

Теория вероятностей и статистика в Машинном Обучении



academy

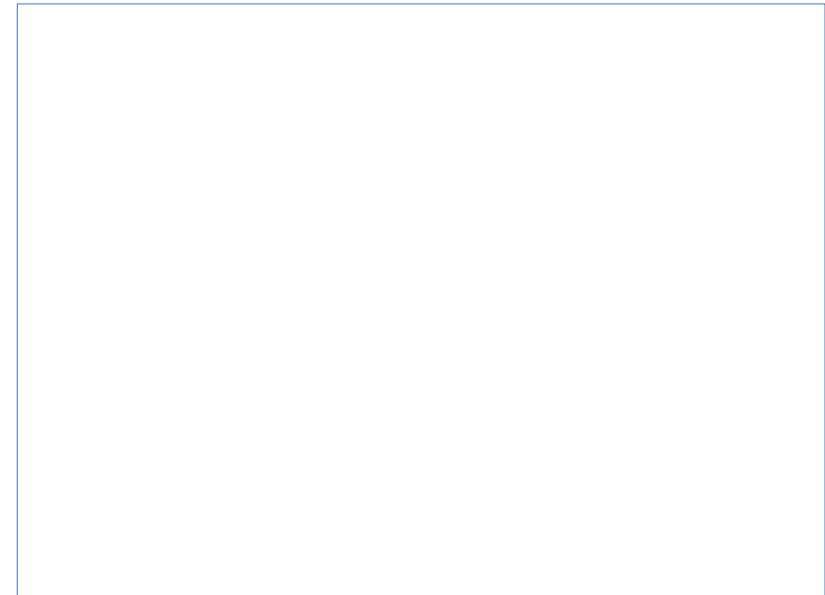


- МОМЕНТЫ И КРИТИЧЕСКИЕ ГРАНИЦЫ СЛУЧАЙНОЙ ВЕЛИЧИНЫ

- Меры связи случайных величин

1. Коэффициент ковариации
2. Коэффициент корреляции Пирсона
3. Коэффициент корреляции Спирмена

- ЦЕНТРАЛЬНАЯ ПРЕДЕЛЬНАЯ ТЕОРЕМА

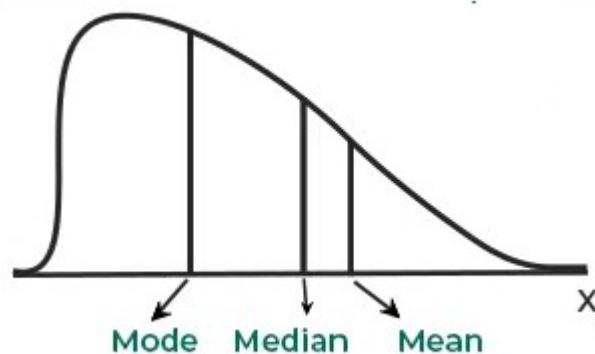




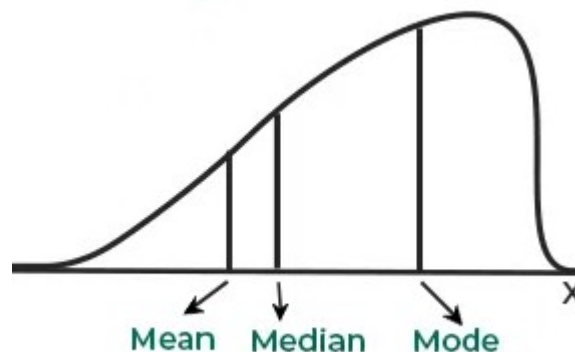
МОМЕНТЫ И КРИТИЧЕСКИЕ ГРАНИЦЫ СЛУЧАЙНОЙ ВЕЛИЧИНЫ

- Математическое ожидание
- Дисперсия
- Асимметрия
- Эксцесс
- Квантили

Медина и мода случайной величины



Ассиметрия > 0



Ассиметрия < 0



Медина и мода случайной величины

В качестве показателя центра группирования значений случайной величины, наряду с математическим ожиданием, используются также **медиана** и **мода**.

Медианой x_{med} случайной величины X называется ее квантиль уровня 0,5 (или, что то же самое, ее 50%-ная точка):

$$x_{\text{med}} = x_{0,5} = \omega_{0,5}.$$

В качестве медианы дискретной случайной величины (поскольку квантиль дискретной случайной величины может быть определена не однозначно) обычно берут значение x_{med} , получаемое линейной аппроксимацией:

$$x_{\text{med}} = x_l + \frac{x_{l+1} - x_l}{p_{l+1}} \left(\frac{1}{2} - \sum_{i=1}^l p_i \right),$$

где **медианный отрезок** $[x_l; x_{l+1}]$ определяется из условий

$$\sum_{i=1}^l p_i \leq 0,5; \quad \sum_{i=1}^{l+1} p_i > 0,5.$$

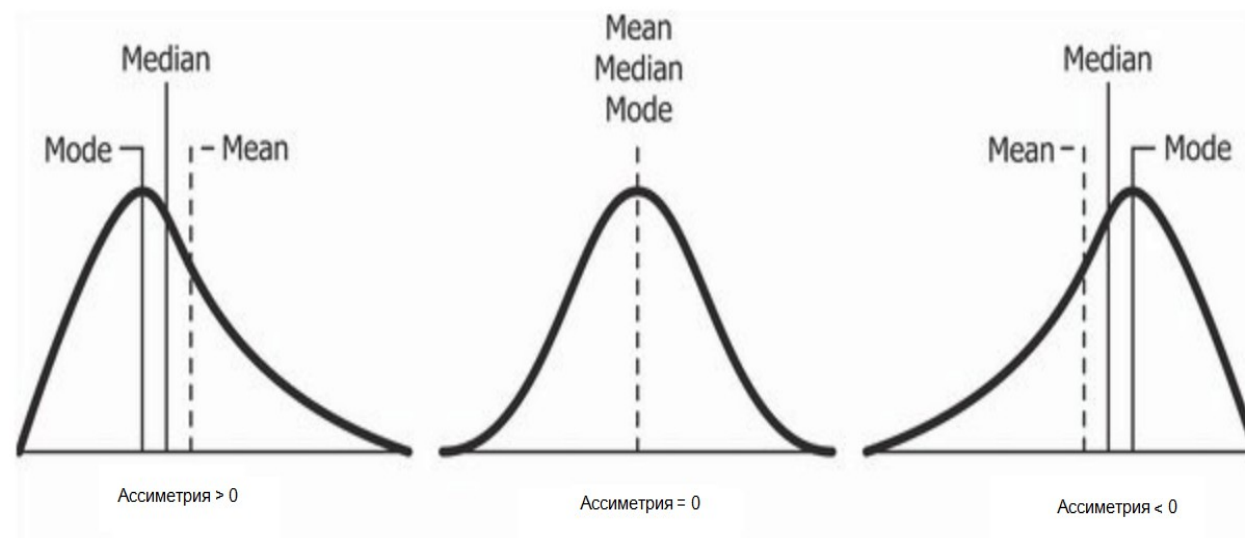
Модой абсолютно непрерывной случайной величины X называется точка локального максимума плотности распределения:

$$f_X(x_{\text{mod}}) = \max_{x \in \mathcal{Y}} f_X(x).$$

Модой дискретной случайной величины X называется значение этой случайной величины, соответствующее наибольшей вероятности:

$$x_{\text{mod}} = x_j, \text{ такое что } p_i = \max_j p_j.$$

Случайные величины, имеющие одну моду, называются **унимодальными**.



Бимодальное



Мультимодальное

The screenshot shows the Microsoft Excel interface. In the top-left corner, the formula bar displays the function `=НОРМОБР(1-0.95;1000;1000)`. Below it, the active worksheet is named "Математика.xlsx". The grid shows columns D, E, and F. Cell E2 contains the numerical value **-644.854**, which is highlighted by a black rectangular selection box.

МЕРЫ СВЯЗИ СЛУЧАЙНЫХ ВЕЛИЧИН



Для измерения «степени зависимости» случайных величин **ковариация** случайных величин X и Y :

$$\text{cov}(X, Y) = E\{[X - E(X)][Y - E(Y)]\}.$$

$$\text{cov}(X, Y) = E(XY) - E(X)E(Y).$$

КОММУТАТИВНОСТЬ КОВАРИАЦИИ. Для любых случайных величин X и Y

$$\text{cov}(X, Y) = \text{cov}(Y, X).$$

СВЯЗЬ КОВАРИАЦИИ И ДИСПЕРСИИ. Для любой случайной величины X

$$\text{cov}(X, X) = \text{Var}(X).$$

ЛИНЕЙНОСТЬ КОВАРИАЦИИ ПО КАЖДОМУ АРГУМЕНТУ. Для любых случайных величин X и Y и любых чисел $a, b \in \mathbb{R}$

$$\text{cov}(aX, bY) = ab \text{cov}(X, Y).$$

АДДИТИВНОСТЬ КОВАРИАЦИИ ПО КАЖДОМУ АРГУМЕНТУ. Для любых случайных величин X, Y и Z

$$\begin{aligned}\text{cov}(X + Y, Z) &= \text{cov}(X, Z) + \text{cov}(Y, Z), \\ \text{cov}(X, Y + Z) &= \text{cov}(X, Y) + \text{cov}(X, Z).\end{aligned}$$

ФОРМУЛА ДЛЯ КОВАРИАЦИИ НЕЗАВИСИМЫХ СЛУЧАЙНЫХ ВЕЛИЧИН. Для независимых случайных величин X и Y

$$\text{cov}(X, Y) = 0.$$

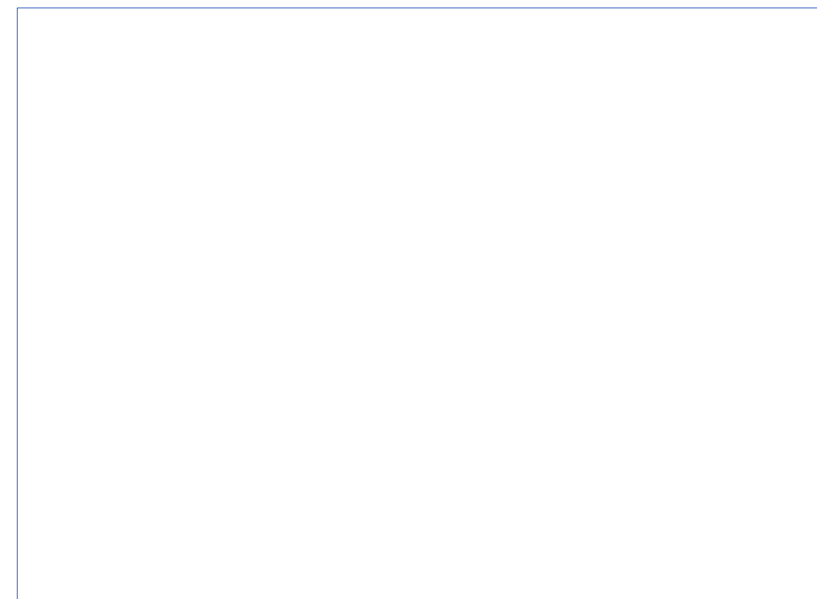
ФОРМУЛА ДЛЯ ДИСПЕРСИИ СУММЫ ПРОИЗВОЛЬНЫХ СЛУЧАЙНЫХ ВЕЛИЧИН. Дисперсия суммы произвольных (зависимых или независимых) случайных величин X и Y рассчитывается по формуле

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{cov}(X, Y).$$

две случайные величины X и Y называются *независимыми*, если для всех $x, y \in \mathbb{R}$

$$P\{(X \leq x) \cap (Y \leq y)\} = P\{X \leq x\}P\{Y \leq y\},$$

т. е. если для всех $x, y \in \mathbb{R}$ события $\{X \leq x\}$ и $\{Y \leq y\}$ независимы.



Коэффициент корреляции



Коэффициент корреляции случайных величин X и Y

$$R(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{E\{[X - E(X)][Y - E(Y)]\}}{\sigma_X \sigma_Y} = \frac{E(XY) - E(X)E(Y)}{\sigma_X \sigma_Y}.$$

ИНВАРИАНТНОСТЬ КОЭФФИЦИЕНТА КОРРЕЛЯЦИИ ОТНОСИТЕЛЬНО ЛИНЕЙНЫХ ПРЕОБРАЗОВАНИЙ АРГУМЕНТОВ. Для любых случайных величин X и Y и любых чисел $a, b, c, d \in \mathbb{R}$

$$R(aX + b, cY + d) = R(X, Y).$$

НОРМИРОВАННОСТЬ КОЭФФИЦИЕНТА КОРРЕЛЯЦИИ. Для любых случайных величин X и Y

$$-1 \leq R(X, Y) \leq 1.$$

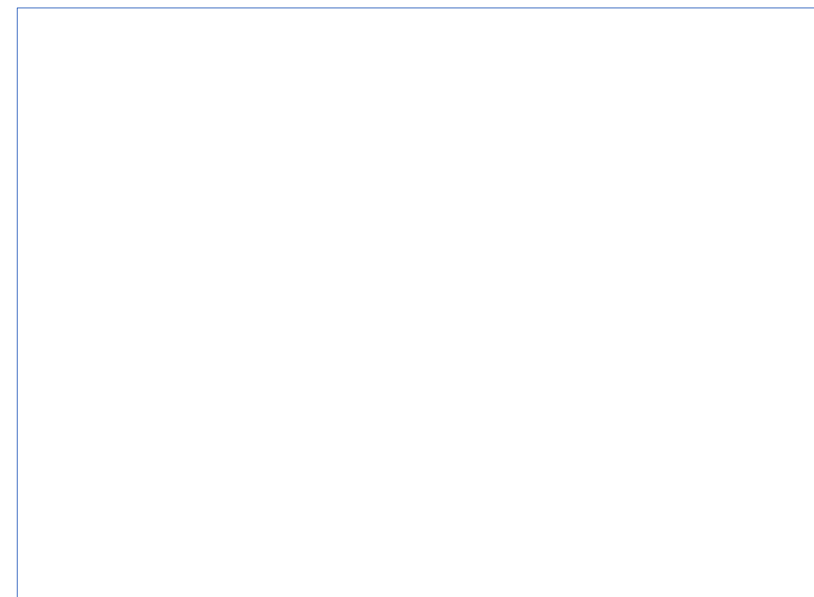
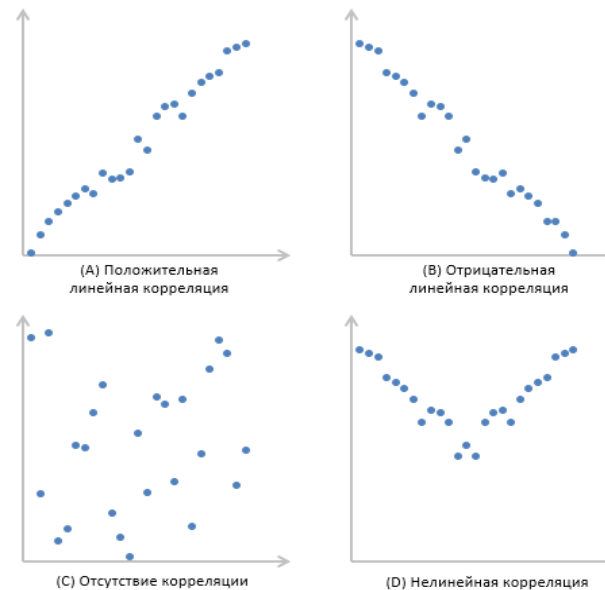
ФОРМУЛА ДЛЯ КОЭФФИЦИЕНТА КОРРЕЛЯЦИИ НЕЗАВИСИМЫХ СЛУЧАЙНЫХ ВЕЛИЧИН. Для независимых случайных величин X и Y

$$R(X, Y) = 0;$$

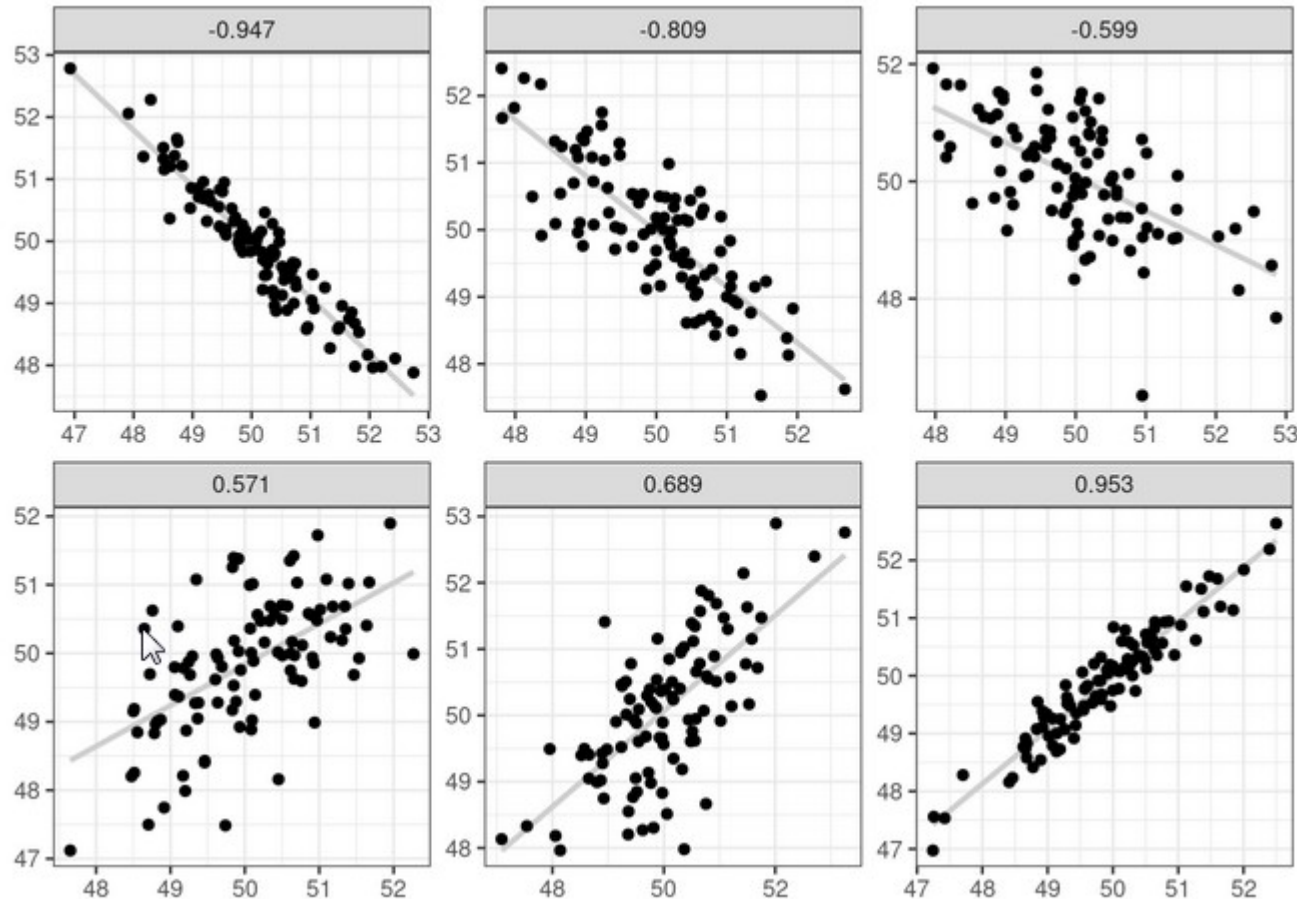
ФОРМУЛА ДЛЯ КОЭФФИЦИЕНТА КОРРЕЛЯЦИИ ЛИНЕЙНО СВЯЗАННЫХ СЛУЧАЙНЫХ ВЕЛИЧИН. Для линейно связанных случайных величин X и $Y = aX + b$ ($a, b \in \mathbb{R}, a \neq 0$)

$$|R(X, Y)| = 1, \quad (\star)$$

если для каких-либо случайных величин X и Y выполняется (\star) , то существуют такие числа $a, b \in \mathbb{R}, a \neq 0$, что $\mathbf{P}\{Y = aX + b\} = 1$.



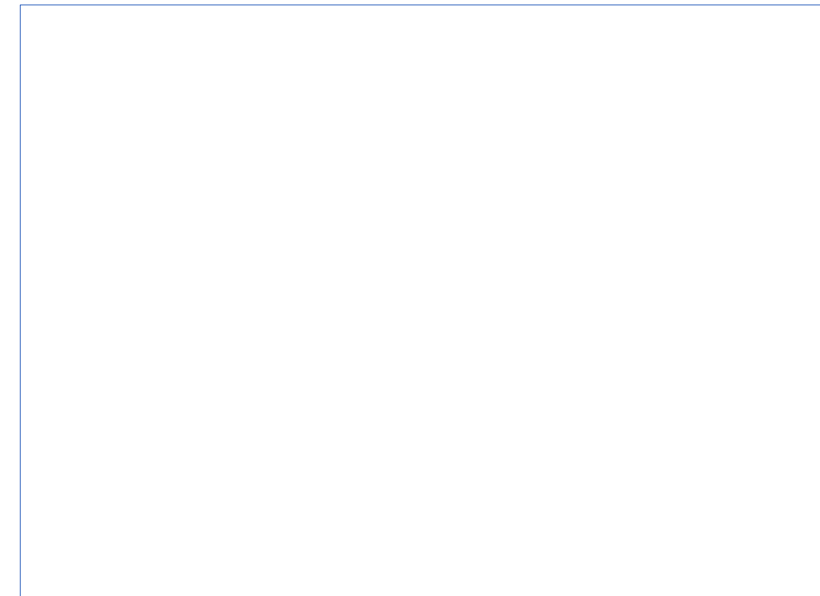
Примеры зависимых случайных величин



Интерпретация коэффициента корреляции производится исходя из уровня силы связи:

$[0,2; 0,3)$ – **слабая** положительная связь,
 $[0,30; 0,7)$ – **умеренная** положительная связь,
 $[0,70; 1,0]$ – **сильная** положительная связь,

$(-0,3; -0,2]$ – **слабая** отрицательная связь,
 $(-0,7; -0,3]$ – **умеренная** отрицательная связь,
 $[-1,0; -0,7]$ – **сильная** отрицательная связь.





Вычисление коэффициента корреляции (Пирсона). Применение

$$r = \frac{n \sum xy - (\sum x)(\sum y)}{\sqrt{[n \sum x^2 - (\sum x)^2][n \sum y^2 - (\sum y)^2]}}$$

Excel: **КОРРЕЛ(массив1;массив2)**

$$\text{Коэффициент детерминации} = R^2 = 1 - \frac{D[y|x]}{D[y]} = 1 - \frac{\sigma^2}{\sigma_y^2}$$

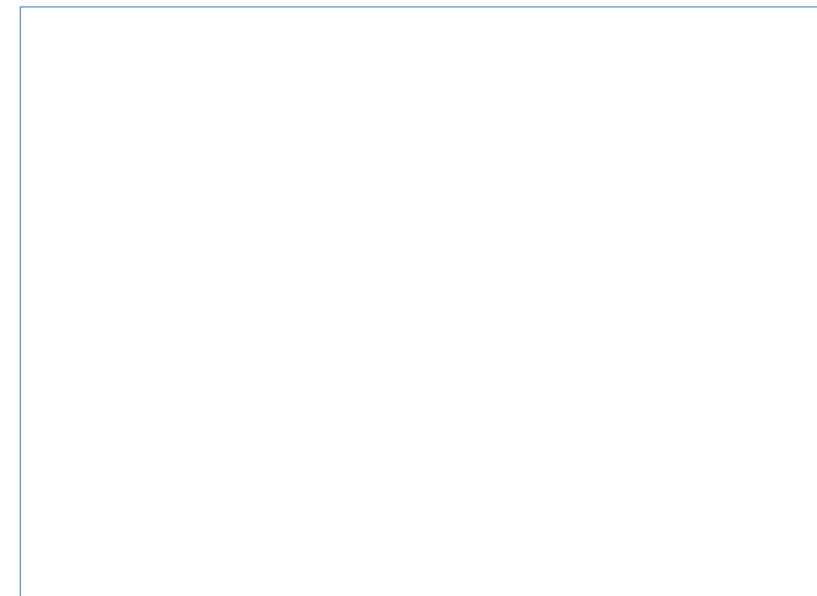
обычно применяется при оценке качества **регрессионных моделей**

Доля дисперсии зависимой переменной, объясняемая рассматриваемой моделью зависимости (объясняющими переменными), или
1- доля необъяснённой дисперсии

1. Коэффициент корреляции Пирсона чувствителен к выбросам. Одно аномальное значение может существенно исказить коэффициент. Поэтому перед проведением анализа следует проверить и при необходимости удалить выбросы. Другой вариант – перейти к ранговому коэффициенту корреляции Спирмена. Рассчитывается также, только не по исходным значениям, а по их рангам

2. Синоним корреляции – это взаимосвязь или совместная вариация. Поэтому наличие корреляции ($r \neq 0$) еще не означает причинно-следственную связь между переменными. Вполне возможно, что совместная вариация обусловлена влиянием третьей переменной. Совместное изменение переменных без причинно-следственной связи называется **ложная корреляция**.

3. Отсутствие линейной корреляции ($r = 0$) не означает отсутствие взаимосвязи. Она может быть нелинейной. Частично эту проблему решает ранговая корреляция Спирмена, которая показывает совместный рост или снижение рангов, независимо от формы взаимосвязи.



Ранговый коэффициент корреляции



Коэффициент корреляции Спирмена:

$$r_s = \rho_{rg_X, rg_Y} = \frac{\text{cov}(rg_X, rg_Y)}{\sigma_{rg_X} \sigma_{rg_Y}},$$

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

Расчет коэффициента Спирмена состоит из следующих этапов:

1) Ранжирование признаков по возрастанию.

Ранг – это порядковый номер. Если встречаются два одинаковых значения, им присваивают одинаковое значение ранга.

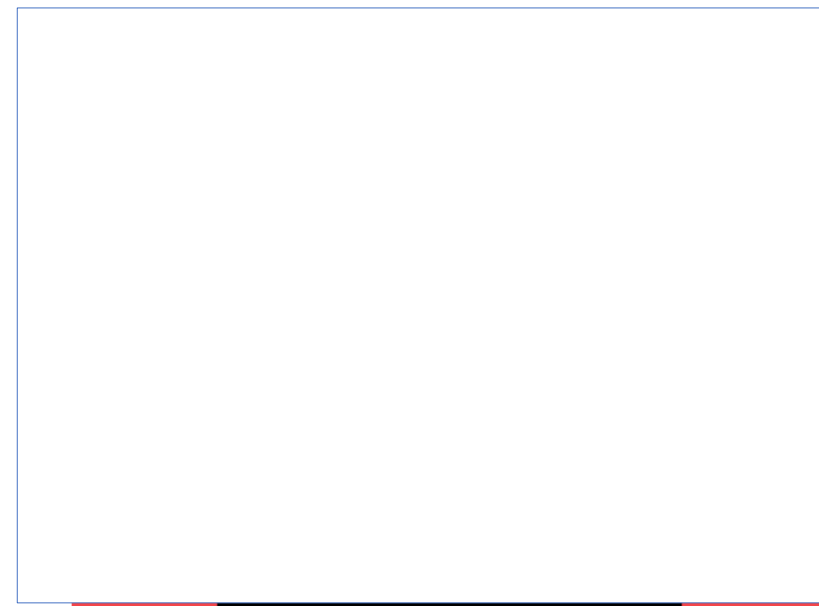
2) Определение разности рангов каждой пары сопоставляемых значений, $d = dx - dy$.

3) Возведение в квадрат разность d_i и нахождение общей суммы, $\sum d^2$.

4) Вычисление коэффициента корреляции рангов по формуле (*)

Шкала Чеддока

Количественная мера тесноты связи	Качественная характеристика силы связи
0,1 - 0,3	Слабая
0,3 - 0,5	Умеренная
0,5 - 0,7	Заметная
0,7 - 0,9	Высокая
0,9 - 0,99	Весьма высокая



ЦЕНТРАЛЬНАЯ ПРЕДЕЛЬНАЯ ТЕОРЕМА



В ряде задач приходится сталкиваться с ситуацией, когда исследуемая случайная величина является суммой большого числа независимых слагаемых, влияние каждого из которых на сумму очень мало. На основании центральной предельной теоремы часто можно до наблюдения того или иного явления сказать, что соответствующая случайная величина должна иметь нормальное распределение или близкое к нему.

ТЕОРЕМА ЛЕВИ. Если независимые случайные величины $X_1, X_2, \dots, X_n, \dots$ распределены по одному и тому же закону с математическим ожиданием a и средним квадратическим отклонением σ , то при $n \rightarrow \infty$ функция распределения случайной величины

$$Z_n = \frac{\sum_{i=1}^n (X_i - a)}{\sigma\sqrt{n}}$$

сходится (в определенном смысле) к функции распределения стандартной нормальной случайной величины $\mathcal{N}(0; 1)$. Это записывается так:

$$Z_n \Rightarrow \mathcal{N}(0; 1).$$

