество различных проблем



## ***КОРРЕЛЯЦИИ (correlation)***

До сих пор нас в выборках интересовала только **одна зависимая переменная**.

Мы изучали, отличается ли распределение этой переменной в одних условиях от распределения той же переменной в других условиях.

Обратимся к ситуации, когда зависимых переменных будет **ДВЕ** и более.

Нас интересует вопрос, в какой степени эти переменные связаны между собой, совместно изменяются.

Это могут быть измерения одной особи или связанных пар.

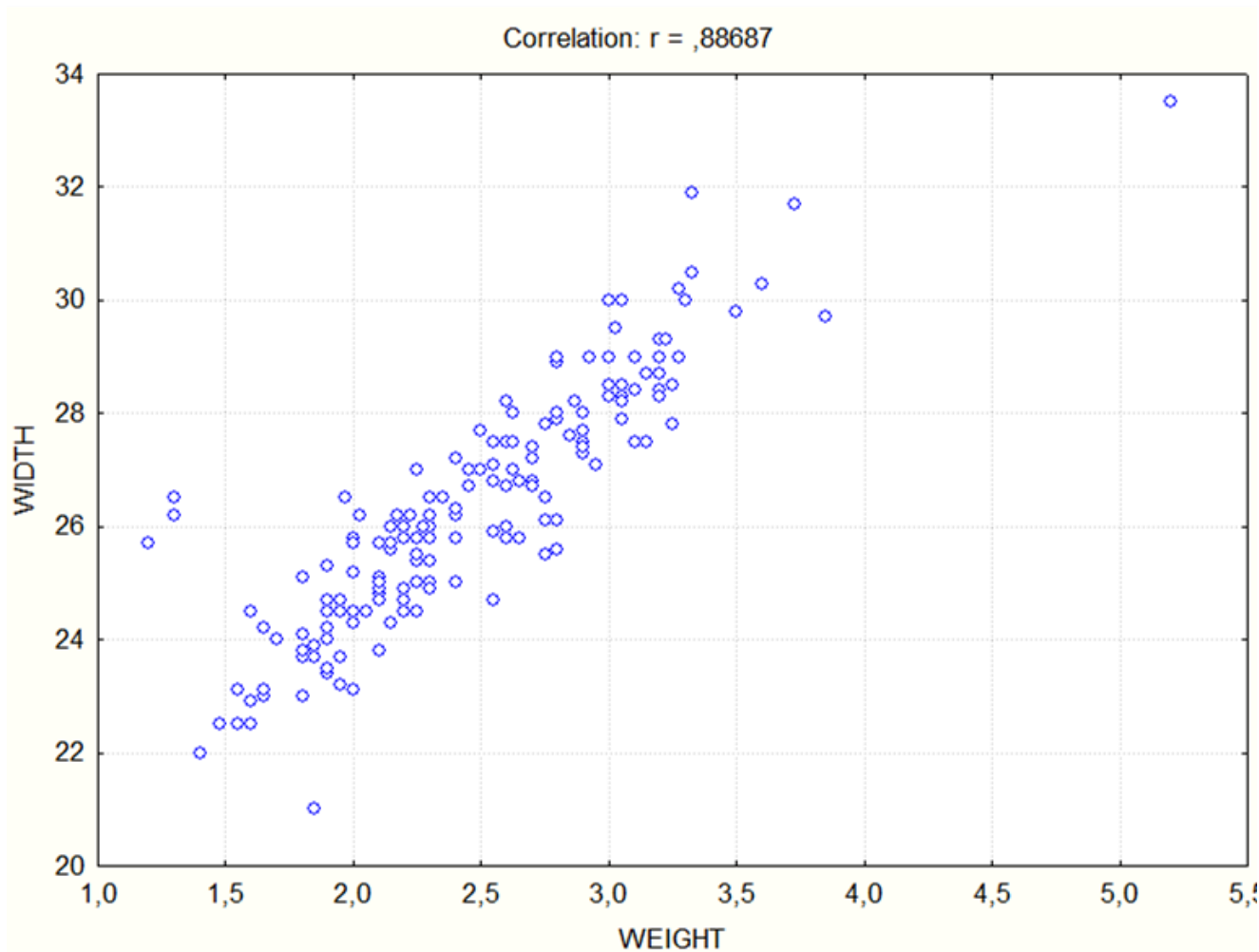
## Коэффициент корреляции

Для оценки тесноты (силы) связи служит коэффициент **корреляции**.

1. Может принимать значения от -1 до +1
2. Знак коэффициента показывает *направление связи* (прямая или обратная)
3. Абсолютная величина показывает *силу* связи
4. всегда основан на парах чисел (измерений 2-х переменных от одной особи или 2-х переменных от разных, но связанных особей)

# Скаттерплот

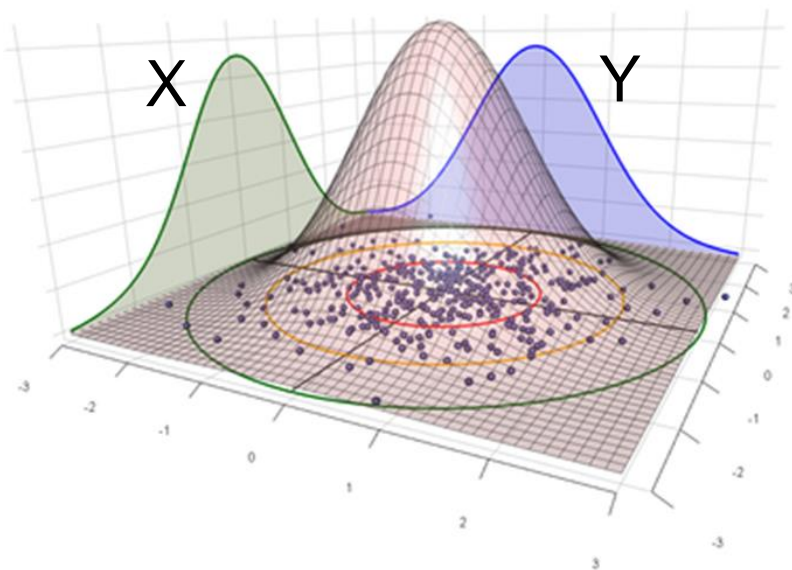
(= диаграмма рассеяния; scatterplot, scatter diagram)



**Две** характеристики: – наклон (направление связи) и ширина (сила связи) воображаемого эллипса

Оценкой *линейной зависимости (связи)* между 2-мя непрерывными переменными, служит *коэффициент корреляции Пирсона* (Подробнее о параметрических методах оценки связи см. лекцию – «Корреляционный и регрессионный анализ»).

Однако существуют **ограничения** для его **применения**.



Значения Y и X должны быть распределены нормально - *двумерное нормальное распределение* (bivariate normal distribution)

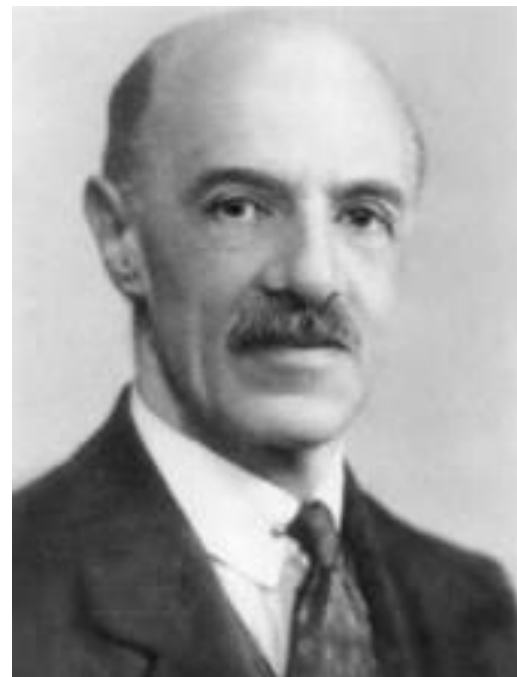
Если хотя бы одна переменная имеет не нормальное распределение или она порядковая, то для оценки зависимости используют **непараметрические коэффициенты корреляции**.

## Коэффициент корреляции **Спирмана** (*Spearman rank order correlation*)

Коэффициент ранговой корреляции Спирмена - непараметрический аналог коэффициента корреляции Пирсона.

Определяется не по величинам переменных признаков, а по рангам - номерам в порядке возрастания величин признаков.

Критерий **разработан и предложен** для проведения корреляционного анализа в **1904 году** Чарльзом Эдвардом Спирменом, английским психологом, профессором Лондонского и Честерфилдского университетов.

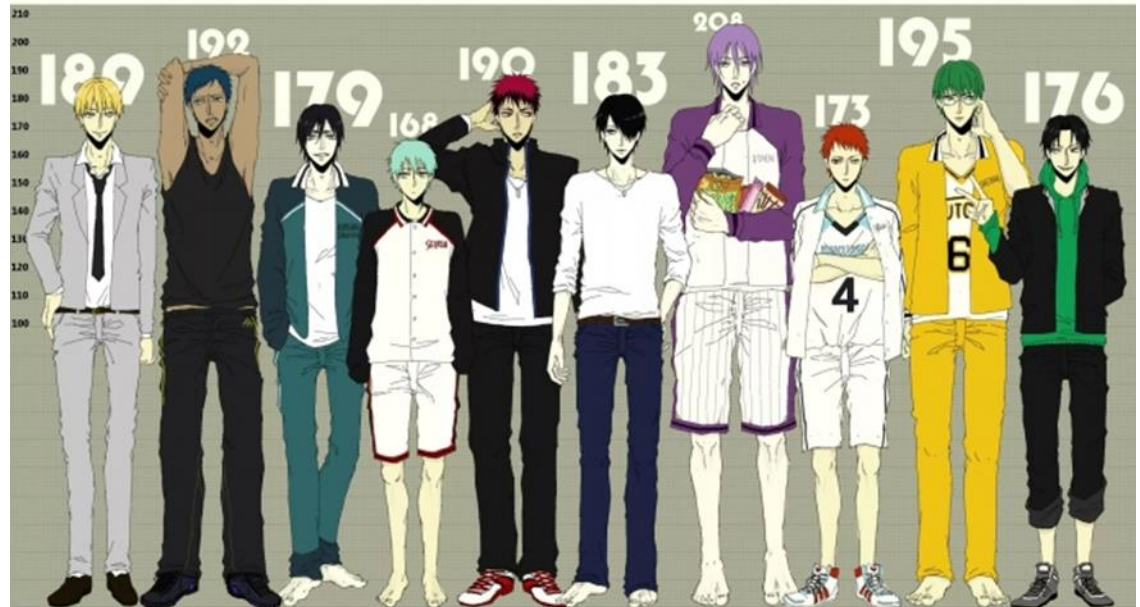


**Чарльз Эдвард Спирмен**  
(*Charles Edward Spearman*  
1863-1945)

Необходимо узнать существует ли зависимость роста сына от роста отца?

**Переменные:**

1. рост сына (Y);
2. рост отца (X)



1. Ранжируем данные для **каждой переменной** от меньшего к большему;
2. Если встретились одинаковые значения (***tied ranks***), присваиваем им средние ранги;
3. Считаем разности рангов в каждой паре данных;

Рост отца		Рост сына		$d_i$	$d_i^2$
Значение, см	ранг	Значение, см	ранг		
167	1	169	2	-1	1
169	2	171	3	-1	1
170	1,5	166	1	2,5	6,25
170	1,5	172	4	-0,5	0,25
172	5	180	7	-2	4
173	6	176	5	1	1
174	7	177	6	1	1
175	8	182	8,5	-0,5	0,25
179	9	182	8,5	0,5	0,25
180	10	186	10	0	0
					$\sum d_i^2 = 15$

#### 4. Считаем коэффициент $r_s$

$$r_s = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

разности рангов

число строк  
(размер выборки)

$$\underline{r_s} = 1 - \frac{6 \times 15}{10^3 - 10} \approx 0,9091$$

При  $\alpha=0,001$  значение  $r_{кр}=0,903$

$r_s > r_{кр}$ , таким образом нулевая гипотеза должна быть отвергнута и наблюдается зависимость роста взрослого сына от роста его отца



$$H_0 : \rho_s = 0$$

$$H_1 : \rho_s \neq 0$$



Статистика критерия – сам коэффициент корреляции Спирмана (имеет t-распределение)

**Коэффициент Спирмана – аналог коэффициента корреляции Пирсона**, стремится к нему в больших выборках. Мощность – около 91% коэффициента Пирсона.

Минимальный объем выборки: Для расчета требуется не менее 5 наблюдений по каждой переменной, лучше  $\geq 10$ .

# Spearman Rank Order Correlations

Зависимость роста сына от роста отца.dat

File Edit Transform Plot **Univariate** Multivariate Model Diversity Timeseries

Show

☐ Row attributes

☐ Column attributes

рост отца

1	• 167	1
2	• 169	1
3	• 170	1
4	• 170	1
5	• 172	1
6	• 173	1
7	• 174	1
8	• 175	1
9	• 179	1
10	• 180	1
11	•	
12	•	
13	•	

Summary statistics

One-sample tests (t, Wilcoxon, single-case)

Two-sample tests

ANOVA etc. (several samples)

**Correlation**

Intraclass correlation

Normality tests

Contingency table (chi<sup>2</sup> etc.)

Mantel-Cochran-Haenszel test

Risk/odds

Single proportion test

Multiple proportion CIs

Ratios of counts CI

Survival analysis

Combine errors

Зависимость роста сына от роста отца.

File Edit Transform Plot **Univariate**

Show

☐ Row attributes

☐ Column attributes

Click mode

☒ Select


☐ Drag rows/c

	рост отца	рост сына
1	• 167	169
2	• 169	171
3	• 170	166
4	• 170	172
5	• 172	180
6	• 173	176
7	• 174	177
8	• 175	182
9	• 179	182
10	• 180	186

# Spearman Rank Order Correlations

Correlation

	рост отца	рост сына
рост отца		0,00027383
рост сына	0,90854	



Correlation statistic

- ☐ Linear r (Pearson)
- ☐ Spearman's D
- ☒ Spearman's rs
- ☐ Kendall's tau
- ☐ Polyserial rho
- ☐ Partial linear

Table format

- ☒ Statistic \ p(uncorr)
- ☐ Statistic
- ☐ p(uncorr)
- ☐ Permutation p

☐ Bonferroni correction

Отвергаем  $H_0$ :

Оказалось, что рост сына положительно связан с ростом его отца.

## Коэффициент корреляции **Кендалла** (*Kendall's coefficient of rank correlation, Kendall- $\tau$* )

В **1937 году** вышла статья **Мориса Кендалла** где описан новый коэффициент ранговой корреляции -  $\tau$  («тау»).

Оценивает разность между вероятностью того, что порядок данных в обеих переменных одинаков, и вероятностью того, что порядки разные.

Только для **ранговых** переменных! Для количественных лучше коэффициент Спирмана, особенно для больших выборок.



**Морис Кендалл**  
(*Sir Maurice George Kendall*  
1907-1983)

## **Пример.**

Обследовано 20 больных серповидноклеточной анемией.

*Оценены:*

тяжесть (в баллах) и коэффициент адгезии эритроцитов.

*Необходимо ответить  
на вопрос:*

Связана ли адгезивность  
эритроцитов и тяжестью  
серповидноклеточной  
анемии?



## Kendall's coefficient of rank correlation, Kendall- $\tau$

The screenshot shows the SPSS Correlation dialog box and the Univariate menu. The Correlation dialog box has two tabs: 'Table' and 'Plot'. The 'Table' tab is active, showing a table with two rows and two columns. The first row is 'Тяжесть забс' and the second row is 'Кoeffициен'. The values are 6,0131E-06 and 0,7342 respectively. The 'Plot' tab is also visible. The Univariate menu is open, showing various statistical tests. The 'Correlation' option is highlighted. The 'Table format' section shows 'Statistic \ p(uncorr)' selected. The 'Table' tab in the Correlation dialog box shows a table with two rows and two columns. The first row is 'Тяжесть забс' and the second row is 'Кoeffициен'. The values are 6,0131E-06 and 0,7342 respectively. The 'Plot' tab is also visible. The Univariate menu is open, showing various statistical tests. The 'Correlation' option is highlighted. The 'Table format' section shows 'Statistic \ p(uncorr)' selected.

Table	Plot
Тяжесть забс	Кoeffициен
Тяжесть забс	6,0131E-06
Кoeffициен	0,7342

серповидноклеточная анемия.dat

File Edit Transform Plot Univariate Multivariate Model Diversity Tir

Show

☐ Row attributes

☐ Column attributes

Тяжесть забол

1 0

2 0

3 1

4 1

5 1

6 1

7 1

8 1

9 2

10 2

11 3

12 3

13 3

6,3

Summary statistics

One-sample tests (t, Wilcoxon, single-case)

Two-sample tests

ANOVA etc. (several samples)

Correlation

Intraclass correlation

Normality tests

Contingency table (chi^2 etc.)

Mantel-Cochran-Haenszel test

Risk/odds

Single proportion test

Multiple proportion CIs

Ratios of counts CI

Survival analysis

Combine errors

Correlation statistic

☐ Linear r (Pearson)

☐ Spearman's D

☐ Spearman's rs

☒ Kendall's tau

☐ Polyserial rho

☐ Partial linear

Table format

☒ Statistic \ p(uncorr)

☐ Statistic

☐ p(uncorr)

☐ Permutation p

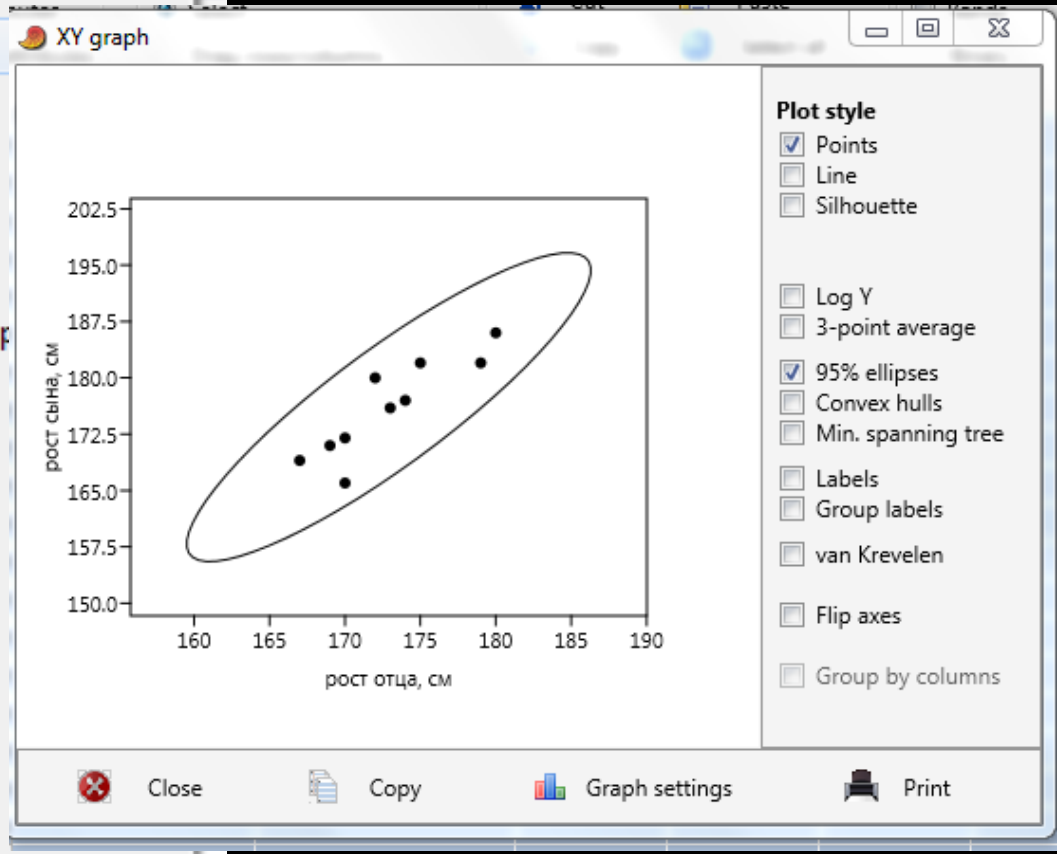
Отвергаем  $H_0$ :  
аггезивность эритроцитов  
положительно связана  
тяжестью  
серповидноклеточной  
анемией

## Kendall's coefficient of rank correlation, Kendall- $\tau$

графика — диаграммы  
рассеяния (scattergram)

Зависимость роста сына от роста отца.dat

	рост отца
1	167
2	169
3	170
4	170
5	172
6	173
7	174
8	175
9	179
10	180
11	
12	
13	
14	
15	
16	



## **Корреляционный анализ: (линейные) соотношения между двумя непрерывными переменными**

✓ **Коэффициент корреляции Пирсона  $r$** , который используется для выявления взаимосвязи между двумя приблизительно нормально распределенными непрерывными переменными. В действительности переменные должны удовлетворять совместно «двумерному нормальному распределению».

✓ **Коэффициент ранговой корреляции Спирмена,  $\rho$  ( $\rho$ )**, применяемый для выявления взаимосвязи между двумя непрерывными переменными, по крайней мере одна из которых распределена не по нормальному закону.

✓ **Коэффициент ранговой корреляции Кендалла,  $\tau$  ( $\tau$ )** применяемый для выявления взаимосвязи между двумя порядковыми переменными или между одной порядковой и одной непрерывной.



Спасибо за внимание!

