

New Interpretation of Principal Components Analysis

Zenon Gniazdowski*

Warsaw School of Computer Science

Abstract

A new look on the principal component analysis has been presented. Firstly, a geometric interpretation of determination coefficient was shown. In turn, the ability to represent the analyzed data and their interdependencies in the form of easy-to-understand basic geometric structures was shown. As a result of the analysis of these structures it was proposed to enrich the classical PCA. In particular, it was proposed a new criterion for the selection of important principal components and a new algorithm for clustering primary variables by their level of similarity to the principal components. Virtual and real data spaces, as well as tensor operations on data, have also been identified. The anisotropy of the data was identified too.

Keywords — determination coefficient, geometric interpretation of PCA, selection of principal components, clustering of variables, tensor data mining, anisotropy of data

1 Introduction

In the method of principal component, a primary set of data consisting of n mutually correlated random variables can be represented by a set of independent hypothetical variables called principal components. A new dataset typically contains fewer variables than the original data. The smaller number of principal components contains almost the same information as the full set of primary variables [1][2].

This work is an attempt to make a new interpretation of results of the classical principal components analysis. Here it should be noted that it is not the purpose of this work full presentation of the principal components method. This can be found in the literature. So, tips on how to formulate a problem in the form of an optimization task can be found in [1]. Examples of the use of this method can be found in [1] and [2]. In this article, the principal component method will be presented in such a way as to be able to demonstrate the new capabilities of this method.

*E-mail: zgniazdowski@wwsi.edu.pl

2 Preliminaries

In this paper vectors and matrices will be used. If a matrix is created from vectors, the vectors will be the columns of that matrix. The tensor concept will be introduced and its rotation will be described too. Only first and second rank tensors¹ will be used in this article [3]. For description of the operations on tensors the matrix notation will be used, instead of dummy indexes [4]. The terms "coordinate system" and "the base" will be used interchangeably.

2.1 Used abbreviations and symbols

It is assumed that if a vector or matrix or any other object has the "prime" symbol ($'$), this object is described in a coordinate system (the base) other than the standard coordinate system. If there is no this symbol, the object is given in the standard coordinate system:

- a – a vector representing a single point in the space of standardized primary random variables (one row in the input data table).
- A – matrix of vectors representing standardized primary variables in a standard base (standard coordinate system) Note: The column in matrix A can not be identified with a single point in the data space (previously denoted as a).
- A' – matrix of vectors representing standardized primary variables in the base of eigenvectors.
- C – correlation coefficient matrix.
- C' – correlation coefficient matrix after diagonalization (the matrix of correlation coefficients in the base of eigenvectors).
- I – a identity matrix containing all axes of the standard base.
- p – a vector representing a single point in the principal components space (one row in the principal components data table).
- p_{ci} – i -th principal component or i -th vector representing this principal component.
- P – matrix of vectors representing principal components in the standard coordinate system (standard base).
- P' – matrix of vectors representing principal components in the base of eigenvectors. Note: The column in matrix P' can not be identified with a single point in the principal components data table (previously denoted as p).
- R – rotation matrix described transition from standard coordinate system.
- u_i – i -th vector representing the axis of the coordinate system after rotation (i -th eigenvector).
- U – matrix containing directional vectors of coordinate system after rotation (matrix of eigenvectors).
- v – vector in standard coordinate system (standard base).
- v' – vector in the base after rotation (in rotated coordinate system).

¹It should be noted that despite the objections, publications interchangeably uses the terms "rank of tensor" and "order of tensor" [3].

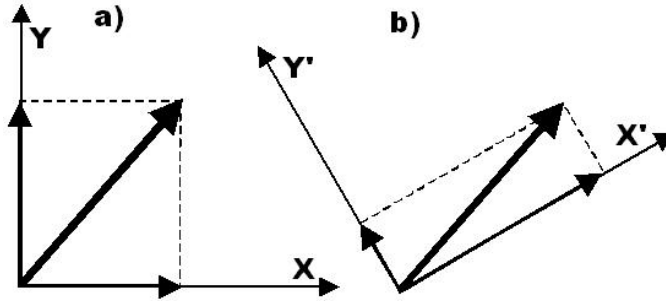


Figure 1: Rotation of Cartesian coordinate system. Components of vector: a) before coordinate system rotation, b) after coordinate system rotation

2.2 The concept of tensor

For description of the tensor concept, Cartesian coordinate system is assumed. Quantities, which are not dependent on this coordinate system, are scalars. Scalar is described by one number. On the contrary, some other quantities are defined with respect to coordinate system. These quantities are tensors. Tensor is a quantity, which is described by a number of components, whose values are depending on the coordinate system. In tensor notation, tensor rank manifests by number of indices:

- Scalar is a tensor of rank zero and it has no index.
- Vector is a first rank tensor and has only one index. For the given coordinate system in $n - D$ space, vector is completely defined by their n components. These elements are perpendicular projections of the vector to the respective axes of the system.
- Tensor of rank two has two indices. It is a quantity, which is defined by $n \times n$ numbers, which form a square matrix. As examples of the second rank tensor, a matrix of quadratic form can be used as well as a matrix of correlation coefficients. At this point it should be noted that second rank tensors are represented by square matrices, but not all square matrices are tensors [3].

Higher rank tensors can be also considered, but they are not the subject of this study.

2.2.1 Rotation of coordinate system

As an example vector on a plane (first rank tensor in $2 - D$) can be considered (see Fig. 1). This vector observed in different coordinate systems (before rotation and after rotation) is the same vector. In these different coordinate systems it has different components (projections on coordinate system axes). For known vector components in the coordinate system before rotation, the vector components in the coordinate system after rotation should be found. For

Table 1: Direction cosines between axes before and after rotation

		Axes before rotation		
		X_1	...	X_n
Axes	X'_1	$\cos(X'_1, X_1)$...	$\cos(X'_1, X_n)$
after	\vdots	\vdots	\vdots	\vdots
rotation	X'_n	$\cos(X'_n, X_1)$...	$\cos(X'_n, X_n)$

solving problem of changes of tensor components, transformation of coordinate system has to be described. Transformation of tensor component can be described after them.

Rotation of coordinate system without change of origin is considered. Axes before rotation are denoted as X_1, X_2, \dots, X_n . After rotation the same axes are denoted as X'_1, X'_2, \dots, X'_n . Table 1 shows direction cosines between coordinate system axes before and after rotation. To find this table, the set of cosines between vectors should be calculated using formula:

$$\cos(X'_i, X_j) = \frac{X'_i \cdot X_j}{\|X'_i\|_2 \cdot \|X_j\|_2}. \quad (1)$$

Since standard base vectors as well as vectors of finale base have unit lengths, the denominator in the above expression is equal to one and the cosine is equal to the scalar product of the corresponding vectors.

The content of Table 1 forms a rotation matrix. Denoting $\cos(X'_i, X_j)$ as r_{ij} , rotation matrix is a square matrix of $n \times n$ dimension, which is denoted as R :

$$R = \begin{bmatrix} r_{11} & \dots & r_{1n} \\ \vdots & \ddots & \vdots \\ r_{n1} & \dots & r_{nn} \end{bmatrix}. \quad (2)$$

Components of rotation matrix R are mutually dependent. Multiplying transposition matrix R^T by matrix R , identity matrix is achieved. It means matrix R is an orthogonal matrix:

$$R^T = R^{-1}. \quad (3)$$

Let the original base is a standard coordinate system represented by the n vectors b_1, b_2, \dots, b_n such that all components of b_i are all zero, apart from the unit element with the index i . As the columns of the first matrix the successive standard base vectors will be inserted. These columns will form the identity matrix :

$$I = \begin{bmatrix} 1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 1 \end{bmatrix}. \quad (4)$$

On the other hand, the final base should be a set of orthogonal vectors u_1, u_2, \dots, u_n such that $u_i = [u_{1i}, \dots, u_{ni}]^T$ and $\|u_i\|_2 = 1$. Let successive vectors of the new base become the successive columns of the second matrix:

Table 2: Transformation of tensors in matrix notation

Rank of tensor	New components expressed by old	Old components expressed by new	Note
0	$\varphi' = \varphi$	$\varphi = \varphi'$	$\varphi, \varphi' - \text{scalar}$
1	$v' = Rv$	$v = R^T v'$	$v, v' - \text{vector}$
2	$C' = RCR^T$	$C = R^T C' R$	$C, C' - \text{second rank tensor expressed as a square matrix}$

$$U = \begin{bmatrix} u_{11} & \dots & u_{1n} \\ \vdots & \ddots & \vdots \\ u_{n1} & \dots & u_{nn} \end{bmatrix}. \quad (5)$$

In the algorithm of finding the product of two matrices, the element with indices i and j in the resulting matrix is the scalar product of the i -th row in the first matrix and the j -th column in the second matrix. Because of the unit lengths of column vectors in matrices (4) and (5), the problem of finding the R rotation matrix is reduced to finding the product of two matrices:

$$R = U^T I. \quad (6)$$

Hence, the transition (rotation) matrix from the primary to the final base is a matrix whose rows are the directional vectors of the final base:

$$R = U^T = \begin{bmatrix} u_{11} & \dots & u_{n1} \\ \vdots & \ddots & \vdots \\ u_{1n} & \dots & u_{nn} \end{bmatrix}. \quad (7)$$

2.2.2 Transformation of tensor component

Vector $V = [v_1, v_2, \dots, v_n]^T$ is considered in some Cartesian coordinate system. If vector components before rotation are given and transformation R is known, then it is possible to find the set of vector components in the new coordinate system. If coordinate system is rotated with the use of matrix (2), then vector v will be observed as a vector v' with new components $v' = [v'_1, v'_2, \dots, v'_n]^T$, in the new coordinate system. The change of vector components is described by the formula [4]:

$$v' = Rv. \quad (8)$$

To return from the new coordinate system to the old one, both sides of equation (8) should be multiplied on the left by the inverse (transposition) of matrix R :

$$v = R^T v'. \quad (9)$$

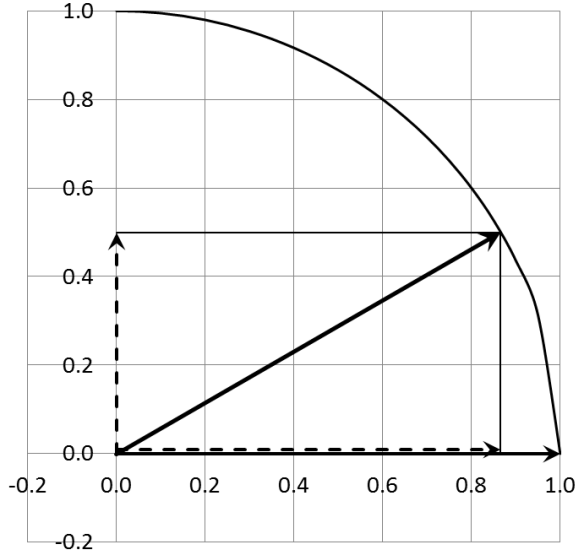


Figure 2: Geometric interpretation of determination coefficient

Since components of vector are dependent on coordinate system, in the same way components of higher rank tensors are also dependent on coordinate system. If coordinate system is rotated then tensor components are changed with this rotation. If tensor components before rotation are given and transformation of coordinate system (2) is known, then it is possible to find the set of tensor components in the new coordinate system. In Table 2, formulas for transformation of tensor components up to rank two are presented [4].

2.3 Geometric interpretation of determination coefficient

A measure of the relation between two random variables X and Y is the correlation coefficient. Denoting the random components of both variables as $x = X - \bar{X}$ and $y = Y - \bar{Y}$, the correlation coefficient can be expressed as the ratio of the scalar product of these vectors to the product of their lengths:

$$\rho_{X,Y} = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}} = \frac{x \cdot y}{\|x\|_2 \cdot \|y\|_2} = \cos(x, y). \quad (10)$$

This coefficient expresses the cosine of the angle between the random components of both variables [5]. Its square is called the coefficient of determination [6]. In order to present the geometric interpretation of the determination coefficient, the following facts should be noted:

- The standard deviation of the random variable is equal to the root of its variance [7]. This means that a random variable with a unit variance also has a unit standard deviation. This

fact applies to the standardized random variables.

- If the random variable is a sum of independent random variables, then its variance is equal to the sum of the variances of the added variables [7].
- Standardization of the variable does not affect the value of the correlation coefficient.

Two standardized random variables with known correlation coefficients are considered. Symbolically, they can be represented as vectors of unit lengths placed on a circle with a radius of one (see Fig. 2). The vector lying along the horizontal axis represents the explanatory variable. The vector set at an angle represents the variable to be explained. The angle between these vectors is chosen so that its cosine is equal to the correlation coefficient between these two variables. Lengths of the vectors represent the standard deviations of (standardized) random variables. Squares of lengths of vectors represent the variances of these variables. Since vectors have unit lengths, the cosine (the correlation coefficient) is represented by the length of the projection of the explained vector per the explanatory vector.

There is a question, what part of the variance of the explained variable is represented by the explanatory variable. The ratio of the square of the length of the projection of explained vector per the explanatory vector (on horizontal axis) to the square of the length of the explanatory vector is the coefficient of determination. This coefficient indicates what part of the variance of the explained variable is explained by the variance of the explanatory variable [6]. It is a number from the range $(0, 1)$, which can also be expressed in percentages. Because of the symmetry of the determination coefficient, the presented reasoning can be reversed by changing the explained variable with the explanatory variable. From here it can also be said that the coefficient of determination describes the level of a common variance of two correlated random variables.

From Pythagorean theorem, the square of the length of the explained vector can be expressed as a sum of the squares of the lengths of its projections per orthogonal axes (Figure 2). One axis is consistent with the direction of the explanatory vector and the second axis is in the orthogonal direction. Orthogonal projections indicated by dashed lines represent independent random variables. Just as the random variable can be represented as the sum of independent random variables, its variance can be represented as the sum of the variances of these (summed) variables. Just as the vector representing an explained variable can be represented as the sum of its orthogonal components, the square of the length of that vector is the sum of the squares of lengths of those components. Like the ratio of the projection of the explained vector to the explanatory vector is the cosine of the angle between these two vectors, the components of explained vector are the correlation coefficients between the explained variable and the independent explanatory variables represented by the orthogonal vectors.

The above conclusion can be generalized to a multidimensional space. If there is an explained variable and there is a set of explanatory variables, then there is an orthogonal base where the vector representation of the explained variable will have components equal to the correlation coefficients between the explained variable and successive explanatory variables.

Table 3: The data for analysis (in centimeters)

Sepal Length	Sepal Width	Petal Length	Petal Width	Class
5.1	3.5	1.4	0.2	Iris-setosa
4.9	3	1.4	0.2	...
4.7	3.2	1.3	0.2	...
⋮	⋮	⋮	⋮	⋮
6.4	3.2	4.5	1.5	Iris-versicolor
6.9	3.1	4.9	1.5	...
5.5	2.3	4	1.3	...
⋮	⋮	⋮	⋮	⋮
6.3	3.3	6	2.5	Iris-virginica
5.8	2.7	5.1	1.9	...
7.1	3	5.9	2.1	...
⋮	⋮	⋮	⋮	⋮
5.845	3.121	3.770	1.199	Average
0.833	0.480	1.773	0.763	Standard deviation
0.693	0.230	3.143	0.582	Variance

3 Principal Component Analysis

The matrix of correlation coefficients is used for finding the principal components. For this correlation matrix, the eigenproblem is solved. Several of the largest eigenvalues can to explain the most of the variation of analysed random variables. The original set of n mutually correlated random variables can be represented by a smaller set of independent hypothetical variables. This new set contains less variables than the original dataset. That means that there is a space in which several hypothetical variables adequately explain the random behavior of the analyzed primary dataset [1][2].

3.1 Data for analysis

As a data for the analysis the flower Iris data proposed by Sir Ronald Fisher in 1930 will be used [8]. On the one hand, these data are not too difficult to analyze. On the other hand, this data are enough to show many aspects of the principal component analysis. This data table contains four columns with correlated numeric data and a fifth column with nominal data. Only numeric columns will be used in this article. Table 3 lists several records of these data. At the bottom of the columns with numbers the mean, standard deviation, and the variance for each column are presented.

Table 4: The matrix of correlation coefficient (cosines)

	Sepal Length	Sepal Width	Petal Length	Petal Width
Sepal Length	1	-0.063	0.866	0.816
Sepal Width	-0.063	1	-0.321	-0.300
Petal Length	0.866	-0.321	1	0.959
Petal Width	0.816	-0.300	0.959	1

Table 5: The matrix of significance levels

	Sepal Length	Sepal Width	Petal Length	Petal Width
Sepal Length	0	0.446	0.000	0.000
Sepal Width	0.446	0	0.000	0.000
Petal Length	0.000	0.000	0	0.000
Petal Width	0.000	0.000	0.000	0

Table 6: The matrix of angles given in degrees

	Sepal Length	Sepal Width	Petal Length	Petal Width
Sepal Length	0	93.59	30.05	35.29
Sepal Width	93.59	0	108.74	107.47
Petal Length	30.05	108.74	0	16.44
Petal Width	35.29	107.47	16.44	0

Table 7: The matrix of determination coefficients

	Sepal Length	Sepal Width	Petal Length	Petal Width
Sepal Length	100.00%	0.39%	74.93%	66.63%
Sepal Width	0.39%	100.00%	10.32%	9.01%
Petal Length	74.93%	10.32%	100.00%	91.99%
Petal Width	66.63%	9.01%	91.99%	100.00%

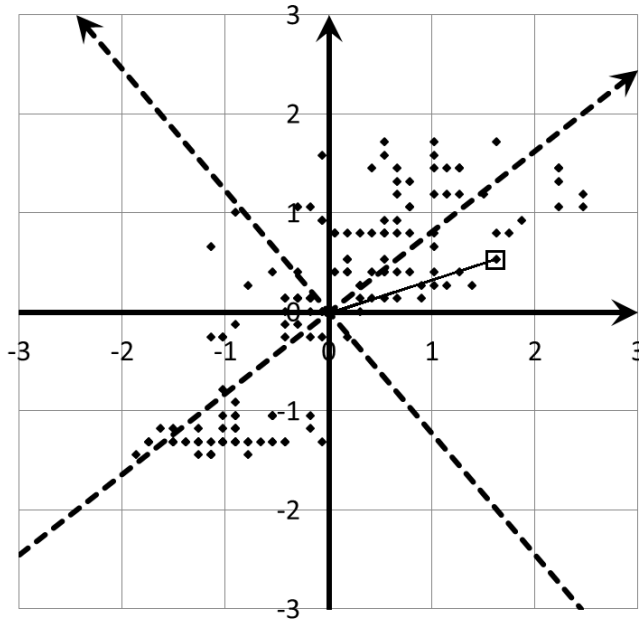


Figure 3: Rotation of the coordinate system. The point surrounded by the square will have different coordinates in the standard coordinate system (the axes marked with solid lines) and in the system of eigenvectors (the axes marked with dashed lines)

For considered data the matrix of correlation coefficients was calculated (Table 4). The significance level [5] for each correlation coefficient was also examined (Table 5). Except in one case, all correlations are significant on the level less than 0.001. Table 6 shows the angles in degrees calculated on the basis of the cosine (correlation coefficients) between the components of the considered random variables [5]. Angles close to the right angle confirms the strong independence of random variables. Angles of close to zero or 180 degrees confirm the strong dependence of random variables. The table shows that the second variable (Sepal Width) is almost orthogonal to the other variables. On the basis of the matrix of correlation coefficients were also found corresponding coefficients of determination [6]. Table 7 shows the results given as a percentage. The coefficients of determination between the second and the other variables do not exceed 11%. On the other hand, the determination coefficients between the first, third and fourth variables are not less than 66%

3.2 Reconstruction of principal components

Based on the matrix of correlation coefficients, the principal components were analyzed. As a result of the solution of the problem, the correlation coefficient matrix was diagonalized.

Table 8: Eigenvectors in columns

0.534	0.317	0.757	0.203
-0.213	0.948	-0.229	-0.066
0.584	0.026	-0.212	-0.783
0.573	0.030	-0.574	0.584

Table 9: Several objects in the principal components space

No.	p_{c1}	p_{c2}	p_{c3}	p_{c4}
1	-2.184	0.393	-0.176	0.048
2	-2.091	-0.674	-0.233	0.068
3	-2.341	-0.356	0.033	0.036
\vdots	\vdots	\vdots	\vdots	\vdots
Average	0.000	0.000	0.000	0.000
Standard deviation	1.694	0.984	0.398	0.181
Variance	2.868	0.967	0.159	0.033

The following eigenvalues were obtained: 2.849, 0.961, 0.158 and 0.033. The corresponding eigenvectors with length reduced to the unity are shown in the columns in Table 8. The order of eigenvectors corresponds to eigenvalues sorted from largest to smallest. Orthogonal eigenvectors represent the new base in which the primary random variables will be represented. Transposed matrix of eigenvectors creates an orthogonal rotation matrix (7). This matrix will be used to find mutually independent principal components.

Each row in the table of standardized primary variables (Table 3) is a vector representing one point in the primary variables space. Its projection to the eigenvectors (new axes of the system) will give a point in the principal components space. Denoting by a^T the row from the standardized primary data (Table 3) and by p^T the unknown row in the principal components space, and using the matrix R , we can accomplish this in the following way:

$$p^T = Ra^T. \tag{11}$$

Transformation (11) is equivalent to the rotation of the standard coordinate system (standard base) to the coordinate system defined by the eigenvectors. Before the rotation the row vector is known in the standard base as the vector a^T . After making the transformation (11), its description is found in the base of eigenvectors. For example, in Figure 3, the point (surrounded by a square) with coordinates (1, 632, 0.528) in the standard coordinate system, after rotation of this system becomes a point with coordinates (1, 899, -0.243). Repeating procedure (11) for all points in the space of standardized primary variables gives all points in the space of principal components. Table 9 shows several objects in the spaces of principal components. At the bottom of this table are given mean values, standard deviations and variances for each

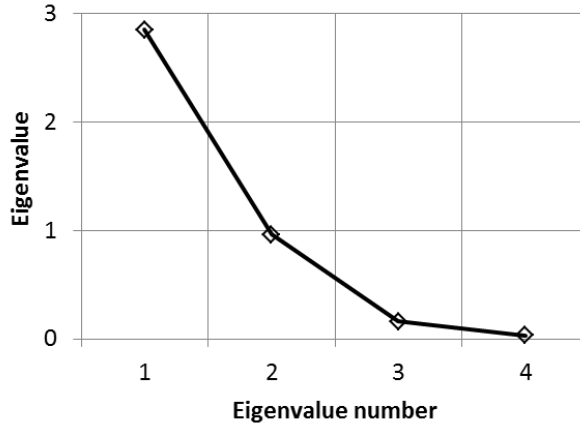


Figure 4: The scree plot

Table 10: The percentage of variance explained by the successive principal components (in brief: PC)

No.	Eigenvalue	Cumulative eigenvalues	Percentage of variance explained by each PC	Cumulative percentage of variance
1	2.849	2.849	71.22%	71.22%
2	0.961	3.810	24.02%	95.24%
3	0.158	3.967	3.94%	99.19%
4	0.033	4.000	0.81%	100.00%

column. It is noted that the variances of each variable p_{c1}, \dots, p_{c4} are (with an accuracy of calculation errors) equal to the previously obtained eigenvalues.

3.3 Choosing a subset of principal components

An important issue is the choice of the amount of the extracted principal components. Various criteria are used for this purpose [2]:

- Percentage criterion of the part of the variance which is explained by the principal components. It is assumed that there are so many principal components that the sum of the eigenvalues associated with the successive principal component is not less than the established threshold relative to the trace of the correlation coefficient matrix.
- Criterion of a scree plot (see Fig. 4). The plot shows the successive eigenvalues, from largest to smallest. Its shape resembles a scree. It is as many principal components as eigenvalues located on the slope of the scree.

Table 11: Correlation coefficients between primary variables and principal components

	Sepal Length	Sepal Width	Petal Length	Petal Width
p_{c1}	0.901	-0.359	0.986	0.968
p_{c2}	0.311	0.929	0.025	0.030
p_{c3}	-0.301	0.091	0.084	0.228
p_{c4}	0.037	-0.012	-0.141	0.105

Table 12: Determination coefficients between primary variables and principal components

	Sepal Length	Sepal Width	Petal Length	Petal Width	Σ
p_{c1}	0.812	0.129	0.972	0.936	2.849
p_{c2}	0.097	0.863	0.001	0.001	0.961
p_{c3}	0.090	0.008	0.007	0.052	0.158
p_{c4}	0.001	0.000	0.020	0.011	0.033
Σ	1.000	1.000	1.000	1.000	4.000

- Criterion of eigenvalue. The number of principal components is equal to the number of eigenvalues of not less than one.

Because each of the above criteria may suggest another number of components, the final decision about their numbers is taken by the human. Table 10 shows eigenvalues, cumulative eigenvalues, the percentage of variance explained by the principal components, and the cumulative percentage of the variance. This table shows that two principal components carry over 95% of the information contained in primary variables. Also, the scree plot shows that there are two values on its slope, while the other two are off the slope. According to the eigenvalue criterion, the second eigenvalue is slightly less than unity. By analyzing all the criteria for selecting the principal components presented above, it can be concluded that two principal components can be selected for representation of statistical behavior of a set of analyzed variables. Such a quantity sufficiently explains at least 95% of their variance.

4 Expandability of PCA capabilities

Formula (11) made it possible to find the principal components based on the standardized primary variables. In order to interpret the obtained results, correlation coefficients between primary variables and principal components were calculated (Table 11). Based on the correlation, appropriate coefficients of determination were calculated. For the given coefficients of determination the sums of the elements in columns and rows were calculated too (Table 12).

4.1 New interpretation of PCA in virtual vector space

The columns in the Tables 11 and 12 refer to standardized primary variables. On the other hand, rows in Tables 11 and 12 refer to the principal components. The consistent mutual analysis of both tables allows to find a new geometric interpretation of the principal component method.

4.1.1 Standardized primary variables - analysis of columns

Elements in the columns of Table 12 sum to the unity, which is the variance of subsequent standardized primary variables. This means that the standardized primary variable is the sum of independent variables, and its variance is equal to the sum of variances of those variables (see Sec. 2.3). It can be said that standardized primary variables divide their variance between independent principal components. For example, a standardized variable Sepal Length divides its variance between four mutually independent principal components p_{c1} , p_{c2} , p_{c3} as well as p_{c4} . About 81% of the variance of the variable Sepal Length is a part of the principal component p_{c1} , almost 10% is a part of p_{c2} , about 9% is passed to the principal component p_{c3} , and the remainder of variance (less than 1%) is a part of p_{c4} .

The following columns of Table 11 are the correlation coefficients between successive primary variables, and independent principal components. Figure 2 shows that the correlation coefficient can be interpreted as a projection of a unit length vector per axis of the coordinate system. This projection is one of the vector components. Consequently, the correlation coefficients in a given column can be treated as components of a vector. If the correlation coefficients in a given column of Table 11 are taken as vector components, then by the Pythagorean theorem the sum of the squares of the vector projections lengths per perpendicular axes equals the square of the length of this vector.

Principal components are associated with eigenvectors. The components of subsequent column vectors in Table 11 are projections for orthogonal eigenvectors. This means that primary variables are represented as vectors described in the base of eigenvectors. These vectors are columns in the matrix A' :

$$A' = \begin{pmatrix} 0.901 & -0.359 & 0.986 & 0.968 \\ 0.311 & 0.929 & 0.025 & 0.030 \\ -0.301 & 0.091 & 0.084 & 0.228 \\ 0.037 & -0.012 & -0.141 & 0.105 \end{pmatrix} \quad (12)$$

4.1.2 Principal components - analysis of rows

The elements in the rows of Table 12 sums to eigenvalues that are the variances of the principal components. The variance of the given principal component consists of mutually independent (orthogonal) parts of the variance of the standardized primary variables. This means that this principal component consists of the sum of independent random variables, and its variance is equal to the sum of the variances of those variables. And so, to the variance of the principal component p_{c1} the variables Sepal Length, Petal Length as well as Petal Width give the most.

To the variance of the principal component p_{c2} contributes most variable Sepal Width. The variances of the principal components p_{c3} and p_{c4} are much less than the variances of the principal components p_{c1} and p_{c2} .

In the subsequent rows of Table 11 there are coefficients of correlation between successive principal components and primary variables. By virtue of Pythagorean theorem, the sum of the squares of these components is equal to the square of the length of the vector (Table 12). These are equal to the eigenvalues, so they are equal to variances of consecutive principal components.

Primary variables are associated with a standard base. Components of consecutive row vectors are projections on the orthogonal axes of the standard base. This means that the principal components are represented as vectors in the standard base. Rows in Table 11 are columns in the matrix P :

$$P = \begin{pmatrix} 0.901 & 0.311 & -0.301 & 0.037 \\ -0.359 & 0.929 & 0.091 & -0.012 \\ 0.986 & 0.025 & 0.084 & -0.141 \\ 0.968 & 0.030 & 0.228 & 0.105 \end{pmatrix} \quad (13)$$

4.1.3 Vector representation of both standardized primary variables and principal components

As a consequence of the analysis, it was shown that both the standardized primary variables and the principal components can be seen as a vectors described in some bases. The columns in Table 11 are vectors representing standardized primary variables in the eigenvectors base. These columns form the matrix A' . On the other hand, the rows in Table 11 are vectors representing the principal components in the standard base. The transposition of the rows of Table 11 forms the matrix P . Between the representation of the standardized primary variables (matrix A' in the base of eigenvectors) and the representation of the principal components (Matrix P in standard base) there is a relation:

$$A' = P^T. \quad (14)$$

For vectors representing standardized primary variables, the natural base is the standard base. For vectors representing principal components, the natural base is the base of the eigenvectors. The question is: What is their character in their natural bases.

Standardised primary variables in the natural base It is known representation A' of the standardized primary variables (12) in the base of eigenvectors. On the other hand, it has been shown that the eigenvectors are the rows of a rotation matrix R . This matrix (or its transposition) describes the transition between the standard base and the eigenvectors base. The equation (9) was used to find description of matrix A in the standard base:

$$A = R^T A'. \quad (15)$$

As a result of this transformation, the symmetric matrix A was obtained:

$$A = \begin{pmatrix} 0.815 & 0.032 & 0.442 & 0.375 \\ 0.032 & 0.978 & -0.157 & -0.132 \\ 0.442 & -0.157 & 0.705 & 0.532 \\ 0.375 & -0.132 & 0.532 & 0.748 \end{pmatrix} \quad (16)$$

Principal components in the natural base For vectors representing principal components, the natural base is the base of eigenvectors. Vectors describing these principal components in the standard base are known. These vectors forms the matrix P . To find the representation of the principal components in the eigenvectors base, formula (8) was used:

$$P' = RP. \quad (17)$$

As a result, the diagonal matrix P' was obtained:

$$P' = \begin{pmatrix} 1.688 & 0.000 & 0.000 & 0.000 \\ 0.000 & 0.980 & 0.000 & 0.000 \\ 0.000 & 0.000 & 0.397 & 0.000 \\ 0.000 & 0.000 & 0.000 & 0.180 \end{pmatrix} \quad (18)$$

When the matrix (18) is raised to the square, on its diagonal appear the variances of the principal components. This means that the squares of the lengths of the vectors representing the principal components are equal to the variances of these principal components. Obtained result is consistent with intuition:

- The principal components are represented by orthogonal vectors conforming to the directions of eigenvectors.
- The values on the diagonal are equal (with precision of numerical errors) to the standard deviations of the principal components.

4.1.4 Ability to reconstruction of correlation matrix

It is assumed that in matrix A containing vectors representing standardized primary variables, complete information is provided about the statistical properties of these variables. Therefore, using this matrix, correlation coefficients between variables can be recreated. Correlation coefficients have the interpretation of cosines of angles. If the coordinates of individual vectors in a certain base are given, the cosines between the vectors can be found. Since vectors representing standardized primary variables have unit lengths, the corresponding cosines are equal to the scalar product of the corresponding vectors. Therefore, using the vector representation of the primary variables A , the matrix C containing the correlation coefficients of the primary variables can be calculated by performing the following matrix multiplication:

$$C = A^T A. \quad (19)$$

Table 13: Level of reconstruction of primary variables

	Sepal Length	Sepal Width	Petal Length	Petal Width	Average in row
p_{c1}	81.17%	12.88%	97.23%	93.61%	71.22%
p_{c2}	9.65%	86.28%	0.06%	0.09%	24.02%
Σ	90.82%	99.16%	97.29%	93.70%	95.24%

This formula is analogous to the formula (10) for the correlation coefficient of two random variables with a precision of skipped product of the lengths of vectors.

It should be noted that the operation (19) can be performed both for vectors represented in the standard base (columns of the matrix (16)) as well as in the base formed by the eigenvectors (columns of the matrix (12)).

4.2 The level of explanation of primary variables

Proposed in section 3.2 two principal components will adequately explain at least 95% of the variance of the set of analyzed variables. The level of explanation of primary variables by the selected two principal components applies to all variables, not to individual variables. On the other hand, it can be expected that the variance of each primary variable will be explained by selected principal components to varying degrees. Since Table 10 refers to all variables rather than to individual variables, additional analysis using determination coefficients is therefore needed. By reducing the determination coefficient table (Table 12) to two rows and giving results in percent (Table 13), you can estimate the level of reproduction of the variance of primary variables by the selected principal components. Thus, the first primary variable will be retained in more than 90 percent, the second variable in over 99 percent, the third variable in over 97 percent, and the fourth variable in over 93 percent. For individual primary variables, the loss of information is less than 10%, less than 1%, less than 3%, and less than 7%, respectively. It can also be seen that average values in the rows (last column in Table 13) are identical to the percentage values in the columns in Table 10. The first two values agree with the percent variance explained by the principal components p_{c1} and p_{c2} , and the third value agrees with the corresponding cumulative value (the cumulative percentage of the variance for two principal components). This means that the result of 95% obtained in Table 10 is only the average result for all variables. The use of this criterion refers to the mean level of explanation of the variables.

In practice it may happen that some of the primary variables will be underrepresented. An additional criterion for selecting the appropriate number of principal components should therefore be used. This criterion should take into account the level of reconstruction of individual primary variables, as in the last row in Table 13. The number of principal components should be such that the representation level of each variable is at least satisfactory.

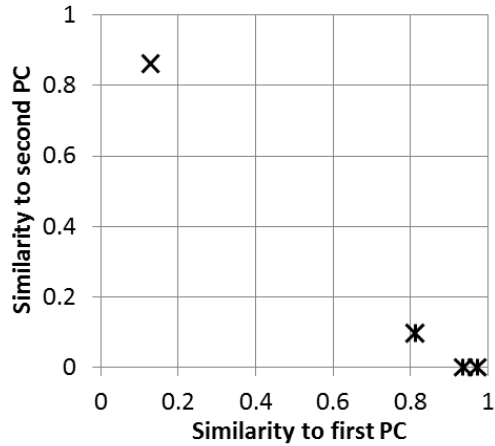


Figure 5: The similarity of the primary variables and selected principal components. Two different clusters are labeled with different markers.

Table 14: The similarity of the primary variables and selected principal components measured by the coefficient of determination

	Sepal Length	Sepal Width	Petal Length	Petal Width
Similarity to p_{c1}	0.812	0.129	0.972	0.936
Similarity to p_{c2}	0.097	0.863	0.001	0.001

Table 15: Operations on tensors

Transition to the new coordinate system	Back to the old coordinate system	Note
$A' = RA$	$A = R^T A'$	Standardized primary variables are represented as vectors (columns in matrices A and A')
$P' = RP$	$P = R^T P'$	Principal components are represented as vectors (columns in matrices P and P')
$C' = RCR^T$	$C = R^T C' R$	The correlation coefficient matrix (C, C') is a tensor of rank two

Table 16: Relationships between tensors

	$A =$	$A' =$	$P =$	$P' =$	$C =$	$C' =$
$A \rightarrow$	\times	RA	AR^T	RAR^T	$A^T A$	$RA^T AR^T$
$A' \rightarrow$	$R^T A'$	\times	$(A')^T$	$R(A')^T$	$(A')^T A'$	$R(A')^T A' R^T$
$P \rightarrow$	PR	P^T	\times	RP	$R^T P^T PR$	$P^T P$
$P' \rightarrow$	$R^T P' R$	$P' R$	$R^T P'$	\times	$R^T (P')^T P' R$	$(P')^T P'$
$C \rightarrow$					\times	RCR^T
$C' \rightarrow$					$R^T C' R$	\times

4.3 The concept of clustering of primary variables based on common variance

The coefficient of determination can be considered as a measure of similarity between two random variables. The greater the common variances, measured by the coefficient of determination, the variables are more similar. Table 12 contains the coefficients of determination between the primary variables and the principal components. Table 14 contains Table 12 reduced to the first two rows. By analyzing Table 14, it can be seen that the three primary variables have considerable common variances with the first principal component p_{c1} and one primary variable has a considerable common variance with the second principal component p_{c2} . The first, third and fourth primary variables will be in a set of variables similar to the first principal component p_{c1} . The second primary variable will remain in a one-element set of primary variables similar to the second principal component p_{c2} . This can be seen in Figure 5. Based on this, primary variables can be clustered according to the strength of similarity to the principal components. A naive method can be used for clustering. If the similarity between a given primary variable and i -th principal component is not less than 50%, then a given variable belongs to the i -th cluster. For variables that are not qualified for any cluster, you must create an additional cluster. A more complex method, such as the k -means method may also be used. For this purpose it can be considered the metric l_p (Euclidean, Manhattan, Chebyshev for $p = 2, 1$ and ∞ , respectively) or the cosine similarity [9]. It should be noted that in the case of a greater number of selected principal components, the clusters obtained by different methods may differ.

5 Discussion

Several new results have appeared in this article. The above propositions suggest some unresolved problems that need to be clarified or at least commented on. If this is not possible then problems should be formulated for a later solution.

5.1 Virtual representation of primary variables and principal components

Real variables Sepal Length, Sepal Width, Petal Length as well as Petal Width are the primary variables. Real variables are obtained from the measurement. For the analysis of the principal

components, standardized random variables are considered. Principal components are the same standardized primary variables described in the transformed coordinate system. A single point in a space is the same point, regardless of whether it is described as a point in the standard base, or a point in the base of the eigenvectors. Both in the case of standardized primary variables and in the case of principal components we have a real data type.

On the other hand, both primary variables and principal components have their vector representation. In this representation only statistical information is contained. It is information about the variance of individual standardized primary variables or principal components, as well as their interdependence. There is no access to single points in measurements space. It can be said that the vector representation is a representation in the virtual vector space.

There are some new questions to which the answer is not yet known. Can the virtual representation obtained in this work be practically used in any way? Can it be useful in data analysis?

5.2 Symmetry of matrix A representing primary variables in standard base

It was shown that the method of principal components can be represented in a virtual vector space. This applies both to the standardized primary variables and to the principal components. As a result of the calculations, a symmetric matrix (16) of vectors representing the primary variables in the standard base was obtained. It can be shown that the result is correct. Left-multiplying the equation (15) by a matrix R , the equation is obtained:

$$A' = RA. \quad (20)$$

Using equations (20), (14) and (17) the P' in the following form can be shown:

$$P' = RA^T R^T. \quad (21)$$

From the formula (18), it can be seen that P' in the base of eigenvectors is a diagonal matrix. Thus, there is symmetry:

$$P' = (P')^T. \quad (22)$$

By comparing the formula (21) with its transposition, the identity is obtained:

$$A = A^T. \quad (23)$$

This means that the A matrix representing the standardized primary variables in the standard base is symmetric. Therefore, the question arises: Can the information contained in this symmetry be useful in data analysis?

5.3 Tensor Data Mining

In the presented analysis, tensor operations were identified. These operations consist in finding a tensor description in a rotated coordinate system. The first example is the transformation (11) executed on a vector (first order tensor). For a given point described in the standard base, this operation finds its coordinates in the base of the eigenvectors. If this transformation is

applied to all points in the space of the standardized primary variables, then all points in the principal component space will be obtained. In virtual representation, vectors (first rank tensors) represent the primary variables as well as the principal components. On the other hand, tensor of the second rank is the matrix of correlation coefficients and its form after diagonalization.

Operations on tensors are summarized in Table 15. Table 16 contains possible to obtain relations between tensors. Six expressions of Table 15 can be identified in cells in the Table 16 located directly above and below the main diagonal.

5.4 Anisotropy of data

In the colloquial meaning, anisotropy consists in the fact that some observed quantity depend on the coordinate system. The opposite of anisotropy is isotropy, which lies in the fact that the quantity in all directions is the same. Tensors are natural carriers of anisotropy. Since tensor operations have been identified, anisotropy has been identified too. Because the tensors represent random data, so the data anisotropy was also identified.

Anisotropy of data can be seen in the fact that a single point in the space of real data can be described in different ways. Once it is visible, as a single row in Table 3 that contains standardized primary variables. After transformation (11) becomes a row in Table 9 containing points in the principal components space. Anisotropy can also be seen in the fact that the vector representation of the standardized primary variables once takes the form of matrix A , another time it takes the form of matrix A' . The situation is similar in the case of the representation of principal components, once in the form of the matrix P , and again in the form of the matrix P' . The above examples refer to vectors (first rank tensors).

In turn, symmetric matrix of correlation coefficients (second rank tensor), depending on the coordinate system, is C matrix or C' matrix. In all these cases, the same objects (vectors, second rank tensors) observed in different coordinate systems (different bases) are described by different components. They are all anisotropic objects.

5.5 Clustering of random variables

The paper proposes the possibility of clustering correlated random variables because of their similarity to the principal components. It should be noted that clustering of points is commonly used in data analysis. For this purpose many different methods are used, from the classical k -means [1] [10], to the spectral methods [11]. On the other hand, the data table consists of columns and rows. The columns contain random variables and the rows contain points in the space of these variables. It can be said that in the data analysis, the rows in the data table are clustered.

So far, the author has not encountered clustering of random variables (columns in the data table). Perhaps to distinguish between these two types of clustering, it is sensible to use different naming conventions. We suggest that clustering of points in the space of random variables was called horizontal clustering, and clustering of random variables (columns) was called vertical clustering.

6 Conclusions

In this paper the method of principal components was analyzed. As a result of this analysis some interesting results have been achieved. These results will be presented synthetically below:

1. A geometric analysis of the determination coefficient was performed. This analysis led to the conclusion that the coefficient of determination describes the level of common variance of two correlated random variables. It was also found that this coefficient is a good measure describing the similarity between the correlated random variables.
2. Geometric interpretation of the principal component method in a vector space was proposed. This interpretation is based on the generalization of the Pythagorean theorem. For this purpose, correlation coefficients and coefficients of determination between the primary variables and the principal components obtained were analyzed:
 - (a) It has been found that standardized primary variables can be decomposed into the sum of independent (orthogonal) random variables. In this way, the variance of each standardized primary variable consists of the sum of the variances of these mutually independent (orthogonal) components. On the one hand it has been noted that the correlation coefficients between the primary variable and the principal components can be interpreted as vector components representing the standardized primary variable in the virtual vector space. On the other hand, the modules of these correlation coefficients can be interpreted as the standard deviation of the above-mentioned independent random variables.
 - (b) A similar interpretation can be applied to the principal components. Each principal component can be decomposed into the sum of mutually independent (orthogonal) random variables. In the same way, the variance of each principal component consists of the variances of these mutually independent (orthogonal) components. Similarly, as for standardized primary variables, there is also a vector representation for the principal components.
3. As a result of the analysis, the principal component method has been enriched with two additional elements:
 - (a) The first element is the proposed algorithm for clustering primary variables by their level of the similarity to the principal components. For the distinction between clustering of points (clustering of rows in a data table) and clustering of variables (clustering columns from a data table), the terms "horizontal clustering" and "vertical clustering" have been proposed. The clustering of the variables proposed in this article is vertical clustering.
 - (b) The second element is an additional criterion for selection of the correct number of principal components because of the level of reconstruction of the primary variables. Using this criterion will select the appropriate number of principal components in such a way that the level of representation of each primary variable was at least satisfactory.

4. A distinction has been made between real data and virtual data. Real data are the primary variables obtained from the measurements as well as the principal components. A vector representation of standardized primary variables and principal components can be called a virtual representation of data or simply virtual data.
5. Operations on tensors have been identified. In this way it is stated that the principal component analysis is a part of the broader concept called Tensor Data Mining. Since tensor operations have been identified and tensors are the natural carrier of anisotropy, the concept of Data Anisotropy has been introduced.

Acknowledgments

The term "Anisotropy of Data" was proposed by Professor Michał Grabowski, when the author of this article stated that in data analysis can be seen tensors as in the case of anisotropic physical phenomena.

References

- [1] David J. Hand, Heikki Mannila, and Padhraic Smyth. *Principles of data mining*. MIT press, 2001.
- [2] Daniel T. Larose. *Data mining methods and models*. John Wiley & Sons, 2006.
- [3] Taha Sochi. Introduction to tensor calculus. *arXiv preprint arXiv:1603.01660*, 2016.
- [4] John F. Nye. *Physical properties of crystals*. Clarendon Press, Oxford, 1957.
- [5] Zenon Gniazdowski. Geometric interpretation of a correlation. *Zeszyty Naukowe WWSI*, 7(9):27–35, 2013.
- [6] Richard Taylor. Interpretation of the correlation coefficient: a basic review. *Journal of diagnostic medical sonography*, 6(1):35–39, 1990.
- [7] Athanasios Papoulis. *Probability, random variables and stochastic processes*. McGraw-Hill, 1965.
- [8] Ronald A. Fisher. The use of multiple measurements in taxonomic problems. *Annals of Eugenics*, 7(2):179–188, 1936.
- [9] Michel Marie Deza and Elena Deza. *Encyclopedia of Distances*. Springer, 2009.
- [10] Daniel T Larose. *Discovering knowledge in data: an introduction to data mining*. John Wiley & Sons, 2005.
- [11] Ulrike Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416, 2007.

