

## **Permutation methods in association analysis**

### **1. Introduction**

The permutation model in hypothesis testing was introduced by R. A. Fisher in 1925. These methods permit us to test hypotheses with as minimal assumptions as possible. These tests require high computing power and therefore they have found the applications in last years. However, the concept of permutation methods is much wider than the issue of permutation testing. In 1923 J. Spława-Neyman introduced a permutation model for the analysis of field experiments [Spława-Neyman 1923, Berry et al. 2014]. This paper was published in Polish language. Neyman's work was translated into English in 1990 and published in „*Statistical Science*” [Ledwina 2012]. Then this article reached the international statistician community [Lehmann and Romano 2005, p. 210]. This paper was recognized as a pioneering achievement in the field of statistical methodology for the analysis of causal relationships.

The aim of the article is to present the possibilities of applying the permutation methods in the analysis of dependencies. The article presents selected possibilities of applying the permutation methods in the analysis of association between variables. The main attention is paid to various possibilities of using permutation methods and possibilities of data labels rearranging. A proposal of a method allowing to investigate the significance of the relationship between two data sets is presented. The presented example illustrate the use of permutation methods to testing significance of the relationship between two sets of variables. The proposed method is compared to well-known canonical correlation analysis.

### **2. Permutation methods**

In mathematics, the term “permutation” is used as the act of rearranging in an ordered fashion objects or values. In statistics the idea of permutation objects is used in several methods, especially in testing hypotheses (permutation tests).

Permutation tests permit us to choose the form of the test statistic. Through sample size reduction, permutation tests can reduce the costs of experiments and surveys. Permutation tests are the most powerful of statistical procedures. There are 5 steps in the process of permutation testing [Good 2005]

1. Identify the null and the alternative hypothesis.
2. Choose the form of the test statistic.
3. Calculate the test statistic for the sample data.
4. Determine the frequency distribution of the test statistic using data permutations..
5. Make a decision using this empirical distribution as a guide.

The basic idea in permutation testing is to generate a reference distribution by recalculating a test statistic for many permutations of the data. The term “permutation methods” should not, however, be limited to the problem of testing hypotheses. In recent years, various statistical methods have been proposed that refer to the idea of the permutation model introduced by J. Sława-Neyman [1923].

L. Corain et al. [2016] present a novel permutation-based nonparametric approach for ranking several multivariate populations. Using data collected from both experimental and observation studies, it covers some of the most useful designs widely applied in research and industry investigations, such as the multivariate analysis of variance and multivariate randomized complete block designs.

K. J. Berry et al. [2018] use rearranging data to generate probability values and measures of effect size for various measures of association. They define association for two interval-level variables, measures of association for two nominal-level variables or two ordinal-level variables, and measures of agreement for two nominal-level or two ordinal-level variables.

K. J. Berry et al [2016] provide a synthesis of a number of statistical tests and measures, which, at first consideration, appear disjoint and unrelated. Numerous comparisons of permutation and classical statistical methods are presented, and the two class are compared via probability values and, where appropriate, measures of effect size.

P. W. Mielke and K. J. Berry Jr. [2007] provide a wide treatment of statistical inference using permutation techniques. Its purpose is to make available to practitioners

a variety of useful and powerful data analytical tools that rely on very few distributional assumptions. Although many of these procedures have appeared in journal articles, they are not readily available to practitioners.

### **3. Permutational methods in hypothesis testing**

Testing of statistical hypotheses is a major branch of study in classical statistical inference. On the basis of a relatively small sample, one can infer about the characteristics of a much larger population. There are two basic kinds of statistical hypotheses: parametric and non-parametric. According to the parametric or non-parametric hypothesis, a parametric statistical test or a nonparametric test is used. Parametric tests require that the sample is taken from a specified distribution, usually the normal distribution. Permutation tests such as nonparametric tests do not require specific population distributions of the variables such as the normal distribution. These tests use data labels rearranging. A typical application of permutation tests is to compare the distributions of two or more populations based on two samples taken independently.

#### **3.1. Methods of data permutations**

The permutation model is based on the data labels permutations. T. W. O’Gorman [2012, p. 78] lists some permutations methods which can be used for testing the significance of the coefficient in the linear regression model. In the linear regression model there is one dependent variable  $Y$  and  $k$  independent variables  $X_1, X_2, \dots, X_k$ . Some methods of rearranging data labels include:

- permute the dependent variable,
- permute the independent variable for the considered variable,
- permute the residuals from the reduced model,
- permute the residuals from the complete model.

To test the significance of a parameter in the linear regression model various methods of rearranging data labels can be used. The results of testing the significance of the coefficient for different permutation methods are not equivalent. T. W. O’Gorman [2012] points that it is not clear which method is superior.

Let us assume that  $\mathbf{y} = [y_1, y_2, \dots, y_k]^T$  and we want to shuffle this vector. To permute this vector we can use the square matrix of dimension  $k \times k$   $\mathbf{P}_k = [a_{ij}]$  where  $a_{ij}$  are only zeros and ones for  $i, j = 1, 2, \dots, k$  and  $\sum_{i=1}^k a_{ij} = 1$  for  $j = 1, 2, \dots, k$  and  $\sum_{j=1}^k a_{ij} = 1$  for  $i = 1, 2, \dots, k$ .

The examples of permutation matrix in the case  $k = 4$  are:

$$\mathbf{P}_1 = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, \mathbf{P}_2 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}, \mathbf{P}_3 = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \end{bmatrix}$$

It is easy to notice that for the first matrix we have  $\mathbf{P}_1 \mathbf{y} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} = \begin{bmatrix} y_2 \\ y_3 \\ y_1 \\ y_4 \end{bmatrix}$

So vector  $[y_2, y_3, y_1, y_4]^T$  is a rearrangement of vector  $[y_1, y_2, y_3, y_4]^T$

For the above given permutation matrices  $\mathbf{P}_1$ ,  $\mathbf{P}_2$  and  $\mathbf{P}_3$  we have

$$\mathbf{P}_1 \mathbf{y} = [y_2, y_3, y_1, y_4]^T$$

$$\mathbf{P}_2 \mathbf{y} = [y_1, y_4, y_2, y_3]^T$$

$$\mathbf{P}_3 \mathbf{y} = [y_2, y_3, y_4, y_1]^T$$

### 3.2. Properties of permutations

Let  $\mathbf{A}$  be the vector given by  $\mathbf{A} = [y_1, y_2, \dots, y_k]^T$ . The permutation of vector  $\mathbf{A}$  is function  $\alpha : \mathbf{A} \rightarrow \mathbf{A}$  that is bijective (i.e. both one-to one and onto).

Identity permutation. Let  $\mathbf{D} = [d_{ij}]$  be the identity matrix such  $d_{ii}=1$ , and  $d_{ij} = 0$  for  $i \neq j$  then matrix  $\mathbf{D}$  leads to the identity permutation. If  $\mathbf{D}$  is the matrix of permutation  $\alpha$ , then for each vector  $\mathbf{A}$ :  $\alpha \mathbf{A} = \mathbf{A}$ .

Composition of permutations. Let  $\mathbf{S}_1$  and  $\mathbf{S}_2$  be two matrices of permutations  $\alpha$  and  $\beta$ , then  $\mathbf{S}_2 \mathbf{S}_1$  is the matrix of composition of permutations  $\gamma = \beta \alpha$

Inverse permutation. Let  $\mathbf{S}$  be the matrix of permutation  $\alpha$ , then  $\mathbf{S}^T$  is the matrix of inverse permutation  $\beta$ , that  $\beta \alpha \mathbf{A} = \alpha \beta \mathbf{A} = \mathbf{A}$  for each vector  $\mathbf{A}$ .

Let  $S_1$ ,  $S_2$  and  $S_3$  be the matrices of permutation  $\alpha$ ,  $\beta$  and  $\gamma$ . Permutation composition is associative, eg.  $S_3 S_2 S_1$  is the matrix of permutation  $\gamma(\beta\alpha) = (\gamma\beta)\alpha$

#### 4. Methods of data labels permutation in correlation and multiple regression models

Pearson correlation coefficient  $\rho$ , tells us about the strength of the linear relationship between two variables  $X$  and  $Y$ . To perform a test of the significance of the correlation coefficient a random sample of size  $n$  should be taken. The sample data is used to compute  $r$ , the correlation coefficient for the sample. Permutation tests use rearranging data labels to obtain the empirical distribution of the coefficient. The data should be permuted  $N$  times ( $N$  usually is greater or equals 1000). The empirical distribution of  $r$  is obtained from  $N$  permuted sets of data. The decision of hypothesis  $H_0$  is based on the value of  $r$  and its empirical distribution. The original two-dimensional data and three examples of random permuted data are presented in fig. 1.

X	Y	X	Y	X	Y	X	Y
x1	y1	x1	y2	x1	y6	x1	y5
x2	y2	x2	y1	x2	y5	x2	y1
x3	y3	x3	y9	x3	y8	x3	y9
x4	y4	x4	y6	x4	y7	x4	y7
x5	y5	x5	y10	x5	y9	x5	y6
x6	y6	x6	y4	x6	y4	x6	y8
x7	y7	x7	y7	x7	y1	x7	y3
x8	y8	x8	y8	x8	y10	x8	y2
x9	y9	x9	y5	x9	y2	x9	y4
x10	y10	x10	y3	x10	y3	x10	y10

Fig. 1. Example of a two-dimensional data set  $(X, Y)$  and three random data sets  $(X, Y)$  with randomly rearranged variable  $Y$

Source: author's own elaboration.

There are more possibilities of rearranging data labels in the case of a multiple regression model than in the case of testing the significance of the correlation coefficient. Let us consider the multiple regression linear model

$$Y = a_1 X_1 + a_2 X_2 + \dots + a_k X_k + e \quad (1)$$

There are many ways of data labels rearranging for one dependent variable  $Y$  and  $k$  independent variables  $X_1, X_2, \dots, X_k$ . There are some examples of possible ways of permutation of one dependent variable and  $k$  ( $k > 2$ ) independent variables:

- permute the dependent variable  $Y$  (it is the same result as permute all independent variables  $X_1, X_2, \dots, X_k$  in the same way),
- permute only one independent variable  $X_i$  ( $i = 1, 2, \dots, k$ )
- permute two independent variables  $X_i$  and  $X_j$  ( $i, j = 1, 2, \dots, k, i \neq j$ ) the same way
- permute two independent variables  $X_i$  and  $X_j$  ( $i, j = 1, 2, \dots, k, i \neq j$ ) independently
- permute two independent variables  $X_i$  and  $X_j$  ( $i, j = 1, 2, \dots, k, i \neq j$ ) in the same way and the variable  $X_s$  ( $s = 1, 2, \dots, k$ ) independently

Labels of original data should be rearranged for testing the significance of the parameters of the linear model. The empirical distribution of the parameter is obtained on the basis of the results for the  $N$  permuted models. Let  $S, S_1, S_2, \dots, S_k$  are the permutation matrices. Some typical methods of permutation for the multiple regression model (1) for  $k = 3$  are:

$$Y = a_1SX_1 + a_2SX_2 + \dots + a_kSX_k + e$$

$$Y = a_1S_1X_1 + a_2XS_2X_2 + \dots + a_kS_kX_k + e$$

$$Y = a_1X_1 + a_2SX_2 + \dots + a_kSX_k + e$$

## 5. The significance of the association of two sets of variables

In addition to the methods of studying the relationship between two variables or a dependent variable and many dependent variables, there are statistical methods that measure the association between two sets of variables  $\mathbf{X}$  and  $\mathbf{Y}$ . One of these methods is the canonical correlations analysis. Let  $\mathbf{X} = (X_1, \dots, X_p)$  and  $\mathbf{Y} = (Y_1, \dots, Y_q)$  be two sets of variables. To determine the strength of the association the permutation method could be used. The use of permutation methods requires data labels to be rearranged. These two sets can be rearranged in many ways. The variables in one set or in two sets can be permuted independently or dependently. The variables in one set (or two sets) could be grouped and permuted within these groups independently or dependently. The original set of variables for  $p = q = 3$  and two examples of permuted sets is presented in fig. 2. In

the first case (on the left) the variables  $X_1$ ,  $X_2$  and  $X_3$  are permuted dependently. In the second case (on the right) all variables are permuted independently.

Y1	Y2	Y3	X1	X2	X3
y11	y21	y31	x11	x21	x31
y12	y22	y32	x12	x22	x32
y13	y23	y33	x13	x23	x33
y14	y24	y34	x14	x24	x34
y15	y25	y35	x15	x25	x35
y16	y26	y36	x16	x26	x36
y17	y27	y37	x17	x27	x37
y18	y28	y38	x18	x28	x38
y19	y29	y39	x19	x29	x39

Y1	Y2	Y3	X1	X2	X3
y11	y21	y31	x12	x22	x32
y12	y22	y32	x14	x24	x34
y13	y23	y33	x16	x26	x36
y14	y24	y34	x15	x25	x35
y15	y25	y35	x11	x21	x31
y16	y26	y36	x18	x28	x38
y17	y27	y37	x19	x29	x39
y18	y28	y38	x13	x23	x33
y19	y29	y39	x17	x27	x37

Y1	Y2	Y3	X1	X2	X3
y11	y23	y35	x12	x21	x31
y12	y22	y34	x14	x27	x32
y13	y21	y39	x13	x28	x36
y14	y25	y37	x18	x24	x33
y15	y24	y31	x15	x26	x34
y16	y26	y33	x17	x22	x38
y17	y29	y36	x16	x29	x39
y18	y28	y32	x11	x25	x37
y19	y27	y38	x19	x23	x35

Fig. 2. The original set of two sets of variables (upper) and two examples of possible permutations of these sets

Source: author's own elaboration.

## 6. Canonical correlations

Pearson correlation coefficient measures the strength of the linear dependency for two variables  $X$  and  $Y$ . The multiple linear regression model could be used to describe the dependency between the dependent variable  $Y$  and a set of independent variables  $X_1, X_2, \dots, X_k$ . Sometimes the association between the two sets of variables  $\mathbf{X}$  and  $\mathbf{Y}$  should be considered. Canonical correlation analysis could be used in this case. This method was proposed by H. Hotelling in 1935-36. Canonical correlation analysis is employed to study the relationships between two variable sets when each variable set consists of at least two variables. The main objectives of canonical correlations are as follows [Thompson 1984]:

- determining the strength of the relationships that may exist between the two sets,
- deriving a set of weights for each set of dependent and independent variables so that the linear combinations of each set are maximally correlated,
- explaining the nature of any relationships existing between the sets of dependent and independent variables.

Canonical correlation analysis develops a number of independent canonical functions that maximize the correlation between the linear composites, also known as canonical variates, which are sets of dependent and independent variables.

The form of the Pearson correlation coefficient:

$$\rho = \frac{Cov(X, Y)}{\sqrt{Var(X)Var(Y)}} \quad (2)$$

The form of canonical correlation coefficient:

$$\rho = \max_{a \in R^p, b \in R^q} \frac{Cov(a'X, b'Y)}{\sqrt{Var(a'X)Var(b'Y)}} \quad (3)$$

Canonical correlation analysis constructs such vectors  $\mathbf{a}$  and  $\mathbf{b}$  based on the following criteria:

1. The first canonical variate  $\mathbf{U}_1 = \mathbf{a}'_1\mathbf{X}, \mathbf{V}_1 = \mathbf{b}'_1\mathbf{Y}$  is constructed from the maximization of (3).
2. The second canonical variate pair  $\mathbf{U}_2 = \mathbf{a}'_2\mathbf{X}, \mathbf{V}_2 = \mathbf{b}'_2\mathbf{Y}$  is constructed from the maximization of (3) with restriction that  $\mathbf{D}^2(\mathbf{U}_2) = \mathbf{D}^2(\mathbf{V}_2) = 1$  and pairs  $(\mathbf{U}_1, \mathbf{V}_1)$  and  $(\mathbf{U}_2, \mathbf{V}_2)$  are uncorrelated.
3. At the  $k$  step, the  $k$ -th canonical variate pair  $\mathbf{U}_k = \mathbf{a}'_k\mathbf{X}, \mathbf{V}_k = \mathbf{b}'_k\mathbf{Y}$  - is obtained from the maximization of (3) with restriction that  $\mathbf{D}^2(\mathbf{U}_k) = \mathbf{D}^2(\mathbf{V}_k) = 1$  and  $(\mathbf{U}_k, \mathbf{V}_k)$  are uncorrelated with the previous  $(k-1)$  canonical variate pairs.
4. Repeat step 3 until the number of canonical variates  $s = \min(p, q)$ .



The first canonical correlation coefficient is denoted by  $r_1$  the second by  $r_2$  and so on. The number of calculated CCA coefficients is equal to the minimum of the number of variables in two considered sets  $s = \min(p, q)$ . For testing the significance of the first canonical correlation coefficient, the following statistics can be used: Wilks' Lambda, Hotelling-Lawley Trace, Pillai-Bartlett Trace, Roy's Largest Root.

## 7. Testing the significance of the association of two sets of variables

One problem to be considered in the canonical correlation analysis is to test the hypothesis that none of the canonical correlation coefficients  $r_1, r_2, \dots, r_s$  is significant. To test such hypothesis we usually use lambda Wilk's statistic, Lawley-Hotelling's trace, Pillai trace or Roy's largest root. The form of these statistics is as follows

Lambda Wilk's statistic:

$$\Lambda_1 = \prod_{i=1}^m (1 - r_i^2) \quad (4)$$

This statistic is distributed as the Wilks  $\Lambda$ -distribution. Rejection of the null hypothesis is for small values of  $\Lambda_1$ .

Pillai trace statistic:

$$V^{(m)} = \sum_{i=1}^m r_i^2 \quad (5)$$

E. L. Rencher and J. P. Christensen [2012, p. 391–395] has tables providing critical values for this statistic.

Lawley-Hotelling's statistic:

$$U^{(m)} = \sum_{i=1}^m \frac{r_i^2}{1 - r_i^2} \quad (6)$$

E. L. Rencher and J. P. Christensen [2012, p. 391–395] has tables providing critical values for this statistic.

Roy's largest root statistic:

$$\theta = r_1^2 \quad (7)$$

E. L. Rencher and J. P. Christensen [2012, p. 391–395] has tables providing critical values for Roy's largest root statistic. This statistic is based only on one largest canonical

correlation coefficient and the other statistics use all canonical correlation coefficients. For this reason, Roy's largest root statistic will not be included in the computer simulation studies.

## 8. Monte Carlo study

The proposed method of testing the significance of dependency between two sets of variables is based on the permutation method. In the proposed method lambda Wilk's statistic (4) is used and the dependent rearrangement of one of the variables set labels is used (see fig. 2). The proposal was compared to the three well known methods of testing the significance of correlations in the canonical analysis in Monte Carlo study. The permutation had  $N = 1000$  random rearrangements of data labels. The significance level  $\alpha = 0.05$  was assumed in computer simulations.

In the first part of a computer simulation two sets of variables  $\mathbf{X}$  and  $\mathbf{Y}$  are considered. All variables were normally distributed. The sample size was  $n = 100$ . The above assumptions are typical for the canonical correlation analysis.

The first set of variables is given by  $\mathbf{X} = (X_1, X_2, X_3)$  where variables  $X_1, X_2, X_3$  are independent and normally distributed with mean 10 and variance 1. The second set of variables is given by  $\mathbf{Y} = (1 - \delta)\mathbf{Z} + \delta\mathbf{X}$  where  $\mathbf{Z} = (Z_1, Z_2, Z_3)$  and variables  $Z_1, Z_2, Z_3$  are independent and normally distributed with mean 10 and variance 1. The parameter  $\delta$  ( $\delta = 0.00, 0.05, \dots, 0.25$ ) describes the strength of association of two sets of variables. For  $\delta = 0$  the sets  $\mathbf{X}$  and  $\mathbf{Y}$  are not associated.

The estimated probabilities of rejection  $H_0$  (there is no association between  $\mathbf{X}$  and  $\mathbf{Y}$ ) are presented in table 1. The probabilities were estimated based on 1000 random samples. In the case of normally distributed variables with the sample of size  $n = 100$  the size of all considered tests is close to the significance level  $\alpha$ . The power of all test is quite similar. There is no advantage of the proposal in this case.

Table 1. The estimated probability of rejection  $H_0$  – normal distribution,  $n = 100$

$\delta$	Method			
	Permutation	Lambda Wilk's	Lawley-Hotelling's	Pillai trace
0.00	0.051	0.053	0.054	0.049
0.05	0.081	0.080	0.082	0.075
0.10	0.199	0.202	0.204	0.193
0.15	0.497	0.495	0.499	0.487
0.20	0.866	0.861	0.861	0.862
0.25	0.993	0.991	0.990	0.989

Source: author's own elaboration.

The advantages of permutation tests are for the applications of small samples and non-normally distributed variables. The next three simulations are based on normally, beta and gamma distributed variables for small samples where  $n = 20$ .

The association of two sets  $\mathbf{X}$  and  $\mathbf{Y}$  of the three dimensional variables were analysed in this part of the study. The considered sets are  $\mathbf{X} = (X_1, X_2, X_3)$  and  $\mathbf{Y} = (1 - \delta)\mathbf{Z} + \delta\mathbf{X}$  where  $\mathbf{Z} = (Z_1, Z_2, Z_3)$ , variables  $Z_1, Z_2, Z_3$  are independently distributed and the parameter  $\delta$  ( $\delta = 0.00, 0.05, \dots, 0.40$ ) describes the strength of association of two sets.

The parameters of random variables  $X_1, X_2, X_3$  and  $Z_1, Z_2, Z_3$  are as follows

- Normal distribution with expected value 10 and variance 1.
- Beta distribution with shape parameters  $s_1 = 2$  and  $s_2 = 2$ .
- Gamma distribution with shape parameter  $s = 2$ .

Table 2. The estimated probability of rejection  $H_0$  – normal distribution,  $n = 20$

$\delta$	Method			
	Permutation	Lambda Wilk's	Lawley-Hotelling's	Pillai trace
0.00	0.046	0.041	0.052	0.033
0.05	0.069	0.064	0.078	0.047
0,10	0.070	0.074	0.087	0.046
0.15	0.088	0.090	0.110	0.056
0,20	0.186	0.173	0.187	0.144
0.25	0.285	0.263	0.278	0.226
0,30	0.496	0.468	0.479	0.438
0.35	0.726	0.688	0.682	0.664
0.40	0.907	0.873	0.858	0.859

Source: author's own elaboration.

The probabilities of rejection of  $H_0$  are presented in tables 2 – 4. The probabilities were estimated based on 1000 random samples. In the permutation method for each sample the data labels were randomly rearranged  $N = 1000$  times. In the first case samples were taken from normal distribution but the sample size was  $n = 20$ . The results are presented in table 2. In the case of normally distributed variables and small samples (table 2) the test based on permutation method has the greatest power. The three other tests have similar powers. The results of computer study for beta distributed variables are presented in table 3.

Table 3. The estimated probability of rejection  $H_0$  – beta distribution,  $n = 20$

$\delta$	Method			
	Permutation	Lambda Wilk's	Lawley-Hotelling's	Pillai trace
0.00	0.056	0.058	0.074	0.042
0.05	0.057	0.055	0.071	0.041
0,10	0.069	0.062	0.073	0.045
0.15	0.101	0.095	0.111	0.066
0.20	0.155	0.149	0.172	0.114
0.25	0.252	0.242	0.253	0.202
0,30	0.495	0.449	0.457	0.411
0.35	0.701	0.653	0.649	0.639
0.40	0.904	0.887	0.870	0.883

Source: author's own elaboration.

The size of Lawley-Hotelling's test is greater than the significance level  $\alpha$  in the case of beta distributed variables. The most powerful test is the test based on permutation method. The results of computer study for gamma distributed variables are presented in table 4.

Table 4. The estimated probability of rejection  $H_0$  – gamma distribution,  $n = 20$

$\delta$	Method			
	Permutation	Lambda Wilk's	Lawley-Hotelling's	Pillai trace
0.00	0.055	0.061	0.068	0.039
0.05	0.059	0.074	0.083	0.043
0,10	0.075	0.080	0.091	0.057
0.15	0.122	0.127	0.141	0.089
0.20	0.170	0.173	0.193	0.136
0.25	0.313	0.291	0.316	0.258
0,30	0.498	0.474	0.480	0.436
0.35	0.725	0.701	0.695	0.681
0.40	0.903	0.887	0.876	0.879

Source: author's own elaboration.

The size of Lawley-Hotelling's test is greater than the significance level  $\alpha$  in the case of gamma distributed variables. The most powerful test is the test based on permutation method.

## 9. Conclusions

Permutation tests are one type of permutation methods. They are very useful for non-normally distributed data and for small samples. However, permutation methods can be used not only for testing hypotheses. They can be used for constructing rankings of multivariate data or measuring the association between variables and sets of variables. In these methods the form of data labels rearrangement is very important. The data labels can be permuted in many ways. Some of these methods of labels rearranging which can be used in the association analysis are described in the paper.

As an example of the use of the permutation method in the association analysis the method of testing the significance of association of two sets of variables is presented in the paper. This method use the lambda Wilk's statistic and is based on the dependent rearranging data labels in variables in one of the considered sets. The properties of the proposal were compared to well-known tests in canonical correlation analysis: lambda Wilk's, Lawley-Hotelling's and Pillai trace. Monte Carlo study shows that for the normally distributed variables and for big samples the power of all tests is quite similar. The proposal is most powerful for non-normally distributed variables and for small samples.

## Bibliography

- Berry K. J., Johnston J. E., Mielke Jr. P. W. [2014] *A Chronicle of Permutation Statistical Methods*, Springer International Publishing, New York.
- Berry K. J., Johnston J. E., Mielke P. W. Jr. [2018] *The Measurement of Association. A Permutation Statistical Approach*. Springer Nature Switzerland. Cham.
- Berry K. J., Mielke P. W. Jr., Johnston J. E. [2016] *Permutation Statistical Methods. An Integrated Approach*. Springer Nature Switzerland. Cham.
- Bonnini S., Corain L., Marozzi M., Salmaso L. [2014] *Nonparametric Hypothesis Testing Rank and Permutation Methods with Applications in R*. John Wiley & Sons, Ltd. Chichester.
- Corain L., Arboretti R., Bonnini S. [2016] *Ranking of Multivariate Populations. A Permutation Approach with Applications*. CRC Press. Boca Raton.
- Good P. [2005] *Permutation, Parametric and Bootstrap Tests of Hypotheses*, Springer Science Business Media, Inc., New York.

- Ledwina T. [2012] *Neyman Jerzy (1894 – 1981)* [w:] Statystycy polscy, Główny Urząd Statystyczny, Polskie Towarzystwo Statystyczne. Warszawa.
- Lehmann E. L., Romano J. P. [2005] *Testing Statistical Hypotheses*, Springer-Verlag New York
- Mielke P. W., Berry K. J. Jr. [2007] *Permutation Methods. A Distance Function Approach*. Springer Science+Business Media, LLC. New York.
- O’Gorman T. W. (2012) *Adaptive Tests of Significance Using Permutations of Residuals with R and SAS*, John Wiley and Sons, New Jersey.
- Rencher, A.C. and Christensen, W.F. (2012) *Methods of Multivariate Analysis*. Wiley, Hoboken, 800.
- Thompson B. [1984] *Canonical Correlation Analysis : Uses and Interpretation*, Sage Publications, Inc. London.

## **Metody permutacyjne w analizie współzależności**

### **Streszczenie**

W 1925 roku R. A. Fisher zaproponował zastosowanie metod permutacyjnych do weryfikacji hipotez statystycznych. Testy permutacyjne pozwalają na weryfikację hipotez bez znacznych wymagań dotyczących założeń. Przeprowadzenie takich testów wymaga dużej mocy obliczeniowej i z tego względu dopiero w ostatnich latach metody te zyskują na praktycznym znaczeniu. Pojęcie metody permutacyjnej jest jednak znacznie szersze niż zagadnienie permutacyjnego testowania hipotez. Już w 1923 J. Sława-Neyman wykorzystał model permutacyjny do analizy wyników doświadczeń polowych.

Celem artykułu jest prezentacja możliwości zastosowania metod permutacyjnych w analizie zależności. W artykule przedstawiono wybrane możliwości zastosowania metod permutacyjnych w testowaniu hipotez dotyczących zależności pomiędzy zmiennymi. Przedstawiono propozycję metody pozwalającej na badanie istotności zależności pomiędzy dwoma zbiorami danych. Rozważania uzupełniono porównaniem rozmiaru i mocy proponowanego testu i testów znanych z analizy kanonicznej. Proponowany test charakteryzuje się większą mocą od pozostałych w przypadku, gdy próby są małe i pobrane z innych rozkładów niż rozkład normalny.

Słowa kluczowe: metody permutacyjne, testy permutacyjne, permutowanie obiektów, analiza zależności, korelacje kanoniczne.

## **Permutation methods in association analysis**

### **Abstract**

The permutation model in hypothesis testing was introduced by R. A. Fisher in 1925. These methods permit us to test hypotheses with as minimal assumptions as possible. These tests require high computing power and therefore they have found the applications in last years. However, the concept of permutation methods is much wider than the issue of permutation testing. In 1923 J. Spława-Neyman introduced a permutation model for the analysis of field experiments.

The aim of the article is to present the possibilities of applying the permutation methods in the analysis of dependencies. The article presents selected possibilities of data rearranging in dependency analysis. A proposal of a method allowing to investigate the significance of the relationship between two data sets was presented. The considerations were supplemented by comparing the size and power of the proposed test and tests known from the canonical analysis. The proposal is most powerful for non-normally distributed variables and for small samples.

Key words: permutation methods, permutation methods, data labels permutation, association analysis, canonical correlations.