# NONPARAMETRIC STATISTICS. PART 3. CORRELATION COEFFICIENTS

**Marina A. Nikitina\*, Irina M. Chernukha**
V. M. Gorbatov Federal Research Center for Food Systems, Moscow, Russia

*Abstract*
*A measure of correlation or strength of association between random variables is the correlation coefficient. In scientific research, correlation analysis is most often carried out using various correlation coefficients without explaining why this particular coefficient was chosen and what the resulting value of this coefficient means. The article discusses Spearman correlation coefficient, Kendall correlation coefficient, phi (Yule) correlation coefficient, Cramér's correlation coefficient, Matthews correlation coefficient, Fechner correlation coefficient, Tschuprow correlation coefficient, rank-biserial correlation coefficient, point-biserial correlation coefficient, as well as association coefficient and contingency coefficient. The criteria for applying each of the coefficients are given. It is shown how to establish the significance (insignificance) of the resulting correlation coefficient. The scales in which the correlated variables should be located for the coefficients under consideration are presented. Spearman rank correlation coefficient and other nonparametric indicators are independent of the distribution law, and that is why they are very useful. They make it possible to measure the contingency between such attributes that cannot be directly measured, but can be expressed by points or other conventional units that allow ranking the sample. The benefit of rank correlation coefficient also lies in the fact that it allows to quickly assess the relationship between attributes regardless of the distribution law. Examples are given and step-by-step application of each coefficient is described. When analyzing scientific research and evaluating the results obtained, the strength of association is most commonly assessed by the correlation coefficient. In this regard, a number of scales are given (Chaddock scale, Cohen scale, Rosenthal scale, Hinkle scale, Evans scale) grading the strength of association for correlation coefficient, both widely recognized and not so well known.*

## Introduction

Correlation (from the Latin *correlatio*), or correlation dependence is a statistical relationship between two or more random variables (or values that may be considered as such with some acceptable degree of accuracy), while changes in the values of one or more of these variables are accompanied by a systematic change in the values of other variable(s) [1].

The term "correlation" was first used by the French paleontologist Jean Cuvier (1769–1832) in 1806: he developed the "law of correlation" for parts and organs of living organisms to restore the appearance of fossil animals. The law of correlation helps to reconstruct the appearance of the entire animal and its place in the system using skulls, bones, etc. from excavations: if the skull has horns, then it was an herbivore, and its legs had hooves; if legs have claws, then it was a carnivore without horns, but with large cuspids [2]. The following story is known about Cuvier and the "law of correlation". During a university holiday, students decided to play a prank on Professor Cuvier. They dressed one of the students in a goatskin with horns and hooves and lifted

him in Cuvier's bedroom window. The student stomped his hooves and yelled: "I'll eat you!" Cuvier woke up, saw a silhouette with horns and calmly answered: "If you have horns and hooves, then according to the law of correlation, you are an herbivore, and you cannot eat me. And for not knowing the law of correlation, you'll get a bad mark!" [3].

However, in statistics, the term "correlation" (in relation to Spearman correlation) was first used by the English biologist and statistician Galton F. (1822–1911) at the end of the 19th century. In 1892, he was the first to propose principles on how to calculate the correlation coefficient. His work was greatly influenced by the papers of Charles Darwin, who was his cousin. At a meeting of the Royal Society in 1888, Galton F. has presented a report "Correlations and their measurement, mainly from anthropometric data", which was devoted to the correlation between the length of arms and legs in a well-proportioned person. An article based on the 1888 report was published next year [4]. "*Two variable organs are considered correlated when a change in one of them is accompanied, in general, by a greater or lesser change in*

*the same direction in the other organ. Thus, the length of the arm is considered to be correlated with the length of the leg, because a person with a long arm usually has a long leg, and vice versa*" [4].

Galton F. calculated the correlation coefficient in anthropometry and in heredity studies. At University College London, Galton F. was the supervisor of Pearson K. (1857–1936), and then they worked together for many years. Pearson K. subsequently became a brilliant mathematician and biographer of Galton F.

Pearson K. is the founder of mathematical statistics, in particular the theory of correlation. He improved mathematical tools for calculating correlation. As a result, widely recognized Pearson correlation coefficient appeared, or analysis using Pearson method. In addition to Pearson K., Francis Ysidro Edgeworth and Walter Frank Raphael Weldon also worked on Pearson correlation coefficient [5]. He also developed nonparametric xi-squared coefficient. These coefficients are widely used in psychodiagnostics studies. Due to them, a tradition of using quantitative methods in the development and use of psychological tests was established.

Along with this, the following scientists made a significant contribution to the development of correlation analysis: Charles Edward Spearman (1863–1945), Maurice George Kendall (1907–1983), Alexander Tschuprow (1874–1926), George Udny Yule (1871–1951) and many others.

There are two types of association between phenomena, i. e. functional and correlation ones.

Correlation relationships between attributes may arise in different ways.

**The first** (most important) **way** is the causal dependence of the resulting attribute (its variation) on the variation of the factor attribute. For example, attribute x is a score for assessing soil fertility, and attribute y is the yield of an agricultural crop. Here it is completely clear which attribute acts as an independent variable (factor) x, and which attribute acts as a dependent variable (result) y.

**The second way** is contingency, which arises in the presence of a common cause. There is a well-known classic example given by the largest statistician in Russia at the beginning of the 20th century, Tschuprow A.: if we take the number of fire brigades in the city as attribute x, and the amount of losses from fires per year in the city as attribute y, then there is a direct correlation between attributes x and y in the Russian cities; on average, the more firefighters in a city, the greater the losses from fires! Did the firefighters set fires for fear of losing their jobs? No, the point is different. This correlation cannot be interpreted as an association between cause and consequence; both attributes are consequences of a common cause, i. e. the size of the city. It is quite logical that in large cities there are more fire departments, but there are more fires and losses from them per year than in small cities.

**The third way** correlation arises is a relationship of attributes, each of which is both a cause and a consequence. This is, for example, the correlation between the levels of labor productivity of workers and the level of wages for 1 hour of labor (rate). On the one hand, the level of wages is a consequence of labor productivity: the higher it is, the higher the payment. But, on the other hand, established rates play a stimulating role: with the right payment system, they act as a factor on which labor productivity depends. In such an attribute system, both formulations of the problem are permissible; each attribute can act as an independent variable x and a dependent variable y.

The publication provides an overview and systematizes information on the conditions for using various correlation coefficients and their grading scales.

**Objects and methods**

The research materials are monographs, manuals, articles, educational documents on statistics.

The authors searched for publications using key phrases: "correlation coefficients", "rank correlation coefficients", "association coefficient and contingency coefficient", "grading scales for correlation coefficients" in Scopus, PubMed, MEDLINE, Web of Knowledge, Google Scholar, IEEE Xplore, Science Direct, eLibrary (RSCI) databases.

The identified publications were preliminarily analyzed in the context of abstracts. The authors selected the following exclusion criteria:
1. works, scientific publications, textbooks devoted to "classical" methods of statistics;
2. publications not related to food and agricultural products.
   Inclusion criteria:
1. scientific articles, textbooks, monographs devoted to nonparametric statistics;
2. publications predominantly in English.

**Main part**
*Correlation coefficients*

Typically, correlation coefficient is a measure of correlation (or strength of association) between random variables. The following correlation coefficients exist: Pearson coefficient, Spearman coefficient, Kendall coefficient, etc. [6]. Table 1 presents the types of correlation coefficients and describes the scales in which variables vary.

*Pearson correlation coefficient (r)*

Pearson coefficient is used quite often by researchers. But before choosing this criterion you need to: 1) know the data type; 2) know the distribution of the studied attributes in the general population, and if this is unknown, you need to check the distribution of both variables in the sample; 3) construct scattergrams in order to make sure that the association between variables is linear, and also to check the homoscedasticity [8].

With skewed distributions, as well as in the presence of true outliers (if researchers decide to include them for analysis), it is better to use nonparametric correlation coefficients, i. e. Spearman coefficient, Kendall coefficient, Cramér's coefficient, etc. It is worth noting that in foreign publications, Spearman correlation coefficient is found much more often [9,10].

**Table 1. Types of correlation coefficients [7]**

| Correlation coefficient | Types of scale | |
|---|---|---|
| | variable X | variable Y |
| Pearson coefficient ($r$) | Interval scale with normal distribution | Interval scale with normal distribution |
| Spearman coefficient ($\rho$) | Interval scale with normal distribution | Ordinal scale |
| | Interval scale with normal distribution | Interval scale with normal distribution |
| Kendall coefficient ($\tau$) | Ordinal scale | Ordinal scale |
| Phi correlation coefficient ($\varphi$) for tables $2 \times 2$ | Nominal scale | Nominal scale |
| Cramér's coefficient (V) for tables more than $2 \times 2$ | Nominal scale | Nominal scale |
| Rank-biserial correlation coefficient ($r_{rb}$) | Nominal scale | Ordinal scale |
| Point-biserial correlation coefficient ($r_{pb}$) | Nominal scale | Interval scale with normal distribution |
| Matthews correlation coefficient (MCC) | Nominal scale | Nominal scale |
| Fechner correlation coefficient ($r_\Phi$) | Interval scale with normal distribution | Interval scale with normal distribution |
| Tschuprow contingency coefficient ($r_{ch}$) | Nominal scale | Nominal scale |

*Spearman rank correlation coefficient ($\rho$)*

Rank correlation coefficients are used to measure relationships between attributes, the values of which may be ordered or ranked according to the decrease (or increase) of a given indicator in the objects under study.

Spearman rank correlation coefficient is a quantitative assessment of the association between phenomena, which is used in nonparametric methods. It determines the strength and direction of the correlation association between two attributes or two profiles of attributes. In this case, the actual degree of parallelism between the two quantitative series of observations being studied is determined and an assessment of the established association strength is given using a quantitatively expressed coefficient [11]. The criterion was developed and proposed for correlation analysis in 1904 by Charles Edward Spearman, an English psychologist, professor at the Universities of London and Chesterfield.

$$\rho = 1 - \frac{6 \sum_{i=1}^{n} d_i^2}{n^3 - n} \tag{1}$$

where $d_i = x_i - y_i$ is the difference between ranks;
$n$ is the number of attribute values observed.

Rank is the position of an element in a variation series. A variation series is a series whose elements are arranged in ascending or descending order.

The significance of Spearman correlation coefficient may be determined by Student's test (t-test) with the number of degrees of freedom equal to $n - 2$.

$$t = \rho \cdot \sqrt{\frac{n-2}{1-\rho^2}} \tag{2}$$

The criteria for applying the nonparametric Spearman rank correlation coefficient are described in detail in [10], and are as follows:

1. Examination of quantitative data distributions for normality is not required. It can be used for samples whose data partially or completely does not follow the law of normal distribution.

2. If the data from one of the samples can be presented on an ordinal scale, the data from the second sample must be quantitative.

3. If the sample size exceeds 5 observations.

4. If there are a large number of identical ranks for one or both compared variables, then Spearman correlation coefficient gives rough values. Ideally, both correlated series should represent two sequences of divergent values.

*Kendall correlation coefficient ($\tau$)*

Kendall rank correlation [12] is an alternative to Spearman correlation in the case of two ordinal scales. This method is a measure of the strength of a nonlinear association and uses an increase or decrease in the resultant attribute as the factor attribute increases. Thus, the calculation of Kendall correlation coefficient involves counting the number of coincidences and inversions. To use Kendall correlation coefficient, there is only one requirement: the scales of the X and Y variables must be ordinal.

$$\tau = \frac{4Q}{n(n-1)} - 1 \tag{3}$$

where $Q$ is the minimum number of exchanges of neighboring elements in one of the rankings for its coincidence with another ranking.

The statistics for the test of significance of this coefficient have a normal distribution $N(0, 1)$.

$$T_{kr} = z_{kr} \cdot \sqrt{\frac{2 \cdot (2n+5)}{9n \cdot (n-1)}} \tag{4}$$

where $n$ is a sample size; $z_{kr}$ is a critical point of the two-sided critical region determined by Laplace function $\Phi(z_{kr}) = \frac{1-\alpha}{2}$ [13,14,15,16].

If $|\tau| < T_{kr}$, rank correlation between attributes is insignificant.

If $|\tau| > T_{kr}$, there is a significant rank correlation between attributes.

Same as for Spearman correlation coefficient: when ranks coincide $\tau = 1$, with opposite ranks $\tau = -1$. Kendall rank correlation coefficient has some advantages over Spearman coefficient. In particular, it may also be used for multivariate analysis. With a sufficiently large number of objects ($n \geq 10$), there is a simple association between the values of rank correlation coefficients $\rho = 1.5 \cdot \tau$ [17].

**Example.** Two panelists conduct a sensory analysis of 10 cooked sausage samples: ranked in descending order.

Ranks by the first panelist: 2, 3, 1, 6, 5, 4, 8, 7, 10, 9

Ranks by the second panelist: 3, 1, 2, 6, 7, 4, 5, 9, 10, 8

Using Spearman rank correlation coefficient and Kendall rank correlation coefficient, it should be determined whether the ratings of the panelist are consistent.

**Solution.**

I.1. To determine Spearman rank correlation coefficient, let's find the difference between the ranks $d_i$, and the square of the difference between the ranks $d_i^2$. The results are presented in Table 2.

**Table 2. Calculation results**

| No. of the cooked sausage sample | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| Panelist 1 | 2 | 3 | 1 | 6 | 5 | 4 | 8 | 7 | 10 | 9 |
| Panelist 2 | 3 | 1 | 2 | 6 | 7 | 4 | 5 | 9 | 10 | 8 |
| $d_i$ | –1 | 2 | –1 | 0 | –2 | 0 | 3 | –2 | 0 | 1 |
| $d_i^2$ | 1 | 4 | 1 | 0 | 4 | 0 | 9 | 4 | 0 | 1 |

2. Let's determine the sum of squares of the rank difference. The number of samples is 10, i.e. n=10.

$$\sum_{i=1}^{n} d_i^2 = 1 + 4 + 1 + 0 + 4 + 0 + 9 + 4 + 0 + 1 = 24$$

3. According to the formula (1), let's calculate Spearman rank correlation coefficient
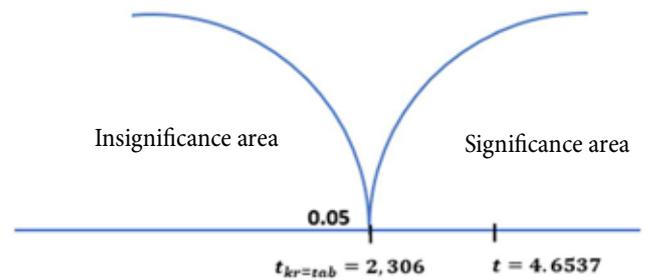
$$\rho = 1 - \frac{6 \cdot 24}{10^3 - 10} = 0.8545$$

4. Test of significance is carried out according to the formula (2). Let's calculate Student's test (t-test)

$$t = 0.8545 \cdot \sqrt{\frac{10 - 2}{1 - 0.8545^2}} = 4.6537$$

5. Let's calculate the critical values of Student's test, significance level $p = 0.05$, the number of degrees of freedom in our case will be equal to $\nu = n - 2 = 10 - 2 = 8$. We can use statistical tables [13,14,15,16] or the function in MS Excel, TINV (Figure 1).

6. Let's plot "the axis of significance" (Figure 2) for our example.



**Figure 2**. The axis of significance

Since 4.6537 > 2.306, the correlation is statistically significant.

II.1. To determine Kendall rank correlation coefficient, let's find the minimum number of exchanges of neighboring elements in one of the rankings for its coincidence with the other ranking. The results are presented in Table 3.

**Table 3. Calculation results**

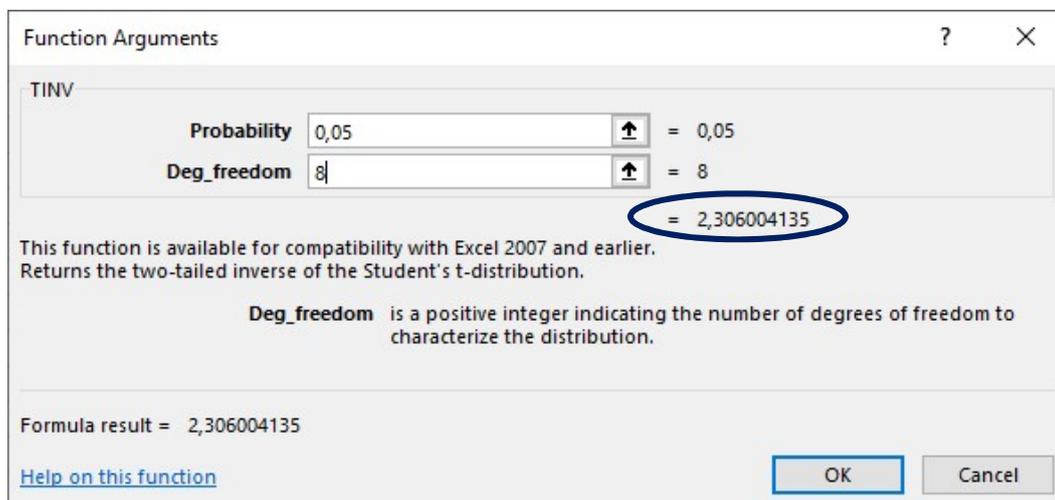| No. of the cooked sausage sample | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| Panelist 1 | 2 | 3 | 1 | 6 | 5 | 4 | 8 | 7 | 10 | 9 |
| Panelist 2 | 3 | 1 | 2 | 6 | 7 | 4 | 5 | 9 | 10 | 8 |
| $Q_i$ | 7 | 8 | 7 | 4 | 3 | 4 | 3 | 1 | 0 | 0 |

$Q_1 = 7$, since in the line "Panelist 2" to the right of 3 (the values of samples that are to the right of the sample under consideration), there are 7 values larger than 3 (samples 4, 5, 6 ,7, 8, 9, 10).

$Q_2 = 8$, since in the line "Panelist 2" to the right of 1, there are 8 values larger than 1 (samples 3, 4, 5, 6, 7, 8, 9, 10).

$Q_3 = 7$, since in the line "Panelist 2" to the right of 2, there are 7 values larger than 2 (samples 4, 5, 6, 7, 8, 9, 10).

$Q_4 = 4$, since in the line "Panelist 2" to the right of 6, there are 4 values larger than 6 (samples 5, 8, 9, 10).

We filled in further in the same way.



**Figure 1.** Calculation of the critical (reference) value of Student's test

2. Let's calculate the sum of the values $Q_i$

$$\sum_{i=1}^{n} Q = 7 + 8 + 7 + 4 + 3 + 4 + 3 + 1 + 0 + 0 = 37$$

3. According to the formula (3), let's calculate Kendall rank correlation coefficient

$$\tau = \frac{4 \cdot 37}{10 \cdot (10 - 1)} - 1 = 0.6444$$

4. Test of significance is carried out according to the formula (4). Let's calculate the criterion $T_{kr}$

$$T_{kr} = z_{kr} \cdot \sqrt{\frac{2 \cdot (2n + 5)}{9n \cdot (n - 1)}}$$

Critical point $z_{kr}$ in Laplace table [13,14,15,16] is equal to 1.96 at $\Phi(z_{kr}) = \frac{1-\alpha}{2} = \frac{1-0.05}{2} = 0.475$

$$T_{kr} = z_{kr} \cdot \sqrt{\frac{2 \cdot (2n + 5)}{9n \cdot (n - 1)}} = 1.96 \cdot \sqrt{\frac{2 \cdot (2 \cdot 10 + 5)}{9 \cdot 10 \cdot (10 - 1)}} =$$

$$= 1.96 \cdot \sqrt{\frac{50}{810}} = 0.017$$

5. Since $|\tau| > T_{kr}$, rank correlation between scores in two tests is significant.

*Phi ($\varphi$) correlation coefficient and Cramér's V-coefficient*

To study the strength of association between variables measured on a nominal scale, phi correlation coefficient and Cramér's coefficient are used.

In statistics, phi correlation coefficient (or root mean square contingency coefficient) is a measure of association between two binary variables. In machine learning, it is known as Matthews correlation coefficient (MCC) introduced by biochemist Brian W. Matthews in 1975 and is used as an indicator of the quality of binary (two-class) classifications [18]. Phi correlation coefficient was introduced by Pearson K. [19], also known as Yule phi coefficient introduced by George Udny Yule in 1912 [20].

The criteria for applying phi ($\varphi$) correlation coefficient:
1. Variables X and Y must be measured on a dichotomous scale.
2. The number of attributes in the compared variables X and Y must be the same.

Two binary variables are considered positively associated if the majority of the data is in the diagonal cells. Otherwise, if most of the data falls off the diagonal, then the binary variables are considered negatively associated.

Phi correlation coefficient may be calculated using fourfold contingency table 2×2 (Table 4).

**Table 4. Fourfold contingency table 2×2**

|  | Y=1 | Y=0 |  |
|---|---|---|---|
| X=1 | a | b | m1=(a+b) |
| X=0 | c | d | m2=(c+d) |
|  | n1=(a+c) | n2=(b+d) | n=(a+b+c+d) |

$$\varphi = \frac{a \cdot d - b \cdot c}{\sqrt{m1 \cdot m2 \cdot n2 \cdot n1}} \qquad (5)$$

where *a, b, c, d* are non-negative values of the number of observations, which add up to *n*, the total number of observations.

Phi ($\varphi$) correlation coefficient is related to point-biserial correlation coefficient and Cohen's *d* and estimates the degree of relationship between two variables (2 × 2) [21].

Matthews correlation coefficient (MCC) is defined identically to phi correlation coefficient and is widely used in the fields of bioinformatics and machine learning. The coefficient is considered as a balanced measure that can be used even if the classes have very different sizes [22].

Matthews correlation coefficient (MCC) returns a value between -1 and +1. Coefficient of "+1" represents a perfect prediction, "0" is no better than a random prediction, and "-1" indicates a complete discrepancy between the prediction and observation.

Cramér's V-coefficient is a modified phi correlation coefficient for tables larger than 2×2. This indicator of association between two nominal variables varies from 0 to + 1 (inclusive). It is based on Pearson chi-square statistics and was published by Harald Cramér in 1946 [23].

The criteria for applying Cramér's V-coefficient:
1. Variables X and Y must be measured on a nominal scale, where the number of codings is more than two (not dichotomous scales).
2. The number of attributes in the compared variables X and Y must be the same.

Like phi correlation coefficient, Cramér's V-coefficient is calculated using contingency tables (larger than 2×2).

$$V = \sqrt{\frac{\chi^2}{n \cdot min(row-1, column-1)}} \qquad (6)$$

Chi-square test is calculated according to the formula:

$$\chi^2 = \sum \frac{(n_i - \hat{n}_i)^2}{\hat{n}_i} \qquad (7)$$

where $n_i$ is the actual number of observations in *ij* cells; $\hat{n}_i$ is the expected number of observations in *ij* cells.

A general overview of the expected values is presented in Table 5.

**Table 5. A general overview of the table of the expected values**

|  | There is an outcome (1) | There is no outcome (0) | Total |
|---|---|---|---|
| There is a risk factor (1) | $\frac{(A+B) \cdot (A+C)}{A+B+C+D}$ | $\frac{(A+B) \cdot (B+D)}{A+B+C+D}$ | A+B |
| There is no risk factor (0) | $\frac{(C+D) \cdot (A+C)}{A+B+C+D}$ | $\frac{(C+D) \cdot (B+D)}{A+B+C+D}$ | C+D |
| Total | A+C | B+D | A+B+C+D |

**Example.** A study is being conducted on the effect of smoking on the risk of developing arterial hypertension. For this purpose, two groups of subjects were selected: the first group included 70 people who smoke at least 1 pack of cigarettes daily, the second group included 80 non-smokers of the same age. In the first group, 40 people had high

blood pressure. In the second group, arterial hypertension was observed in 32 people. Thus, in the group of smokers, normal blood pressure was in 30 people (70 – 40 = 30), and in the group of non-smokers, normal blood pressure was in 48 people (80 – 32 = 48)."

**Solution.** Let's generate a contingency table (Table 6).

**Table 6. Fourfold contingency table 2×2**

| | There is an arterial hypertension (1) | There is no arterial hypertension (0) | |
|---|---|---|---|
| Smokers (1) | A=40 | B=30 | A+B=70 |
| Non-smokers (0) | C=32 | D=48 | C+D=80 |
| | A+C=72 | B+D=78 | A+B+C+D=150 |

A general overview of the expected values is presented in Table 7 according to the formulas in Table 6.

$$\frac{(A+B)\cdot(A+C)}{A+B+C+D} = \frac{70\cdot72}{150} = 33.6$$

$$\frac{(C+D)\cdot(A+C)}{A+B+C+D} = \frac{80\cdot72}{150} = 38.4$$

$$\frac{(A+B)\cdot(B+D)}{A+B+C+D} = \frac{70\cdot78}{150} = 36.4$$

$$\frac{(C+D)\cdot(B+D)}{A+B+C+D} = \frac{80\cdot78}{150} = 41.6$$

**Table 7. A general view of the expected values**

| | There is an arterial hypertension (1) | There is no arterial hypertension (0) | |
|---|---|---|---|
| Smokers (1) | 33.6 | 36.4 | 70 |
| Non-smokers (0) | 38.4 | 41.6 | 80 |
| | 72 | 78 | 150 |

Let's calculate chi-square test:

$$\chi^2 = \sum \frac{(n_i - \hat{n}_i)^2}{\hat{n}_i} = 4.3956$$

We determine the number of degrees of freedom $v = (row-1)\cdot(column-1) = (2-1)\cdot(2-1) = 1$,

where row is the number of rows (in our example row=2), column is the number of columns (in our example column=2).

We find the critical value of Pearson chi-square test at significance level of p=0.05. We can use statistical tables [13,14,15,16] or the MS Excel function, CHIINV (Figure 3).

At significance level of p=0.05 and number of degrees of freedom equal to 1, $\chi^2_{kr(tab)} = 3.8415$.
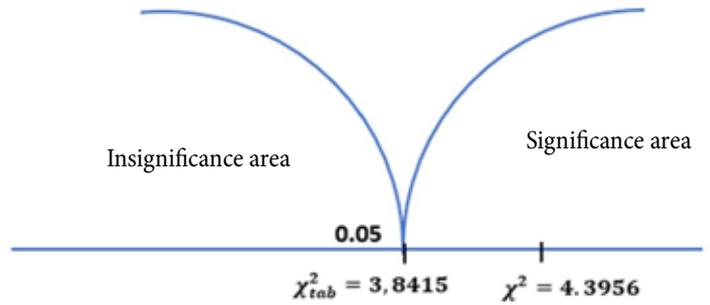
Let's plot "the axis of significance" (Figure 4).



**Figure 4.** The axis of significance

Since 4.396 > 3.841, the dependence of the arterial hypertension incidence on smoking is statistically significant. The significance level of this relationship corresponds to p<0.05.

Cramér's V-coefficient:

$$V = \sqrt{\frac{\chi^2}{n\cdot min(row-1, column-1)}} = 0.1712$$

*Fechner correlation coefficient ($r_\phi$)*

The simplest indicators of strength of association include the sign correlation coefficient, which was proposed by the German physicist, philosopher and psychologist, founder of psychophysics, Gustav Theodor Fechner (1801-1887). In his posthumously published collective measurement theory (Kollektivmasslehre, 1897) [24], Fechner introduced the concept of the two-sided Gauss' law (Zweiseitige Gauss'sche Gesetz) or two-part normal distribution to account for the asymmetries he observed in empirical frequency distributions in many areas.
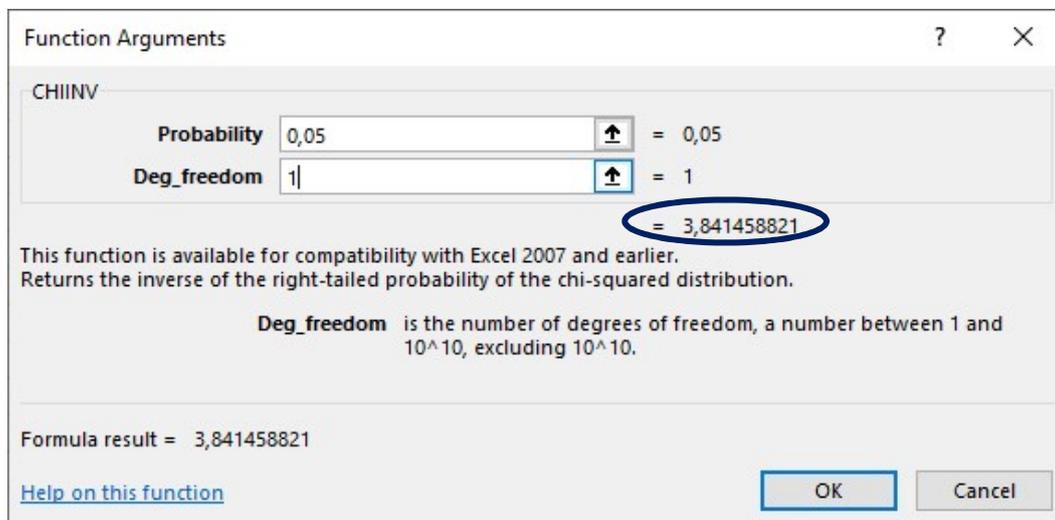


**Figure 3.** Calculation of the critical (reference) value of Pearson chi-square test

Fechner correlation coefficient is based on assessing the degree of consistency in the directions of deviations in individual values of factor attribute and resulting attribute from the corresponding averages. To calculate it, the average values of the resulting attribute and factor attribute are calculated, and then deviation signs for all values of correlated pairs of attributes are assigned.

$$r_\phi = \frac{C-H}{C+H} \qquad (8)$$

where $C$ is the number of coincidences of identical difference signs, both positive and negative $(x_i - \bar{x})$ and $(y_i - \bar{y})$;
$H$ is the number of non-coincided difference signs $(x_i - \bar{x})$ and $(y_i - \bar{y})$;
$\bar{x}, \bar{y}$ are the average values of vectors (samples) $x_i, y_i$.

Like Pearson correlation coefficient, Fechner correlation coefficient may be in the range from –1 to +1. With a positive correlation, it has a positive sign, and with a negative correlation, it has a negative sign.

When using Fechner correlation coefficient, it should be noted that the distribution law of Fechner coefficient is unknown. Therefore, the question of assessing reliability remains.

**Example.** Based on the data accumulated on the milk fat content for cows at the farm and their 12 daughters of the same age (Table 8), we need to determine the relationship between the milk fat content for cows of the maternal generation and their offspring of the same age [25].

**Table 8. Given data**

| Fat percentage in milk | | | |
|---|---|---|---|
| cows $x_i$ | daughters $y_i$ | cows $x_i$ | daughters $y_i$ |
| 3.10 | 3.65 | 3.80 | 3.61 |
| 3.17 | 3.11 | 3.65 | 3.98 |
| 3.76 | 3.57 | 3.34 | 3.36 |
| 3.61 | 3.61 | 3.65 | 3.89 |
| 3.27 | 3.44 | 3.45 | 3.45 |
| 3.61 | 3.71 | 4.05 | 3.79 |

**Solution.**

1. First let's calculate arithmetic averages of the vectors of milk fat content for cows ($x_i$) and their daughters ($y_i$).

$\bar{x}$ = (3.10 + 3.17 + 3.76 + 3.61 + 3.27 + 3.61 + 3.80 +

+3.65 + 3.34 + 3.65 + 3.45 + 4.05)/12 = 3.5383

$\bar{y}$ = (3.65 + 3.11 + 3.57 + 3.61 + 3.44 + 3.71 + 3.61 +

+3.98 + 3.36 + 3.89 + 3.45 + 3.79)/12 = 3.5975

2. Then let's calculate the difference $(x_i - \bar{x})$ and $(y_i - \bar{y})$; data are presented in Table 9.

3. We calculate the number of coincided signs C = 10, and the number of non-coincided signs H = 2 (highlighted in green in Table 9).

4. Let's calculate Fechner correlation coefficient

$$r_\phi = \frac{10 - 2}{10 + 2} = 0.6667$$

**Table 9. Deviations from the average values**

| Fat percentage in milk | | Deviations from the average values | | Signs of deviations from the average values | |
|---|---|---|---|---|---|
| cows $x_i$ | daughters $y_i$ | $(x_i - \bar{x})$ | $(y_i - \bar{y})$ | $(x_i - \bar{x})$ | $(y_i - \bar{y})$ |
| 3.10 | 3.65 | –0.4383 | 0.0525 | – | + |
| 3.17 | 3.11 | –0.3683 | –0.4875 | – | – |
| 3.76 | 3.57 | 0.2217 | –0.0275 | + | – |
| 3.61 | 3.61 | 0.0717 | 0.0125 | + | + |
| 3.27 | 3.44 | –0.2683 | –0.1575 | – | – |
| 3.61 | 3.71 | 0.0717 | 0.1125 | + | + |
| 3.80 | 3.61 | 0.2617 | 0.0125 | + | + |
| 3.65 | 3.98 | 0.1117 | 0.3825 | + | + |
| 3.34 | 3.36 | –0.1983 | –0.2375 | – | – |
| 3.65 | 3.89 | 0.1117 | 0.2925 | + | + |
| 3.45 | 3.45 | –0.0883 | –0.1475 | – | – |
| 4.05 | 3.79 | 0.5117 | 0.1925 | + | + |

Thus, it can be stated that there is a moderate association between the milk fat content from cows of the maternal generation and their offspring of the same age.

*Rank-biserial correlation coefficient ($R_{rb}$)*

In cases where one variable is measured on a dichotomous scale (variable X), and the other variable is measured on a rank scale (variable Y), rank-biserial correlation coefficient is used. Variable X measured on a dichotomous scale have only two values (codes), 0 and 1. It should be especially emphasized: despite the fact that this coefficient varies in the range from –1 to +1, its sign does not matter for the interpretation of the results. This is another exception to the general rule.

This coefficient is calculated according to the formula:

$$R_{rb} = \frac{(\bar{x}_1 - \bar{x}_0) \cdot 2}{N} \qquad (9)$$

where $\bar{x}_1$ is the average rank for those elements of variable Y that correspond to code (attribute) 1 in variable X;
$\bar{x}_0$ is the average rank for those elements of variable Y that correspond to code (attribute) 0 in variable X;
$N$ is the total number of elements in variable X.

**Example.** A psychologist tests a hypothesis about whether there are gender differences in verbal ability.

**Solution.** To solve this problem, 15 teenagers of different genders were ranked by a literature teacher according to the degree of expression of verbal abilities. The data obtained are presented in the form of a table (Table 10).

In this case, the correctness of the ranking need not be checked, since there are no coincided ranks and the ranking is carried out in order. In Table 10, boys are designated by code 1 (green), and girls are designated by code 0. In our case, there are 9 boys and 6 girls.

1. Let's calculate the average rank values separately for boys and girls.

$$\bar{x}_1 = \frac{1 + 6 + 9 + 7 + 4 + 3 + 5 + 12 + 2}{9} = \frac{49}{9} = 5.44$$

$$\bar{x}_0 = \frac{10 + 15 + 8 + 13 + 11 + 14}{6} = \frac{71}{6} = 11.83$$

**Table 10. Verbal abilities of teenagers**

| No. of the subject | Gender | Verbal ability ranks |
|---|---|---|
| 1 | 1 | 1 |
| 2 | 0 | 10 |
| 3 | 1 | 6 |
| 4 | 1 | 9 |
| 5 | 0 | 15 |
| 6 | 1 | 7 |
| 7 | 0 | 8 |
| 8 | 0 | 13 |
| 9 | 1 | 4 |
| 10 | 1 | 3 |
| 11 | 1 | 5 |
| 12 | 0 | 11 |
| 13 | 1 | 12 |
| 14 | 1 | 2 |
| 15 | 0 | 14 |

2. Let's calculate rank-biserial correlation coefficient $R_{rb}$ according to the formula (9):

$$R_{rb} = \frac{(\bar{x}_1 - \bar{x}_0) \cdot 2}{N} = \frac{(5.44 - 11.83) \cdot 2}{15} = -0.852$$
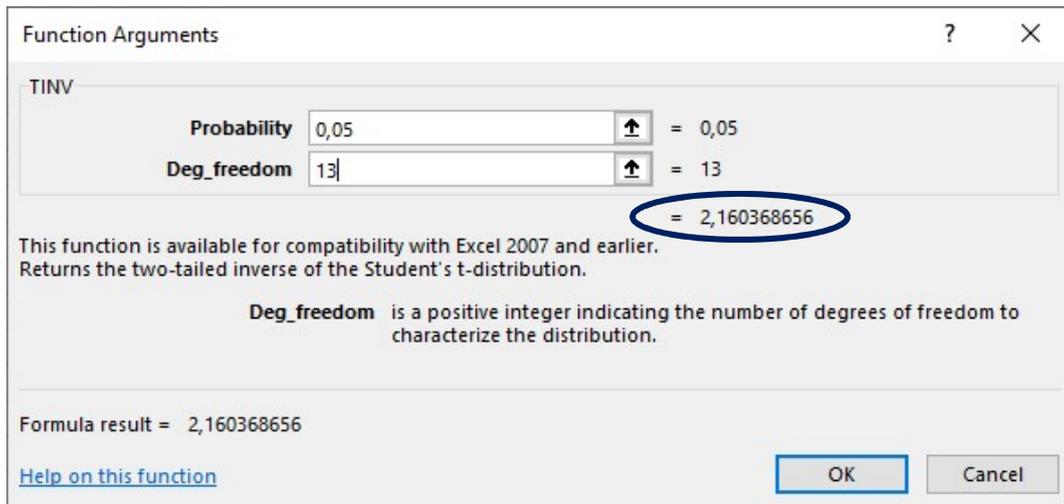
3. Let's check the significance of the resulting correlation coefficient using the formula

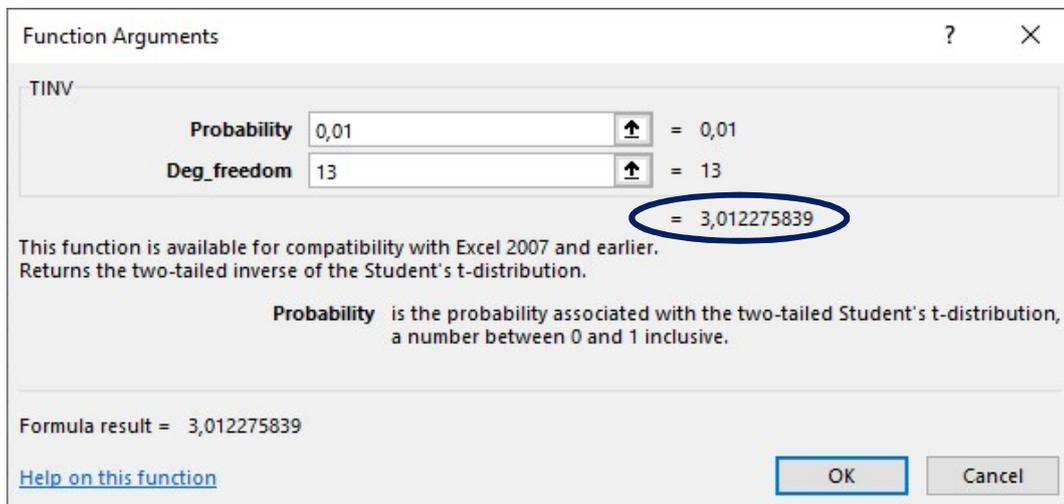$$T_\Phi = |R_{rb}| \cdot \sqrt{\frac{N-2}{1-R_{rb}^2}} \qquad (10)$$

at $v = N - 2 = 15 - 2 = 13$ ($v$ is the degree of freedom by which the reference (critical) value is found and compared with the calculated value obtained according to the formula (10).

$$T_\Phi = |R_{rb}| \cdot \sqrt{\frac{N-2}{1-R_{rb}^2}} = |-0.852| \cdot \sqrt{\frac{15-2}{1-(-0.852)^2}} =$$

$$= 0.852 \cdot \sqrt{\frac{13}{1-0.725904}} = 0.852 \cdot \sqrt{\frac{13}{0.274096}} =$$

$$= 0.852 * 6.88684529674 = 5.87$$

In our case, the number of degrees of freedom will be equal to $v=13$. To calculate the critical (reference) value of Student's test, we can use statistical tables [13,14,15,16] or the MS Excel function, TINV (Figures 5 and 6).



**Figure 5.** Calculation of the critical (reference) value of Student's test at significance level of $p<0.05$



**Figure 6.** Calculation of the critical (reference) value of Student's test at significance level of $p<0.01$

4. The critical (reference) value of Student's test for $P < 0.05$ is equal to $t_{kr=tab} = 2.16$ and for $P < 0,01$ is equal to $t_{kr=tab} = 3.01$. In the accepted notation form it looks like this:

$$t_{kr=tab} = \begin{cases} 2.16 \ at\ P \leq 0.05 \\ 3.01 \ at\ P \leq 0.01 \end{cases}$$

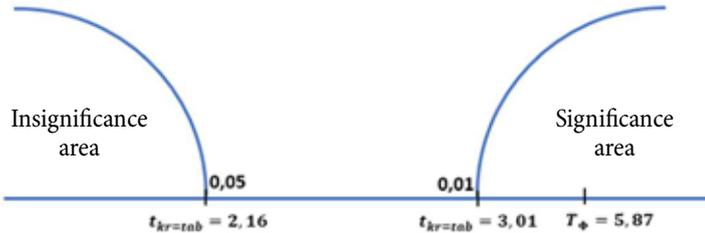5. Let's plot "the axis of significance" (Figure 7):



**Figure 7.** The axis of significance

The result is in significance area. Therefore, hypothesis $H_1$, according to which the resulting rank-biserial correlation coefficient is significantly different from zero, is accepted. In other words, in this sample of teenagers, significant gender differences were found in the degree of expression of verbal abilities.

To use the rank-biserial correlation coefficient, the following criteria must be met:

1. The variables being compared must be measured on different scales: X on a dichotomous scale; Y on a rank scale.

2. The number of varying attributes in the compared variables X and Y must be the same.

3. To assess the level of reliability of the rank-biserial correlation coefficient, we need to use the formula to determine $T\phi$ and a table of critical values for Student's t-test at $\upsilon = N - 2$.

*Tschuprow contingency coefficient* ($K_\text{ч}$) is calculated according to the formula

$$K_\text{ч} = \sqrt{\frac{\varphi^2}{(K_1-1)\cdot(K_2-1)}} \qquad (11)$$

where $K_1$ and $K_2$ the number of groups in the columns and the number of groups in the rows.

The result of assessing the strength of association obtained by Tschuprow contingency coefficient is more accurate, since it takes into account the number of groups for each of the studied attributes.

It is also beneficial to use when there is a greater division of population into groups according to correlated attributes. Pearson contingency coefficient is used mainly in the case of a square table, while Tschuprow contingency coefficient is suitable for measuring association in rectangular tables also.

It is believed that already with contingency coefficient value of 0.3, we can talk about a strong association between the variation of the studied attributes.

**Example.** Using Tschuprow contingency coefficient, it is necessary to determine the strength of association between the grain yield in agricultural enterprises of the region and their legal form according to Table 11.

**Table 11. Grouping of agricultural enterprises with different grain yields according to legal form**

| Grain yield, dt/ha ($x_i$) | Number of enterprises, units ($f_{i0}$) | State enterprises ($f_{i1}$) | collective enterprises ($f_{i2}$) | Farming enterprises ($f_{i3}$) |
|---|---|---|---|---|
| | | including | | |
| 15.80-18.97 | 3 | 2 | 1 | – |
| 18.97-22.14 | 4 | – | 4 | – |
| 22.14-25.31 | 11 | 3 | 8 | – |
| 25.31-28.48 | 7 | 1 | 3 | 3 |
| 28.48-31.65 | 4 | – | 1 | 3 |
| 31.65-34.82 | 1 | – | – | 1 |
| Total: | 30 | 6 | 17 | 7 |

Let's transform the table into more convenient form for calculating Tschuprow contingency coefficient (Table 12).

**Table 12. Distribution of agricultural enterprises in the region by their legal form and level of grain yield**

| Group of enterprises | By grain yield (dt/ha) | | | | | | Total | Average yield by group, dt/ha |
|---|---|---|---|---|---|---|---|---|
| By legal form | 15.80-18.97 | 18.97-22.14 | 22.14-25.31 | 25.31-28.48 | 28.48-31.65 | 31.65-34.82 | | |
| Average value of the range | 17.4* | 20.6 | 23.7 | 26.9 | 30.1 | 33.2 | | |
| State enterprises | 2 | – | 3 | 1 | – | – | 6 | 22.14** |
| Collective enterprises | 1 | 4 | 8 | 3 | 1 | – | 17 | 23.54 |
| Farming enterprises | – | – | – | 3 | 3 | 1 | 7 | 29.16 |
| TOTAL: | 3 | 4 | 11 | 7 | 4 | 1 | 30 | 24.57 |

* $\frac{15.8+18.97}{2} = 17.38 \approx 17.4$; $\frac{18.97+2.14}{2} = 20.55 \approx 20.6$, etc.

** $\frac{17.4\cdot2+23.7\cdot3+26.9\cdot1}{6} \approx 22.14$; $\frac{17.4\cdot1+20.6\cdot4+23.7\cdot8+26.9\cdot3+30.1*1}{17} \approx 23.54$, etc.

According to the formula

$$\varphi^2 = \sum \frac{f_{ij}^2}{F_i F_j} - 1 \qquad (12)$$

where $F_i = \sum_i f_{ij}, F_j = \sum_j f_{ij}$

the mean square contingency is equal to

$$\varphi^2 = \left( \frac{2^2}{3\cdot6} + \frac{1^2}{3\cdot17} + \frac{4^2}{4\cdot17} + \frac{3^2}{11\cdot6} + \frac{8^2}{11\cdot17} + \frac{1^2}{7\cdot6} + \frac{3^2}{7\cdot17} + \right.$$
$$\left. + \frac{3^2}{7\cdot7} + \frac{1^2}{4\cdot17} + \frac{3^2}{4\cdot7} + \frac{1^2}{1\cdot7} \right) - 1 =$$
$$= \left( \frac{4}{18} + \frac{1}{51} + \frac{16}{68} + \frac{9}{66} + \frac{64}{187} + \frac{1}{42} + \frac{9}{119} + \frac{9}{49} + \frac{1}{68} + \frac{9}{28} + \frac{1}{7} \right)$$
$$- 1 = 0.718$$

According to the formula

$$K_\text{ч} = \sqrt{\frac{\varphi^2}{(K_1 - 1) \cdot (K_2 - 1)}} = \sqrt{\frac{0.718}{(6 - 1) \cdot (3 - 1)}} = 0.268$$

Tschuprow contingency coefficient is = 0.268. Since this value verges towards 0.3, we can talk about the presence of a fairly strong association between the yield of grain crops and the legal form of enterprises.

*Tschuprow correlation coefficient ($r_{ch}$)*

Tschuprow correlation coefficient ($r_{ch}$) is calculated using the following formula:

$$r_{ch} = \pm\sqrt{\frac{\chi^2}{N\sqrt{(a-1)\cdot(b-1)}}} \qquad (13)$$

where $\chi^2$ is the empirical value of the chi-square test;
$N$ is the sample size (number of objects for which both attributes were taken into account);
$a, b$ is the number of modalities of both attributes.

The reliability of Tschuprow correlation coefficient is assessed by the value of the chi-square test. Chi-square test is calculated according to the formula:

$$\chi^2 = \sum \frac{(n_i - \hat{n}_i)^2}{\hat{n}_i} \qquad (14)$$

Number of components added when calculating $\chi^2$ is equal to the product $a \times b$.

The null hypothesis is that there is no reliable association between the variables. If $\chi^2 > \chi^2_{0.05}$, the null hypothesis is rejected (the association between variables is significant); if $\chi^2 < \chi^2_{0.05}$, the null hypothesis is accepted (the association between variables is insignificant).

If it is proven that the association is insignificant, Tschuprow correlation coefficient is not calculated and is set to **0**.

**Example.** To establish the association between the shape of the glands on the leaf petioles and the degree (score) of powdery mildew damage to peach, 1319 cultivars were studied. The frequencies of occurrence of peach cultivars by combination of modalities of these attributes are as follows (Table 13). What is the correlation between the shape of the glands on the petioles and powdery mildew in peach?

**Table 13. Given data: powdery mildew damage on petioles**

| Powdery mildew damage | Shape of the glands | |
|---|---|---|
| | reform | rounded |
| Absent or minor | 453 | 40 |
| Medium or severe | 46 | 780 |

**Solution.** The attribute "shape of the glands" is nominal, since the modalities "reniform" and "rounded" cannot be ranked. The attribute "powdery mildew damage" can be considered as an ordinal attribute, since its states, i.e. "absent or minor" and "medium or severe" are easily ranked. If at least one of the attributes is nominal, then Tschuprow

correlation coefficient is used to estimate the correlation between it and other attributes.

1. First, we generate a table for frequencies of occurrence of cultivars based on the two studied attributes (Table 14) and calculate the theoretically expected frequencies, provided that there is no correlation between these attributes:

Empirical and theoretically expected frequencies of occurrence of peach cultivars based on the combination of modalities "shape of the glands" and "powdery mildew damage", provided that there is no correlation between these attributes.

**Table 14. Frequency of cultivar occurrence by two attributes**

| Powdery mildew damage | Shape of the glands | | | | Sum |
|---|---|---|---|---|---|
| | reniform | | rounded | | |
| | $n_i$ | $\hat{n}_i$ | $n_i$ | $\hat{n}_i$ | |
| Absent or minor | 453 | 186.51 | 40 | 306.49 | **493** |
| Medium or severe | 46 | 312.49 | 780 | 513.51 | **826** |
| TOTAL: | **499** | | **820** | | |

$$\hat{n}_{11} = \frac{499 \cdot 493}{1319} = 186.51$$

$$\hat{n}_{12} = \frac{820 \cdot 493}{1319} = 306.49$$

$$\hat{n}_{21} = \frac{499 \cdot 826}{1319} = 312.49$$

$$\hat{n}_{22} = \frac{820 \cdot 826}{1319} = 513.51$$

2. Let's calculate chi-square test value:

$$\chi^2 = \sum \frac{(n_i - \hat{n}_i)^2}{\hat{n}_i} = \frac{(453 - 186.51)^2}{186.51} + \frac{(46 - 312.49)^2}{312.49} +$$
$$+ \frac{(40 - 306.49)^2}{306.49} + \frac{(780 - 513.51)^2}{513.51} = 978.04$$

3. We find the critical value of Pearson chi-square test at significance level of p=0.05 and the degrees of freedom equal to $\upsilon = 2 - 1 = 1$. To calculate the critical (reference) value of Pearson chi-square test, we can use statistical tables [13,14,15,16] or the MS Excel function, CHIINV.

The critical (reference) value of Pearson chi-square test at significance level of $p = 0.05$ and the degrees of freedom equal to $\upsilon = 1$ is 3.84 ($\chi^2_{0.05} = 3.84$).
$\chi^2 = 978.03 > \chi^2_{0.05} = 3.84$

Statistical conclusion: the correlation between powdery mildew damage and the shape of the glands is significant.

4. Let's calculate Tschuprow correlation coefficient:

$$r_{ch} = \pm\sqrt{\frac{\chi^2}{N\sqrt{(a-1)\cdot(b-1)}}} = \pm\sqrt{\frac{978.04}{1319\sqrt{(2-1)\cdot(2-1)}}}$$
$$= \pm\sqrt{0.7415} = \pm0.86$$

**Conclusion:** The correlation between the powdery mildew score and the type of glands is reliable and strong. However, it is impossible to establish which variable is an argument and which is a function. Though, it may be logically assumed that the shape of the glands is an independent variable (argument), and the powdery mildew damage

is a dependent variable (function). After all, it is difficult to say that the degree of powdery mildew damage changes the cultivar of peaches and the shape of the glands on their leaves. Conversely, the assumption that the degree of damage to peach leaves by powdery mildew depends on the shape of the leaf glands is reasonable.

Spearman rank correlation coefficient and other non-parametric indicators are independent of the distribution law, and that is why they are very useful. They make it possible to measure the contingency between such attributes that cannot be directly measured, but can be expressed by points or other conventional units that allow ranking the sample. The benefit of rank correlation coefficient also lies in the fact that it allows to quickly assess the relationship between attributes regardless of the distribution law.

To determine the strength of association between two attributes, each of which consists of only two groups, *association coefficient and contingency coefficient are used*.

If there is a relationship between the variation of attributes, this means their association, or relationship. If the association was formed randomly, this means contingency. To evaluate association in this case, a number of indicators are used.

To calculate them, Table 15 is generated, which shows the association between two phenomena, each of which must be alternative, i.e. consisting of two different attribute values (for example, a product is good or defective).

**Table 15. For calculation of association coefficient and contingency coefficient**

| a | c | a+c |
|---|---|---|
| b | d | b+d |
| a+b | c+d | a+b+c+d |

The coefficients are calculated using the formulas:
Association coefficient:

$$K_a = \frac{ad-bc}{ad+bc} \quad (15)$$

Contingency coefficient:

$$K_k = \frac{ad-bc}{\sqrt{(a+b)\cdot(b+d)\cdot(a+c)\cdot(c+d)}} \quad (16)$$

Contingency coefficient is always less than association coefficient.

Association is considered confirmed if

$$K_a \geq 0.5 \text{ or } K_k \geq 0.3$$

**Example.** We study the association between the participation of the population of one of the cities in environmental actions and their level of education. The survey results are characterized by the following data (Table 16). Let's define: 1) association coefficient; 2) contingent coefficient.

**Solution.**
1. Example calculation of association coefficient

$$K_a = \frac{ad-bc}{ad+bc} = \frac{78\cdot68 - 22\cdot32}{78\cdot68 + 22\cdot32} = \frac{5304-704}{5304+704} = \frac{4600}{6008}$$
$$= 0.7656$$

**Table 16. Dependence of the participation of the city population in environmental actions on educational level**

| Groups of workers | Population of the city, persons | among them | |
|---|---|---|---|
| | | Participants in the actions, persons | Not participants in the actions, persons |
| With secondary education | 100 | 78 | 22 |
| Without secondary education | 100 | 32 | 68 |
| TOTAL: | 200 | 110 | 90 |

2. Example calculation of contingency coefficient

$$K_k = \frac{ad-bc}{\sqrt{(a+b)\cdot(b+d)\cdot(a+c)\cdot(c+d)}}$$

$$= \frac{78\cdot68 - 22\cdot32}{\sqrt{(78+22)\cdot(22+68)\cdot(78+32)\cdot(32+68)}}$$

$$= \frac{4600}{\sqrt{100\cdot90\cdot110\cdot100}} = \frac{4600}{\sqrt{99000000}} = \frac{4600}{9949.87437106}$$

$$= 0.4623$$

Thus, there is an association between the participation of the city population in environmental actions and its educational level.

When measuring the strength of association between qualitative alternative attributes and a continuously varying quantitative attribute, *biserial correlation coefficient* ($r_{bs}$) is used. The coefficient is calculated according to the formula:

$$r_{bs} = \frac{\bar{x}_1 - \bar{x}_2}{s} \cdot \sqrt{\frac{n_1 \cdot n_2}{N\cdot(N-1)}} \quad (17)$$

where $\bar{x}_1$ and $\bar{x}_2$ are the average values for alternative groups;
$s$ is the standard deviation;
$n_1$ and $n_2$ are sizes of alternative groups;
$N = (n_1 + n_2)$ is the total number of observations.

*Biserial correlation coefficient* varies from **–1** to **+1**; at $x_1 = x_2, r_{bs} = 0$. As for association coefficient, the sign of biserial coefficient has no meaning.

**Example.** We study the effect of tops affected by buck eye rot on the yield of "Priekulsky ranny" potato (Table 17). It is necessary to determine whether there is a correlation between potato yield and tops affected by buck eye rot.

**Table 17. Given data**

| Yield, kg per bush ($X$) | | 0.7 | 0.6 | 0.5 | 0.4 | 0.3 | 0.2 |
|---|---|---|---|---|---|---|---|
| Number of bushes, pcs. | total ($f$) | 12 | 15 | 18 | 13 | 9 | 6 |
| | incl. affected ($f_1$) | 0 | 4 | 9 | 10 | 7 | 6 |

**Solution.**
1. We generate calculation table (Table 18).
2. Let's calculate average values for alternative groups:

$$\bar{x}_1 = \frac{\sum_{i=1}^{n_1} f_{1i}X}{n_1} = \frac{14.2}{36} = 0.3944;$$

$$\bar{x}_2 = \frac{\sum_{i=1}^{n_2} f_{2i}X}{n_2} = \frac{21.3}{37} = 0.5757;$$

**Table 18. Calculation table**

| X | $f_1$ | $f_2$ | $f = f_1 + f_2$ | $f_1 X$ | $f_2 X$ | $fX$ | $X^2$ | $fX^2$ |
|------|------|------|------|------|------|------|------|------|
| 0,7 | 0 | 12 | 12 | 0 | 8.4 | 8.4 | 0.49 | 5.88 |
| 0,6 | 4 | 11 | 15 | 2.4 | 6.6 | 9.0 | 0.36 | 5.40 |
| 0,5 | 9 | 9 | 18 | 4.5 | 4.5 | 9.0 | 0.25 | 4.50 |
| 0,4 | 10 | 3 | 13 | 4.0 | 1.2 | 5.2 | 0.16 | 2.08 |
| 0,3 | 7 | 2 | 9 | 2.1 | 0.6 | 2.7 | 0.09 | 0.81 |
| 0,2 | 6 | 0 | 6 | 1.2 | 0 | 1.2 | 0.04 | 0.24 |
| **Sum** | **36** | **37** | **73** | **14.2** | **21.3** | **35.5** | **1.39** | **18.91** |

3. Let's calculate standard deviation:

$$s = \sqrt{\frac{\sum fX^2 - \frac{(\sum fX)^2}{N}}{N-1}} = \sqrt{\frac{18.91 - \frac{35.5^2}{73}}{73-1}}$$

$$= \sqrt{\frac{18.91 - 17.26}{72}} = \sqrt{0.0229} = 0.15$$

4. Let's calculate biserial correlation coefficient:

$$r_{bs} = \frac{\bar{x}_1 - \bar{x}_2}{s} \cdot \sqrt{\frac{n_1 \cdot n_2}{N \cdot (N-1)}} =$$

$$= \frac{0.3944 - 0.5757}{0.15} \cdot \sqrt{\frac{36 \cdot 37}{73 \cdot (73-1)}} = \frac{-0.1813}{0.15} \cdot \sqrt{\frac{1332}{5256}}$$

$$= \frac{-0.1813}{0.15} \cdot \sqrt{\frac{1332}{5256}} = -1.21 \cdot \sqrt{0.2534} = -0.6091$$

5. Let's calculate biserial correlation coefficient error:

$$s_{r_{bs}} = \sqrt{\frac{1 - r_{bs}^2}{N-2}} = \sqrt{\frac{1 - (-0.6091)^2}{73-2}} = \sqrt{\frac{1 - 0.3710}{71}} =$$

$$= \sqrt{0.0089} = 0.094$$

6. Let's calculate the criterion for the significance of biserial correlation coefficient:

$$t_r = \frac{r_{bs}}{s_{r_{bs}}} = \frac{-0.6091}{0.094} = -6.48$$

Since the sign of the criterion does not have any meaning, we discard it.

Using a statistical table or using the MS Excel TINV function, we find the value of Student's test at a 5% significance level and the number of degrees of freedom equal to $v = N - 2 = 73 - 2 = 71$. The critical (reference) value of Student's test is $t_{0.05} = 1.99$.

The criterion is greater than Student's test, therefore there is a significant correlation between the attributes.

**Conclusion:** With an increase in the incidence of buck eye rot on tops, the yield of "Priekulsky ranny" potatoes decreases significantly.

We presented a description of correlation coefficients and demonstrated examples of their application, so it is interesting to further discuss what grading scales exist for interpreting these coefficients.

Thus, we know that Pearson correlation coefficient is in the range from –1 to 1. The closer the resulting correlation coefficient to –1 or 1, the stronger the association between the studied indicators. When assessing the strength of association for correlation coefficients, various scales are used.

*Chaddock scale*
In 1925, the American statistician Robert Emmet Chaddock (1879–1940) introduced a scale for Pearson correlation coefficient [26]. This scale is the first gradation of correlation strength: 1) 0.1–0.3, poor association; 2) 0.3–0.5, fair association; 3) 0.5–0.7, good association; 4) 0.7–0.9, very good association.

*Cohen scale 1960–1988*
In the 1960s, statistician in the field of psychology and sociology Jacob Cohen (1923–1998, USA) proposed his "statistical power" scale for use in cases where the effects were small [27].

The power (of a test or research) is influenced by: 1) effect size, i. e. the degree of its manifestation; 2) the selected level of statistical significance (α, the probability of erroneously rejecting the null hypothesis; for us, usually at p <0.05); 3) size of sample from the general population [28,29].

According to Cohen scale, Pearson correlation coefficient has the following gradation: 1) 0.1, small association; 2) 0.3, medium association; 3) more than 0.5, large association.

Later, "Cohen's subjective standards" were brought to the logical form of ranges in very few sources [30, 31]: 1) 0.1–0.3, small association; 2) 0.3–0.5, medium association; 3) more than 0.5, large association.

However, in most sources, Cohen scale is quoted in its original form of three values.

*Rosenthal scale*
In the work by Rosenthal J. A. [32] published in 1996, Cohen scale was supplemented with a range of very strong association: 1) 0.1 (–0.1), weak association; 2) 0.3 (–0.3), moderate association; 3) 0.5 (–0.5), strong association; 4) 0.7 (–0.7), very strong association.

In modern publications, when using Cohen scale, Rosenthal gradation is used [33,34].

*Hinkle scale 1979–2003 (versions)*
Scale by D. E. Hinkle appears in publications dated 2011 to 2018 [35,36,37,38]. These publications contain references to monographs by Dennis E. Hinkle published by him in collaboration with other scientists [39,40] in the period of 1979–2003.

The following gradings are used in publications: 1) 0–0.3, little if any or negligible association; 2) 0.3–0.5, low association; 3) 0.5–0.7, moderate association; 4) 0.7–0.9, high or strong association; 5) 0.9–1.0, very high or very strong association.

A similar, but somewhat expanded scale is given on the website of Andrews University (USA, Michigan) [41]. To the

listed gradations, another association has been added: little association, <0.3. Thus, in [41] there are both 'Little' (<0.3) and 'Low' (0.3–0.5) correlation coefficient r values.

The manual "The Basic Practice of Statistics" [42] proposes the following gradation: 1) less than 0.3, very weak association; 2) 0.3–0.5, weak association; 3) 0.5–0.7, moderate association; 4) more than 0.7, strong association. Scale truncated at both ends by D. E. Hinkle et al. are presented in the manual "Statistics for Dummies" [43]: 1) 0.3–0.5, weak association; 2) 0.5–0.7, moderate association; 3) more than 0.7, strong association.

*Evans scale*

In 1996, the monograph by James D. Evans "Straightforward statistics for the behavioral sciences" [44] was published in the USA, in which another effect size scale was proposed. The scale is made by dividing the range of 0 to 1.0 into equal segments and does not provide for an insignificant correlation. This scale is used in publications (2012–2019) on psychology [35,45,46,47,48], programming [49], and a textbook on statistics [50]. The gradation of this scale is as follows: 1) 0–0.19, very weak association; 2) 0.2–0.39, weak association; 3) 0.40–0.59, moderate association; 4) 0.6–0.79, strong association; 5) 0.80–1.0, very strong association.

All given scales are used for grading Pearson correlation coefficient. To grade other coefficients (Spearman coefficient, Kendall coefficient, Cramér's coefficient, etc.), a search for publications in the ScienceDirect and PubMed systems gave the following information. The manual "Statistics without Maths for Psychology" [51] uses an original scale for grading Spearman correlation coefficient. The article [36] uses

Hinkle scale for Spearman correlation coefficient. The review article [52] presents the original grading scale for Spearman coefficient, Kendall coefficient, Phi coefficient, Cramer's V-coefficient, and concordance correlation coefficient (CCC).

The study [53] presents a detailed overview of the effect size grading for Hill yield criterion "strength of association" according to the correlation coefficient value parameter. Koterov et al. [53] analyzed 121 sources and collected information on 19 scales. They note that Chaddock scale from 1925 is not currently used abroad, but is widely represented in the countries of the former USSR. The most well-recognized grading scales for the correlation coefficient, to which there are many references, are Cohen scale, scale by D. E. Hinkle et al., Evans scale. Along with this, it is noted that there are a number of scales by other authors published once both in educational material (including on-line), in publications, and even in manuals or monographs. Quotations from such sources were rare, and in most cases simply absent.

**Conclusion**

In the third part of the article "Nonparametric Statistics", Spearman correlation coefficient, Kendall correlation coefficient, phi (Yule) correlation coefficient, Cramér's coefficient, Matthews coefficient, Fechner coefficient, Tschuprow coefficient, rank-biserial correlation coefficient, point-biserial correlation coefficient, as well as association coefficient and contingent coefficient were reviewed. Scales for grading the strength of association for correlation coefficients are given, both widely known and widely used, and those found in individual publications. Examples of calculating correlation coefficients and explanations are given.

## REFERENCES

1. Shmoilova, R.A., Minashkin, V.G., Sadjvnikova, N.A., Shuvalova, E.B. (2014). Theory of statistics. Moscow: Finance and Statistics, 2014. (In Russian)
2. Cuvier, G. (1805). Leçons d'anatomie comparée (volumes in-8). Paris: Baudouin. (In French)
3. Eliseeva, I.I., Yuzbashev, M.M. (2004). General theory of statistics. Moscow: Finance and Statistics, 2004. (In Russian)
4. Galton, F. (1888). Co-relations and their measurement, chiefly from Anthropometric Data. *Proceedings of the Royal Society of London*, 45, 135–145. https://doi.org/10.1098/rspl.1888.0082
5. Pearson, K. (1895). Notes on regression and inheritance in the case of Two Parents. *Proceedings of the Royal Society of London*, 58, 240–242. https://doi.org/10.1098/rspl.1895.0041
6. Kraemer, H.C. (2006). Correlation coefficients in medical research: from product moment correlation to the odds ratio. *Statistical Methods in Medical Research*, 15(6), 525–544. https://doi.org/10.1177/0962280206070650
7. Bavrina, A.P., Borisov, I.B. (2021). Modern rules of the application of correlation analysis. *Medicinskij al'manah*, 3(68), 70–79. (In Russian)
8. Grjibovski, A.M. (2008). Correlation analysis. *Human Ecology*, 9, 50–60. (In Russian)
9. Leach, C. (1979). Introduction to statistics: a nonparametric approach for the social sciences. Chichester: Wiley, 1979.
10. Noether, G.E. (1981). Why Kendall Tau? *Teaching Statistics*, 3(2), 41–43. https://doi.org/10.1111/j.1467-9639.1981.tb00422.x
11. Bonett, D.G., Wright, T.A. (2000). Sample size requirements for estimating Pearson, Kendall and Spearman correlations. *Psychometrica*, 65, 23–28. https://doi.org/10.1007/BF02294183
12. Kendall, M.G. (1938). A new measure of rank correlation. *Biometrika*, 30(1–2), 81–93. https://doi.org/10.1093/biomet/30.1-2.81
13. Gubler, E.V., Genkin, A.A. (1973). Application of nonparametric statistical criteria in biomedical research. Leningrad: Medicine, 1973. (In Russian)
14. Rosenbaum, S. (1954). Tables for a nonparametric test of location. *Annals of Mathematical Statistics*, 25(1), 146–150. https://doi.org/10.1214/aoms/1177728854
15. Stepanov, V.G. (2019). Application of nonparametric statistical methods in agricultural biology and veterinary medicine research. St-Petersburg: Lan. 2019. (In Russian)
16. Edelbaeva, N.A., Lebedinskaya, O.G., Kovanova, E.S., Tenetova, E.P., Timofeev, A.G. (2019). Fundamentals of nonparametric statistics. Moscow: YUNITY-DANA. 2019. (In Russian)
17. Klyachkin, V.N. (2021). Expert methods. Chapter in a book: Statistical methods in quality management: computer technologies. Moscow: Finance and Statistics, 2021. (In Russian)
18. Matthews, B.W. (1975). Comparison of the predicted and observed secondary structure of T4 phage lysozyme. *Biochimica et Biophysica Acta (BBA) — Protein Structure*, 405(2), 442–451. https://doi.org/10.1016/0005-2795(75)90109-9

19. Cramer, H. (1946). Variables and Distributions. Chapters in a book: Mathematical Methods of Statistics. Princeton: Princeton University Press, 1946.

20. Yule, G.U. (1912). On the methods of measuring association between two attributes. *Journal of the Royal Statistical Society*, 75(6), 579–652. https://doi.org/10.2307/23401263

21. Aaron, B., Kromrey, J.D., Ferron, J.M. (1998, 2–4 November). *Equating r-based and d-based effect-size indices: Problems with a commonly recommended formula*. Paper presented at the annual meeting of the Florida Educational Research Association. Orlando, FL. Retrieved from https://files.eric.ed.gov/fulltext/ED433353.pdf Accessed June 20, 2023.

22. Boughorbel, S., Jarray, F., El-Anbari, M. (2017). Optimal classifier for imbalanced data using Matthews coefficient metric. *PLOS ONE*, 12(6), Article e0177678. https://doi.org/10.1371/journal.pone.0177678

23. Cramer, H. (1946). The two-dimensional case. Chapter in a book: Mathematical Methods of Statistics. Princeton: Princeton University Press, 1946.

24. Fechner, G.T. (1897). Kollektivmasslehre. Leipzig: Verlag von Wilhelm Engelmann, 1897. (In German)

25. Lakin, G.F. (1990). Biometrics. Moscow: High School. 1990. (In Russian)

26. Chaddock, R.E. (1925). Principles and methods of statistics. Boston, New York, 1925.

27. Cohen, J. (1973). Statistical power analysis and research results. *American Educational Research Journal*, 10(3), 225–229. https://doi.org/10.2307/1161884

28. Cohen, J. (1988). Statistical power analysis for the behavioral sciences. Hillsdale. Mahwah, NJ: Lawrence Erlbaum Associates, 1988.

29. Cohen, J. (1992). A power primer. *Psychological Bulletin*, 112(1), 155–159. https://doi.org/10.1037/0033-2909.112.1.155

30. Divaris, K., Vann, W.F., Baker, A.D., Lee, J.Y. (2012). Examining the accuracy of caregivers' assessments of young children's oral health status. *The Journal of the American Dental Association*, 143(11), 1237–1247. https://doi.org/10.14219/jada.archive.2012.0071

31. Lomax, R.G., Hahs-Vaughn, D.L. (2012). Statistical Concepts. A Second Course. New-York: Taylor and Francis Group, LLC, 2012.

32. Rosenthal, J.A. (1996). Qualitative descriptors of strength of association and effect size. *Journal of Social Service Research*, 21(4), 37–59. https://doi.org/10.1300/j079v21n04_02

33. Berry, K.J., Johnston, J.E., Mielke, P.W. (2018). The Measurement of Association. A Permutation Statistical Approach. Cham: Springer Nature Switzerland AG, 2018. https://doi.org/10.1007/978-3-319-98926-6

34. de Menezes, R.F., Bergmann, A., Thuler, L.C.S. (2013). Alcohol consumption and risk of cancer: a systematic literature review. *Asian Pacific Journal of Cancer Prevention*, 14(9), 4965–4972. https://doi.org/10.7314/apjcp.2013.14.9.4965

35. Bourne, P.A., Hudson-Davis, A. (2016). Psychiatric induced births in Jamaica: homicide and death effects on pregnancy. *Psychology and Behavioral Science International Journal*, 1(2), Article 555558. https://doi.org/10.19080/PBSIJ.2016.01.555558

36. Mukaka, M.M. (2012). Statistics Corner: A guide to appropriate use of correlation coefficient in medical research. *Malawi Medical Journal*, 24(3), 69–71.

37. Schober, P., Boer, C., Schwarte, L.A. (2018). Correlation coefficients: appropriate use and interpretation. *Anesthesia and Analgesia*, 126(5), 1763–1768. https://doi.org/10.1213/ANE.0000000000002864

38. Kotrlik, J.W., Williams, H.A., Jabor, M.K. (2011). Reporting and interpreting effect size in quantitative agricultural education research. *Journal of Agricultural Education*, 52(1), 132–142. https://doi.org/10.5032/jae.2011.01132

39. Hinkle, D.E., Wiersma, W., Jurs, S.G. (1979). Applied Statistics for the Behavioral Sciences. Chicago: Rand McNally College Pub. Co., 1979.

40. Hinkle, D.E., Wiersma, W., Jurs, S.G. (2003). Applied Statistics for the Behavioral Sciences. 5th Ed. Boston: Houghton Mifflin, 2003.

41. *Correlation Coefficients. Applied Statistics. Lesson 5.* (2005). Andrews University (Michigan). Retrieved from https://www.andrews.edu/~calkins/math/edrm611/edrm05.htm. Accessed June 20, 2023.

42. Moore, D., Notz, W.I., Fligher, M.A. (2012). The Basic Practice of Statistics. Publisher: W. H. Freeman, 2012.

43. Rumsey, D.J. (2016). Statistics For Dummies. 2nd Ed. New York: For Dummies, 2016.

44. Evans, J.D. (1996). Straightforward statistics for the behavioral sciences. Pacific Grove, Calif.: Brooks/Cole Publ. Co: An International Thomson Publ. Co., 1996.

45. Neill, J. (2018). *Survey research and design in psychology. Lecture 4*. Retrieved from https://upload.wikimedia.org/wikiversity/en/f/fd/SRDP_Lecture04Handout_Correlation_6slidesperpage.pdf. Accessed June 20, 2023.

46. Yavna, D.V., Kupriyanov, I.V., Bunyaeva, M.V. (2016). Sensory and perceptual processes: a tutorial. Rostov-on-Don: Publishing House of the Southern Federal University. 140. (In Russian)

47. Chakkapark, J, Vinitwatanakun, W. (2017). The relationship between division heads' leadership styles and teacher satisfaction at Siam Commercial College of Technology. *Scholar: Hum Sciences*, 9(1), 36–47.

48. Gerguri, D. (2018). Leader-staff relationships in Kosovo customs: leadership and its impact on customs effectiveness. *Styles of Communication*, 10(1), 108–124.

49. Miletic, M., Vukusic, M., Mausa, G., Galinac, T. (2017, 11–13 September). *Relationship between design and defects for software in evolution*. Proceedings of the Workshop of Software Quality, Analysis, Monitoring, Improvement, and Applications. Belgrade, Serbia, 2017.

50. Pearson's correlation. (2011). Site Statstutor. Statistics support for students. UK. Retrieved from www.statstutor.ac.uk/resources/uploaded/pearsons.pdf. Accessed June 20, 2023.

51. Dancey, C.P., Reidy, J. (2007). Statistics without Maths for Psychology. Harlow: Pearson Education Limited, 2007.

52. Akoglu, H. (2018). User's guide to correlation coefficients. *Turkish Journal of Emergency Medicine*, 18(3), 91–93. https://doi.org/10.1016/j.tjem.2018.08.001

53. Koterov, A.N., Ushenkova, L.N., Zubenkova, E.S., Kalinina, M.V., Biryukov, A.P., Lastochkina, E.M. et al. (2019). Strength of association. Report 2. Graduations of correlation size. *Medical Radiology and Radiation Safety*, 64(6), 12–24. https://doi.org/10.12737/1024-6177-2019-64-6-12-24

## AUTHOR INFORMATION

**Marina A. Nikitina,** Doctor of Technical Sciences, Docent, Leading Scientific Worker, the Head of the Direction of Information Technologies of the Center of Economic and Analytical Research and Information Technologies, V. M. Gorbatov Federal Research Center for Food Systems. 26, Talalikhina, 109316, Moscow, Russia. Tel: +7–495–676–95–11 (297), E-mail: m.nikitina@fncps.ru
ORCID: https://orcid.org/0000–0002–8313–4105
* corresponding author

**Irina M. Chernukha,** Doctor of Technical Sciences, Professor, Academician of the Russian Academy of Sciences, Head of the Department for Coordination of Initiative and International Projects, V. M. Gorbatov Federal Research Center for Food Systems. 26, Talalikhina, 109316, Moscow, Russia. Tel: +7–495–676–95–11 (109), E-mail: imcher@inbox.ru
ORCID: https://orcid.org/0000–0003–4298–0927

All authors bear responsibility for the work and presented data.

All authors made an equal contribution to the work.

The authors were equally involved in writing the manuscript and bear the equal responsibility for plagiarism.

The authors declare no conflict of interest.