

Pattern Normalization – a New Tool for Dynamic Comparisons

Iwona Müller-Fraćzek¹ | Nicolaus Copernicus University, Toruń, Poland

Abstract

The article presents a new method of normalization – normalization with respect to pattern (or pattern normalization in short). It has properties expected for this type of transformation: preserves skewness, kurtosis and the Pearson correlation coefficients. Although pattern normalization uses only observations from the current unit of time, it can be used in dynamic research. An additional advantage of new normalization is the ability to reflect different analysis environments. The effects of pattern normalization are illustrated by an empirical example. Indicators monitoring the implementation of the Europe 2020 Strategy are used. Normalizations are carried out for two reference groups: the entire EU and countries that joined the EU in 2004. The results for two years are compared. The example of Poland shows that the “dynamic image” of the country is affected by the use of pattern normalization itself as well as by the choice of the environment. In this context pattern normalization is similar to dynamic standardization, and different than dynamic scaling.

Keywords

Normalization, standardization, transformation of variables

JEL code

C43, C19, C38

INTRODUCTION

We understand normalization as procedure of pre-treatment of data in order to allow for their mutual comparison and further analysis. Such a procedure is used, for example, in a study of a complex phenomenon, i.e. a qualitative phenomenon that is characterized by a collection of quantitative variables. Without losing generality, we assume that this is the phenomenon observed for objects in space, such as socio-economic development of countries. In this case, normalization deprives variables of their units and unifies their ranges. After normalization we can compare variables separately or construct a composite indicator. The composite indicator is one-dimensional image of multidimensional phenomenon (compare Saisana and Saltelli, 2011; Saltelli, 2007).

There are many normalization formulas (see Jajuga and Walesiak, 2000; Milligan and Cooper, 1988; Młodak, 2006; Steinley, 2004). Most often they are given for a static analysis, i.e. for a fixed point in time. Normalization problems appear when we want to compare a given phenomenon at several time points. In this case diagnostic variables should also be comparable over time.

To achieve this effect we can use two approaches. In first of them, we exploit all values of variable (both in space and time) to determine the parameters needed for normalization (compare Nardo et al., 2005). We can call this approach the stochastic one, because we treat observations for a given time point as randomly selected sample of population. But it is rather controversial in regional comparisons where

¹ Faculty of Economic Sciences and Management, Gagarina 13 a, 87-100 Toruń, Poland. E-mail: muller@umk.pl, phone (+48)566114718.

we work with the whole population of objects in space, but not with a sample (compare Zeliaś, 2002). In addition, a practical disadvantage of this solution is the need to recalculate all results with the appearance of observations for next unit of time.

In the second approach, parameters needed for normalization do not result directly from variable distributions. They are taken in advance, the same for all the units of time (also future). This solution is used, for example, in the very popular Human Development Index (HDI), as well as in a newer proposal i.e. the Adjusted Mazziotta-Pareto Index, called AMPI in short (compare Mazziotta and Pareto, 2015, 2016).

The article proposes another way of solution to dynamic problems. We introduce a new method of feature normalization – normalization with respect to the pattern (or pattern normalization for short). The method is consistent with the static approach (only current observation are taken), but it can be used to compare objects at different time points. The method meets requirements of normalization that are suggested in literature (compare e.g. Jajuga and Walesiak, 2000; Młodak, 2006). It preserves skewness and kurtosis. Moreover, the absolute values of the Pearson correlation coefficients are not changed after normalization.

An additional advantage of pattern normalization is the possibility to reflect different environments in research. This is the same as in standardization and on the contrary to scaling (or min-max normalization) used in the mentioned HDI or AMPI. This property is illustrated by an empirical example. Indicators monitoring the implementation of the Europe 2020 Strategy (see European Commission, 2010) are normalized in two environments, one is the whole European Union, and the second – a group of countries that joined the EU in 2004. The example of Poland shows differences for both environments.

The article is divided into 6 parts. Section 1 introduces the normalization with respect to pattern. Section 2 presents properties of the pattern normalization. Section 3 discusses advantages of new proposal. Section 4 illustrates theoretical consideration. The article ends with conclusions.

1 DEFINITION OF PATTERN TRANSFORMATION

Consider a set of $n \in N$ objects in space. For these objects, we analyze a phenomenon which is not directly measurable and it is composed of many aspects (a complex phenomenon). Various aspects of this phenomenon are characterized by measurable diagnostic variables, that is, variables for which a connection with a certain aspect of the complex phenomenon is not in doubt and the direction of this relationship can be determined (a stimulant is a diagnostic variable that has a positive impact on the analyzed complex phenomenon, while a destimulant negative).² An example of a complex phenomenon is socio-economic development of the European Union countries, and diagnostic variables for this phenomenon are, among others, the indicators monitoring implementation of the Europe 2020 Strategy (considered in Section 4).

The analyzed objects can be ordered due to individual diagnostic variables, i.e. in relation to particular aspects of the complex phenomenon. To order objects due to all aspects of this phenomenon, we can construct a synthetic variable (a composite indicator).³ One of the stages of such construction is normalization of variables.

For a given unit of time consider one diagnostic variable $x = (x_1, x_2, \dots, x_n) \in R^n$. This variable is a stimulant (then we write $x \in S$, where S denotes the set of stimulants) or a destimulant ($x \in D$ respectively). We choose a pattern – the most beneficial of all values of the variable x . This name was inspired by the Hellwig's paper (Hellwig, 1968). The pattern is unique for all objects and is described by the formula:

² Other types of variables are not considered. If they must be used in the study, they should be transformed into stimulants.

³ In this case, the diagnostic variables must meet additional statistical requirements such as sufficient variability or weak correlation. This is beyond the scope of the article, for more details we refer, for example, to Zeliaś (2002).

$$x^+ = \begin{cases} \max_{i=1, \dots, n} x_i & \text{if } x \in S, \\ \min_{i=1, \dots, n} x_i & \text{if } x \in D. \end{cases} \quad (1)$$

After specifying the pattern x^+ we can consider a new variable u instead of the variable x given by:

$$u_i = \frac{|x_i - x^+|}{\sum_{j=1}^n |x_j - x^+|} = \begin{cases} \frac{x^+ - x_i}{\sum_{j=1}^n (x^+ - x_j)} & \text{if } x \in S, \\ \frac{x_i - x^+}{\sum_{j=1}^n (x_j - x^+)} & \text{if } x \in D. \end{cases} \quad \text{for } i = 1, \dots, n. \quad (2)$$

The Formula (2) determines a transformation of initial variable $x = (x_1, x_2, \dots, x_n)$ into a new variable $u = (u_1, u_2, \dots, u_n)$. After this transformation (transformation with respect to pattern) the new variable describes the same aspect of complex phenomenon as x describes. So u is a diagnostic variable of this phenomenon.

2 STATIC PROPERTIES OF PATTERN TRANSFORMATION

2.1 Basic properties

1. All diagnostic variables after pattern transformation are unitless, non-negative and limited to interval $[0,1]$. Because of that, the new set of diagnostic variables contains comparable elements.
2. The lower is the value u_i the better is the situation of the i -th object. It means that the variable after the pattern transformation becomes destimulant irrespective of its initial nature. So, the pattern transformation unifies the nature of the diagnostic variables.
3. Transforming of variables does not affect the ordering of objects.

2.2 Extreme values after pattern normalization

1. The variable u can take the zero value only for the pattern object:

$$u_i = 0 \Leftrightarrow x_i = x^+ \quad \text{for } i = 1, \dots, n. \quad (3)$$

2. Since the pattern is chosen among the values of the variable x , the zero value is taken:

$$\min_{i=1, \dots, n} u_i = 0. \quad (4)$$

3. The value u_i equals 1 when all objects except the i -th one are patterns:

$$u_i = 1 \Leftrightarrow \forall_{j \neq i} x_j = x^+ \quad \text{for } i = 1, \dots, n. \quad (5)$$

4. The maximum value of u depends on the nature of variable x :

$$\max_{i=1, \dots, n} u_i = \begin{cases} \frac{\max_{i=1, \dots, n} x_i - \min_{i=1, \dots, n} x_i}{\sum_{j=1}^n (\max_{i=1, \dots, n} x_i - x_j)} & \text{if } x \in S, \\ \frac{\max_{i=1, \dots, n} x_i - \min_{i=1, \dots, n} x_i}{\sum_{j=1}^n (x_j - \min_{i=1, \dots, n} x_i)} & \text{if } x \in D. \end{cases} \quad (6)$$

2.3 Descriptive characteristics of transformed variables

1. The mean value of u depends only on the number of objects:

$$\bar{u} \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n u_i = \frac{1}{n}. \quad (7)$$

2. The variance of u is described by:

$$S^2(u) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n (u_i - \bar{u})^2 = \frac{S^2(x)}{n^2 (x^+ - \bar{x})^2}. \quad (8)$$

3. The standard deviation of u depends on the nature of variable x and it is expressed by:

$$S(u) \stackrel{\text{def}}{=} \sqrt{S^2(u)} = \begin{cases} \frac{S(x)}{n(x^+ - \bar{x})} & \text{if } x \in S, \\ \frac{S(x)}{n(\bar{x} - x^+)} & \text{if } x \in D. \end{cases} \quad (9)$$

4. The coefficient of variation of u is given by:

$$CV(u) \stackrel{\text{def}}{=} \frac{S(u)}{\bar{u}} = \begin{cases} \frac{S(x)}{x^+ - \bar{x}} & \text{if } x \in S, \\ \frac{S(x)}{x - x^+} & \text{if } x \in D. \end{cases} \quad (10)$$

5. The 3rd central moment of u is expressed by:

$$\mu_3(u) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n (u_i - \bar{u})^3 = \frac{\mu_3(x)}{n^3 (\bar{x} - x^+)^3}. \quad (11)$$

6. The absolute value of the coefficient of skewness is preserved:

$$A(u) = \frac{\mu_3(u)}{S^3(u)} = \begin{cases} -A(x) & \text{if } x \in S, \\ A(x) & \text{if } x \in D. \end{cases} \quad (12)$$

7. The 4th central moment of u is given by:

$$\mu_4(u) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n (u_i - \bar{u})^4 = \frac{\mu_4(x)}{n^4 (x^+ - \bar{x})^4}. \quad (13)$$

8. The kurtosis of u does not change after the pattern transformation:

$$K(u) \stackrel{\text{def}}{=} \frac{\mu_4(u)}{S^4(u)} - 3 = K(x). \quad (14)$$

2.4 Linear relation between variables after transformation

Denote by $u_1 = (u_{11}, u_{12}, \dots, u_{1n})$ and $u_2 = (u_{21}, u_{22}, \dots, u_{2n})$ two diagnostic variables after pattern transformation.

1. The covariance between u_1 and u_2 equals:

$$\text{cov}(u_1, u_2) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n (u_{1i} - \bar{u}_1) (u_{2i} - \bar{u}_2) = \begin{cases} \frac{\text{cov}(x_1, x_2)}{n^2 (x_1^* - \bar{x}_1) (x_2^* - \bar{x}_2)} & \text{if } x_1, x_2 \in S \text{ or } x_1, x_2 \in D, \\ \frac{-\text{cov}(x_1, x_2)}{n^2 (x_1^* - \bar{x}_1) (x_2^* - \bar{x}_2)} & \text{otherwise.} \end{cases} \tag{15}$$

2. The absolute value of the Pearson correlation coefficient is preserved:

$$\rho(u_1, u_2) \stackrel{\text{def}}{=} \frac{\text{cov}(u_1, u_2)}{S(u_1) \cdot S(u_2)} = \begin{cases} \rho(x_1, x_2) & \text{if } x_1, x_2 \in S \text{ or } x_1, x_2 \in D, \\ -\rho(x_1, x_2) & \text{otherwise.} \end{cases} \tag{16}$$

3 DISCUSSION ON PATTERN TRANSFORMATION

The transformation described by Formula (2) can be called normalization, because it makes variables comparable (1) and has expected properties. First, it preserves two important characteristics of variable distribution – skewness (6) and kurtosis (8). Second, this conversion does not disrupt linear relation between variables – the absolute value of the Pearson correlation coefficient does not change (2).

Unlike other methods the pattern normalization is not just a technical procedure, it has clear interpretation. u_i specifies the share of distance between the i -th object and the pattern in the total distance of all objects from the pattern. We can say that by pattern normalizing we get a relative assessment of the objects situations.

The values of the variable u characterize the positions of objects in the whole system of objects. The value u_i is influenced by all the values of the variable x , so it is important in which environment (a reference group) the normalization is carried out. This is particularly important when analyzing changes of u over time. For one reference group, normalized values for i -th object can increase, and for another group they can decrease. Such situations are presented in an empirical example described in Section 4.

Similar property occurs for standardization:

$$z_i = \frac{x_i - \bar{x}}{S(x)}, \quad \text{for } i = 1, \dots, n, \tag{17}$$

where the reference group is represented by the arithmetic mean \bar{x} and standard deviation $S(x)$, calculated on the basis of the values for all objects. However, it is different for scaling (or min-max normalization):

$$s_i = \frac{x_i - \min_i x_i}{\max_i x_i - \min_i x_i}, \quad \text{for } i = 1, \dots, n. \tag{18}$$

In this case only the maximum and minimum values represent the environment and influence the values of the variables after transformation.

A major advantage of normalization with respect to pattern appears in dynamic approach. In the case of other types of normalization, if we transform the variable for each time unit separately (i.e. we use

a static approach for each unit of time), the results are not comparable over time. To achieve comparability over time, two ways are possible. Firstly, the parameters needed for normalization (e.g. average, deviation, extreme values) can be determined on the basis of all observations (in space and time). Secondly, some reference values for these parameters can be established, that are common to all objects and all (also future) units of time (this can be done on the basis of expert knowledge).

In the case of pattern normalization, we obtain comparability over time using a static approach, because for each unit of time, we distribute the same "mass" (equal to 1) between the same number of objects. For a given object if value of a normalized variable increases, it means that this object increases its share in the total distance from the pattern. So in comparison to other objects, it moves away from "the best" object, so its relative situation is getting worse. Although current data are the sole data used to convert variables, after normalization variables are naturally comparable over time.

The property mentioned above is very advantageous when creating dynamic synthetic variables (composite indicators). The results obtained for a certain time interval are permanent and do not require recalculation after the appearance of observations for the next time period.

4 EMPIRICAL EXAMPLE

To illustrate the effects of pattern normalization, indicators monitoring implementation of the Europe 2020 Strategy (European Commission, 2010) are used. Data come from the statistical office of Poland (Statistics Poland, 2018). 4 stimulants and 7 destimulants are transformed. They are:

- x_1 – Gross domestic expenditure on R&D (% of GDP; $x_1 \in S$),
- x_2 – Early leavers from education and training (%; $x_2 \in D$),
- x_3 – Tertiary educational attainment of persons aged 30–34 (%; $x_3 \in S$),
- x_4 – Greenhouse gas emissions (1990 = 100; $x_4 \in D$),
- x_5 – Share of renewables in gross final energy consumption (%; $x_5 \in S$),
- x_6 – Consumption of primary energy (kg of oil equivalent per 1 000 EUR of GDP; $x_6 \in D$),
- x_7 – Employment rate of persons aged 20–64 (%; $x_7 \in S$),
- x_8 – Share of people at risk-of-poverty or social exclusion (%; $x_8 \in D$),
- x_9 – People living in households with very low work intensity (%; $x_9 \in D$),
- x_{10} – People at risk-of-poverty rate (after social transfers) (%; $x_{10} \in D$),
- x_{11} – Severely materially deprived people (%; $x_{11} \in D$).

The pattern normalization is carried out for two years: 2010 and 2015, as well as in two environments: the entire European Union (abbr. EU28) and 10 countries that joined the EU in 2004 (abbr. EU10). Table 1 and Table 2 show the characteristics of indicators before and after normalization.

Table 1 Characteristics of indicators before (abbr. raw) and after (abbr. norm) normalization in both environments (reference groups) EU28 and EU10 – year 2010

| Indicator | | Reference group | Max | Min | Mean | Standard deviation | Skewness | Kurtosis |
|-------------|------|-----------------|-------|-------|-------|--------------------|----------|----------|
| $x_1 \in S$ | raw | EU28 | 3.730 | 0.450 | 1.514 | 0.893 | 0.759 | -0.285 |
| | norm | | 0.053 | 0.000 | 0.036 | 0.014 | -0.759 | -0.285 |
| | raw | EU10 | 2.060 | 0.450 | 0.991 | 0.497 | 0.913 | -0.389 |
| | norm | | 0.151 | 0.000 | 0.100 | 0.047 | -0.913 | -0.389 |

| Table 1 | | (continuation) | | | | | | |
|-------------|------|-----------------|---------|---------|---------|--------------------|----------|----------|
| Indicator | | Reference group | Max | Min | Mean | Standard deviation | Skewness | Kurtosis |
| $x_2 \in D$ | raw | EU28 | 28.300 | 4.700 | 12.168 | 6.306 | 1.176 | 0.878 |
| | norm | | 0.113 | 0.000 | 0.036 | 0.030 | 1.176 | 0.878 |
| | raw | EU10 | 23.800 | 4.700 | 9.910 | 5.589 | 1.286 | 1.073 |
| | norm | | 0.367 | 0.000 | 0.100 | 0.107 | 1.286 | 1.073 |
| $x_3 \in S$ | raw | EU28 | 50.100 | 18.300 | 34.325 | 9.894 | -0.135 | -1.506 |
| | norm | | 0.072 | 0.000 | 0.036 | 0.022 | 0.135 | -1.506 |
| | raw | EU10 | 45.300 | 20.400 | 32.220 | 8.743 | 0.057 | -1.403 |
| | norm | | 0.190 | 0.000 | 0.100 | 0.067 | -0.057 | -1.403 |
| $x_4 \in D$ | raw | EU28 | 163.770 | 43.200 | 90.637 | 27.493 | 0.271 | 0.072 |
| | norm | | 0.091 | 0.000 | 0.036 | 0.021 | 0.271 | 0.072 |
| | raw | EU10 | 163.770 | 43.200 | 83.126 | 36.835 | 0.953 | -0.185 |
| | norm | | 0.302 | 0.000 | 0.100 | 0.092 | 0.953 | -0.185 |
| $x_5 \in S$ | raw | EU28 | 47.200 | 1.000 | 15.857 | 10.765 | 0.889 | 0.485 |
| | norm | | 0.053 | 0.000 | 0.036 | 0.012 | -0.889 | 0.485 |
| | raw | EU10 | 30.400 | 1.000 | 14.370 | 8.632 | 0.348 | -0.882 |
| | norm | | 0.183 | 0.000 | 0.100 | 0.054 | -0.348 | -0.882 |
| $x_6 \in D$ | raw | EU28 | 464.900 | 82.400 | 191.943 | 93.795 | 1.269 | 1.134 |
| | norm | | 0.125 | 0.000 | 0.036 | 0.031 | 1.269 | 1.134 |
| | raw | EU10 | 417.900 | 142.000 | 250.140 | 75.287 | 0.512 | 0.340 |
| | norm | | 0.255 | 0.000 | 0.100 | 0.070 | 0.512 | 0.340 |
| $x_7 \in S$ | raw | EU28 | 78.100 | 59.900 | 68.136 | 5.334 | 0.239 | -1.188 |
| | norm | | 0.065 | 0.000 | 0.036 | 0.019 | -0.239 | -1.188 |
| | raw | EU10 | 75.000 | 59.900 | 66.000 | 4.496 | 0.489 | -0.576 |
| | norm | | 0.168 | 0.000 | 0.100 | 0.050 | -0.489 | -0.576 |
| $x_8 \in D$ | raw | EU28 | 49.200 | 14.400 | 24.575 | 8.202 | 1.222 | 1.221 |
| | norm | | 0.122 | 0.000 | 0.036 | 0.029 | 1.222 | 1.221 |
| | raw | EU10 | 38.200 | 14.400 | 25.070 | 6.994 | 0.409 | -0.806 |
| | norm | | 0.223 | 0.000 | 0.100 | 0.066 | 0.409 | -0.806 |

| Indicator | | Reference group | Max | Min | Mean | Standard deviation | Skewness | Kurtosis |
|----------------|------|-----------------|--------|-------|--------|--------------------|----------|----------|
| $x_9 \in D$ | raw | EU28 | 22.900 | 4.900 | 9.657 | 3.442 | 1.913 | 5.314 |
| | norm | | 0.135 | 0.000 | 0.036 | 0.026 | 1.913 | 5.314 |
| | raw | EU10 | 12.600 | 4.900 | 8.570 | 2.269 | 0.322 | -0.736 |
| | norm | | 0.210 | 0.000 | 0.100 | 0.062 | 0.322 | -0.736 |
| $x_{10} \in D$ | raw | EU28 | 21.600 | 9.000 | 15.957 | 3.455 | 0.045 | -0.981 |
| | norm | | 0.065 | 0.000 | 0.036 | 0.018 | 0.045 | -0.981 |
| | raw | EU10 | 20.900 | 9.000 | 15.190 | 3.610 | 0.087 | -0.896 |
| | norm | | 0.192 | 0.000 | 0.100 | 0.058 | 0.087 | -0.896 |
| $x_{11} \in D$ | raw | EU28 | 45.700 | 0.500 | 10.621 | 10.038 | 1.879 | 3.383 |
| | norm | | 0.159 | 0.000 | 0.036 | 0.035 | 1.879 | 3.383 |
| | raw | EU10 | 27.600 | 5.900 | 13.350 | 7.040 | 0.727 | -0.721 |
| | norm | | 0.291 | 0.000 | 0.100 | 0.094 | 0.727 | -0.721 |

Note: After pattern normalization, the minimum value is always zero, while the mean is always $1/n$, i.e. 0.036 for EU28 and 0.1 for EU10 (compare properties 2.2.2, 2.3.1).

Source: Own calculation

| Indicator | | Reference group | Max | Min | Mean | Standard deviation | Skewness | Kurtosis |
|-------------|------|-----------------|--------|--------|--------|--------------------|----------|----------|
| $x_1 \in S$ | raw | EU28 | 3.270 | 0.480 | 1.610 | 0.823 | 0.606 | -0.864 |
| | norm | | 0.060 | 0.000 | 0.036 | 0.018 | -0.606 | -0.864 |
| | raw | EU10 | 2.200 | 0.480 | 1.208 | 0.523 | 0.490 | -0.757 |
| | norm | | 0.173 | 0.000 | 0.100 | 0.053 | -0.490 | -0.757 |
| $x_2 \in D$ | raw | EU28 | 20.000 | 2.700 | 9.821 | 4.397 | 0.919 | 0.236 |
| | norm | | 0.087 | 0.000 | 0.036 | 0.022 | 0.919 | 0.236 |
| | raw | EU10 | 19.800 | 5.000 | 8.760 | 4.500 | 1.331 | 0.837 |
| | norm | | 0.394 | 0.000 | 0.100 | 0.120 | 1.331 | 0.837 |
| $x_3 \in S$ | raw | EU28 | 57.600 | 25.300 | 40.496 | 9.077 | -0.066 | -1.069 |
| | norm | | 0.067 | 0.000 | 0.036 | 0.019 | 0.066 | -1.069 |
| | raw | EU10 | 57.600 | 27.800 | 40.610 | 9.915 | 0.246 | -1.095 |

| Table 2 | | (continuation) | | | | | | |
|----------------|------|-----------------|---------|--------|---------|--------------------|----------|----------|
| Indicator | | Reference group | Max | Min | Mean | Standard deviation | Skewness | Kurtosis |
| $x_3 \in S$ | norm | EU10 | 0.175 | 0.000 | 0.100 | 0.058 | -0.246 | -1.095 |
| $x_4 \in D$ | raw | EU28 | 144.450 | 41.990 | 80.583 | 24.071 | 0.489 | 0.135 |
| | norm | | 0.095 | 0.000 | 0.036 | 0.022 | 0.489 | 0.135 |
| | raw | EU10 | 144.450 | 41.990 | 73.372 | 30.371 | 1.057 | 0.356 |
| | norm | | 0.326 | 0.000 | 0.100 | 0.097 | 1.057 | 0.356 |
| $x_5 \in S$ | raw | EU28 | 53.900 | 5.000 | 19.811 | 11.697 | 0.944 | 0.560 |
| | norm | | 0.051 | 0.000 | 0.036 | 0.012 | -0.944 | 0.560 |
| | raw | EU10 | 37.600 | 5.000 | 18.270 | 9.491 | 0.615 | -0.603 |
| | norm | | 0.169 | 0.000 | 0.100 | 0.049 | -0.615 | -0.603 |
| $x_6 \in D$ | raw | EU28 | 448.500 | 62.000 | 165.468 | 85.366 | 1.538 | 2.563 |
| | norm | | 0.133 | 0.000 | 0.036 | 0.029 | 1.538 | 2.563 |
| | raw | EU10 | 358.000 | 90.500 | 208.490 | 67.861 | 0.400 | 0.482 |
| | norm | | 0.227 | 0.000 | 0.100 | 0.058 | 0.400 | 0.482 |
| $x_7 \in S$ | raw | EU28 | 80.500 | 54.900 | 69.936 | 5.796 | -0.469 | 0.056 |
| | norm | | 0.087 | 0.000 | 0.036 | 0.020 | 0.469 | 0.056 |
| | raw | EU10 | 76.500 | 67.700 | 70.630 | 3.160 | 0.636 | -1.170 |
| | norm | | 0.150 | 0.000 | 0.100 | 0.054 | -0.636 | -1.170 |
| $x_8 \in D$ | raw | EU28 | 41.300 | 14.000 | 24.318 | 6.743 | 0.675 | -0.204 |
| | norm | | 0.094 | 0.000 | 0.036 | 0.023 | 0.675 | -0.204 |
| | raw | EU10 | 30.900 | 14.000 | 23.890 | 5.240 | -0.369 | -0.983 |
| | norm | | 0.171 | 0.000 | 0.100 | 0.053 | -0.369 | -0.983 |
| $x_9 \in D$ | raw | EU28 | 19.200 | 5.700 | 10.343 | 3.256 | 0.945 | 0.343 |
| | norm | | 0.104 | 0.000 | 0.036 | 0.025 | 0.945 | 0.343 |
| | raw | EU10 | 10.900 | 6.600 | 8.130 | 1.375 | 0.633 | -0.861 |
| | norm | | 0.281 | 0.000 | 0.100 | 0.090 | 0.633 | -0.861 |
| $x_{10} \in D$ | raw | EU28 | 25.400 | 9.700 | 17.061 | 3.940 | 0.208 | -0.876 |
| | norm | | 0.076 | 0.000 | 0.036 | 0.019 | 0.208 | -0.876 |

| Indicator | | Reference group | Max | Min | Mean | Standard deviation | Skewness | Kurtosis |
|----------------|------|-----------------|--------|-------|--------|--------------------|----------|----------|
| $x_{10} \in D$ | raw | EU10 | 22.500 | 9.700 | 16.760 | 4.080 | 0.003 | -1.037 |
| | norm | | 0.181 | 0.000 | 0.100 | 0.058 | 0.003 | -1.037 |
| $x_{11} \in D$ | raw | EU28 | 45.700 | 0.500 | 10.621 | 10.038 | 1.879 | 3.383 |
| | norm | | 0.159 | 0.000 | 0.036 | 0.035 | 1.879 | 3.383 |
| | raw | EU10 | 27.600 | 5.900 | 13.350 | 7.040 | 0.727 | -0.721 |
| | norm | | 0.291 | 0.000 | 0.100 | 0.094 | 0.727 | -0.721 |

Note: After pattern normalization, the minimum value is always zero, while the mean is always 1/n, i.e. 0.036 for EU28 and 0.1 for EU10 (compare properties 2.2.2, 2.3.1).

Source: Own calculation

Poland is selected as an example. Table 3 compares the results of normalization for both years. An influence of normalization itself and normalization environment on the dynamic image of the country are examined. That is, for a given object (Poland) we analyze what happens to the normalized value of indicator if the raw value improves (or gets worse). These aspects are important when comparing the pattern normalization with scaling (min-max normalization). Scaling, the most popular method of dynamic normalization, in this context can be called neutral. It does not affect the dynamic image of objects, moreover, the environment of scaling does not matter.

| Indicator | | 2010 | 2015 | 2015 to 2010 |
|-------------|------|--------|--------|--------------|
| $x_1 \in S$ | raw | 0.720 | 1.000 | + |
| | EU28 | 0.049 | 0.049 | - |
| | EU10 | 0.125 | 0.121 | + |
| | rank | 5 | 4 | + |
| $x_2 \in D$ | raw | 5.400 | 5.300 | + |
| | EU28 | 0.003 | 0.013 | - |
| | EU10 | 0.013 | 0.008 | + |
| | rank | 4 | 3 | + |
| $x_3 \in S$ | raw | 34.800 | 43.400 | + |
| | EU28 | 0.035 | 0.030 | + |
| | EU10 | 0.080 | 0.084 | - |
| | rank | 6 | 6 | 0 |

Table 3

(continuation)

| Indicator | | 2010 | 2015 | 2015 to 2010 |
|----------------|------|---------|---------|--------------|
| $x_4 \in D$ | raw | 87.170 | 82.760 | + |
| | EU28 | 0.033 | 0.038 | - |
| | EU10 | 0.110 | 0.130 | - |
| | rank | 7 | 7 | 0 |
| $x_5 \in S$ | raw | 9.300 | 11.800 | + |
| | EU28 | 0.043 | 0.044 | - |
| | EU10 | 0.132 | 0.133 | - |
| | rank | 4 | 3 | + |
| $x_6 \in D$ | raw | 278.300 | 227.300 | + |
| | EU28 | 0.064 | 0.057 | + |
| | EU10 | 0.126 | 0.116 | + |
| | rank | 8 | 8 | 0 |
| $x_7 \in S$ | raw | 64.300 | 67.800 | + |
| | EU28 | 0.049 | 0.043 | + |
| | EU10 | 0.119 | 0.148 | - |
| | rank | 3 | 2 | + |
| $x_8 \in D$ | raw | 27.800 | 23.400 | + |
| | EU28 | 0.047 | 0.033 | + |
| | EU10 | 0.126 | 0.095 | + |
| | rank | 7 | 5 | + |
| $x_9 \in D$ | raw | 7.300 | 6.900 | + |
| | EU28 | 0.018 | 0.009 | + |
| | EU10 | 0.065 | 0.020 | + |
| | rank | 4 | 3 | + |
| $x_{10} \in D$ | raw | 17.600 | 17.600 | 0 |
| | EU28 | 0.044 | 0.038 | + |
| | EU10 | 0.139 | 0.112 | + |

Table 3

(continuation)

| Indicator | | 2010 | 2015 | 2015 to 2010 |
|----------------|------|--------|-------|--------------|
| $x_{10} \in D$ | rank | 8 | 7 | + |
| | raw | 14.200 | 8.100 | + |
| $x_{11} \in D$ | EU28 | 0.048 | 0.029 | + |
| | EU10 | 0.111 | 0.059 | + |
| | rank | 7 | 4 | + |

Note: + improvement, – deterioration, 0 no changes.

Source: Own calculation

In the analyzed period in Poland, raw values of all indicators except x_{10} improve, i.e. values of stimulants increase, values of destimulants decrease. Changes would be the same after dynamic scaling, but after pattern normalization the changes over time are not so uniform.

From this point of view, the indicators monitoring the implementation of the Europe 2020 Strategy can be divided into three groups. In the first one there are x_6 , x_8 , x_9 , x_{11} . In their case, pattern normalization does not change the dynamics of variables. Raw indicators are improved as well as normalized indicators (for both environments).

The second group are indicators for which normalization changes the "dynamic image" of Poland, but the normalization environment is irrelevant. This group includes x_{10} . The raw value of this indicator does not change, but after normalization it improves in both environments. This means that Poland's objective situation has not improved, but the relative one has (because the situation of other countries in this period has deteriorated). The next in this group are x_4 and x_5 . For them, the impact of normalization is more evident. Although the raw values of these indicators improve, the situation of Poland in both considered environments has got worse. All three indicators after normalization increase.

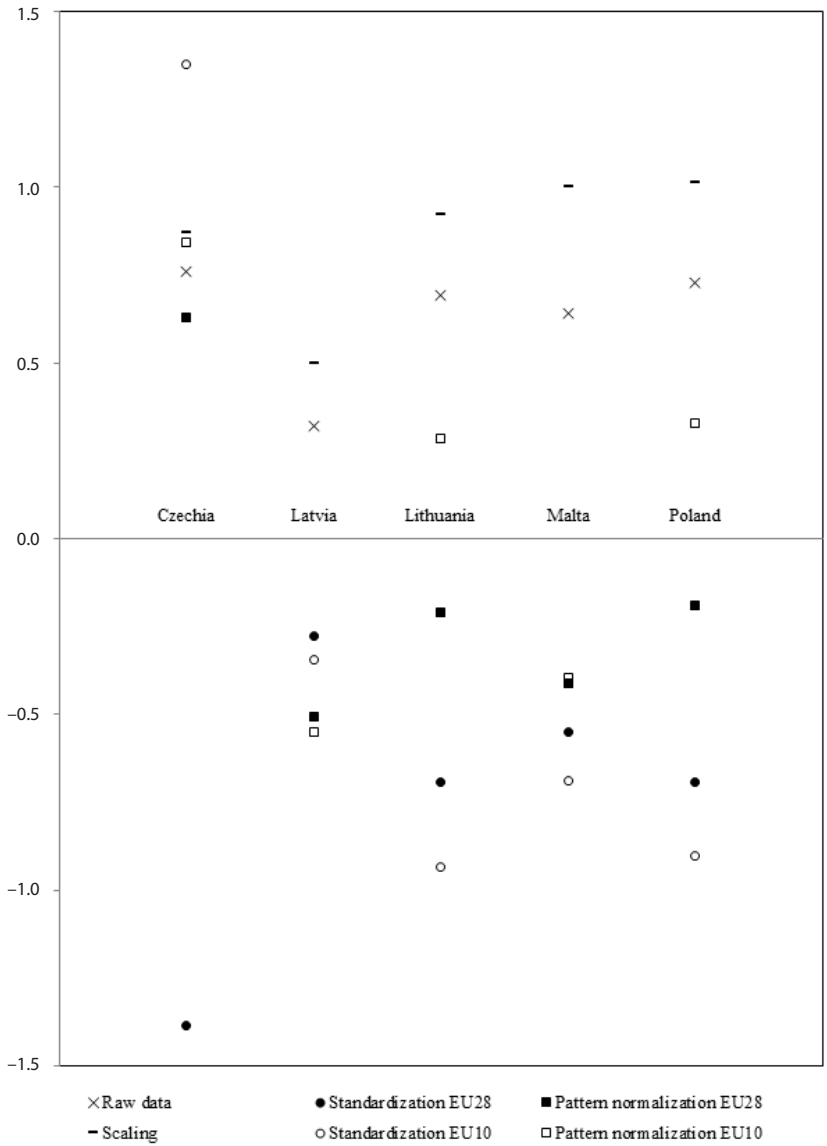
The last group are indicators for which the normalization environment is important. For x_3 and x_7 , Poland has improved against the background of a bigger environment, and it has declined against a smaller one. For x_1 , x_2 the situation is reversed.

It is interesting to confront the above considerations with the analysis of changes in Poland's position in the ranking. For some indicators the relative improvement is so great that the position of Poland in the ranking also improves (e.g. x_8), however the relative change can be insufficient to improve the position (e.g. x_6). There is also a situation (x_5) in which the direction of changes in the normalized variables and positions is reversed.

Next, the pattern normalization is compared with the most popular methods of normalization: standardization (17) and scaling (18). The direction of changes in the values of the normalized variable in 2015 as compared to 2010 is analyzed. Dynamic standardization and dynamic scaling is performed based on data from both years. For selected countries Figure 1 shows a relative increase in the value of variable $x_1 \in S$ in 6 versions: without normalization, after scaling (in this case, the reference group does not matter), after standardization and pattern normalization for both E28 and E10 environments.

For all presented countries, the raw values of the variable x_1 increase, and thus its values after scaling increase as well. This differs the scaling from the pattern normalization and the standardization. For the last two normalizations the direction of changes in transformed values does not necessarily coincide with the direction of changes in raw data. Moreover, for a certain country, the normalized value for one reference group may increase, and for another group it may decrease.

Figure 1 Relative increments of variable x_1 before and after normalization



Note: To make the graph more transparent, the increments are transformed with the cube root. If a country is located above the axis, its situation in 2015 improved compared to 2010, that is, the value of the stimulant increased (x, z, s), and the value of the destimulant decreased (u).

Source: Own calculation

CONCLUSIONS

The article presents a new transformation of diagnostic variables, that plays a double role in analyses of complex phenomenon: it unifies the nature of variables and makes variables comparable. The transformation

is called normalization with respect to the pattern (or pattern normalization in short). The pattern normalization has properties expected for this type of transformation.

The values of variables after normalization with respect to pattern characterize the relative situation of the objects, i.e. the situation on the background of the environment in which the research is carried out. Changing the environment can change the research results. This feature is both an advantage and the biggest disadvantage of the proposed method. The pattern normalization can only be used in research in which the context of the environment is important. "Objective" changes may be distorted during this transformation.

A main advantage of new normalization is the possibility of use in dynamic analysis (i.e. for different time units). However, it is not necessary to re-calculate results with the appearance of observations for next unit of time, as, for example, in the case of dynamic standardization.

The effects of normalization with respect to pattern are illustrated by an empirical example. Indicators monitoring the implementation of the Europe 2020 Strategy are normalized. Normalizations are carried out for two environments: the entire EU and countries that joined the EU in 2004. The results for two years are compared. The example of Poland shows that the "dynamic image" of the country is affected by the use of normalization itself as well as by the choice of the environment in normalization.

Pattern normalization can be used in common construction of composite indicators instead of other methods of normalization. A possible applications are shown in Müller-Fraćzek (2017, 2018).

The proposed construction can have various modifications. First of all, we can change the mass distributed between objects (for example to n). We can also change the measure of distance or the method of choosing the pattern.

References

- EUROPEAN COMMISSION. *Europe 2020. A strategy for smart, sustainable and inclusive growth* [online]. Brussels: European Commission, 2010. [cit. 20.12.2018] <<https://eur-lex.europa.eu/LexUriServ/LexUriServ.do?uri=COM:2010:2020:FIN:EN:PDF>>.
- HELLWIG, Z. Zastosowanie metody taksonomicznej do typologicznego podziału krajów ze względu na poziom ich rozwoju i strukturę kwalifikowanych kadr. *Przegląd Statystyczny*, 1968, 15(4), pp. 307–327.
- JAJUGA, K. AND WALESIAK, M. Standardisation of Data Set under Different Measurement Scales. In: DECKER, W. AND GAUL, W. eds. *Classification and Information Processing at the Turn of the Millennium*, Berlin, Heidelberg: Springer, 2000.
- MAZZIOTTA, M. AND PARETO, A. Comparing Two Non-Compensatory Composite Indices to Measure Changes over Time: a Case Study [online]. *Statistika: Statistics and Economy Journal*, 2015, 95(2), pp. 44–53.
- MAZZIOTTA, M., PARETO, A. On a generalized non-compensatory composite index for measuring socio-economic phenomena. *Social Indicators Research*, 2016, 127(3), pp. 983–1003.
- MILLIGAN, G. AND COOPER, M. A study of standardization of variables in cluster analysis. *Journal of Classification*, 1988, 5(2), pp. 181–204.
- MŁODAK, A. Multilateral normalizations of diagnostic features. *Statistics in Transition*, 2006, 7(5), pp. 1125–1139.
- MÜLLER-FRĄCZEK, I. Propozycja miary syntetycznej. *Przegląd Statystyczny*, 2017, 64(4), pp. 413–428.
- MÜLLER-FRĄCZEK, I. Dynamic measure of development. In: PAPIEŻ, M. AND ŚMIECH, S. eds. *The 12th Professor Aleksander Zelias International Conference on Modelling and Forecasting of Socio-Economic Phenomena. Conference Proceedings*, Cracow: Foundation of the Cracow University of Economics, 2018.
- NARDO, M., SAISANA, M., SALTELLI, A., TARANTOLA, S., HOFFMAN, A., GIOVANNINI, E. *Handbook on Constructing Composite Indicators*. OECD publishing, 2005.
- SAISANA, M. AND SALTELLI, A. Rankings and Ratings: Instructions for use. *Hague Journal on the Rule of Law*, 2011, 3.2, pp. 247–268.
- STATISTICS POLAND. *Europe 2020 indicators* [online]. Warsaw, 2018. [cit. 20.12.2018] <<https://stat.gov.pl/en/international-statistics/international-comparisons/tables-about-countries-by-subject/europe-2020-indicators>>.

ANNEX

Proof of (3):

$$u_i = 0 \Leftrightarrow \frac{|x_i - x^+|}{\sum_{j=1}^n |x_j - x^+|} = 0 \Leftrightarrow |x_i - x^+| = 0 \Leftrightarrow x_i = x^+.$$

Proof of (5):

$$u_i = 1 \Leftrightarrow \frac{|x_i - x^+|}{\sum_{j=1}^n |x_j - x^+|} = 1 \Leftrightarrow |x_i - x^+| = \sum_{j=1}^n |x_j - x^+| \Leftrightarrow \forall_{j \neq i} x_j = x^+.$$

Proof of (6):

$$\text{If } x \in S, \text{ then: } \max_i u = \frac{\min_i(x^+ - x_i)}{\sum_{j=1}^n (x^+ - x_j)} = \frac{x^+ - \min_i x_i}{\sum_{j=1}^n (x^+ - x_j)} = \frac{\max_i x_i - \min_i x_i}{\sum_{j=1}^n (\max_i x_i - x_j)}.$$

$$\text{If } x \in D, \text{ then: } \max_i u = \frac{\max_i(x_i - x^+)}{\sum_{j=1}^n (x_j - x^+)} = \frac{\max_i x_i - x^+}{\sum_{j=1}^n (x_j - x^+)} = \frac{\max_i x_i - \min_i x_i}{\sum_{j=1}^n (x_j - \max_i x_i)}.$$

Proof of (7):

$$\bar{u} = \frac{1}{n} \sum_{i=1}^n \frac{|x_i - x^+|}{\sum_{j=1}^n |x_j - x^+|} = \frac{1}{n} \frac{\sum_{i=1}^n |x_i - x^+|}{\sum_{j=1}^n |x_j - x^+|} = \frac{1}{n}.$$

Proof of (8): Assume that $x \in S$, but the proof is similar when $x \in D$.

$$\begin{aligned} S^2(u) &= \frac{1}{n} \sum_{i=1}^n \left(\frac{x^+ - x_i}{\sum_{j=1}^n (x^+ - x_j)} - \frac{1}{n} \right)^2 = \frac{1}{n^3} \sum_{i=1}^n \left(\frac{x^+ - x_i}{x^+ - \frac{1}{n} \sum_{j=1}^n x_j} - 1 \right)^2 = \frac{1}{n^3} \sum_{i=1}^n \left(\frac{x^+ - x_i}{x^+ - \bar{x}} - 1 \right)^2 \\ &= \frac{1}{n^3} \sum_{i=1}^n \left(\frac{\bar{x} - x_i}{x^+ - \bar{x}} \right)^2 = \frac{\frac{1}{n} \sum_{i=1}^n (\bar{x} - x_i)^2}{n^2 (x^+ - \bar{x})^2} = \frac{S^2(x)}{n^2 (x^+ - \bar{x})^2}. \end{aligned}$$

Proof of (11): Assume that $x \in S$, but the proof is similar when $x \in D$.

$$\begin{aligned} \mu_3(u) &= \frac{1}{n} \sum_{i=1}^n \left(\frac{x^+ - x_i}{\sum_{j=1}^n (x^+ - x_j)} - \frac{1}{n} \right)^3 = \frac{1}{n^4} \sum_{i=1}^n \left(\frac{x^+ - x_i}{x^+ - \frac{1}{n} \sum_{j=1}^n x_j} - 1 \right)^3 \\ &= \frac{1}{n^4} \sum_{i=1}^n \left(\frac{x^+ - x_i}{x^+ - \bar{x}} - 1 \right)^3 = \frac{1}{n^4} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\bar{x} - x^+} \right)^3 = \frac{\mu_3(x)}{n^3 (\bar{x} - x^+)^3}. \end{aligned}$$

Proof of (13): Assume that $x \in S$, but the proof is similar when $x \in D$.

$$\begin{aligned} \mu_4(u) &= \frac{1}{n} \sum_{i=1}^n \left(\frac{x^+ - x_i}{\sum_{j=1}^n (x^+ - x_j)} - \frac{1}{n} \right)^4 = \frac{1}{n^5} \sum_{i=1}^n \left(\frac{x^+ - x_i}{x^+ - \frac{1}{n} \sum_{j=1}^n x_j} - 1 \right)^4 \\ &= \frac{1}{n^5} \sum_{i=1}^n \left(\frac{x^+ - x_i}{x^+ - \bar{x}} - 1 \right)^4 = \frac{1}{n^5} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{x^+ - \bar{x}} \right)^4 = \frac{\mu_4(x)}{n^4(x^+ - \bar{x})^4}. \end{aligned}$$

Proof of (15): Assume that x_1 and x_2 are stimulants. The proof in other cases is similar.

$$\begin{aligned} \text{cov}(u_1, u_2) &= \frac{1}{n} \sum_{i=1}^n \left(\frac{x^+ - x_{1i}}{\sum_{j=1}^n (x_1^+ - x_{1j})} - \frac{1}{n} \right) \left(\frac{x_2^+ - x_{2i}}{\sum_{j=1}^n (x_2^+ - x_{2j})} - \frac{1}{n} \right) \\ &= \frac{1}{n^3} \sum_{i=1}^n \left(\frac{x_1^+ - x_{1i}}{x_1^+ - \frac{1}{n} \sum_{j=1}^n x_{1j}} - 1 \right) \left(\frac{x_2^+ - x_{2i}}{x_2^+ - \frac{1}{n} \sum_{j=1}^n x_{2j}} - 1 \right) \\ &= \frac{1}{n^3} \sum_{i=1}^n \left(\frac{\bar{x}_1 - x_{1i}}{x_1^+ - \bar{x}_1} \cdot \frac{\bar{x}_2 - x_{2i}}{x_2^+ - \bar{x}_2} \right) = \frac{\frac{1}{n} \sum_{i=1}^n (\bar{x}_1 - x_{1i})(\bar{x}_2 - x_{2i})}{n^2(x_1^+ - \bar{x}_1)(x_2^+ - \bar{x}_2)} = \frac{\text{cov}(x_1, x_2)}{n^2(x_1^+ - \bar{x}_1)(x_2^+ - \bar{x}_2)}. \end{aligned}$$