# THE IMPACT OF THE TRAINING SET SIZE ON THE CLASSIFICATION OF REAL ESTATE WITH AN INCREASED FISCAL BURDEN[1]

**Sebastian Gnat, PhD**
*Faculty of Economics and Management*
*University of Szczecin*
*e-mail: sebastian.gnat@usz.edu.pl*

## Abstract

The introduction of an ad valorem tax can lead to an increase in the tax burden on real estate. There are concerns that this increase will be large and widespread. Before undertaking any actual actions related to the real estate tax reform, pilot studies and statistical analyses need to be conducted in order to verify the validity of those concerns and other aspects regarding the replacement of a real estate tax, agricultural tax and forest tax with an ad valorem tax. The article presents results of research on the effectiveness of the classification of real estate into a group at risk of an increase of tax burden with the use of the k-nearest neighbors method. The main focus was to determine the size of a real estate set (training data set) on the basis of which classification is conducted, as well as on the efficiency of that classification, depending on the size of such data set.

## 1. Introduction

Ad valorem tax, i.e. a real estate tax calculated on the basis of a real estate value (in line with the Polish regulations the cadastral value) exists in many countries in Europe and in the world. Polish fiscal solutions regarding real estate are based on property area as the tax basis. There are many studies regarding communes' area-based revenue from real estate (TROJANEK, KISIAŁA 2016). The area-based system is criticized solution as being ineffective (ETEL, DOWGIER 2013). What is more, real estate tax fails to perform non-fiscal functions; it is not used as an instrument of a rational space management policy (WÓJTOWICZ 2006; GNAT, SKOTARCZAK 2006). Despite the discussions carried out in the spheres of science and politics and in the press, as well as numerous attempts to demonstrate the advantages of taxing real estate value, it seems its implementation in Poland does not currently seem possible. There are a number of reasons for such a situation. They are of social and political as well as economic and organizational nature. Public opinion is quite explicitly negative with regard to the tax based on real estate value (ad valorem tax). There is a multitude of concerns as to the tax amounts; it is argued that the process of determining a cadastral tax will be complex and costly.

Legislative work was started in the past regarding the introduction of universal real estate taxation, the aim of which is determining real estate value in a mass fashion. On account of the above-

mentioned negative social perception, the real estate taxation reform has been suspended. One of the reasons for such a situation is the social conviction of a substantial increase in fiscal burden. The studies conducted so far focused on pilot appraisals and statistical analyses of the possible effects of introducing value-based real estate taxation. Comprehensive pilot studies, though highly precise, are long-lasting and costly. On account of the above, the article constitutes an attempt at employing another approach to the process of appraising the effects of introducing an ad valorem tax. An assumption was made that it is possible to draw a sample of real properties, to appraise them and to predict a tax burden increase on the basis of the information about the properties of the selected real estate located in the territory of the entire municipality. The objective of the article is to define the minimum size of such a sample (training data set) for which the predicted accuracy (classification of such real properties into data sets featuring either an increase or decrease of taxation) will be at an acceptably high lev-el. Correct assessment of the effects of implementing the real estate tax reform carried out solely on the selected (possibly small) sample will be significantly cheaper and quicker than comprehensive studies.

The issue of determining the optimal size of a sample is important not only in terms of the analysis of the effects of introducing an ad valorem tax, but it is also of vital significance for *SAREMA* (the Szczecin Algorithm of Real Estate Mass Appraisal), since one of its stages involves an individual appraisal of a selected sample of representative real properties. The size of the sample both determines the accuracy of mass appraisals as well as significantly affecting the costs of implementing the procedure of mass appraisal.

The *K*-nearest neighbor method was employed for conducting land plot classification on account of a change in fiscal burden. The *KNN* method is an example of a statistical learning algorithm. Some others, like probit models can be successfully utilized on the real estate market (see PLUMMER, 2014). The models mentioned above, along with many others, belong to a group of statistical learning algorithms of supervised learning that can be used for classification in many fields, i.e. fraud detection (see SONG et al. 2014), mortgage analysis (FELDMAN, GROSS 2005) or tenure choice analysis (ENSTRÖM-ÖST, SÖDERBERG, WILHELMSSON 2017).

The subject matter of the study concerns urban land plots located in the territory of one of the municipalities of the West Pomeranian Voivodeship.

## 2. *K* nearest neighbors method

The *KNN* classifier (*k* nearest neighbors), is recognized as a lazy learning classifier based on examples (BOSCHETTI, MASSARON 2017). *KNN* is an algorithm existing in the machine learning realm, but there are also many other statistical techniques used in real estate market analysis (GŁUSZCZAK 2015). In the case of the *KNN* classifier it is not necessary to estimate a model (HASTIE et al. 2009). The *KNN* algorithm can also be used in regression problems (see PACE 1996), though it is mainly applied is classification problems. The operation of the classifier comes down to two steps. In the first step for point $x_0$ we find $k$ training points $x(r), r = 1,…,k$ located closest to $x_0$. In the second step a classification is made based on the so-called majority voting. Point $x_0$ is recognized as belonging to the class which is most numerously represented in a training point set. In the event of the same number of votes during majority voting, the *KNN* algorithm implemented in the scikit-learn package of Python programming language (PEDREGOSA *et al.* 2011) used in the study favours the neighbors located closest to the analyzed point. In a situation in which distances from neighbours are similar, the algorithm selects the first class label found in a training data set.

Despite its simplicity, the *KNN* method is successfully applied in a large number of problems related to classification, including handwriting identification, satellite image analysis or *ECG* records (cf. BAO, ISHII 2002; LAGERHOLM et al. 2000). Its effectiveness becomes evident particularly in such classification problems in which the decision-making borders are highly irregular.

The application of $k$ nearest neighbors method requires a data set to be normalized. Owing to the fact that a *KNN* classifier classifies a given observation by identifying points located closest to it, numerical scales of variables are of significant importance. All big scale variables will have a far greater influence on the distance between observations, and thereby on classification results, than small scale variables will (cf. JAMES et al. 2013, p. 165). In this context the need for the standardization of variable values describing the analyzed objects ought to be mentioned. In turn, the question of which points should be recognized as the ones closest to point $x_0$ depends on the adopted metric. The Euclidean distance was applied in this study.
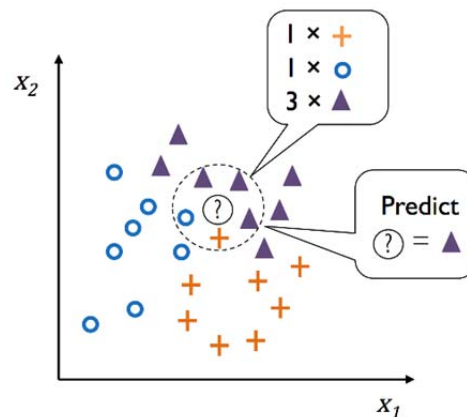
**Fig. 1.** Example of classification with *KNN* method. *Source*: RASCHKA (2018).

### 3. Mass real estate appraisal

The process of real estate mass appraisal (i.e. universal real estate taxation), which will need to precede the introduction of a cadastral tax in Poland, is regulated in the Management of Real Estate Act of 21 August 1997 (hereinafter referred to as: the *MREA*) and the regulation of the Council of Ministers dated 29 June 2005 on universal real estate taxation. Article 162.1 of the *MREA* states: "Cadastral values, determined in the process of universal real estate taxation, ought to take into account the differences existing between individual real properties and an approximation to the market value that can be obtained with the application of the rules assumed for mass appraisal". It involves two major elements. Firstly, the legislator specifies that the cadastral value ought to be close to the market value. Secondly, mass appraisal must be specified as a valuation method. Real estate mass appraisal is a process in which cadastral value of a number of real properties is determined simultaneously. Such an approach requires the use of a specific procedure based on mathematical models. Proposals of real estate mass appraisal algorithms can be found in the works of many authors (cf. inter alia HOZER et al. 1999; CZAJA 2001; SAWIŁOW 2009). The study was based on a mass appraisal procedure developed by J. Hozer:

$$W_{ji} = WWR_j \times pow_i \times C_{baz} \times \prod_{k=1}^{K}(1 + A_{kp}) \qquad (1)$$

$$WWR_j = \sqrt[n_j]{\prod_{i=1}^{n_j} \frac{w_{ji}^{rz}}{w_{ji}^{h}}} \qquad (2)$$

where:
$W_{ji}$ – market (or cadastral) value $i$–th real property in $j$–th elementary terrain,
$WWR_j$ – market value coefficient in $j$–th elementary terrain $(j = 1, 2, …, J)$,
$pow_i$ – surface of $i$–th real property,
$C_{baz}$ – price of 1m² of the cheapest land (without the technical infrastructure) in the appraised area,
$A_{kp}$ – influence of $p$–th category of $k$–th attribute $(k = 1, 2, …, K; p = 1,2, …, k_p)$,
$K$ – number of attributes,
$k_p$ – number of categories of $k$–th attributes,
$w_{ji}^{rz}$ – value of $i$-th real property determined by a real estate appraiser in $j$–th elementary terrain,
$n_j$ – number of representative properties assessed by real estate appraiser in $j$–th elementary terrain,
$w_{ji}^{h}$ – hypothetical value of $i$-th property in $j$–th elementary terrain determined with the following formula:

$$w_{ji}^{h} = pow_i \cdot C_{baz} \prod_{k=1}^{K}(1 + A_{kp}) \qquad (3)$$

The amounts of real estate taxes were determined for urbanized lands located within the territory of Kolbaskowo municipality of the West Pomeranian Voivodeship in line with a resolution of its Municipality Council in force in 2015. Hence, no real tax burden amounts were used, only their approximations. Furthermore, buildings and structures were not considered. Even if a land plot was

equipped with technical utilities, that fact was not taken into account. The focus of the analysis concerned lands exclusively. Secondly, employing the described Szczecin Real Estate Mass Appraisal Algorithm, universal (within the scale of a singular municipality) taxation of real estate was conducted, the subject of which were land plots. The values obtained as a result of applying the algorithm constituted the grounds for assessing the ad valorem tax amount. The subject of appraisal in the study involved 2337 land plots. The market value of a real property is determined by various properties and market conditions. These are elements that affect the value in the opinion of potential market participants. On the basis of data concerning real properties, gathered in the process of land mass appraisal, the following properties were defined:

- area – large (over 5000 m²), medium (between 1000 and 5000 m²), small (less than 1000 m²),
- location – unfavorable (plots located in the smallest towns and on the outskirts of larger towns), average (plots located in larger towns), favorable (plots located in areas considered attractive by potential market participants),
- utility infrastructure – lacking, incomplete (usually without sewage system), complete,
- shape – wrong (difficult to build on, narrow plots of irregular shape), good (close to a square or rectangular shape, easy to build on and use),
- manner of use – farming, industrial, multi-family residential, single-family residential, commercial. The manner of use determined the main character of use of a given land plot.

In order to define various effects of the real estate taxation reform, a percentage rate of an ad valorem tax was used, which will ensure income into a municipality's budget equal to the one determined for a real estate tax (in this case the rate of cadastral tax was 0.36% of cadastral value). Nevertheless, the analyses of the impact exerted by a cadastral tax can be conducted with the gathered data for any rate of cadastral tax.

In Figure 2, a fragment of the analyzed municipality terrain is presented. The darker color marks the plots for which fiscal burden increased with the assumed rate of a cadastral tax. The lighter color designates the plots which, owing to a low real estate value, feature an ad valorem tax amount that is lower than the current real estate tax amount. The plots represented in the figure in white constitute non-urbanized plots, and thereby they are not subject to a real estate tax.

In order to fulfil the study assumptions, it was necessary to appraise the value of all the land plots constituting its subject. Thanks to this, it was possible to estimate an ad valorem tax amount in the case of each plot. Since the statistical analysis of land plots does not constitute the subject matter of this study, a synthetic summary of the simulation for introducing value-based taxation of urbanized land plots was presented in Figures 3 and 4. Figure 3 presents the distribution of absolute changes in the tax amounts. For the majority (but not for all) land plots, this difference is positive, which means an increase of tax burden. More detailed results concerning replacing the real estate tax with an ad valorem tax can be found in the works of GNAT (2010, 2018).
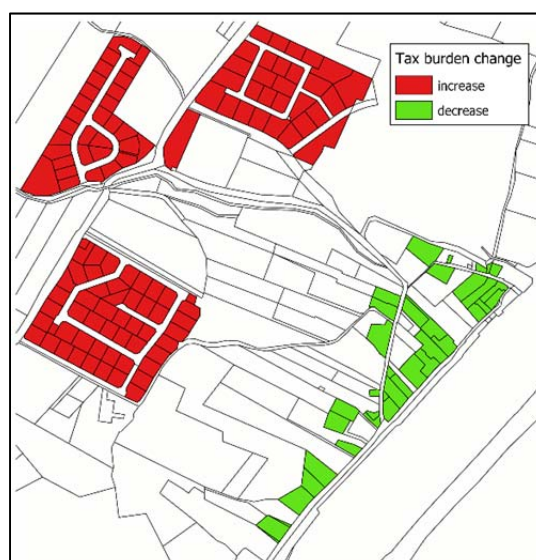


**Fig. 2**. Fragment of the examined municipality taking into account urbanized plots in terms of change in tax burden. *Source:* own study.
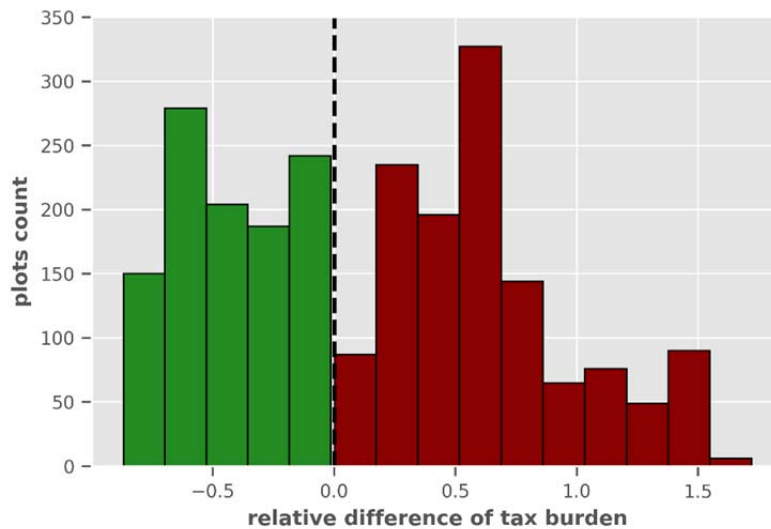
**Fig. 3.** Distribution of relative changes in tax burden amounts for the analyzed land plots. *Source:* own study.

Figure 4 represents the structure of plots of various relative changes in taxation depending on one of the features assumed in the mass appraisal, namely the manner of land use. From the information presented in the figure, it follows that in the case of land used for single-family residential purposes, the majority of plots feature a tax increase, while a portion of land plots will be subject to a tax decrease. An analogous structure of changes occurs for plots designated for multi-family residential purposes. In the case of all the land on which farm buildings have been erected (the analyzed municipality has a rural character) a decrease of tax amount was observed. A similar situation occurs for industrial plots. Plots designated for commercial use are found at the opposite end of the spectrum. In the case of such plots, each one of them would be encumbered with a higher value-based tax than in the case of a tax assessed on the grounds of plot area.
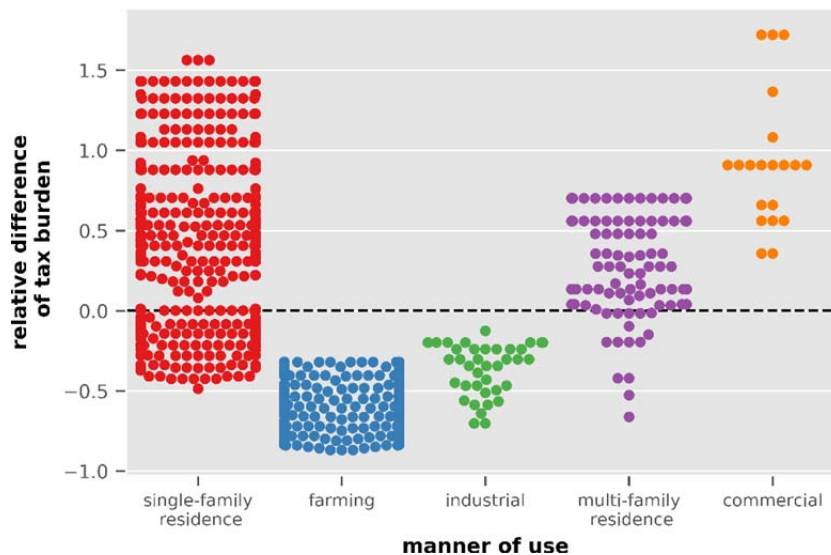


**Fig. 4.** Relative changes of tax burden taking into account the manner of land use. *Source:* own study.

## 4. Classification of real estate based on change in fiscal burden

As had already been mentioned, the study covered 2,337 urbanized land plots located within the area of a single municipality of West Pomeranian Voivodeship. The objective of the study was to define how big a sample from the analyzed data set is required in order to effectively determine the change in fiscal burden for all the urbanized land plots. The sample should be large enough to ensure high accuracy of classification, but small enough for the cost and time of conducting the study to be acceptable to decision-makers. The course of the study first involved proposing several sizes of

training data sets on the grounds of which several learning algorithms, which used 'k nearest neighbors' method, were created. The default k value in the package used is 5, and this was the value applied. Then, on the basis of each algorithm, classification was conducted for the remaining land plots and classification accuracy evaluated. The effectiveness of an algorithm is best defined in the case of such data sets in which there is no definitive predominance of any of the classes contained in that data set. In this study there are two classes designating land plots for which a change in the tax burden, after replacing surface taxation with value taxation, is positive or negative. In Figure 5, the structure of plots was presented taking into account the differentiated classes.
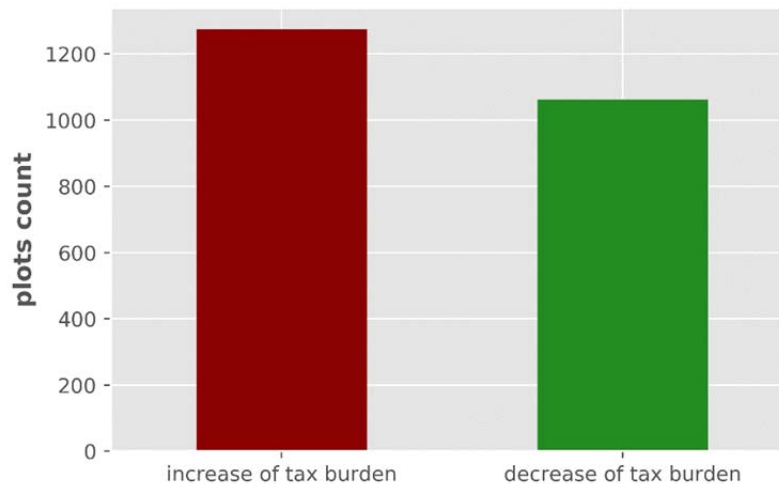


**Fig. 5.** Distribution of land plots on account of the change in tax burdens. *Source:* Own study.

A greater number of plots demonstrate an increase in fiscal burden. The percentage of plots characterized by a rise in tax burden is equal to 54.6%. The land plots for which the simulation demonstrates a decrease in tax amount to 45.5%.

The following sizes of training data sets were proposed in the study as a percentage of the entire dataset: 3%, 5%, 10%, 25%, 50% and 75%. The adopted range encompasses both very small and very large samples, which will enable an effective evaluation of the influence of a training data set size on classification accuracy. Accuracy defines what part of the predicted class labels is consistent with real results:

$$ACC = \frac{TP + TN}{TP+TN+FP+FN},$$ (4)

where:
*ACC*   – classification accuracy,
*TP*   – number of true positive predictions (both actual and predicted values indicated an increased tax burden),
*TN*   – number of true negative predictions (both actual and predicted values indicated decreased tax burden),
*FP*   – number of false positive predictions (prediction indicated increased tax burden but actual tax burden was lower),
*FN*   – number of false negative predictions (prediction indicated decreased tax burden but actual tax burden was higher).

Land plots for individual training datasets were selected in a random fashion. A set of properties characteristics used in the process of real estate mass appraisal was adopted as the characteristics describing individual plots in the *KNN* algorithm. Those characteristics included: surface, location, technical utilities, plot shape and the manner of plot use. The properties were converted into binary variables.

There are a number of measures intended for evaluating the accuracy of classification. Those include, inter alia, confusion matrices or an accuracy coefficient. Accuracy defines what part of the predicted class labels is consistent with real results. It means a percentage of correctly classified labels.

The obtained results of classification accuracy are presented for the determined sizes of training samples in Table 1.

**Table 1**

The accuracy of land plot classification depending on a learning dataset size ($k$=5)

| train set size | classification accuracy |
|---|---|
| 3% | 81.7% |
| 5% | 89.6% |
| 10% | 91.9% |
| 25% | 90.9% |
| 50% | 92.6% |
| 75% | 91.3% |

*Source:* own study.

From the presented data it arises that the accuracy of predictions falls within the range between 81.7% and 92.6%. The poorest results were obtained for a sample containing 3% of land plots. Maximum prediction accuracy occurred for a sample of 50%. However, such a large sample is not satisfactory time- and cost-wise. The second best accuracy was obtained for a sample of 10%, and this was the sample that was deemed to offer the best compromise between classification accuracy and economic conditions. The area under the *ROC* curve was used as an example of a graphic method for evaluating classification accuracy. "The *ROC* curve constitutes a graphic method of representing changes in classifier effectiveness for various possible thresholds (thus, the curve illustrates changes in results depending on the values of parameters). The effectiveness comes from the coefficient of accurate positive predictions and false positive predictions. The first value describes the percentage of correctly predicted positive outcomes, while the second one describes a percentage of incorrectly predicted positive outcomes. The area under the curve informs us how well the classifier works in comparison to a random classifier (with an area under the curve equal to 0.5)" (BOSCHETTI, MASSARON 2017, p. 173). The area under the *ROC* curve is equal to 0.92, being much higher than the area of the random classifier (dashed line in Figure 6).

As had previously mentioned, the number of neighbors equal to 5 was used to build the classifier. For the sample deemed to be the best, a further analysis was conducted and it was examined how a change of $k$ parameter will affect the classification accuracy. In Figure 7, the development of classification accuracy was presented on the basis of a sample equal to 10% for various $k$ values. The lowest classification error was recorded for $k$=6, thus for the value nearest to the initial value. The use of $k$=6 did not produce any other data regarding the evaluation of a sample size. Still, the sample of 10% of the entire plot set must be considered optimal.
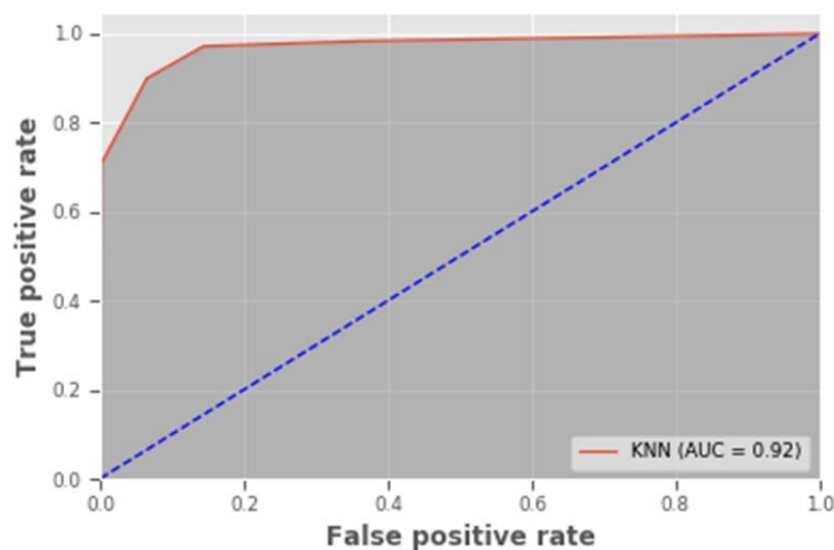


**Fig. 6.** Area under the ROC curve for a training set equal to 10%. *Source*: own study.
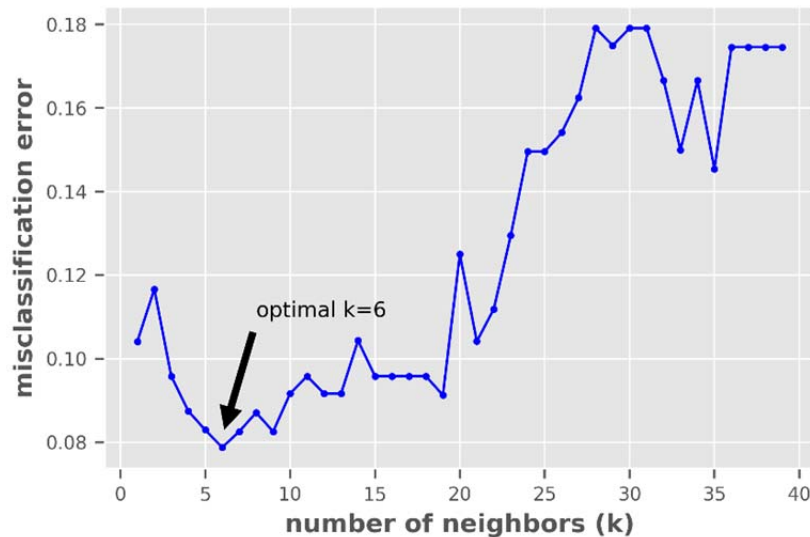
**Fig. 7.** Optimal value of *k* parameter in terms of classification accuracy. *Source:* own study.

**Table 2**

Accuracy of land plot classification depending on a training set size (*k*=6)

| Training set size | Classification accuracy |
|---|---|
| 3% | 82.8% |
| 5% | 85.7% |
| 10% | 91.9% |
| 25% | 91.8% |
| 50% | 88.0% |
| 75% | 90.9% |

*Source*: own work.

The last stage of the study involved evaluating whether the model can be deemed as being overfitted. This is a highly important stage of the evaluation of a classification algorithm. A characteristic of overfitted algorithms is that they reach higher accuracy in a training data set and a one in a testing set. Such a model does not offer good chances for accurate classifications. Figure 8 represents validation curves of the applied KNN model depending on the number of neighbors. It shows that, starting with *k*=2, the classification accuracy was higher in the testing set than in the training one. This means that the model is not overfitted, i.e. it provides a greater chance for effective real estate classification on account of a change in tax burden.

## 5. Conclusions

The article presents the use of the *k* nearest neighbors algorithm for the classification of land plots on account of an increase or decrease in the fiscal burden caused by the replacement of a real estate tax (in which a real estate surface constitutes the basis for the tax assessment) with an ad valorem tax. The objective of the study was determining how the impact of the size of a training set will affect the effectiveness of classification in a testing set, i.e. determining the choice of the optimal size of a training sample. The conducted research demonstrates that, for the gathered data, the *KNN* method produces results of high classification accuracy. For the population of 10% of land plots, deemed to be an optimal sample size, located in one of the municipalities of the Western Pomeranian Voivodeship, the classification accuracy obtained in the set of the remaining plots (a testing set) was nearly 92%. This means that conducting a pilot study for one tenth of the population allowed for performing a highly probable evaluation of over 2 thousand land plots in terms of whether their owners can expect an increase in tax burden. This is an important conclusion in the context of a cost analysis for pilot studies preceding the introduction of a major and necessary real estate tax reform.
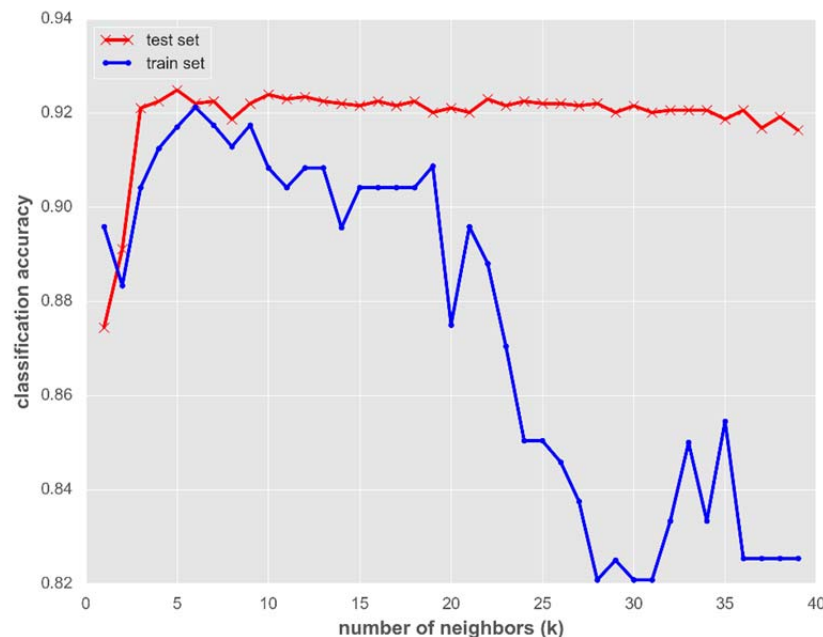
**Fig. 8.** Relative accuracy of a training and testing set depending on the number of neighbors (training set size of 10%). *Source:* Own study.

## 6. References:

BAO Y., ISHII N., 2002, *Combining Multiple K-Nearest Neighbor Classifiers for Text Classification by Reducts*, In: Lange S., Satoh K., Smith C.H. (eds) Discovery Science. DS 2002. Lecture Notes in Computer Science, vol 2534. Springer, Berlin, Heidelberg.

BOSCHETTI A., MASSARON L., 2017, *Python Data Science Essentials* in Polish: *Python, Podstawy nauki o danych*, Helion, Gliwice.

CZAJA J., 2001, *Methods of appraising real property market and cadastral value* in Polish: *Metody szacowania wartości rynkowej i katastralnej nieruchomości*, Komp-system, Kraków.

ENSTRÖM-ÖST C., SÖDERBERG B., WILHELMSSON M., 2017, *Homeownership rates of financially constrained households*, Journal of European Real Estate Research, Vol. 10 Issue: 2, pp.111-123.

ETEL L., DOWGIER R., 2013, *Local taxes and charges – time for a change* in Polish: *Podatki i opłaty lokalne – czas na zmiany,* Białystok: Temida 2.

FELDMAN D., GROSS S., 2005, *Mortgage Default: Classification Trees Analysis*, The Journal of Real Estate Finance and Economics, Volume 30, Issue 4, pp. 369–396.

GŁUSZAK M., 2015, *Multinomial Logit Model of Housing Demand in Poland*, Real Estate Management and Valuation, Vol. 23, No. 1, pp. 84-89.

GNAT S., 2010, *Analysis of the Effects of Replacing Current Property Tax with ad Valorem Property Tax in a Sample Municipality*, Folia Oeconomica Stetinensia, 8 16, pp. 82-98.

GNAT S., 2018, *Analysis of Communes' Potential Fall in Revenue Following Introduction of ad Valorem Property Tax*, Real Estate Management and Valuation, vol. 26, no. 1, pp. 63-72.

GNAT S., SKOTARCZAK M., 2006, *Analysis of local tax rates distributions in the West Pomeranian Voivodeship municipalities in the period of 2002–2004* in Polish: *Analiza rozkładów stawek podatków lokalnych w gminach województwa zachodniopomorskiego w latach 2002–2004*, in: J. Hozer red., *Economic situation vs. real estate market* in Polish: *Koniunktura gospodarcza a rynek nieruchomości*, Szczecin: Uniwersytet Szczeciński, Instytut Analiz, Diagnoz i Prognoz Gospodarczych, pp. 74-82.

HASTIE T., TIBSHIRANI R., FRIEDMAN J., 2009, *The elements of statistical learning*, Springer, New York.

HOZER J., FORYŚ I., ZWOLANKOWSKA M., KOKOT S., KUŹMIŃSKI W., 1999, *Econometric algorithm of land real estate mass appraisal* in Polish: *Ekonometryczny algorytm masowej wyceny nieruchomości gruntowych*, Uniwersytet Szczeciński, Stowarzyszenie Pomoc i Rozwój, Szczecin.

LAGERHOLM M., PETERSON C., BRACCINI G., EDHENBRANDT L., SORNMO L., 2000, *Clustering ECG Complexes Using Hermite Functions and Self-Organizing Maps*, IEEE Trans. Biomed. Eng., vol. 47, No. 7, pp. 838-848.

PACE P.K., 1996, *Relative performance of the grid, nearest neighbor, and OLS estimators*, The Journal of Real Estate Finance and Economics, Volume 13, Issue 3, pp. 203–218.

PEDREGOSA F., VAROQUAUX G., GRAMFORT A., MICHEL V., THIRION B., GRISEL O., BLONDEL M., PRETTENHOFER P., WEISS R., DUBOURG V., VANDERPLAS J., PASSOS A., COURNAPEAU D., BRUCHER M., PERROT M., DUCHESNAY É., 2011, *Scikit-learn: Machine Learning in Python*, JMLR 12, pp 2825-2830.

PLUMMER E., 2014, *The Effects of Property Tax Protests on the Assessment Uniformity of Residential Properties,* Real Estate Economics, Volume 4, Issue 4, pp. 900-937.

RASCHKA S., 2018, *Python. Machine learning* in Polish: *Python. Uczenie maszynowe*, Wydawnictwo Helion, Gliwice.

SAWIŁOW E., 2009, *The application of methods of multi-dimensional comparative analysis for the purpose of cadastral value determination* in Polish: *Zastosowanie metod wielowymiarowej analizy porównawczej dla potrzeb ustalania wartości katastralnych*, Studia i Materiały Towarzystwa Naukowego Nieruchomości, Vol. 17, no. 1.

Song Xin-Ping, Hu Zhi-Hua, Du Jian-Guo, Sheng Zhao-Han, 2014, *Application of Machine Learning Methods to Risk Assessment of Financial Statement Fraud: Evidence from China,* Journal of Forecasting, Volume 33, Issue 8, pp. 611-626.

TROJANEK M., KISIAŁA W., 2016, *The Diversification of Communes' Revenue from Real Estate Across Provinces*, Real Estate Management and Valuation, Vol. 24, No. 2, pp. 36-49.

WÓJTOWICZ K., 2006, *Analysis of potential effects of real estate tax system reform in Poland* in Polish: *Analiza potencjalnych skutków reformy systemu opodatkowania nieruchomości w Polsce*. Public finances. Lublin: Wyd. UMCS.

## Acknowledgments