

# CLUSTER ANALYSIS OF THE CZECH REGIONS BY WAGE LEVEL

Diana Bílková

---

## Abstract

Income levels of the population have been constantly researched by economists in the developed countries mainly due to their connection with the living standards of the population. Knowledge of the wage distribution and its comparison from various socio-economic and time-spatial aspects is a precondition for the assessment of living standards, social security and equality in the division of material values produced by the society. Statistical analysis of the wage distribution also forms the basis for government social policy, taxation, budgetary and other decisions. Moreover, the direct connection between wages and the purchasing power of the population brings tracking the level, structure and development of the wage distribution to the foreground when identifying sales opportunities for the products of both long- and short-term consumption.

The present paper deals with a comparison of wage levels of fourteen regions in the Czech Republic. Similar wage-level clusters were created using the methods of cluster analysis. Three regions with the highest and lowest wage levels, respectively, were selected. For these six regions, the model wage distribution was presented to enable the comparison of wage development over the past seven years. Three-parameter lognormal curves represent the basis of the theoretical wage distribution.

**Key words:** Cluster analysis, method of the furthest neighbour, Euclidean distance matrix, wage models, Akaike and Bayesian information criterions

**JEL Code:** J31, D31, E24

---

## Introduction

Income levels of the population have been constantly researched by economists in the developed countries mainly due to their connection with the living standards of the population. Knowledge of the wage distribution and its comparison from various socio-economic and time-spatial aspects is a precondition for the assessment of living standards,

social security and equality in the division of material values produced by the society. Statistical analysis of the wage distribution also forms the basis for government social policy, taxation and budgetary decisions.

The primary aim of this paper is to create the clusters of Czech Republic regions based on wage level similarities. The farthest neighbour clustering method and the Euclidean distance metric were used for the creation of those clusters. The key objective lies in modelling the wage distribution in individual regions so that its development since 2009 can be traced. Three-parameter lognormal curves represent the basis of the model distribution with parameters estimated employing the maximum likelihood method. The accuracy of the models obtained is assessed applying Akaike and Bayesian Information criteria.

## **1 Review of literature**

Labour market analysis, wages and incomes of the population as well as the gender wage gap remain the current focus of extensive research interest. The following selection of publications provides some examples of the relevant research output. Albelda, and Carr (2014) track the proportion of low-income workers who benefit from employee and welfare programmes in the United States over the period 1979–2011, examining the dependence on workers' gender and marital status. Bárány (2016) demonstrates the relationship between the minimum wage development and income inequality. Bartošová, and Želinský (2013) give an overview of attempts to measure poverty in former Czechoslovakia, they further analysing monetary poverty, relative material deprivation and subjective perception of poverty in both the successor states using EU-SILC microdata. Domínguez-Villalobos, and Brown-Grossman (2010) investigate the effects of industrial restructuring projects on male-female wage inequality in Mexico between 2001 and 2005, pointing to the negative impact of export orientation on the remuneration of both sexes, women, however, being disadvantaged in both absolute and relative proportions. Fisher, Johnson, and Smeeding (2015) explore the distribution of both income and consumption in the United States using a representative sample of individuals yielded from a single data set. Gobillon, Meurs, and Roux (2015) propose a new measurement rate of the employment discrimination on the grounds of gender based on the assignment work plan, suggesting the probability male-female ratio of employment at each step of the income ladder. Jenderny (2016) carries out an analysis of upper-income mobility in Germany using panel microdata from personal tax returns. Kukk, and Staehr (2014) estimate the income range of business households in relation to those of

manual workers in Estonia. Malá (2015) constructs multidimensional probabilistic models of the income distribution of Czech households.

## **2 Database**

Data for the present study are collected from the official website of the Czech Statistical Office (CSO).

The CSO database contains the total wage distribution for the period 2009–2015 covering all employees in the Czech Republic broken down by regions. In individual regions, employees were divided into groups of men and women and the gender wage gap was examined. Annual data are related to gross monthly nominal wages in the respective years, the average (median) wage, for instance, representing average (median) gross monthly wage over the year.

There are also data in the form of the interval frequency distribution with uneven and extreme open intervals. Neither more detailed nor individual data have been currently available. Since only nominal wage data are provided by the CSO, the obtained average and median nominal wages had to be converted to average and median real wages using the CSO-reported inflation rate data.

Only the data on nominal wages having been available, inflation rates had to be used for the conversion to a real wage that reflects purchasing power allowing for a comparison of the wage development without inflation effects in the research period. The rate of inflation is derived from the consumer price index (CPI), based on the Laspeyres price index. The real wage was calculated using the real wage index, the nominal wage index being divided by the CPI (living cost index). The data were processed utilizing SAS and Statgraphics statistical programme packages and Microsoft Excel spreadsheets.

The research data include wages and salaries paid to employees for work performed in the private (business) and public (state budget, non-business) sectors, respectively. In terms of the data presented on the CSO website, “wages” cover remuneration for work done in both the sectors.

# 1 Theory and methods

## 3.1 Cluster analysis

Cluster analysis was used to divide the Czech regions into relatively homogeneous groups according to their respective gross monthly wage levels. Multivariate data analysis, which is often done to process economic data (see, e.g. Lukáš Malec (2016)), may include other approaches to statistical data analysis, namely that of canonical correlation. (Particular aspects of cluster analysis are dealt with in Elena Makhalova, and Iva Pecáková (2015) or Hana Řezanková, and Tomáš Löster (2013).)

Multidimensional observations can be used when classifying a set of objects into several relatively homogeneous clusters. We have a data matrix  $X$  of  $n \times p$  type, where  $n$  is the number of objects and  $p$  is the number of variables. Assuming various decompositions  $S^{(k)}$  of the set of  $n$  objects into  $k$  clusters, we look for the most appropriate decompositions. The aim is to find the objects within certain clusters that are as similar as possible to those from other clusters. Only decompositions with disjunctive clusters and tasks with a specified number of classes are conceded.

### 3.1.1 Criteria for assessing the quality of decomposition

The general task is to assess to what extent the cluster analysis aim has been achieved in a given situation, while applying a specific algorithm. Several criteria – decomposition functions – are proposed for this purpose. The most frequently used ones exhibit the following characteristics. They are the matrices of internal cluster variance

$$\mathbf{E} = \sum_{h=1}^k \sum_{i=1}^{n_h} (\mathbf{x}_{hi} - \bar{\mathbf{x}}_h) \cdot (\mathbf{x}_{hi} - \bar{\mathbf{x}}_h)' \quad (1)$$

and between-cluster variance

$$\mathbf{B} = \sum_{h=1}^k n_h \cdot (\bar{\mathbf{x}}_h - \bar{\mathbf{x}}) \cdot (\bar{\mathbf{x}}_h - \bar{\mathbf{x}})', \quad (2)$$

whose sum is the matrix of total variation

$$\mathbf{T} = \sum_{h=1}^k \sum_{i=1}^{n_h} (\mathbf{x}_{hi} - \bar{\mathbf{x}}) \cdot (\mathbf{x}_{hi} - \bar{\mathbf{x}})'. \quad (3)$$

There are vectors of the observations for the  $i^{\text{th}}$  object and  $h^{\text{th}}$  cluster  $x_{hi}$ , the averages for the  $h^{\text{th}}$  cluster  $\bar{x}_h$  and those for the total set  $\bar{x}$ . There are  $p^{\text{th}}$ -membered vectors,  $\mathbf{E}$ ,  $\mathbf{B}$  and  $\mathbf{T}$  being symmetric square matrices of the  $p^{\text{th}}$  order. The principal aim, consisting in the creation of mutually distant compact clusters, is fulfilled by reaching the minimum of the total sum of the deviation squares of all values of corresponding cluster averages

$$C_1 = \text{st } \mathbf{E} = \sum_{h=1}^k \sum_{i=1}^{n_h} \sum_{j=1}^p (x_{hij} - \bar{x}_{hj})^2, \quad (4)$$

i.e. the Ward criterion. Since the  $\text{st } \mathbf{T}$  is the same for all decompositions, the minimization of the  $\text{st } \mathbf{E}$  means the same as that of the  $\text{st } \mathbf{B}$ . In order to become independent on the used units of measurement (or, more generally, the invariance to the linear transformations), it is recommended to minimize the determinant of the matrix of the internal cluster variance

$$C_2 = |\mathbf{E}|$$

or to maximize the trace criterion

$$C_3 = \text{st}(\mathbf{B}\mathbf{E}^{-1}) \text{ or else } C_4 = \text{st}(\mathbf{B}\mathbf{T}^{-1}).$$

The criteria mentioned above are employed not only retrospectively to assess the decomposition quality accomplished, changes in criterion values also guiding the creation of clusters. Since the criteria ultimately reach the limits ( $C_1$  and  $C_2$  the minimum,  $C_3$  and  $C_4$  the maximum) at  $k = n$ , it is necessary to find the extreme of the purpose function that properly includes the loss following from the growth in the number of clusters. The Ward criterion, for instance, is proposed to move towards the minimization of the quantity

$$Z_1 = C_1 + z \cdot k, \quad (5)$$

where constant  $z$  represents the loss resulting from an increase in the number of clusters by one.

### 3.1.2 Distance and similarity of objects

Having selected the variables characterizing the properties of the clustered objects and found their values, we decided on the method of the evaluation of distance or similarity of objects, the calculation of appropriate measures for all pairs of objects often being the initial stage of clustering algorithm implementation. The symmetric square matrix of  $n \times n$  type has zeros or

ones on the diagonal, depending on whether it is the matrix of distance  $D$  measures or that of similarity  $A$  measures, respectively.

Let us now focus on measuring the distance of the objects described by quantitative variables. The Hemming distance can be used when individual variables are roughly on the same level or at least expressed in the same units of measurement

$$D_H(x_i, x_{i'}) = \sum_{j=1}^p |x_{ij} - x_{i'j}|. \quad (6)$$

The Euclidean distance can be applied in the same case

$$D_E(x_i, x_{i'}) = \sqrt{\sum_{j=1}^p (x_{ij} - x_{i'j})^2} \quad (7)$$

as well as the Chebyshev distance

$$D_C(x_i, x_{i'}) = \max_j |x_{ij} - x_{i'j}|. \quad (8)$$

All the above mentioned measurements have some common drawbacks – the dependence on the used measuring units that sometimes hinders the meaningful acquisition of any sum for different variables and the fact that if the variables are considered in sum with the same weights, the strongly correlated variables have a disproportionately large effect on the outcome. The starting point is the transformation of variables. The adverse effect of the measuring units can be removed by dividing all the values by the balancing factor, which can be presented with the corresponding average  $\bar{x}_j$ , standard deviation  $s_j$  or the range after deletion of extremes

$$\max_i x_{ij} - \min_i x_{ij}.$$

Particular variables can be also assigned more weight – having decided subjectively or on the basis of relevant information – their values then appearing in the formulas for the calculation of distance.

Other measurements of distance and similarity of objects for numerical, ordinal, nominal and alternative variables are described in the professional literature. When dealing with variables of a different type, the Lance-Williams distance is recommended

$$D_{LW}(x_i, x_{i'}) = \frac{\sum_{j=1}^p |x_{ij} - x_{i'j}|}{\sum_{j=1}^p (x_{ij} + x_{i'j})}. \quad (9)$$

### 3.1.3 Algorithm for the creation of hierarchical sequence of decompositions

The creation of a hierarchical sequence of decompositions belongs to the most widely used techniques applied in the cluster analysis, occurring sequentially in the following steps:

- 1)  $\mathbf{D}$  matrix calculation of appropriate measurements of distances;
- 2) the start of the decomposition process  $\mathcal{S}^{(n)}$  from  $n$  clusters, each of them containing one object;
- 3) the assessment of the symmetric matrix  $\mathbf{D}$  (a lower or upper triangle), finding two clusters (the  $h^{\text{th}}$  and  $h'^{\text{th}}$  ones) whose distance  $D_{hh'}$  is minimal;
- 4) the connection of the  $h^{\text{th}}$  and  $h'^{\text{th}}$  clusters into a new  $g^{\text{th}}$  cluster, the replacement of the  $h^{\text{th}}$  and  $h'^{\text{th}}$  row and column in the matrix  $\mathbf{D}$  with those of the new cluster, the order of the matrix being reduced by one;
- 5) renumbering of the order of the cycle  $l = 1, 2, \dots, n - 1$ , the identification of the connected objects  $h, h'$  and the level of the connection  $d_l = D_{hh'}$ ;
- 6) returning to step (3) if the creation of decompositions has not been completed by connecting all objects into a single cluster  $\mathcal{S}^{(1)}$ .

A divisive hierarchical procedure, contrary to the agglomerative hierarchical one, is less-used, starting from a single cluster  $\mathcal{S}^{(1)}$ , splitting one of the clusters into two in each step and obtaining  $\mathcal{S}^{(n)}$  at the end of the process. The results of hierarchical cluster procedures can be effectively displayed in the form of a graphical tree dendrogram.

Given the choice of variables  $x_1, x_2, \dots, x_p$  and the matrix of distances  $\mathbf{D}$ , the results of applying the described algorithm vary according to the way the distance between clusters is evaluated.

### Nearest Neighbour Method

Within the nearest neighbour method, both clusters, whose connection is considered, are represented by objects that are the closest to each other. The  $D_{hh'}$  distance between the  $h^{\text{th}}$  and  $h'^{\text{th}}$  clusters therefore represents the minimum of all  $q = n_h n_{h'}$  distances between their objects, the procedure of the third phase of the above algorithm thus being specified. In the fourth

step, the  $h^{\text{th}}$  and  $h'^{\text{th}}$  rows and columns in the distance matrix are replaced with the new  $g^{\text{th}}$  cluster's row and column of distances. In the  $l^{\text{th}}$  cycle, total  $n - l - 1$  distances determined by

$$D_{gg'} = \min(D_{g'h}, D_{g'h'}) \quad (10)$$

can be written.

If the way of evaluation of the proximity or similarity of clusters is given, which also determines the conversion of the distance matrix in each cycle, the above algorithm allows for the creation of a hierarchical sequence of decompositions and construction of the dendrogram.

When using this method, even considerably distant objects can get together in the same cluster if a large number of other objects create a kind of bridge between them. This typical chaining of objects is considered as a drawback, especially if there is a reason for the clusters to acquire the usual elliptical shape with a compact core. This method, however, possesses many positive features that outweigh the above disadvantage.

### **Farthest Neighbour Method**

The method of the farthest neighbour is based on the opposite principle. The criterion for the connection of clusters is the maximum of  $q$  possible between-cluster distances of objects. When editing the matrix of distances, we proceed according to

$$D_{gg'} = \max(D_{g'h}, D_{g'h'}) . \quad (11)$$

An adverse chain effect does not occur in this case. On the contrary, there is a tendency towards the formation of compact clusters, not extraordinarily large, though.

### **Average Linkage Method (Sokal-Sneath Method)**

As a criterion for the connection of clusters, this method applies an average of the  $q$  possible between-cluster distances of objects. When recalculating the distance matrix, we use

$$D_{gg'} = \frac{n_h \cdot D_{g'h} + n_{h'} \cdot D_{g'h'}}{n_h + n_{h'}} . \quad (12)$$

The method often leads to similar results as the farthest neighbour one.



### Centroid method (Gower method)

Unlike the above methods, this one is not based on summarizing the information on between-cluster distances of objects, the criterion being the Euclidean distance of centroids

$$D_E(\bar{x}_h, \bar{x}_{h'}) = \sum_{j=1}^p (\bar{x}_{hj} - \bar{x}_{h'j})^2 . \quad (13)$$

The recalculation of the distance matrix is done as follows

$$D_{gg'} = \frac{1}{n_h + n_{h'}} \left( n_h \cdot D_{g'h} + n_{h'} \cdot D_{g'h'} - \frac{n_h \cdot n_{h'}}{n_h + n_{h'}} \cdot D_{hh'} \right) . \quad (14)$$

### Ward Method

The method uses a functional of the decomposition quality  $C_1$  in formula (4). The criterion for the cluster connection is an increment to the total intra-group sum of the squares of observation deviations from the cluster average, thus

$$\Delta C_1 = \sum_{i=1}^g \sum_{j=1}^p (x_{gij} - \bar{x}_{gi})^2 - \sum_{i=1}^h \sum_{j=1}^p (x_{hij} - \bar{x}_{hj})^2 - \sum_{i=1}^{h'} \sum_{j=1}^p (x_{h'ij} - \bar{x}_{h'j})^2 . \quad (15)$$

The increment is expressed as a sum of squares in an emerging cluster which is reduced by the sums of squares in both vanishing clusters. Using arithmetic modifications, the expression can be simplified into the form

$$\Delta C_1 = \frac{n_h \cdot n_{h'}}{n_h + n_{h'}} \cdot \sum_{j=1}^p (\bar{x}_{hj} - \bar{x}_{h'j})^2 . \quad (16)$$

This equation is a product of the Euclidean distance between the centroids of clusters considered for the connection and a coefficient depending on the cluster size. The value of this coefficient grows with an increasing size of clusters, and for fixed  $n_h + n_{h'}$  it represents the maximum in the case of same-size ( $n_h = n_{h'}$ ) clusters. Since we create the connections to ensure the minimization of the criterion  $\Delta C_1$ , the Ward method tends to eliminate small clusters, i.e. to form those of roughly the same size, which is often a desirable property. Starting from the matrix of Euclidean distances between objects in the process of its modification, we can use the formula

$$D_{gg'} = \frac{1}{n_h + n_{h'} + n_{g'}} \cdot [(n_h + n_{g'}) \cdot D_{hg'} + (n_{h'} + n_{g'}) \cdot D_{h'g'} - n_{g'} \cdot D_{hh'}] . \quad (17)$$

### 3.2 Other Methods Used

Three-parameter lognormal curves represent the basis of the theoretical wage distribution. The minimum wage in the year was considered as the beginning of the distribution. Remaining two parameters were estimated using maximum likelihood method. The accuracy of the models obtained was compared applying the Akaike and Bayesian information criteria, both of which take a number of the corresponding wage model parameters into account.

Let us consider  $L$  as the maximum value of the likelihood function for an assumed model of data,  $k$  and  $n$  denoting the number of parameters estimated and the sample size, respectively. The Akaike information criterion (AIC) has the form

$$AIC = 2k - 2\ln L \quad (18)$$

and the Bayesian information criterion (BIC) is defined as

$$BIC = k \ln n - 2\ln L. \quad (19)$$

The model with minimal AIC or BIC values is preferred over other alternatives, AIC and BIC criteria also including a penalty which is an increasing function of the number of estimated parameters.

## 4 Results

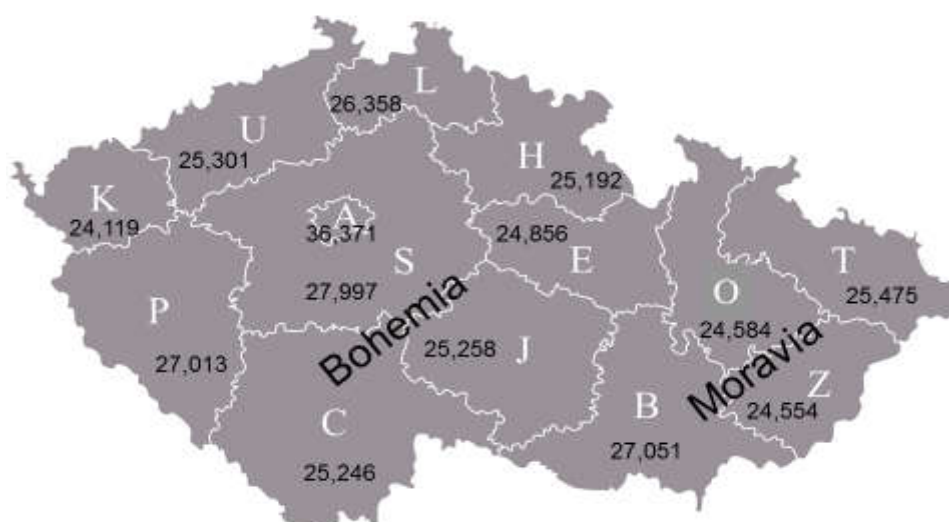
Figures 1 and 2 provide information on the geographical location of each region of the Czech Republic (their official names are presented in Table 1) and the respective level of the gross monthly wage. The figures clearly show a substantially higher wage level in the region of the capital Prague. A relatively high level of wages in Central Bohemian and Pilsen regions is noticeable, low levels, on the other hand, being reported in Karlovy Vary, Zlin and Olomouc regions.

**Tab. 1: Official names<sup>1</sup> of regions of the Czech Republic<sup>2</sup>**

Region	Code	Region	Code
Capital Prague Region	A	Hradec Kralove Region	H
Central Bohemian Region	S	Pardubice Region	E
South Bohemian Region	C	Vysocina Region	J
Pilsen Region	P	South Moravian Region	B
Karlovy Vary Region	K	Olomouc Region	O
Usti Region	U	Zlin Region	Z
Liberec Region	L	Moravian-Silesian Region	T

Source: www.mdcz.cz

**Fig. 1: Average gross monthly wages (in CZK) in respective regions of the Czech Republic in 2015**



Source: www.czso.cz; own customization

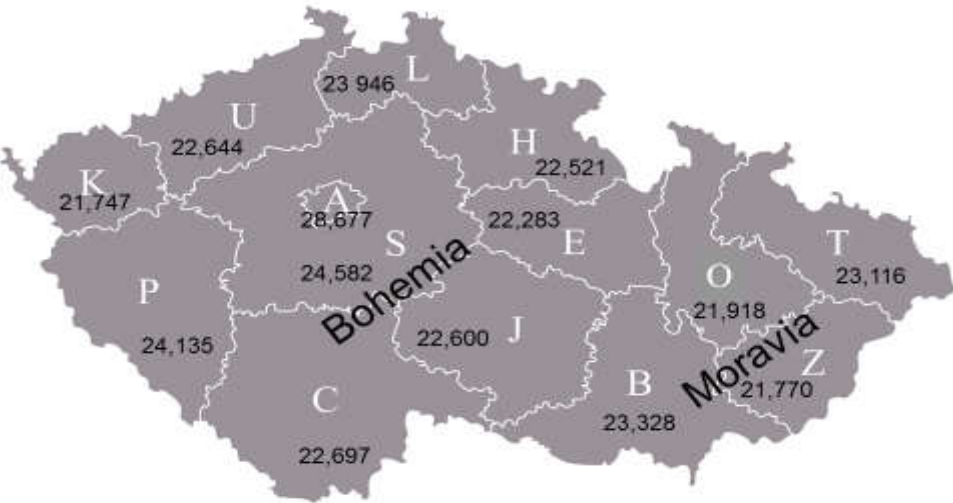
The region of the capital Prague is not food and energy self-sufficient. Its macroeconomic statistics, however, considerably exceed those of other regions, wages earned in Prague resembling those received in more developed countries. One of the reasons is the concentration of industries with higher labour productivity such as finance and informatics. Corporate policies also play an important role, since Prague-based firms often produce values outside the capital but divide the profits at a place where they are headquartered. As expected,

<sup>1</sup> The names of most regions match those of their respective capitals.

<sup>2</sup> Different backgrounds distinguish the Bohemian regions (grey) from Moravian ones (white).

the three regions with the highest wage levels are identical to those with the lowest unemployment rates in the same order. However, the order of the regions at the bottom wage levels is not the same as that of the regions with top rates of unemployment, neither the Moravian-Silesian Region nor Usti Region belonging to the three lowest wage level areas; for details, see Table 2.

**Fig. 2: Median gross monthly wages (in CZK) in respective regions of the Czech Republic in 2015**



Source: www.czso.cz; own customization

**Tab. 2: Average unemployment rates *U* (in %) in regions of the Czech Republic in 2015**

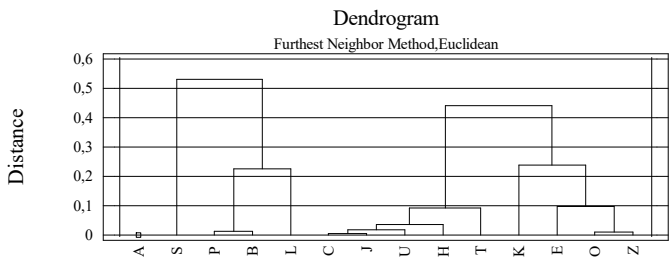
	Region													
	Central Prague Region	Central Bohemian Region	South Bohemian Region	Pilsen Region	Karlovy Vary Region	Usti Region	Liberec Region	Hradec Kralove Region	Pardubice Region	Vysocina Region	South Moravian Region	Olomouc Region	Zlin Region	Moravian-Silesian Region
<i>U</i>	2.8	3.5	4.0	3.8	6.7	7.6	5.5	5.6	4.6	4.7	5.0	5.9	4.7	8.1

Source: www.czso.cz

Figures 3–6 provide an overview of the results of regional cluster analysis according to the wage level employing the method of the farthest neighbour and Euclidean distance metric.

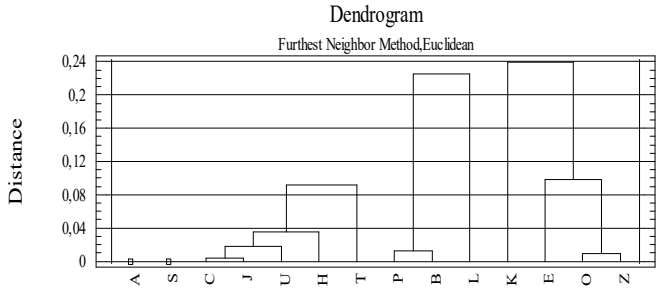
The first cluster always contains only one element – the Capital Prague Region – in the case of both the average and median monthly wage (i.e. three- and five-cluster analysis, respectively), due to markedly higher wage levels in this respective region.

**Fig. 3: Cluster analysis using three clusters, farthest neighbour method and Euclidean distance metric; 2015 average wage**



Source: Own research

**Fig. 4: Cluster analysis using five clusters, farthest neighbour method and Euclidean distance metric; 2015 average wage**

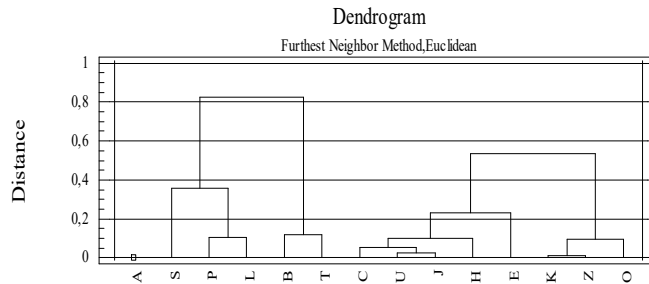


Source: Own research

Within the division of the regions into three clusters by the average wage, the second cluster includes four elements – Central Bohemian, Pilsen, Liberec and South-Moravian regions. According to the median wage division, however, the second cluster has five elements; along with those four mentioned above, this cluster contains the Moravian-Silesian Region, which seems to be rather surprising since this region’s general unemployment rate

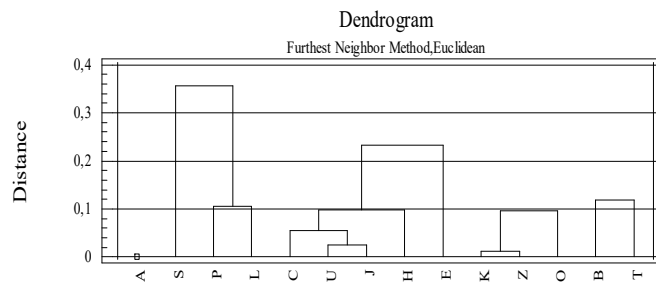
reaches the highest value of the whole Czech Republic. The remaining regions form the third clusters.

**Fig. 5: Cluster analysis using three clusters, farthest neighbour method and Euclidean distance metric; 2015 median wage**



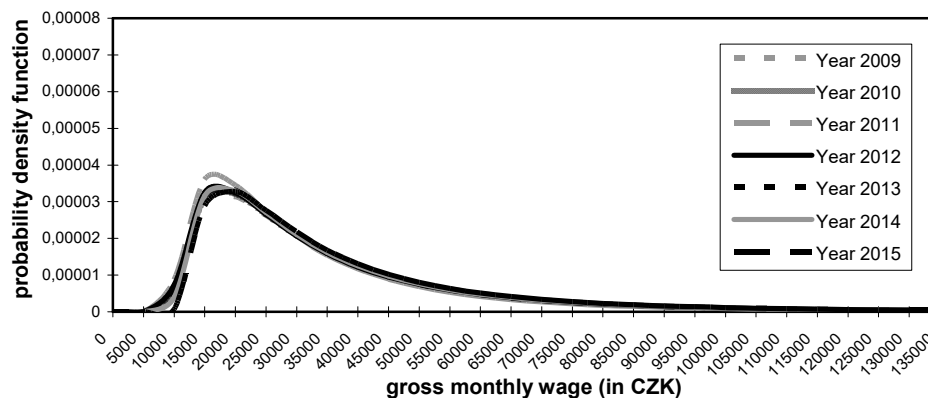
Source: Own research

**Fig. 6: Cluster analysis using five clusters, farthest neighbour method and Euclidean distance metric; 2015 median wage**



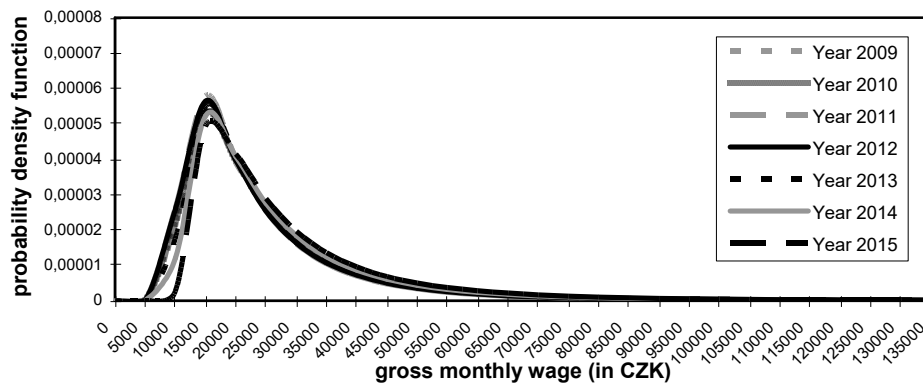
Source: Own research

**Fig. 7: Development of model wage distributions – Capital Prague Region**



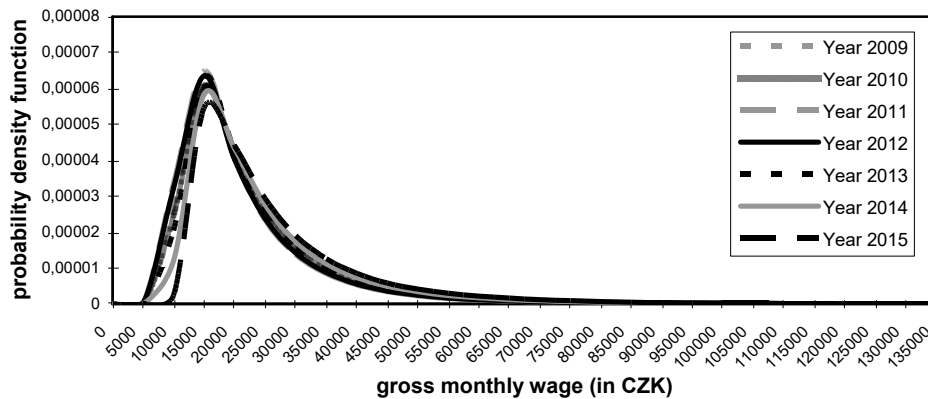
Source: Own research

**Fig. 8: Development of model wage distributions – Central Bohemian Region**



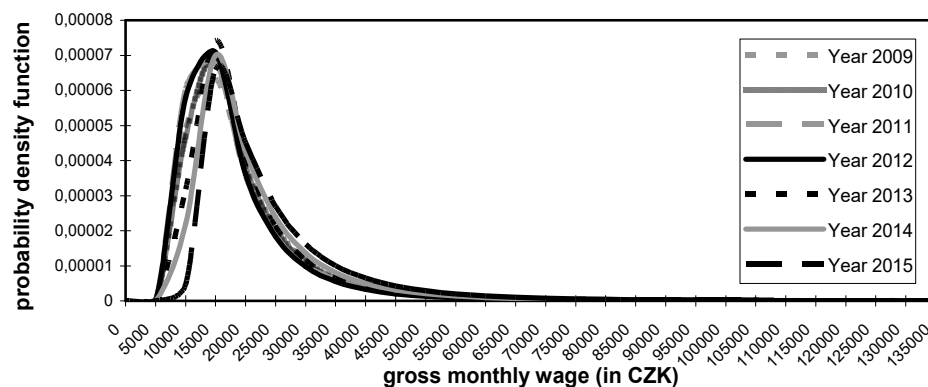
Source: Own research

**Fig. 9: Development of model wage distributions – Pilsen Region**



Source: Own research

**Fig. 10: Development of model wage distributions – Karlovy Vary Region**

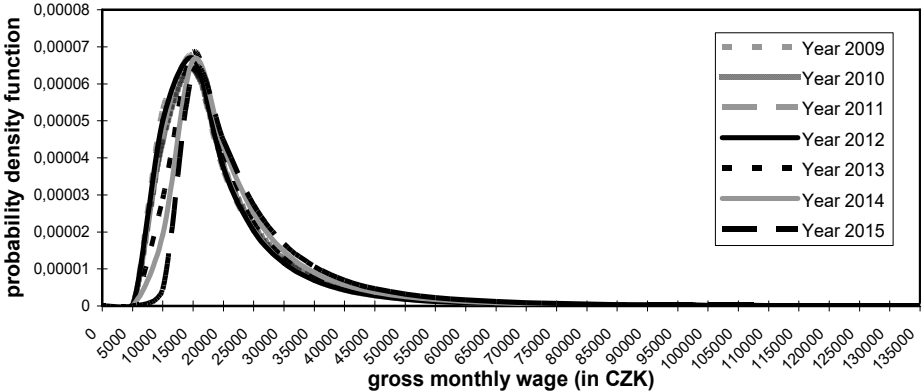


Source: Own research

Within the division of the regions into five clusters by the average wage, the second cluster contains only one element, namely the Central Bohemian Region. However, according to the median wage, the second cluster consists of three elements – Central

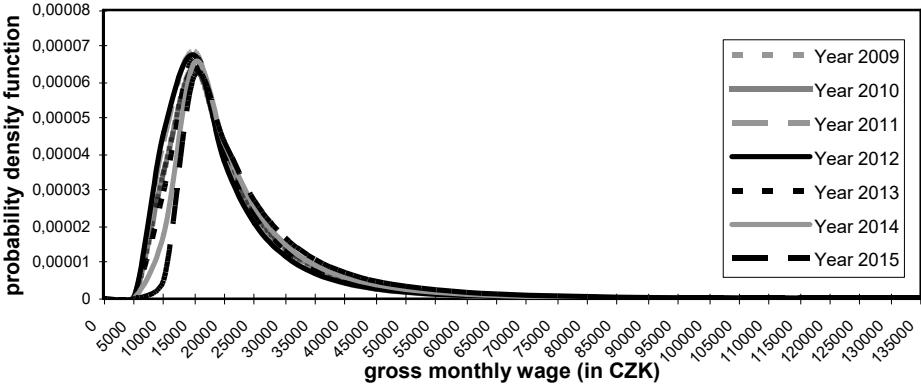
Bohemian, Pilsen and Liberec regions. The third cluster always comprises five elements – South Bohemian, Usti, Hradec Kralove, Vysocina and Moravian-Silesian regions by the average wage and South Bohemian, Usti, Hradec Kralove, Pardubice and Vysocina regions, respectively, according to the median wage. The fourth cluster has only three elements in both cases – Pilsen, Liberec and South Moravian regions according to the average wage, and Karlovy Vary, Olomouc and Zlin regions by the median wage, the latter being those with the lowest wage levels. The fifth clusters formed by the average and median wage contain the four and two remaining regions, respectively; for details, see Figures 5–8).

**Fig. 11: Development of model wage distributions – Zlin Region**



Source: Own research

**Fig. 12: Development of model wage distributions – Olomouc Region**



Source: Own research

Theoretical models for the wage distribution of each region from 2009 onwards have been constructed. They are based on the use of the probability density function of three-parameter lognormal curves, whose parameters were estimated using the maximum likelihood



method. The beginning of these curves is represented by the value of the minimum wage in respective years; see Table 3.

**Tab. 3: Minimum wage development in the Czech Republic (in CZK) since 2009**

Year	2009	2010	2011	2012	2013 <sup>3</sup>	2014	2015	2016	2017
Minimum wage	8,000	8,000	8,000	8,000	8,000 <sup>4</sup> 8,500 <sup>5</sup>	8,500	9,200	9,900	11,000

Source: www.mpsv.cz

The accuracy of the models obtained was compared applying the Akaike and Bayesian information criteria, both of which take a number of the corresponding wage model parameters into account.

Theoretical wage models using three-parameter lognormal curves and the maximum likelihood method of parameter estimation were created. Table 4 indicates the values of parameters estimated, Table 5 presenting the values of the Akaike and Bayesian information criteria, which enable to assess the estimation accuracy. Using Figures 7–12, the theoretical wage models capture the three regions with the highest (Prague, Central Bohemian and Pilsen regions) and the other three with the lowest wage levels (Karlovy Vary, Zlin and Olomouc regions).

The above figures allow for comparison of the wage distribution development of the regions with the highest and lowest wages over the last seven years. As observed in the figures, distributions with a high wage level are characterized by higher variability than those with a low wage level. Moreover, distributions with a low level of wages are more skewed and have higher kurtosis than those with high wage levels. A probability model, usually representing a simple approximation of a rather complex empirical distribution and the knowledge of the development trend of its parameters allow for the estimation of the whole wage distribution for future research purposes. As we can see from Table 5, wage models for the capital Prague region show the lowest accuracy, while Karlovy Vary region with its lowest wage level indicates the best accuracy of wage models, the number of model parameters being implicated in both information criteria (AIC and BIC). The relationship

<sup>3</sup> In 2013, the beginning of lognormal curves was determined proportionally, i.e.

$$8\,000 + 7 * \frac{8\,500 - 8\,000}{12} = 8\,292.$$

<sup>4</sup> From 1<sup>st</sup> January 2013 to 31<sup>st</sup> July 2013.

<sup>5</sup> From 1<sup>st</sup> August 2013 to 31<sup>st</sup> December 2013.

between the model accuracy and the wage level in the respective region is obvious. It holds that low model accuracy corresponds to a high wage level and vice versa.

**Tab. 4: Parameter estimates of three-parameter lognormal distribution using maximum likelihood method (parameter  $\theta$  equalling respective annual minimum wage)**

Region	Est.	Year						
		2009	2010	2011	2012	2013	2014	2015
Capital Prague Region	$\mu$	9,937	9,911	9,803	9,885	9,875	9,885	9,907
	$\sigma^2$	0,765	0,746	0,740	0,761	0,752	0,765	0,777
Central Bohemian Region	$\mu$	9,426	9,443	9,342	9,366	9,412	9,467	9,512
	$\sigma^2$	0,732	0,707	0,683	0,702	0,699	0,711	0,742
South Bohemian Region	$\mu$	9,112	9,156	9,124	9,092	9,176	9,250	9,312
	$\sigma^2$	0,701	0,688	0,633	0,633	0,627	0,632	0,668
Pilsen Region	$\mu$	9,298	9,316	9,202	9,226	9,275	9,371	9,426
	$\sigma^2$	0,653	0,634	0,639	0,649	0,630	0,639	0,682
Karlovy Vary Region	$\mu$	9,113	9,054	8,978	8,963	9,060	9,149	9,239
	$\sigma^2$	0,737	0,627	0,686	0,635	0,577	0,616	0,670
Usti Region	$\mu$	9,264	9,280	9,130	9,169	9,216	9,266	9,337
	$\sigma^2$	0,711	0,674	0,671	0,664	0,637	0,654	0,707
Liberec Region	$\mu$	9,326	9,292	9,146	9,182	9,247	9,308	9,392
	$\sigma^2$	0,950	0,656	0,636	0,625	0,618	0,626	0,689
Hradec Kralove Region	$\mu$	9,144	9,201	9,087	9,139	9,190	9,259	9,304
	$\sigma^2$	0,640	0,671	0,623	0,625	0,620	0,636	0,661
Pardubice Region	$\mu$	9,225	9,168	9,105	9,118	9,160	9,231	9,306
	$\sigma^2$	0,920	0,697	0,642	0,660	0,642	0,650	0,704
Vysocina Region	$\mu$	9,195	9,204	9,093	9,133	9,193	9,255	9,322
	$\sigma^2$	0,754	0,691	0,616	0,640	0,605	0,627	0,685
South Moravian Region	$\mu$	9,358	9,394	9,268	9,311	9,361	9,404	9,460
	$\sigma^2$	1,028	1,025	1,028	1,027	1,018	1,014	1,009
Olomouc Region	$\mu$	9,204	9,203	9,088	9,086	9,164	9,239	9,296
	$\sigma^2$	0,661	0,659	0,632	0,657	0,651	0,633	0,719
Zlin Region	$\mu$	9,077	9,139	9,075	9,065	9,149	9,215	9,277
	$\sigma^2$	0,739	0,704	0,656	0,680	0,626	0,633	0,686
Moravian-Silesian Region	$\mu$	9,196	9,251	9,220	9,233	9,262	9,290	9,345
	$\sigma^2$	0,681	0,664	0,665	0,662	0,656	0,654	0,702

Source: Own research

**Tab. 5: Akaike and Bayesian information criteria values**

Region	Crit.	Year						
		2009	2010	2011	2012	2013	2014	2015
Capital Prague Region	<i>AIC</i>	1,715,847	1,605,510	1,514,319	1,456,900	1,362,167	1,291,967	1,219,753
	<i>BIC</i>	1,715,870	1,605,533	1,514,342	1,456,923	1,362,191	1,291,990	1,219,776
Central Bohemian Region	<i>AIC</i>	544,441	536,317	528,511	543,060	546,067	556,645	577,484
	<i>BIC</i>	544,462	536,338	528,533	543,082	546,089	556,666	577,505
South Bohemian Region	<i>AIC</i>	297,991	296,670	282,364	284,622	285,256	289,311	303,276
	<i>BIC</i>	298,011	296,691	282,384	284,643	285,277	289,331	303,296
Pilsen Region	<i>AIC</i>	269,898	267,230	272,194	278,447	275,599	281,651	298,363
	<i>BIC</i>	269,919	267,250	272,215	278,468	275,619	281,671	298,383
Karlovy Vary Region	<i>AIC</i>	127,711	115,055	123,387	117,626	110,162	116,655	124,583
	<i>BIC</i>	127,730	115,074	123,406	117,644	110,180	116,674	124,602
Usti Region	<i>AIC</i>	348,582	336,708	336,174	334,146	324,956	331,734	350,898
	<i>BIC</i>	348,603	336,728	336,195	334,167	324,977	331,755	350,919
Liberec Region	<i>AIC</i>	231,143	185,448	183,454	183,282	183,845	187,689	203,032
	<i>BIC</i>	231,162	185,468	183,474	183,302	183,864	187,709	203,052
Hradec Kralove Region	<i>AIC</i>	240,158	250,588	239,702	242,212	242,731	249,249	258,341
	<i>BIC</i>	240,178	250,609	239,722	242,233	242,751	249,270	258,361
Pardubice Region	<i>AIC</i>	275,403	235,724	226,401	234,568	233,737	239,552	257,057
	<i>BIC</i>	275,423	235,744	226,421	234,588	233,757	239,572	257,077
Vysocina Region	<i>AIC</i>	238,726	228,425	213,311	222,088	215,872	224,336	241,838
	<i>BIC</i>	238,746	228,445	213,331	222,108	215,892	224,356	241,858
South Moravian Region	<i>AIC</i>	786,921	788,918	793,299	796,116	795,440	796,875	797,641
	<i>BIC</i>	786,943	788,940	793,321	796,138	795,462	796,897	797,663
Olomouc Region	<i>AIC</i>	255,008	261,376	260,425	274,584	279,522	280,684	313,825
	<i>BIC</i>	255,028	261,396	260,445	274,604	279,543	280,705	313,846
Zlin Region	<i>AIC</i>	289,572	282,403	271,177	279,996	265,977	270,105	288,029
	<i>BIC</i>	289,593	282,423	271,197	280,017	265,997	270,125	288,050
Moravian-Silesian Region	<i>AIC</i>	581,796	574,637	578,439	579,575	578,706	580,905	613,021
	<i>BIC</i>	581,817	574,659	578,461	579,597	578,728	580,927	613,042

Source: Own research

## **Conclusion**

Currently, the wage level is increasing in all regions of the Czech Republic, the wages of both men and women growing. Having slowed down between 2010 and 2013, the annual wage growth has accelerated in the whole of the Czech Republic since then, the average wage increasing by more than 2,300 CZK in most regions. The slowest growth was recorded in Moravian-Silesian and Usti regions (by 1,859 and 2,125 CZK, respectively), the fastest in the Central Bohemian region (by 3,260 CZK). The lowest wage level has long been reported in Karlovy Vary region – the average wage was a third lower than in the capital of Prague in 2015, the difference between the two regions gradually diminishing, the greatest one having been in 2011.

The highest and lowest wages are reported in Prague and Karlovy Vary regions, respectively, the average gross monthly wage amounting to 36,371 CZK in the former, compared to only 24,119 CZK in the latter region in 2015. Residents of Central Bohemian, Pilsen and South Moravian regions receive relatively high wages, averaging 27,997, 27,013 and 27,051 CZK, respectively, in the same year. High-income regions, however, are also characterized by relatively wide gender wage gaps.

It follows from the results of the 2015 Living Conditions Survey that net household income per capita increased by 4,400 CZK on average in 2014 (latest data available), compared to the previous year. The annual fall in income occurred only in Vysocina region, the sharpest increase – more than 10 per cent – in Karlovy Vary region, representing an annual growth of more than 14,000 CZK per person in cash. Household income in individual regions varies considerably, depending on local economic conditions. Best-off households are in Prague region, while those in the Moravian-Silesian region report the lowest net income.

## **Acknowledgment**

This paper was subsidized by the funds of institutional support of a long-term conceptual advancement of science and research number IP400040 at the Faculty of Informatics and Statistics, University of Economics, Prague, Czech Republic.

## References

- Albelda, R., & Carr, M. (2014). Double Trouble: US Low-Wage and Low-Income Workers, 1979–2011. *Feminist Economics*, 20(2), 1–28.
- Bárány, Z. L. (2016). The Minimum Wage and Inequality: The Effects of Education and Technology. *Journal of Labor Economics*, 34(1), 237–274.
- Bartošová, J., & Želinský, T. (2013). The Extent of Poverty in the Czech and Slovak Republics 15 Years after the Split. *Post-Communist Economies*, 25(1), 119–131.
- Domínguez-Villalobos, L., & Brown-Grossman, F. (2010). Trade Liberalization and Gender Wage Inequality in Mexico. *Feminist Economics*, 16(4), 53–79.
- Fisher, J., Johnson, D. S., & Smeeding, T. M. (2015). Inequality of Income and Consumption in the U.S.: Measuring the Trends in Inequality from 1984 to 2011 for the Same Individuals. *Review of Income and Wealth*, 61(4), 630–650.
- Gobillon, L., Meurs, D., & Roux, S. (2015). Estimating Gender Differences in Access to Jobs. *Journal of Labor Economics*, 33(2), 317–363.
- Jenderny, K. (2016). Mobility of Top Incomes in Germany. *Review of Income and Wealth*, 62(2), 245–265.
- Kukk, M., & Staehr, K. (2014). Income Underreporting by Households with Business Income: Evidence from Estonia. *Post-Communist Economies*, 26(2), 257–276.
- Makhalova, E., & Pecáková, I. (2015). The Fuzzy Clustering Problems and Possible Solutions. In: *The 9<sup>th</sup> International Days of Statistics and Economics, Conference Proceedings [online]*, Prague, 10.09.2015–12.09.2015, 1052–1061.
- Malá, I. (2015). Multivariate Probability Model for Incomes of the Czech Households. *Politická ekonomie*, 63(7), 895–908.
- Malec, L. (2016). Some Remarks on the Functional Relation Between Canonical Correlation Analysis and Partial Least Squares. *Journal of Statistical Computation and Simulation*, 86(12), 2379–2391.
- Řezanková, H., & Löster, T. (2013). Shluková analýza domácností charakterizovaných kategoriálními ukazateli (Cluster Analysis of Households Characterized by Categorical Indicators). *E+M Ekonomie a Management*, 16(3), 139–147.

**Contact**

Diana Bílková

University of Economics, Prague  
Faculty of Information and Statistics  
Department of Statistics and Probability  
Sq. W. Churchill 1938/4  
130 67 Prague 3  
Czech Republic

Mail: [bilkova@vse.cz](mailto:bilkova@vse.cz)

University of Finance and Administration  
Faculty of Economic Studies  
Department of Informatics and Mathematics  
Estonian Street 500/3  
101 00 Prague 10  
Czech Republic

Mail: [diana.bilkova@vsfs.cz](mailto:diana.bilkova@vsfs.cz)