

Automated Valuation Models: Improving Model Performance by Choosing the Optimal Spatial Training Level

Bastian Krämer^a, Moritz Stang^a, Vanja Doskoč^b, Wolfgang Schäfers^a and Tobias Friedrich^b

^a *International Real Estate Business School, University of Regensburg, Regensburg, Germany*

^b *Hasso Plattner Institute, University of Potsdam, Potsdam, Germany*

Abstract

The use of Automated Valuation Models (AVMs) in the context of traditional real estate valuations and their performance has been discussed in the academic community for several decades. Most studies focus on finding which method is best suited for estimating property values. One aspect that has not yet been studied scientifically is the appropriate choice of the spatial training level. The published research on AVMs usually deals with a manually defined region and fails to test the methods used on different spatial levels. The aim of our research is thus to investigate the impact of training AVM algorithms at different spatial levels in terms of valuation accuracy. We use a dataset with about 1.2 million residential properties from Germany and test four different methods, namely Ordinary Least Square, Generalized Additive Models, eXtreme Gradient Boosting and Deep Neural Network. Our results show that the right choice of spatial training level can have a major impact on the model performance, and that this impact varies across the different methods.

Keywords: Automated Valuation Models (AVMs), Machine Learning, Spatial Training Level, Model Performance, Valuation Accuracy

*Corresponding author: bastian.kraemer@ur.de

Introduction

The use of Automated Valuation Models (AVMs) in the context of traditional real estate valuations and their performance have been discussed in the scientific community for several decades, and increasingly scrutinized by practitioners in recent years. Most studies focus on the comparison of different statistical methods. Accordingly, there is a large body of literature comparing traditional hedonic models with more modern approaches of machine learning (ML), or approaches from the field of spatial econometrics (see e.g. Pace & Hayunga, 2020). The aim of these studies is to find out which method is best suited for estimating real estate values or prices.

The use of Automated Valuation Models (AVMs) in the context of traditional real estate valuations and their performance have been discussed in the scientific community for several decades, and increasingly scrutinized by practitioners in recent years. Most studies focus on the comparison of different statistical methods. Accordingly, there is a large body of literature comparing traditional hedonic models with more modern approaches of machine learning (ML), or approaches from the field of spatial econometrics (see e.g. Pace & Hayunga, 2020). The aim of these studies is to find out which method is best suited for estimating real estate values or prices.

The use of Automated Valuation Models (AVMs) in the context of traditional real estate valuations and their performance have been discussed in the scientific community for several decades, and increasingly scrutinized by practitioners in recent years. Most studies focus on the comparison of different statistical methods. Accordingly, there is a large body of literature comparing traditional hedonic models with more modern approaches of machine learning (ML), or approaches from the field of spatial econometrics (see e.g. Pace & Hayunga, 2020). The aim of these studies is to find out which method is best suited for estimating real estate values or prices.

Apart from the method selection, AVMs can also be optimized in many other areas. For example, the data selection and the cleaning or preparation of the selected data play an important role for the performance of the AVM. Another aspect is the choice of spatial level on which to train the selected methods. This is decisive for determining which data are ultimately included in the estimation of the AVM and thus what information is used and what is ignored. Thanks to georeferencing, models can in principle be trained at any level. For example, a model can be trained either at the level of a city, the associated commuter belt, or even at a nationwide level. However, this aspect has received little to no attention from the academic community until now. The published research on AVMs usually deals only with a manually defined region and fails to test of the methods used on different spatial levels. One reason for this might be

that historically, the availability of suitable real estate data¹ for academic purposes has been limited and therefore, analyses could only be conducted in the limited area where the data was available. However, data availability has improved massively in recent years, which is why this has now become less of a factor (Mortgage Bankers Association, 2019). In the meantime, there are providers of real-estate-related data in almost every country, which centrally force a collection of existing data and make them available for further analysis. Another reason could be the usually assumed heterogeneity of real estate markets. Traditionally, real estate markets are assumed to have a certain regionality, which in turn would mean that data from other diverging regions would not provide further explanatory power. However, the fundamental question arises as to whether this heterogeneity is generally present or whether there are not also basic characteristics that apply consistently to all markets. If this is the case, then it may be possible to achieve a higher degree of valuation accuracy by adding further data from different markets.

Therefore, the question arises as to whether the right choice of spatial level for training the models does not also represent an important, and so far underestimated role in improving the performance of AVMs. The aim of our research is to answer this question and to investigate the influence of training statistical models used for AVMs on different spatial levels.

For this purpose, we compare a total of four different methods trained on four differing spatial levels each, and compare the overall performance of the models. Our objective is not primarily a comparison of the methods used, but a specific comparison within the individual methods with respect to their performance on different spatial levels. We are interested in whether different methods deliver different results, and whether there are any specific patterns to be that emerge. The methods we select for this purpose represent a collection of regularly used ones in academic studies related to AVMs. In addition to parametric Ordinary Least Square (OLS) regressions, we analyse semi-parametric Generalized Additive Models (GAM) as well as eXtreme Gradient Boosting (XGBoost) algorithms and Deep Neural Networks (DNN) from the field of modern ML. Our analysis is based on a dataset of about 1.2 million residential properties across Germany provided by professional real estate appraisers. The four spatial levels are based on the NUTS nomenclature of the European Union. The NUTS (Nomenclature of territorial units for statistics) classification is a hierarchical system for dividing up the economic territory of the EU and the UK. In total, there are four different subdivision levels, called NUTS-

¹ In order to avoid a structural break within the dataset, the data should ideally come from one source or have been collected according to the same criteria.

0, NUTS-1, NUTS-2 and NUTS-3 which we use to train our models on a country, state, cross-regional, and county level, respectively².

Our research has various theoretical and practical implications that collectively help to improve the valuation accuracy of AVMs. Our findings show that the right choice of spatial training level can have a significant influence on the model performance of different AVM algorithms, and that this influence varies considerably, depending on the type of method. The results indicate that for parametric and semi-parametric approaches, it is advisable to choose a training level that is relatively small. This shows that the trained OLS and GAM are not able to draw further explanatory power from observations that lie outside a certain region. The results for the two modern ML algorithms are quite different. We observe that they are able to gain a higher degree of explanatory power by adding further observations, and that this effect outweighs that of local heterogeneity. Therefore, we recommend for modern ML algorithms choosing, a generally higher training level.

Literature Review

AVMs are computer-based applications, which use various statistical and algorithmic approaches to assess the value or price of a property in an automated manner. Used correctly, they can be a cost-effective and rapid alternative to traditional valuation procedures (Schulz et al., 2014). AVMs emerged mainly from the results of research in the area of hedonic price models (HPM). HPMs were developed to estimate the effects of individual characteristics, so-called marginal prices, of a good on its value or price. By aggregating these marginal prices, the overall value of a good can subsequently be calculated (Chau & Chin, 2002). HPMs were first brought into a real estate context by Lancaster (1966) and Rosen (1974). As Malpezzi (2003) and Sirmans et al. (2005) show, a diverse and dynamic field of research has emerged since then, addressing a wide variety of real-estate-specific issues.

To improve the quality of automated real estate appraisals, the focus of the research community in recent years has been almost exclusively on finding the best-fitting method. For this purpose, a large number of different approaches were either newly designed, or methods from other areas were adapted and applied. The applied methods cover the complete bandwidth of statistical methods and can be classified either as parametric, semi-parametric or non-parametric approaches. Regarding parametric approaches, most common multiple linear regression (MLR) models are applied and tested. Schulz et al. (2014), for example, use a flexible parametric

² Due to data availability, the models could not be analysed at an even smaller spatial level.

hedonic regression introduced by Bunke et al. (1999) to measure the potential predictive performance of an AVM applied to the housing market of Berlin in Germany. Other examples of parametric approaches can be found at Tse (2002), Pace and LeSage (2004), Páez et al. (2008), Bourassa et al. (2008) Osland (2010) and Zurada et al. (2011). Semi-parametric approaches can come in a variety of different forms. An often-used semi-parametric approach is the GAM, first introduced by Hastie and Tibshirani (1986). In contrast to traditional MLR models, the GAM is able to automatically control for non-linear relations between the dependent and independent variables. An early and prominent application within a real estate context is the study of Pace (1998). The author applies a GAM to a dataset for residential properties in Memphis (Tennessee) and finds that the GAM is able to outperform parametric and polynomial methods in terms of predictive behaviour. A more recent example of the GAM can be found in Dąbrowski and Adamczyk (2010). Non-parametric approaches are a category of methods which do not need an a-priori specified functional form regarding the predictor. Instead, the form is learned by the information derived from the data itself. Given this flexibility non-parametric approaches are usually able to account for non-linearities and interactions within datasets and able to outperform parametric and semi-parametric approaches (Stang et al., 2022). Prominent examples of non-parametric approaches include modern machine learning methods like Support Vector Machines, Artificial Neural Networks or Tree-Based Models. A real-estate-specific application of ML methods can be found in Mayer et al. (2019). The authors apply three commonly used basic techniques of modern ML (Random Forrest Regression, Gradient Boosting, Neural Networks) and compare their performance against some more traditional parametric approaches. Their findings show that the non-parametric methods are clearly able to outperform the more stricter parametric approaches. Other real-estate-specific applications of non-parametric modelling techniques can be found at Chun Lin and Mohan (2011), Yoo et al. (2012), Antipov and Pokryshevskaya (2012), W. J. McCluskey et al. (2013), Kok et al. (2017) and Yilmazer and Kocaman (2020).

Another aspect with regard to the optimization of the valuation accuracy of AVMs is, besides the method selection, the choice of spatial level for training the models. The level at which the models are trained implies for which data, and thus ultimately also which information is considered in the context of the valuation and which is not. This could have a strong influence on the results of the models and is therefore a factor that should not be neglected. AVM-related studies currently always focus on a predefined region. The region to which the analyses are limited is in most cases the city level or the immediate surroundings of a city. Yao et al. (2018), for example, focus on the city level of Shenzhen (China), and W. McCluskey et al. (2012)

choose the Lisburn District Council area around Belfast (North Ireland) to test their hypotheses. Other authors go a step further and conduct their analysis at the city district level. Baldominos et al. (2018), for example, focused on the Salamanca district of Madrid (Spain), Hong et al. (2020) run their analysis for the Gangnam district of Seoul (South Korea), and Yilmazer and Kocaman (2020) run their model at the Mamak district of Ankara (Turkey). However, none of the authors investigates whether the chosen level is also the best one for training the models.

To the best of our knowledge, there is currently no study that deals with the optimal spatial level for training AVMs. Our aim is therefore to close this gap in the literature and to determine the influence of the choice of spatial level on the quality of statistical valuation results. In particular, we are interested in whether this influence is the same for different types of methods (parametric, semi-parametric, non-parametric) or whether there are fundamental differences. In our analysis, we hence calculate the valuation accuracy of four different statistical methods (OLS, GAM, XGBoost, DNN) each trained at four different spatial levels, and compare their results subsequently.

Data

We base our analysis on a dataset consisting of 1,212,546 residential properties. These observations are distributed across Germany and were collected between 2014 and 2020. The dataset originates from the valuation department of one of Germany's largest mortgage lenders. Table 1 shows the distribution of the data over the observation period.

Table 1: Observations per year

	2014	2015	2016	2017	2018	2019	2020
n	196318	196403	176238	163365	165106	165996	149120
(%)	0.1619	0.1620	0.1453	0.1347	0.1362	0.1369	0.1230

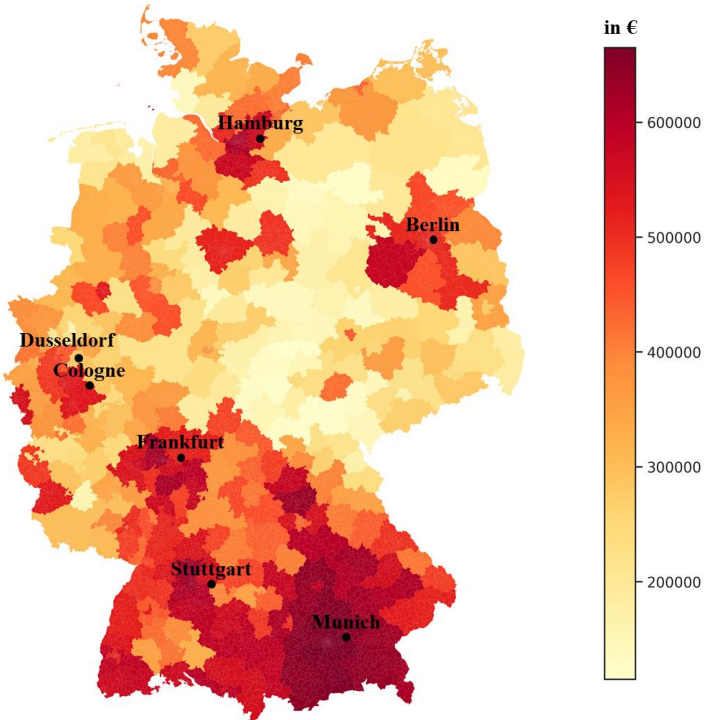
Notes: This table reports the number of observations available for each year. Over the years, the trend is slightly downward. Especially in 2020, the number of observations is lower, due to the COVID restrictions prevailing at that time. Due to the contact restrictions in place, on-site visits by appraisers were limited.

The data are actual valuation data collected by professional appraisers. We use the assessed market value as our target variable. An overview of the average market values across Germany is provided in Figure 1. The areas with the highest market values can be found in the so-called Top-7³ cities and their commuter belts. Furthermore, the market values are by far the highest in the south of Germany and tend to be lower in the east.

³ The Top-7 are the most important cities in Germany, namely Berlin, Munich, Hamburg, Frankfurt, Cologne, Dusseldorf and Stuttgart.

As hedonic characteristics, we use a set of features describing the structural characteristics, the micro-location and the macro-location of the properties. In addition, the year and quarter of the valuation are used to capture a temporal trend and seasonality. An overview of all the features used and their univariate distribution can be seen in Table 2. Before being used, the dataset was cleaned to account for duplicates, incompleteness, and erroneous data points. There are no correlations of concern within the dataset, so that all variables can be integrated accordingly.⁴

Figure 1: Average market value per district



Notes: This figure shows the average market values per NUTS-3 district. The average was calculated using all available observations within the individual districts. The highest market values can be found in the vicinity of the major metropolitan regions and in the south of Germany. The strong discrepancy between the west and east of Germany is striking. The market values observed here are also consistent with other studies (see e.g., Just & Maennig, 2012), so that it can be assumed that the observations used are representative.

⁴ The correlation matrix is available on request.

Table 2: Descriptive statistics

Variable	Unit	Mean	Median	Standard Deviation	Maximum	Minimum
Market value	Integer	228157.10	200000.00	141717.54	3860000.00	20100.00
Modernization year	Integer	1989.10	1988.00	17.19	2020.00	1950.00
Year of construction	Integer	1978.48	1981.00	29.77	2023.00	1900.00
Year of valuation	Integer	2016.82	2017.00	2.00	2020.00	2014.00
Quarter of valuation	Integer	2.45	2.00	1.12	4.00	1.00
Quality grade	Integer	3.12	3.00	0.51	5.00	1.00
Living area	Float	120.31	114.68	51.69	440.00	15.00
Lot size	Float	436.48	323.00	541.66	10000.00	0.00
Longitude	Float	9.25	8.94	1.90	19.25	5.87
Latitude	Float	50.62	50.74	1.85	55.02	47.40
Micro score	Float	72.73	74.20	14.44	99.85	0.00
Unemployment ratio	Float	4.96	4.17	2.89	26.89	0.04
Time on market	Float	12.27	10.90	4.80	106.00	0.20
Basement condominium	Binary	0.38	0.00	0.48	1.00	0.00
No basement	Binary	0.19	0.00	0.39	1.00	0.00
Basement	Binary	0.44	0.00	0.50	1.00	0.00
Owner-occupied & Non-owner-occupied	Binary	0.09	0.00	0.29	1.00	0.00
Owner-occupied	Binary	0.70	1.00	0.46	1.00	0.00
Non-owner-occupied	Binary	0.21	0.00	0.41	1.00	0.00
Object subtype condominium	Binary	0.38	0.00	0.48	1.00	0.00
Object subtype detached house	Binary	0.42	0.00	0.49	1.00	0.00
Object subtype no detached house	Binary	0.20	0.00	0.40	1.00	0.00
Condition good	Binary	0.38	0.00	0.49	1.00	0.00
Condition disastrous	Binary	0.00	0.00	0.02	1.00	0.00
Condition middle	Binary	0.45	0.00	0.50	1.00	0.00
Condition moderate	Binary	0.02	0.00	0.14	1.00	0.00
Condition bad	Binary	0.00	0.00	0.05	1.00	0.00
Condition very good	Binary	0.15	0.00	0.36	1.00	0.00
Regiotype aggro commuter belt	Binary	0.15	0.00	0.36	1.00	0.00
Regiotype aggro cbd	Binary	0.13	0.00	0.34	1.00	0.00
Regiotype aggro middle order centre	Binary	0.13	0.00	0.34	1.00	0.00
Regiotype aggro upper order centre	Binary	0.04	0.00	0.19	1.00	0.00
Regiotype rural commuter belt	Binary	0.15	0.00	0.36	1.00	0.00
Regiotype rural middle order centre	Binary	0.07	0.00	0.26	1.00	0.00
Regiotype rural upper order centre	Binary	0.01	0.00	0.07	1.00	0.00
Regiotype urban commuter belt	Binary	0.15	0.00	0.36	1.00	0.00
Regiotype urban middle order centre	Binary	0.10	0.00	0.29	1.00	0.00
Regiotype urban upper order centre	Binary	0.07	0.00	0.26	1.00	0.00
NUTS-1	String	-	-	-	-	-
NUTS-2	String	-	-	-	-	-
NUTS-3	String	-	-	-	-	-

Notes: This table reports the descriptive statistics of the dataset. Polytomous variables are one-hot encoded to binary variables to account for the requirements of modern machine learning methods. For the rather traditional methods – OLS and GAM – these polytomous variables are dummy encoded. The numbers were determined on the basis of all available observations. Overall, both structural features and location-describing features were used for model estimation. The selection of the parameters was in accordance with other publications in the AVM

literature (see e.g., Metzner & Kindt, 2018). The parameter “market value” is the dependent variable in the model estimation.

Features describing the structural characteristics of properties include the subtype of property, year of construction, modernization year, living area, lot size, use of the property, quality grade, condition and variable denoting whether the property has a basement or not. The subtype of a property can be either a “Condominium”, “Detached house” or “Not a detached house”. The year of modernization represents when the last major refurbishment took place. The use of the building describes the possible uses, whereby the characteristics are either “Owner-occupied & Non-owner-occupied”⁵, “Owner-Occupied” or “Non-owner-occupied”. Basically, the variable describes whether a property can be rented to a third-party or not. The quality of the property is measured via a grade, on a scale ranging from 1 (very poor) to 5 (very good). The general condition of the property is represented by a categorical variable with 5 different categories ranging from disastrous to very good. The variable “Basement condominium” measures whether an apartment has an extra cellar compartment or not, whereas the “Basement” and “No Basement” variables are only valid for detached and non-detached houses.

The features describing the micro-location of the properties are the latitude and longitude, the different regiotypes, and the micro score. The regiotype is provided by Acxiom⁶, and clusters Germany into ten different area types. In general, Acxiom defines four different spatial types: “Central-Business-District”, “Agglomeration Area”, “Urban Area” and “Rural Area”. The last three can be divided further into three sub-categories each (“Upper Centers”, “Middle Centers”, “Commuter Belt”). All addresses in Germany can be allocated to one of the ten area types. The individual area types are determined according to the respective settlement structure and population density within the municipality and its surrounding area. The micro score of a location is calculated via a gravity model and reflects the accessibility in the sense of proximity to selected everyday destinations. A more detailed description of the construction of the micro score of a location can be found in Appendix I. In addition, the two socio-economic variables “unemployment ratio” and “Time-on-Market”, are included to represent the macro-location of the properties. All are available at the ZIP code level.

The spatial breakdown of our dataset is based on the NUTS nomenclature of the European Union. To account for local fixed effects, three features, namely “NUTS-1”, “NUTS-2” and

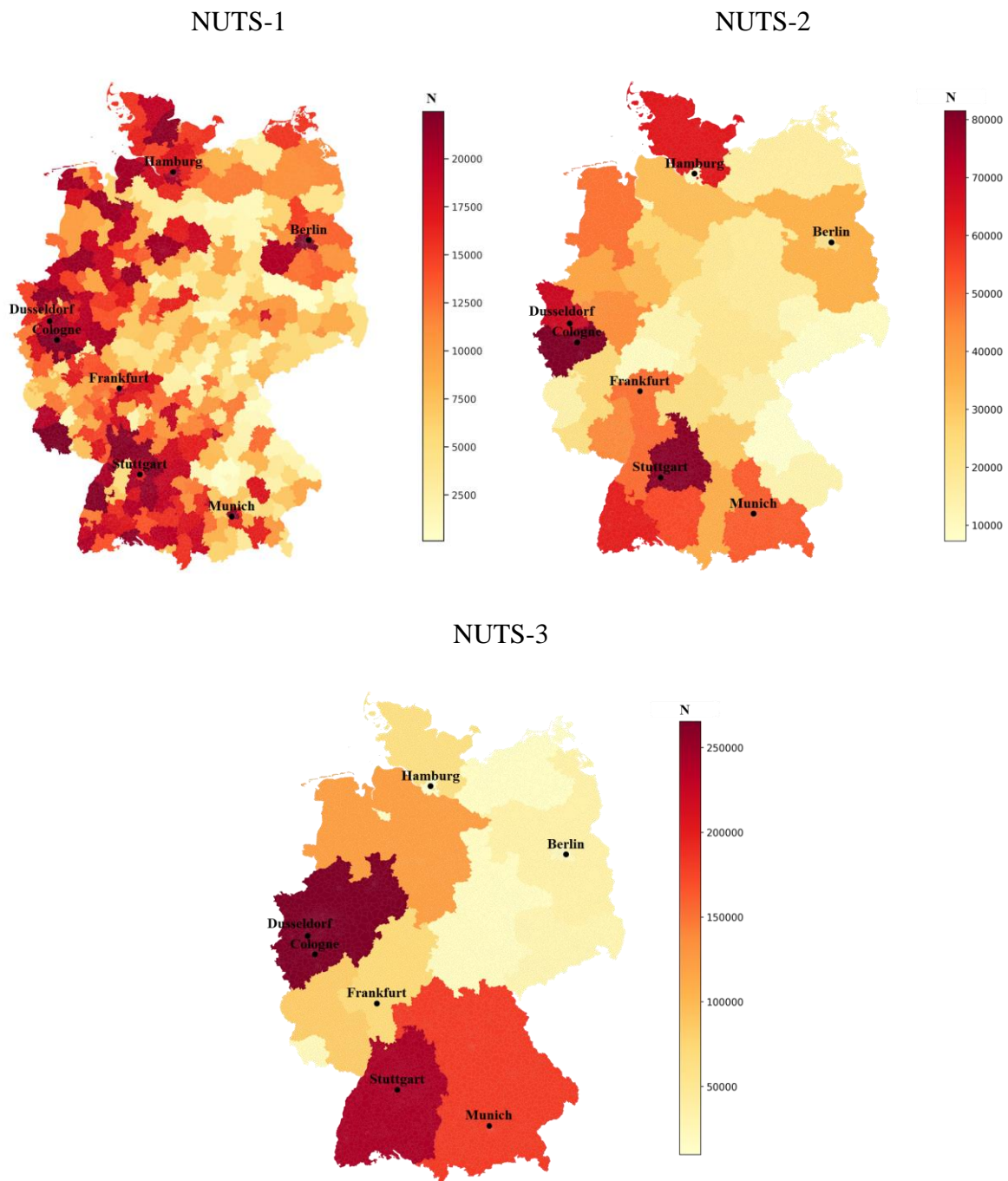
⁵ Applies if the property is both partly owner-occupied and partly non-owner-occupied (e.g., single-family home with attached rental unit).

⁶ Acxiom is an American provider of international macroeconomic and microeconomic data. Further information can be found at: <https://www.acxiom.com/>.

“NUTS-3”, are included in the dataset. NUTS-0 describes the country level, NUTS-1 describes major socio-economic regions within the country, NUTS-2 divides the corresponding NUTS-1 region into smaller basic regions for the application of regional policies, and NUTS-3 again divides the individual NUTS-2 regions into small regions for specific diagnoses.⁷ In general, Germany can be divided into a single NUTS-0, 16 NUTS-1, 38 NUTS-2 and 401 NUTS-3 regions. Since only a few observations were available in some NUTS-3 regions, we had to combine these regions and ended up with a total of 327 NUTS-3 regions for our analysis. Figure 2 provides an overview of the different NUTS regions and the number of observations available for the specific regions. Analysing the NUTS-3 level, most observations are located around the largest German metropolitan areas like Berlin, Hamburg and Munich. In addition, the NUTS-2 and NUTS-1 levels show that a difference can be observed between west and east Germany, with the east tending to have fewer observations. This is consistent with the widely diverging population figures between these regions. A comprehensive introduction to the structure of the German regions can be found at Just and Schäfer (2017), and a more detailed overview of the German real estate markets is given by Just and Maennig (2012).

⁷ Further information about the NUTS nomenclature can be found at: <https://ec.europa.eu/eurostat/web/nuts/background>.

Figure 2: Number of observations per NUTS region



Notes: This figure highlights the observations available for the individual NUTS regions. It can be seen that fewer observations are available, especially in the eastern part of Germany. This can be explained by the generally lower market activity in these regions. Structurally, these regions are mostly rural and characterized by a high level of out-migration and vacancies. The data distribution is therefore not a dataset-specific distortion, but rather a representative reflection of the German residential real estate market.

Methodology

Ordinary Least Square Regression - OLS

The first method applied is an Ordinary Least Square Regression (OLS). The main advantage of the OLS is that it is easy both to understand and to interpret. Therefore, it is the most commonly used machine learning method and often considered as a benchmark. The aim of an OLS is to explain a dependent variable y_i , $i \in \{1, \dots, n\}$, with independent variables $x_{i,1}, \dots, x_{i,k}$, a-priori unknown parameters $\beta_0, \beta_1, \dots, \beta_k$ and an error term ε_i :

$$y_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_k x_{i,k} + \varepsilon_i,$$

for all observations with

$$\mu_i = E[y_i] = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_k x_{i,k}.$$

Thereby, the relationship between the dependent variable and independent variables is assumed to be linear in parameters, and the error terms ε_i are considered to be independent and to have a constant variance. For further information, we recommend a look at Fahrmeir et al. (2013).

Several optimizations were performed to account for locational differences and to achieve the best model performance, including backward stepwise regression, interaction terms and variable transformations.

Generalized Additive Model – GAM

Our second method is a Generalized Additive Model (GAM). It is a further development of the OLS and based essentially on the concept of the Generalized Linear Model. A monotonic link function is used to model the relationship between the expected value of the dependent variable and the independent variables. The main advantage of the GAM over the OLS is that unspecified, non-parametric smoothing functions s_j , $j \in \{1, \dots, k\}$, of the covariates can be included in the model:

$$g(\mu_i) = \beta_0 + s_1(x_{i,1}) + \dots + s_k(x_{i,k}).$$

For a more extensive description of the GAM, we recommend Wood (2017).

Again, multiple model optimizations were carried out. Additional to the above-mentioned methods, this time, different penalized spline types like cubic and thin plate splines were considered. As in the OLS, these optimizations must be implemented manually.

Extreme Gradient Boosting – XGBoost

The third method is a so-called Extreme Gradient Boosting (XGBoost) algorithm, is a tree-based ensemble learning method. Ensemble learning algorithms train many weak learners h_m , in our case, single decision trees and combine them to form one strong learner h :

$$h(\mathbf{y}|\mathbf{x}) = \sum_{i=1}^M u_m h_m(\mathbf{y}|\mathbf{x}),$$

with u_m being used to weight the weak learners. M denotes the number of single trees, \mathbf{x} is the features space and \mathbf{y} the response variable. In boosting, the weak learners h_m are trained sequentially. The algorithm starts with one model and uses the errors made to improve the subsequent trees. In Gradient boosting, the so-called gradient descent algorithm is used to add new trees in order to minimize the loss of the model. The eXtreme Gradient Boosting is a computationally effective and highly efficient version of Gradient Boosting. The advantage of XGBoost is that it can recognise very complex patterns within a large amount of data. However, it is not clear from the model structure why a certain result occurs. The eXtreme Gradient Boosting is a computationally effective and highly efficient version of Gradient Boosting. The XGBoost can automatically detect complex non-linearities or higher-order interactions within a large dataset, with fewer manual optimizations needed, compared to the OLS and GAM. A detailed description of tree-based methods, ensemble learning and gradient boosting can be found in Hastie et al. (2001).

Deep Neural Network – DNN

Lastly, we consider deep neural networks (DNN), a popular and performant machine learning technique. DNNs are designed from biological neural networks (Pham, 1970), like the human brain, and consist of multiple layers, which are typically densely connected. In turn, each layer consists of numerous neurons, each processing the weighted output of all (hence the term dense) neurons of the previous layer, combined with a bias value, and applies a so-called activation function onto this linear combination. To capture this formally, let y be a neuron in the current layer, and let n be the number of neurons in the previous layer. For $i \in \{1, \dots, n\}$, let x_i be the output of the i -th neuron in the previous layer and let w_i be the according weight. Furthermore, let f be the activation function of the current neuron and b the bias term. Then, the output of the neuron y is

$$y = f(b + \sum_{i=1}^n x_i w_i).$$

A DNN then consists of multiple such neurons and layers.

To train a DNN for a specific task and data, the weights and biases are adapted. In a forward-propagation step, the data is passed through the DNN in batches. For each datum in a batch, a prediction is calculated and the predictions are then evaluated with regard to a loss function. The weights and biases are then adjusted to minimize the loss function using gradient descent. After all the data is passed through the DNN once, we say that one epoch has passed. After sufficiently many epochs, the DNN is trained and predictions for a new object can be obtained by passing the object through the DNN again.

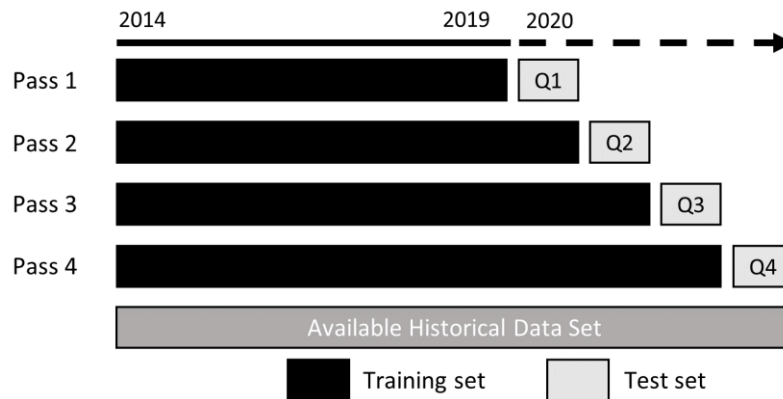
Finding the right architecture of a DNN for the task at hand is an important yet tedious task. We use the hyperparameter optimization framework Optuna (Akiba et al., 2019) to find suitable architectures for each considered region. In particular, we allow Optuna to choose the number of layers, the number of neurons per layer and the activation function per layer. Furthermore, we allow Optuna to choose the dropout rate per layer, which controls how many neurons per layer are actually activated.

The advantages of deep neural networks are that they are very flexible and adapt automatically to all data. Therefore, they can capture complex non-linearities and higher order interactions by themselves. Besides that, compared to other modern machine learning approaches, deep neural networks require less computation power in order to produce reliable results. For more information about DNNs, see Goodfellow et al. (2016).

Testing concept

To evaluate the predictive performance of the models, an extending window approach is implemented according to Mayer et al. (2019). Figure 3 illustrates the testing concept.

Figure 1: Extending window approach



Notes: This figure visualizes the applied extending window approach-testing strategy. The strategy is the right choice for the purposes of this study, as it best reflects the test procedure of conventional AVM providers and thus provides a strong reference to reality. AVM providers usually update their models on a quarterly basis as well. The results obtained in this way therefore represent an extract that is in all probability also achievable in a real-life situation.

The first iteration divides the dataset into a training set with observations from Q1/2014 to Q4/2019 and a test set from Q1/2020. In the next steps, the data of the tested quarter is added to the training set, and the models are retrained and tested on data of the next quarter. The advantages of this approach are that all algorithms are tested on unseen data and thus produce unbiased, robust results. Furthermore, the testing approach provides a realistic testing scenario. In Table 3, the number of training and test observations for each iteration are presented.

Table 3: Training and test observations

Data split	Q1	Q2	Q3	Q4
Training	1,063,426	1,106,866	1,141,612	1,180,741
Test	43,440	34,746	39,129	31,805

Notes: This table shows the number of training and test observations over the four quarters of 2020. The number of training data increases over the quarters by the number of test data from the previous quarter. With regard to the test data, it can be seen in particular that fewer observations are available in Q2 and Q4. This can be attributed to COVID restrictions which made it difficult to conduct assessment visits, especially shortly after the pandemic outbreak (Q2) and during the winter (Q4).

Evaluation metrics

For each model, we compute the Mean Absolute Percentage Error (MAPE) and the Median Absolute Percentage Error (MdAPE) as accuracy measures. Unlike Mayer et al. (2019), we use the relative rather than the absolute measures of error to enable a better comparison between the different spatial levels. In order to obtain an overall picture of the strength and weaknesses of the algorithms, we additionally provide the proportion of predictions within 10 and 20 percent (PE(x)), as well as the coefficient of determination R^2 , following Cajias et al. (2019) and Stang et al. (2022). A detailed description of all metrics can be found in Table 4.

Table 4: Evaluation metrics

Error	Formula	Description
Mean Absolute Percentage Error (MAPE)	$MAPE(y, \hat{y}) = \frac{1}{n} \sum_{i=1}^n \left \frac{y_i - \hat{y}_i}{y_i} \right $	Mean of all absolute percentage errors. A lower MAPE signals higher overall prediction accuracy in percent.
Median Absolute Percentage Error (MdAPE)	$MdAPE(y, \hat{y}) = \text{median}\left(\left \frac{y_i - \hat{y}_i}{y_i} \right \right)$	Median of all absolute percentage errors. A lower MdAPE denotes a higher precision in percent without being sensitive to outliers.
Error buckets (PE(x))	$PE(x) = 100 \left \frac{y_i - \hat{y}_i}{y_i} \right < x$	Percentage of predictions where the relative deviation is less than $x\%$, with x being 10 and 20. A larger PE(x) signals a lower variation in the predictions.
R^2	$R^2(y, \hat{y}) = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$	Coefficient of determination. A high R^2 is an indication of better goodness of fit of the model.

Notes: This table reports the evaluation metrics used to determine the valuation accuracy of the different algorithms. All four metrics are regularly used to assess the quality of AVMS. The choice of several metrics in total allows a more differentiated statement to be made than would be the case with just one metric.

Results

The aim of our study is to find out whether the choice of spatial level for training statistical models has an influence on their performance, and whether this influence is the same for all methods, or whether there are differences between more traditional and modern ML methods. In contrast to other publications, the main focus is not on which method performs best overall, but on an intra-method comparison, in order to determine which spatial level seems best suited for which method. Essentially, this enables finding out whether the assumed local heterogeneity of real estate markets is also reflected in the results of the valuation methods, or whether greater valuation accuracy can be achieved by adding further observations from other submarkets. For this purpose, two traditional approaches (OLS & GAM) as well as two modern ML approaches (XGBoost & DNN) are each trained for different spatial levels (NUTS-0, NUTS-1, NUTS-2, NUTS-3).

Below, we show the results for all four methods. In order to achieve comparability and to be able to make a valid statement, we evaluate the results on an aggregated level. For this purpose, we first provide a table for each method that shows the individual evaluation metrics for the four spatial levels of all test observations. For the metrics in the “NUTS-3” row, for example, all test data is predicted with the various different models calculated at the NUTS-3 level and finally, the metrics are calculated for the nationwide aggregated residuals. For the other three levels, the procedure is then the same. Furthermore, four maps are shown for each method. The maps are essentially a cartographic representation of the results of the MAPE on a NUTS-3 level from the tables presented earlier. The representation allows for more detailed interpretations with respect to regional performance. For example, it allows us to determine whether the results differ across different regions and whether general data availability plays a role.

Results of the ordinary least squares regression

The OLS results presented in Table 5 yield a clear pattern: The smaller the spatial level, the better the performance. In terms of the MAPE, the NUTS-3 models, which divide Germany into a total of 327 submarkets, are more than 3 percentage points better than the NUTS-0 model, which considers Germany as one overall market. In relative terms, this represents a performance increase of 18.0%. R^2 also shows that the NUTS-3 models are far superior to the NUTS-0 model. They are able to explain 4.8% more of the variance of the dependent variable.

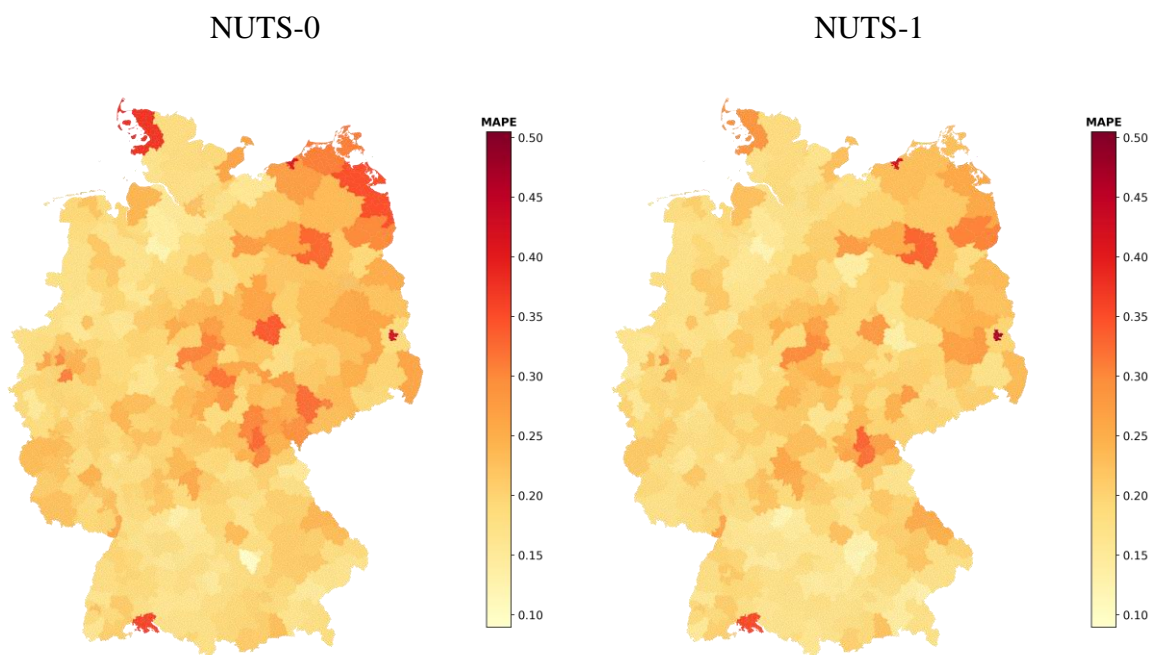
Table 5: OLS – model prediction errors 2020 throughout Germany

Models	MAPE	MdAPE	PE(10)	PE(20)	R ²
OLSNUTS-0	0.2023	0.1521	0.3423	0.6236	0.8214
OLSNUTS-1	0.1914	0.1454	0.3577	0.6473	0.8389
OLSNUTS-2	0.1852	0.1407	0.3688	0.6612	0.8487
OLSNUTS-3	0.1714	0.1294	0.3985	0.7004	0.8692

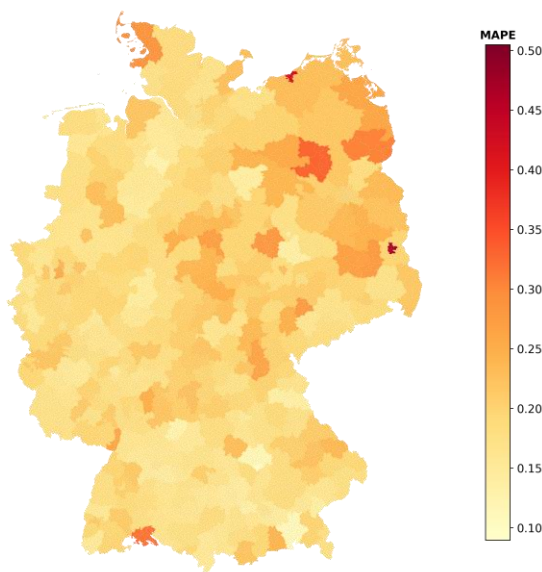
Notes: This table reports the model prediction errors for the OLS. The results are clear across all metrics and show that model performance improves with a decreasing spatial training level. This result confirms the correctness of the proceeding that in parametric approaches, a data selection that is as granular as possible must be conducted in each case.

The cartographic representation in Figure 4 illustrates once more the results from Table 5. It can be clearly seen that the lower the spatial training level, the better the MAPE for each region. The maps further show that the increased performance at the aggregate level can essentially be attributed to improved performance in the eastern parts of Germany. In addition, the German North Sea Island group around Sylt stands out on the top left of the maps. Here, it can be seen that the performance in the NUTS-3 models is much better than in the NUTS-0 model. The real estate market on Sylt and the surrounding islands is characterized by very distinct peculiarities. Residential properties are traded there only at top prices and there is a strong dependency between the specific location of the property and its value.

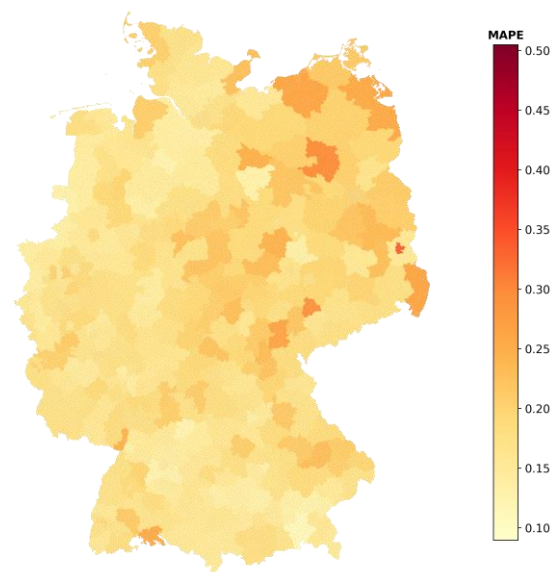
Figure 4: MAPE of the different OLS models



NUTS-2



NUTS-3



Notes: This figure visualizes the MAPE of the four different OLS models. The maps show the average absolute percentage error obtained when applying the individual models within a given region. For a granular representation, the 327 NUTS-3 regions were selected as the corresponding levels of representation. The representation of the scale is chosen so that the minimum and maximum are the largest and smallest errors respectively of all four methods.

In summary, the OLS is only able to capture local effects of the German residential real estate market, when trained on a small spatial level. Therefore, it is advisable to use the smallest possible spatial level, in our case NUTS-3, for training the OLS. These results also seem to make sense in theory, since the OLS is very generalizing in its structure and therefore can hardly (or not at all) take into account local characteristics of individual regions, if training is done on a global level. For the NUTS-0 model, the coefficients of the OLS are smoothed by too many individual regional effects and this leads to a significant deterioration in performance. In the case of an OLS, it should therefore always be ensured that only regional data are used to determine the coefficients, and ideally that different submarkets are delimited from one another in advance.

Results of the generalized additive model

The results for the GAM, shown in Table 6, are also clear and similar to those for the OLS. The more granular the spatial level for training, the better the estimation accuracy. This is true for all five evaluation metrics used. This time, the MAPE at NUTS-3 level is even 23.8% better than the NUTS-0 model. The R^2 also shows that the NUTS-3 models can explain 4.4% more of the variance of the dependent variable than the NUTS-0 model. If we look at the general

performance of the GAM and compare it with the results of the OLS, we see that the GAM is generally better able to correctly estimate the market values of the properties. It can be seen that the use of non-linear functions, which characterizes the GAM, results in a boost in performance. It is interesting to note, however, that this effect only seems to really come into play at a granular level. While the relative difference between the MAPEs of the NUTS-0 models of the OLS and the GAM is only 2.6%, it increases continuously and amounts to 7.7% at the level of the NUTS-3 models.

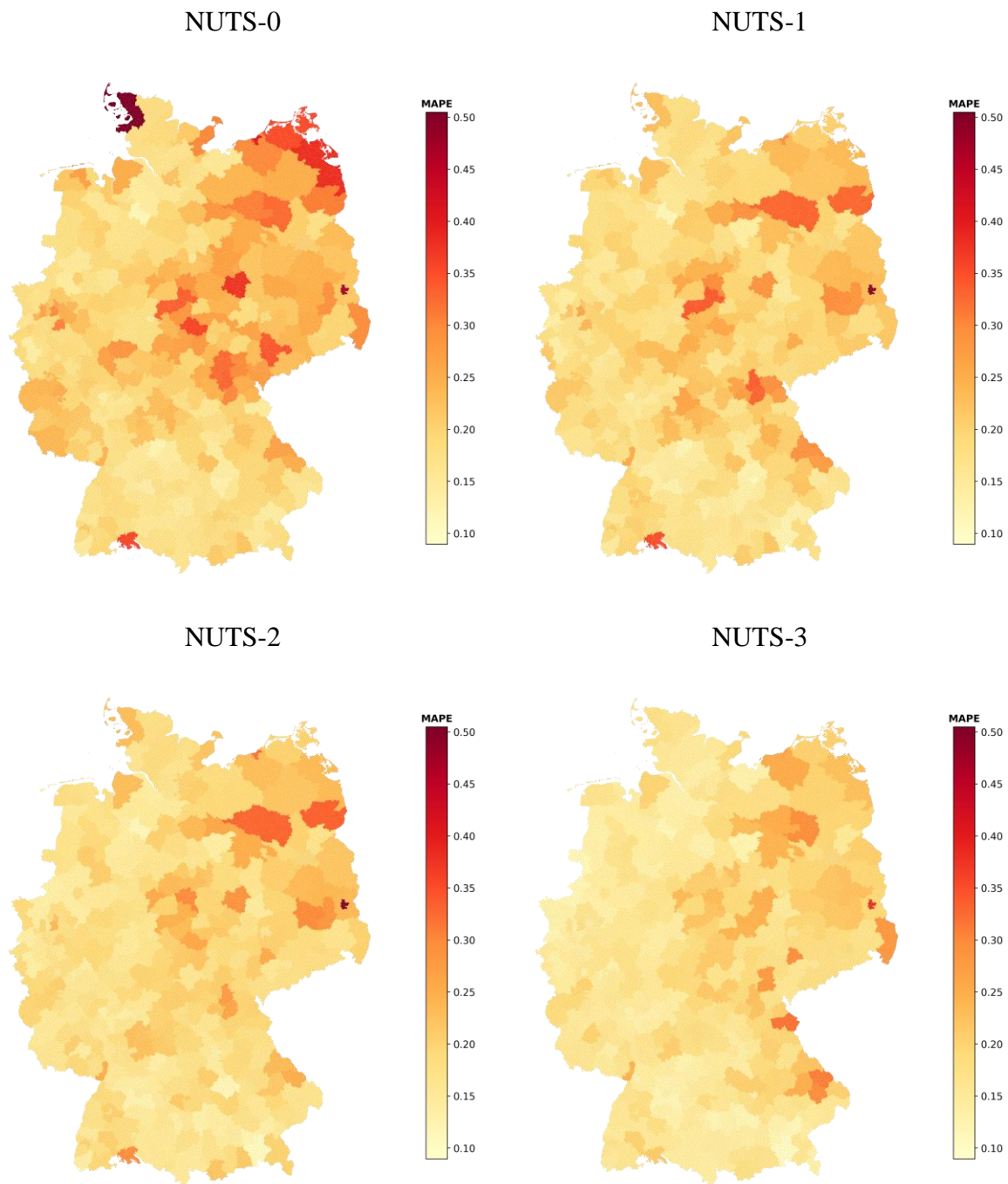
Table 6: GAM – model prediction errors 2020 throughout Germany

Models	MAPE	MdAPE	PE(10)	PE(20)	R ²
GAM _{NUTS-0}	0.1971	0.1423	0.3641	0.6504	0.8299
GAM _{NUTS-1}	0.1832	0.1339	0.3852	0.6800	0.8478
GAM _{NUTS-2}	0.1734	0.1273	0.4044	0.7028	0.8656
GAM _{NUTS-3}	0.1592	0.1160	0.4398	0.7426	0.8737

Notes: This table reports the model prediction errors for the GAM. The results are also clear across all metrics and similar to the results of the OLS. They show that model performance improves with a decreasing spatial training level. Again, the implication is that the smallest spatial level should be chosen to achieve the best model performance.

The cartographic representation in Figure 5 again shows the same picture as the OLS. Once again, it is noticeable that especially the estimation accuracy in the eastern part of Germany can be improved by implementing the method on a granular level. Furthermore, the group of islands around Sylt stands out again and again it applies that the smaller the spatial level for training the model, the better the performance.

Figure 5: MAPE of the different GAM models



Notes: This figure depicts the MAPE of the four different GAM models. The maps show the average absolute percentage error obtained when applying the individual models within a given region. For a granular representation, the 327 NUTS-3 regions were selected as the corresponding levels of representation. The representation of the scale is chosen so that the minimum and maximum are the largest and smallest errors respectively of all four methods.

In summary, the same feedback as for the OLS can be formed for the GAM. On a higher spatial level, the GAM does not manage to represent the heterogeneity of the individual residential real estate markets as well as it does on a granular level. Therefore, when using a GAM for estimating residential property values, the smallest possible level should be used for training.

Results of eXtreme gradient boosting

Compared to the first two methods, the results of the XGBoost yield a different picture. The evaluation metrics from Table 7 show that the performance is very similar on all four NUTS levels, and the greatest accuracy is achieved this time on the NUTS-1 level and not, as with the OLS and the GAM, on the NUTS-3 level. This is interesting in that, as shown in the literature review, most of the academic studies dealing with ML algorithms focus on the NUTS-3 level, and this spatial level even yield the worst performance in our case. Relative to the NUTS-1 level, the NUTS-3 level based on MAPE is 2.9% worse in terms of valuation accuracy. Although the differences between the individual metrics are only small in absolute values, if these are considered in relative terms, then a small performance boost is shown by the correct choice of the spatial level.

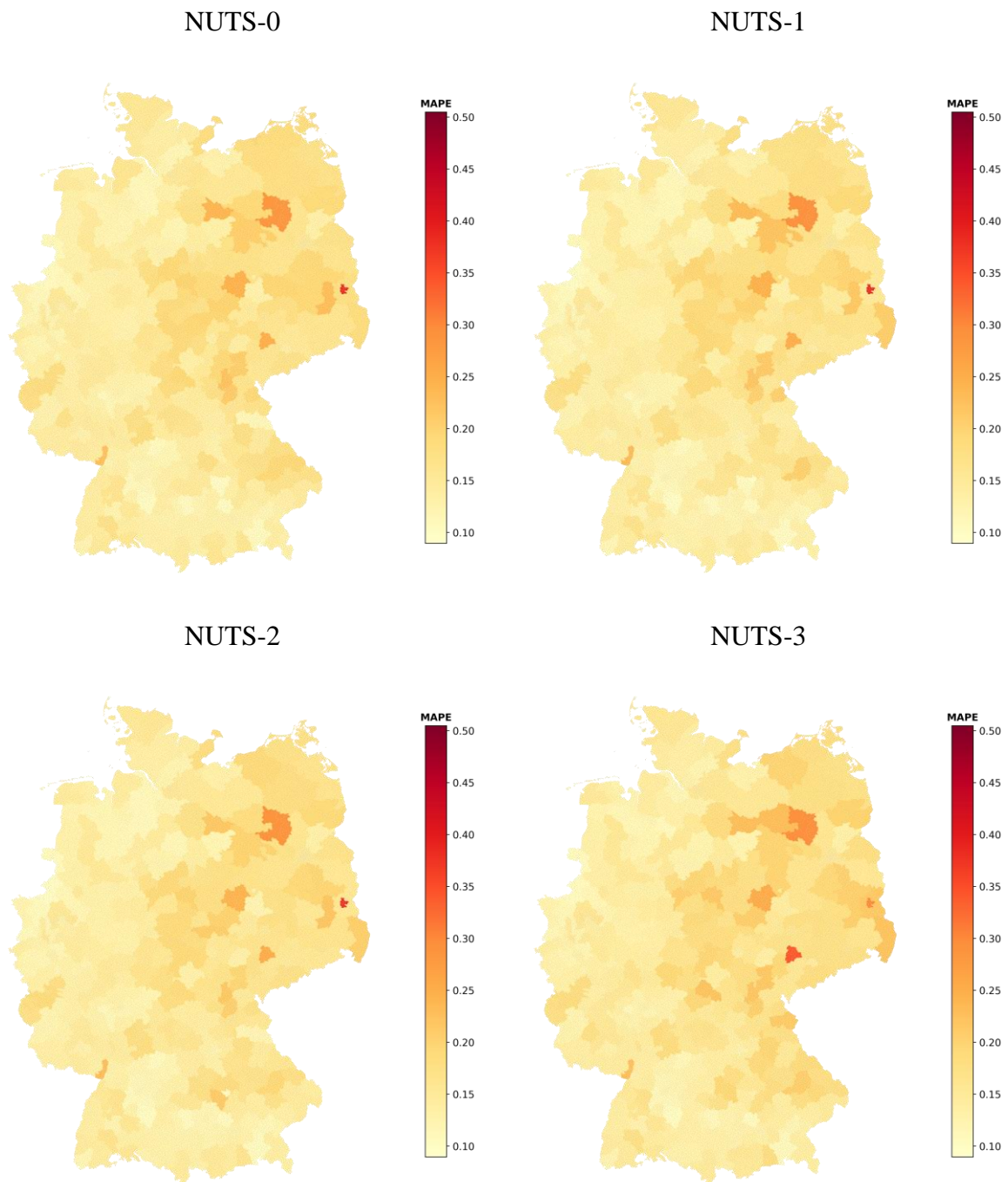
Table 7: XGBoost – model prediction errors 2020 throughout Germany

Models	MAPE	MdAPE	PE(10)	PE(20)	R ²
XGB _{NUTS-0}	0.1426	0.1077	0.4693	0.7780	0.9051
XGB _{NUTS-1}	0.1402	0.1064	0.4739	0.7869	0.9078
XGB _{NUTS-2}	0.1407	0.1071	0.4719	0.7850	0.9074
XGB _{NUTS-3}	0.1442	0.1107	0.4578	0.7733	0.9036

Notes: This table reports the model prediction errors for the XGBoost. Here, too, the results are the same across all evaluation metrics. Unlike for the first two methods, however, the model performance of the XGBoost does not improve with a decreasing spatial training level, but is relatively constant across all levels. The best performance is achieved at the NUTS-1 level, indicating that the XGBoost can gain a higher degree of explanatory power by adding more data.

The analysis of the maps from Figure 6 shows in particular that in the parts of Germany where few observations are available (see Figure 2), the choice of a higher spatial level for training the models leads to a performance improvement. It is an important implication that in regions where little data are available, it can be useful in the case of the XGBoost to include data from other surrounding districts. This represents an essential difference to the results of the GAM and the OLS. For them, especially in the parts of Germany with low data availability, the results deteriorates with the choice of a higher spatial level for training the models.

Figure 6: MAPE of the different XGB models



Notes: This figure visualizes the MAPE of the four different XGBoost models. The maps show the average absolute percentage error obtained when applying the individual models within a given region. For a granular representation, the 327 NUTS-3 regions were selected as the corresponding levels of representation. The representation of the scale is chosen so that the minimum and maximum are the largest and smallest errors respectively of all four methods.

In summary, it can be seen that the heterogeneity of local real estate markets can still be detected by the XGBoost when it is trained at a higher spatial level. In some cases, the use of additional data even leads to a further improvement of the estimation accuracy. Therefore, unlike OLS and GAM, the NUTS-3 level is not the optimal spatial level for training the XGBoost, but the NUTS-1 level. However, the results of Table 7 also show that there seems to be a limit regarding the optimal size of the spatial level. The results at NUTS-0 level are still better than those at NUTS-3 level, but not as good as on NUTS-1 and NUTS-2 level.

Results of the neural network

Finally, in analysing the results of the DNN, we again see a different picture. The evaluation metrics presented in Table 8, show that the DNN can clearly improve its valuation accuracy as the spatial training level increases. This is the exact opposite of the OLS and GAM results, and also a different result compared to the XGBoost. Although the results of the MAPE indicate that the NUTS-1 level performs best here as well, the four other metrics yield a slightly different picture for this specific algorithm. They evaluate the NUTS-0 level as the best suited. In principle, therefore, the situation between the NUTS-0 and NUTS-1 levels is quite similar, influenced only by marginal changes. Compared to the other modern ML algorithm, the XGBoost, it seems that the number of observations used to optimize the algorithm plays an more important role. This is also logical from the point of view of the complexity of the method. The DNN can only show its strength in recognising non-linear relationships and multi-layer interactions, if a sufficiently large number of observations is available. This finding is also in line with those of Nghiep and Al (2001), which show on the basis of a dataset for Rutherford County, Tennessee, that neural networks perform better than multiple regression analysis only with increasing dataset size.⁸

Table 8: DNN – model prediction errors 2020 throughout Germany

Models	MAPE	MdAPE	PE(10)	PE(20)	R ²
DNN _{NUTS-0}	0.1551	0.1080	0.4700	0.7620	0.8705
DNN _{NUTS-1}	0.1542	0.1090	0.4648	0.7595	0.8664
DNN _{NUTS-2}	0.1595	0.1142	0.4471	0.7448	0.8606
DNN _{NUTS-3}	0.1656	0.1176	0.4356	0.7281	0.8461

Notes: This table reports the model prediction errors for the DNN. The results clearly show that model performance can be increased by choosing the highest possible spatial training level. Unlike the XGBoost results, the increase

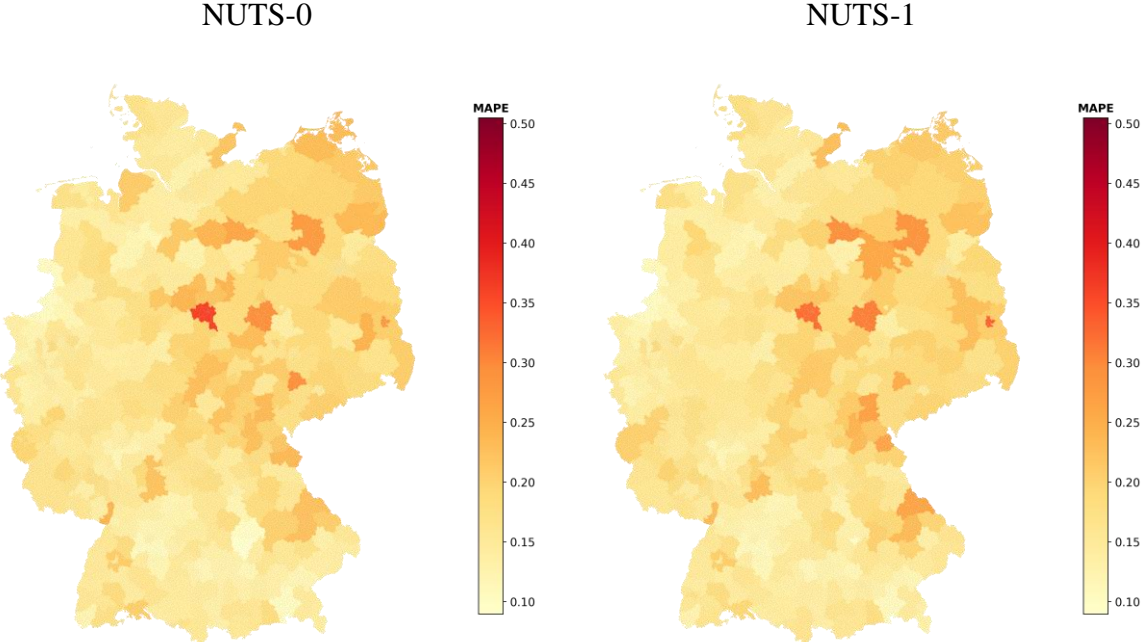
⁸ However, unlike in our study, these authors only work at a county level and only vary the amount of data available within the county. In our case, the amount of data is varied by adding observations from other spatial levels.

in performance is much more significant. The results indicate that the DNN is only able to show its strength from a certain amount of data.

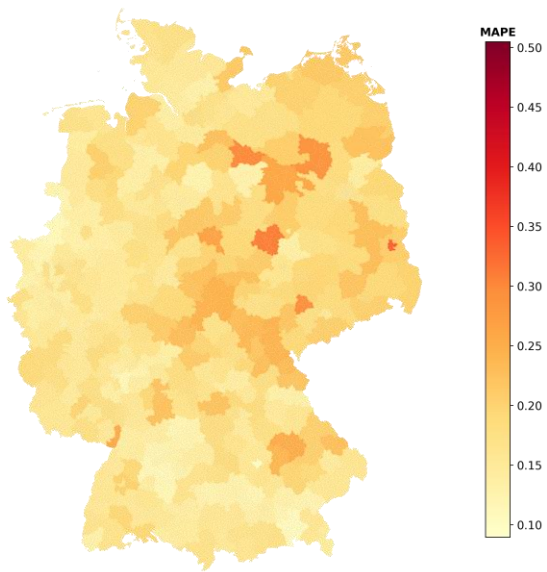
The visual representation of the results in Figure 7 yield a very similar picture to the XGBoost results. It can be seen that by choosing a higher training level for the DNN, the valuation performance can be increased, especially in areas with few observations. Again, the same implication emerges as with the XGBoost, that in regions where few data are available, it can be useful to include data from other surrounding districts. With respect to the four algorithms used for the analysis, this can only be empirically proven for the modern ML algorithms, which represents a significant contribution to the literature of this study.

In summary, the DNN is only able to estimate property values as accurately as possible once a certain number of observations has been used. The effect of adding more observations therefore outweighs the effect of local heterogeneity. Thus, the DNN is independently able to generate further explanatory power for a specific real estate market, even from data outside the specific market. Regional effects can thus be more effectively detected, extracted and extrapolated by modern ML algorithms. With respect to the DNN, it is therefore advisable to choose as high as possible a training level or to maximize the available observations for training the algorithm.

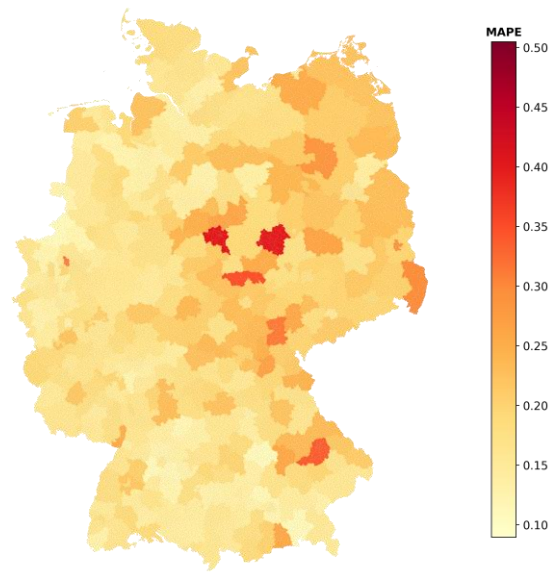
Figure 7: MAPE of the different DNN models



NUTS-2



NUTS-3



Notes: This figure visualizes the MAPE of the four different DNN models. The maps show the average absolute percentage error obtained when applying the individual models within a given region. For a granular representation, the 327 NUTS-3 regions were selected as the corresponding levels of representation. The representation of the scale is chosen so that the minimum and maximum are the largest and smallest errors respectively of all four methods.

Conclusion

This study is intended to answer the question of whether the right choice of the appropriate spatial level for training AVM algorithms also plays an important and so far underestimated role in improving the valuation accuracy of AVMs. We use a dataset consisting of about 1.2 million residential properties across Germany to test our hypotheses for a total of four different but typical AVM algorithms (OLS, GAM, XGBoost, DNN). All four algorithms are each trained on four different spatial levels after which the results are evaluated. The four spatial levels are based on the NUTS nomenclature of the European Union. We use the NUTS-0, NUTS-1, NUTS-2 and NUTS-3 levels to train our models on a country, state, cross-regional, and county level, respectively

Our results indicate that the correct choice of spatial training level can exert a significant influence on the model performance, and that this can vary considerably, depending on the type of method. With respect to the OLS results, it is advisable to select a training level that is as granular as possible, since this is the only way to ensure that the most accurate valuations are made. It can be seen that there are regional differences and thus certain heterogeneities, which the OLS can only recognise as accurately as possible if they are locally limited. The results for the GAM yield a similar picture to the OLS. Here, too, the model performance correlates positively with a smaller spatial training level. Accordingly, the same findings can be generated for the parametric and the semi-parametric approaches. These confirm the correctness of the current trend in academic publications and in practice, of choosing the most granular analysis level possible for traditional econometric methods. These two methods are not able to draw further explanatory power from observations that lie outside a certain region. On the contrary, they even suffer from it. The results of the two applied modern ML algorithms are quite different. With respect to the XGBoost, the evaluation metrics show that the choice of the most suitable spatial level can be made relatively indifferently. Although there are marginal differences with respect to the evaluation accuracy, these are only minor compared to OLS and GAM. In contrast to the parametric and semi-parametric approaches, the results of the non-parametric XGBoost show that the performance actually increases slightly with increasing spatial training level, and the NUTS-1 level seems to be the most appropriate. This trend can be observed even more clearly for the results of the DNN. Here, it can be seen that the performance does not decrease with an increasing training level, as is the case with the OLS and the GAM, but clearly improves. With respect to the two modern ML algorithms, it can be seen that they are able to gain a higher degree of explanatory power by adding further observations, and that this effect outweighs that of local heterogeneity. In particular, their ability

to recognise and map non-linear relationships and multi-layered interactions allows them to exploit overlapping effects of different regions to achieve more accurate real estate valuations. This is particularly evident in regions, so as where there are few observations. In these cases, it is advisable to train a modern ML algorithm with additional regions in order to benefit from their basic commonalities.

In summary, the choice of the right training level should always depend on the method. For parametric and semi-parametric methods, we recommend using a spatial level which is as granular as possible for training the models, since these are only able to separate local heterogeneities from each other to a limited extent. For non-parametric modern ML methods, however, we generally recommend a higher training level. These complex methods are able to detect regional differences independently and to separate them. Furthermore, they benefit from the fact that there are basic commonalities in the functioning of local real estate markets, and that these can be used to increase their explanatory power. With respect to the practical application and implementation of AVM algorithms, this offers the additional advantage that the higher training level means that fewer models have to be trained and calibrated overall. For example, less effort is required with regard to data preparation and processing. Thus, efficiencies can be increased for AVM providers operating nationwide, and significant economic advantages can be achieved.

References

- Akiba, T., Sano, S., Yanase, T., Ohta, T., & Koyama, M. (2019, July 25). *Optuna: A Next-generation Hyperparameter Optimization Framework*.
- Antipov, E. A., & Pokryshevskaya, E. B. (2012). Mass appraisal of residential apartments: An application of Random forest for valuation and a CART-based approach for model diagnostics. *Expert Systems with Applications*, *39*(2), 1772–1778. <https://doi.org/10.1016/j.eswa.2011.08.077>
- Baldominos, A., Blanco, I., Moreno, A., Iturrarte, R., Bernárdez, Ó., & Afonso, C. (2018). Identifying Real Estate Opportunities Using Machine Learning. *Applied Sciences*, *8*(11), 2321. <https://doi.org/10.3390/app8112321>
- Bourassa, S. C., Cantoni, E., & Hoesli, M. (2008). *Predicting House Prices with Spatial Dependence: Impacts of Alternative Submarket Definitions*. <https://doi.org/10.2139/ssrn.1090147>
- Bunke, O., Droge, B., & Polzehl, J. (1999). Model Selection, Transformations and Variance Estimation in Nonlinear Regression. *Statistics*, *33*(3), 197–240. <https://doi.org/10.1080/02331889908802692>
- Cajias, M., Willwersch, J., & Lorenz, F. (2019). *I know where you will invest in the next year – Forecasting real estate investments with machine learning methods*. European Real Estate Society (ERES). ERES. https://ideas.repec.org/p/arz/wpaper/eres2019_171.html
- Chau, K. W., & Chin, T. L. (2002). *A Critical Review of Literature on the Hedonic Price Model*.
- Chun Lin, C., & Mohan, S. B. (2011). Effectiveness comparison of the residential property mass appraisal methodologies in the USA. *International Journal of Housing Markets and Analysis*, *4*(3), 224–243. <https://doi.org/10.1108/17538271111153013>
- Dąbrowski, J., & Adamczyk, T. (2010). Application of GAM additive non-linear models to estimate real estate market value. *Geomatics and Environmental Engineering*, *4*(2), 55-62.
- Fahrmeir, L., Kneib, T., Lang, S., & Marx, B. (2013). *Regression: Models, Methods and Applications*. Springer-Verlag.
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.
- Handy, S. L., & Clifton, K. J. (2001). Evaluating Neighborhood Accessibility: Possibilities and Practicalities. *Journal of Transportation and Statistics*, *4*(2).
- Hastie, T., Friedman, J., & Tibshirani, R. (2001). *The Elements of Statistical Learning*. Springer New York. <https://doi.org/10.1007/978-0-387-21606-5>
- Hong, J., Choi, H., & Kim, W. (2020). A House Price Valuation based on Random Forest Approach: The Mass Appraisal of Residential Property in South Korea. *International Journal of Strategic Property Management*, *24*(3), 140–152. <https://doi.org/10.3846/ijspm.2020.11544>
- Huang, Y., & Dall’erba, S. (2021). Does Proximity to School Still Matter Once Access to Your Preferred School Zone Has Already Been Secured? *The Journal of Real Estate Finance and Economics*, *62*(4), 548–577. <https://doi.org/10.1007/s11146-020-09761-w>
- Just, T., & Maennig, W. (Eds.). (2012). *Understanding German Real Estate Markets*. Springer Berlin Heidelberg. <https://doi.org/10.1007/978-3-642-23611-2>
- Kok, N., Koponen, E.-L., & Martínez-Barbosa, C. A. (2017). Big Data in Real Estate? From Manual Appraisal to Automated Valuation. *The Journal of Portfolio Management*, *43*(6), 202–211. <https://doi.org/10.3905/jpm.2017.43.6.202>

- Lancaster, K. J. (1966). A New Approach to Consumer Theory. *Journal of Political Economy*, 74(2), 132–157. <https://doi.org/10.1086/259131>
- Malpezzi, S. (2003). Hedonic Pricing Models: A Selective and Applied Review. In *Housing Economics and Public Policy* (pp. 67–89). <https://doi.org/10.1002/9780470690680.ch5> (Original work published 2003)
- Mayer, M., Bourassa, S. C., Hoesli, M., & Scognamiglio, D. (2019). Estimation and updating methods for hedonic valuation. *Journal of European Real Estate Research*, 12(1), 134–150. <https://doi.org/10.1108/JERER-08-2018-0035>
- McCluskey, W. J., McCord, M [M.], Davis, P. T., Haran, M [M.], & McIlhatton, D [D.] (2013). Prediction accuracy in mass appraisal: A comparison of modern approaches. *Journal of Property Research*, 30(4), 239–265. <https://doi.org/10.1080/09599916.2013.781204>
- McCluskey, W., Davis, P., Haran, M [Martin], McCord, M [Michael], & McIlhatton, D [David] (2012). The potential of artificial neural networks in mass appraisal: The case revisited. *Journal of Financial Management of Property and Construction*, 17(3), 274–292. <https://doi.org/10.1108/13664381211274371>
- Metzner, S., & Kindt, A. (2018). Determination of the parameters of automated valuation models for the hedonic property valuation of residential properties. *International Journal of Housing Markets and Analysis*, 11(1), 73–100. <https://doi.org/10.1108/IJHMA-02-2017-0018>
- Mortgage Bankers Association. (2019). *The State of Automated Valuation Models in the Age of Big Data*.
- Nghiep, N., & Al, C. (2001). Predicting Housing Value: A Comparison of Multiple Regression Analysis and Artificial Neural Networks. *Journal of Real Estate Research*, 22(3), 313–336. <https://doi.org/10.1080/10835547.2001.12091068>
- Nobis, C., & Kuhnimhof, T. (2018). *Mobilität in Deutschland – MiD: Ergebnisbericht*.
- Osland, L. (2010). An Application of Spatial Econometrics in Relation to Hedonic House Price Modeling. *Journal of Real Estate Research*, 32(3), 289–320. <https://doi.org/10.1080/10835547.2010.12091282>
- Pace, R. K., & Hayunga, D. (2020). Examining the Information Content of Residuals from Hedonic and Spatial Models Using Trees and Forests. *The Journal of Real Estate Finance and Economics*, 60(1-2), 170–180. <https://doi.org/10.1007/s11146-019-09724-w>
- Pace, R. K (1998). Appraisal Using Generalized Additive Models. *Journal of Real Estate Research*, 15(1), 77–99. <https://doi.org/10.1080/10835547.1998.12090916>
- Pace, R. K, & LeSage, J. (2004). Spatial Statistics and Real Estate. *The Journal of Real Estate Finance and Economics*, 29(2), 147–148. <https://doi.org/10.1023/b:real.0000035307.99686.fb>
- Páez, A., Long, F., & Farber, S. (2008). Moving Window Approaches for Hedonic Price Estimation: An Empirical Comparison of Modelling Techniques. *Urban Studies*, 45(8), 1565–1581. <https://doi.org/10.1177/0042098008091491>
- Pham, D. T. (1970). Neural Networks In Engineering. *WIT Transactions on Information and Communication Technologies*, 6, 3–36. <https://doi.org/10.2495/AI940011>
- Powe, N. A., Garrod, G. D., & Willis, K. G. (1995). Valuation of urban amenities using an hedonic price model. *Journal of Property Research*, 12(2), 137–147. <https://doi.org/10.1080/09599919508724137>

- Rosen, S. (1974). Hedonic Prices and Implicit Markets: Product Differentiation in Pure Competition. *Journal of Political Economy*, 82(1), 34–55.
<https://doi.org/10.1086/260169>
- Schulz, R., Wersing, M., & Werwatz, A. (2014). Automated valuation modelling: A specification exercise. *Journal of Property Research*, 31(2), 131–153.
<https://doi.org/10.1080/09599916.2013.846930>
- Sirmans, S., Macpherson, D., & Zietz, E. (2005). The Composition of Hedonic Pricing Models. *Journal of Real Estate Literature*, 13(1), 1–44.
<https://doi.org/10.1080/10835547.2005.12090154>
- Stang, M., Krämer, B., Nagl, C., & Schäfers, W. (2022). From human business to machine learning—Methods for automating real estate appraisals and their practical implications. *Zeitschrift Für Immobilienökonomie*, 1–28.
<https://doi.org/10.1365/s41056-022-00063-1>
- Tse, R. Y. C. (2002). Estimating Neighbourhood Effects in House Prices: Towards a New Hedonic Model Approach. *Urban Studies*, 39(7), 1165–1180.
<https://doi.org/10.1080/00420980220135545>
- Wood, S. N. (2017). *Generalized Additive Models: An Introduction with R, Second Edition*. CRC Press.
- Yang, J., Bao, Y., Zhang, Y., Li, X., & Ge, Q. (2018). Impact of Accessibility on Housing Prices in Dalian City of China Based on a Geographically Weighted Regression Model. *Chinese Geographical Science*, 28(3), 505–515.
<https://doi.org/10.1007/s11769-018-0954-6>
- Yao, Y., Zhang, J., Hong, Y., Liang, H., & He, J. (2018). Mapping fine-scale urban housing prices by fusing remotely sensed imagery and social media data. *Transactions in GIS*, 22(2), 561–581. <https://doi.org/10.1111/tgis.12330>
- Yilmazer, S., & Kocaman, S. (2020). A mass appraisal assessment study using machine learning based on multiple regression and random forest. *Land Use Policy*, 99, 104889. <https://doi.org/10.1016/j.landusepol.2020.104889>
- Yoo, S., Im, J., & Wagner, J. E. (2012). Variable selection for hedonic model using machine learning approaches: A case study in Onondaga County, NY. *Landscape and Urban Planning*, 107(3), 293–306. <https://doi.org/10.1016/j.landurbplan.2012.06.009>
- Zurada, J., Levitan, A., & Guan, J. (2011). A Comparison of Regression and Artificial Intelligence Methods in a Mass Appraisal Context. *Journal of Real Estate Research*, 33(3), 349–388. <https://doi.org/10.1080/10835547.2011.12091311>

Appendix

Appendix 1 – Micro Score

The micro score of a location is calculated via a gravity model, and reflects the accessibility in the sense of proximity to selected everyday destinations. A gravity model is a common method for approximating the accessibility of a location and is based on the assumption that nearby destinations play a greater role in everyday life than more distant destinations (Handy and Clifton (2001)). The relevant points-of-interest (POIs) are selected from the findings of Powe et al. (1995), Metzner and Kindt (2018), Yang et al. (2018), Nobis and Kuhnimhof (2018) and Huang and Dall'erna (2021) and are provided in Table 3.

Table 3: Features of the micro score of a location

Points-of-Interests	Category	Description
University	Education & Work	University campus: institute of higher education
School	Education & Work	Facility for education
Kindergarten	Education & Work	Facility for early childhood care
CBD	Education & Work	Center of the next city
Supermarket	Local Supply	Supermarket – a large shop with groceries
Marketplace	Local Supply	A marketplace where goods are traded daily or weekly
Chemist	Local Supply	Shop focused on selling articles for personal hygiene, cosmetics, and household cleaning products
Bakery	Local Supply	Place for fresh bakery items
ATM	Local Supply	ATM or cash point
Hospital	Local Supply	Facility providing in-patient medical treatment
Doctors	Local Supply	Doctor's practice / surgery
Pharmacy	Local Supply	Shop where a pharmacist sells medications
Restaurant	Leisure & Food	Facility to go out to eat
Café	Leisure & Food	Place that offers casual meals and beverages
Park	Leisure & Food	A park, usually urban (municipal)
Fitness Centre	Leisure & Food	Fitness Centre, health club or gym
Movie Theater	Leisure & Food	Place where films are shown
Theater	Leisure & Food	Theatre where live performances take place
Shopping Mall	Leisure & Food	Shopping Centre– multiple shops under one roof
Department Store	Leisure & Food	Single large shop selling a large variety of goods
Subway Station	Transportation	City passenger rail service
Tram Station	Transportation	City passenger rail service
Railway Station	Transportation	Railway passenger only station.
Bus Stop	Transportation	Bus stops of local bus lines.
E-Charging Station	Transportation	Charging facility for electric vehicles

Note: The descriptions of the selected Points-of-Interest is based on the explanations of Open Street Map.⁹

⁹ See https://wiki.openstreetmap.org/wiki/Map_features.

Our gravity model can be described using an activity function $f(A_p)$ and a distance function $f(D_{i,p})$:

$$A_{i,p} = \sum f(A_p)f(D_{i,p}).$$

Here, $A_{i,p} \in [0,100]$ denotes the accessibility of point i for the POI p , whereby the activity function $f(A_p)$ specifies the relative importance of POI p , with $f(A_p) \in [0,1]$. The function $f(D_{i,p})$ measures the travel time from point i to the POI p by using a non-symmetric sigmoidal distance function. The travel time was obtained for the selected POIs via Open Street Map¹⁰ and normalized using the following function:

$$L(x) = \frac{K}{(1 + Qe^{0.5x})^{\frac{1}{v}}},$$

where $K, Q \in \mathbb{R}$ and $v \in \mathbb{R}^+$ are defined for all possible distances $x \in \mathbb{R}$. Furthermore, we have:

$$\begin{aligned} K &= (1 + Q)^{1+v}, \\ Q &= v \cdot \exp(B \cdot x^*), \\ v &= \frac{\exp(B \cdot x^*) - 1}{\ln(y_i) - 1}, \end{aligned}$$

where x^* denotes a feature specific point of inflection and y^* is 0.5.

¹⁰ <https://www.openstreetmap.org/>