

STACKING-BASED UNCERTAINTY MODELLING OF STATISTICAL AND MACHINE LEARNING METHODS FOR RESIDENTIAL PROPERTY VALUATION

A. Jafari¹, M.R. Delavar^{2*}, A. Stein³

1- MSc. Graduate, Department of GIS, School of Surveying and Geospatial Eng. College of Engineering, University of Tehran, Tehran, Iran - alijafari7525@ut.ac.ir

2- Centre of Excellence in Geomatic Eng. in Disaster Management, School of Surveying and Geospatial Eng. College of Engineering, University of Tehran, Tehran, Iran – mdelavar@ut.ac.ir

3- Faculty of Geo-Information Science and Earth Observation (ITC), University of Twente, The Netherlands – a.stein@utwente.nl

Commission IV, WG IV/3

KEY WORDS: Uncertainty Modelling, Support Vector Regression, Weighted K-Nearest Neighbours, Random Forest, Ordinary Least Squares, Stacking, Residential Property Valuation

ABSTRACT:

Estimating real estate prices helps to adapt informed policies to regulate the real estate market and assist sellers and buyers to have a fair business. This study aims to estimate the price of residential properties in District 5 of Tehran, Capital of Iran, and model its associated uncertainty. The study implements the Stacking technique to model uncertainties by integrating the outputs of basic models. Basic models must have a good performance for their combinations to have acceptable results. This study employs four statistical and machine learning models as basic models: Random Forest (RF), Ordinary Least Squares (OLS), Weighted K-Nearest Neighbour (WKNN), and Support Vector Regression (SVR) to estimate the price of residential properties. The results show that the integrated output is more accurate for the quadruple combination mode than for any of the binary and triple combinations of the basic models. Comparing the Stacking technique with the Voting technique, it is shown that the Mean Absolute Percentage Error (MAPE) reduces from 10.18% to 9.81%. Hence we conclude that our method performs better than the Voting technique.

1. INTRODUCTION

The real estate market is recognized as a fundamental factor to support national economies (Pai and Wang, 2020; Wiles, 1974). From the Iranian people's perspective, there is an impression that housing is one of the best investment options. When there is an opportunity to invest, the real estate market is preferred to other investments such as in the gold or car industries. There is also a belief among the Iranian people that an investment in real estate will never lose its capital. Therefore, most of investors consider real estate a suitable investment. Hence, accurate estimation of real estate prices is helpful for buyers and sellers, and local governments formulate effective and informed policies to regulate the real estate market (Dong et al., 2020). Estimating the price of residential real estate, is challenging, due to several uncertainties in the process (Liu et al., 2011).

Nur et al. (2017) employed regression analysis and particle swarm optimization (PSO) to predict housing prices in Malang, East Java, Indonesia. PSO is used for feature selection, followed by a regression analysis to estimate the price. This study showed that the combination of regression analysis and PSO was a suitable method obtaining an MAPE of 14.19%. Čeh et al. (2018) investigated a random forest (RF) machine learning model to estimate residential property prices and compare it with the Multi Regression Analysis (MRA) model in Ljubljana, Slovenia. The result of this study showed that the RF algorithm could better detect changes in apartment prices and predict them more effectively than MRA in complex urban forms. The MAPE for RF and MRA were equal to 7.27% and 17.48%, respectively. Pai and Wang, (2020) used genetic algorithms and machine learning models to estimate real estate prices in Taiwan, including least squares support vector regression (LSSVR), classification and regression trees (CART), general

regression neural networks (GRNN) and backpropagation neural networks (BPNN). Employing genetic algorithms for feature selection, they showed that the LSSVR model with MAPE = 0.23% has the lowest error and the BPNN model with MAPE = 14.42% has the highest error. Jafari and Delavar (2021) integrated a support vector machine (SVM), a genetic algorithm (GA), and PSO to estimate the price of residential properties in Tehran. The MAPE values for SVR integrated with PSO and GA were equal to 10.13% and 10.14%, respectively. Koohpayma and Argany (2020) estimated prices of residential apartments in Tehran using ordinary least squares (OLS) and geographically weighted regression (GWR). Their performance was compared with different structural and spatial parameters as well as their effect on house price estimation. They concluded that GWR has a good performance for estimating the price of a residential apartment.

Previous research showed that methods used to combine the output from the models that consider the uncertainty of their results have not been extensively considered. Furthermore, a model may not solve a problem or solve it optimally in some cases alone. Therefore, instead of using one model to solve a problem, several models are combined to enhance the accuracy of the output estimation (Dixit, 2017). Stacking is a combination method for the outputs of the models that is used in this study. We assume that using several models and combining their outputs, increases the accuracy of estimating residential property prices.

2. BASIC CONCEPTS

In this section, a theoretical framework for the employed methods in this paper is discussed.

* Corresponding author

2.1 Ordinary Least Squares

Ordinary least squares regression (OLS), also called linear regression, is a widely used method for fitting linear statistical models. It depends upon a number of k independent variables. If $k = 1$, it is called simple linear regression, while for $k > 1$, it is called multiple linear regression. OLS applied to linear regression was introduced by Gauss and presents a straightforward type of statistical prediction (Weisberg, 2005). In OLS modelling, the unknown coefficients are assumed to be universal, i.e., they are assumed to be fixed in the study area and are location independent. The model can be written as Equation (1) (Weisberg, 2005):

$$y_i = \beta_0 + \beta_1 x_{i,1} + \beta_2 x_{i,2} + \beta_3 x_{i,3} + \dots + \beta_k x_{i,k} + \varepsilon_i$$

$$\rightarrow y_i = \beta_0 + \sum_{j=1}^k \beta_j x_{i,j} + \varepsilon_i \quad (1)$$

$i = 1, 2, 3, \dots, n > k$

where, β_1, \dots, β_k , are the unknown coefficients of the independent variables x_1, \dots, x_k and y_i is the dependent variable.

2.2 Weighted K-Nearest Neighbours

Nearest neighbour classification is one of the oldest and simplest non-parametric algorithms for classifying data. It uses a simple and efficient way on the basis of neighbour samples to classify new sample data, i.e. sample data that were not used in the learning process (Fan et al., 2019; Gou et al., 2012). This classifier is used in two ways including k-nearest neighbour (KNN) and weighted k-nearest neighbour (WKNN). In the KNN algorithm, the nearest neighbour can make different decisions for different values of k , while the optimal k improves model performance (Gou et al., 2012). When k is low, the classifier may be affected by the noise samples which leads to a wrong decision. While if k is large, it may affect the values at sampled positions, also leading to wrong decisions (Gou et al., 2012). To solve this problem, weights for each neighbour are determined based upon their distance and similarity; closer neighbours are given greater weight, while farther neighbours give a lower weight. Finally, a weighted average is produced (Fan et al., 2019; Gou et al., 2012). In this study, Equation (5) has been used to obtain the weights (Gou et al., 2012):

$$d_1 = d(\bar{x}, x_1^{NN}) \quad (2)$$

$$d_i = d(\bar{x}, x_i^{NN}) \quad (3)$$

$$d_k = d(\bar{x}, x_k^{NN}) \quad (4)$$

$$\bar{w}_i = \begin{cases} \frac{d_k - d_i}{d_k - d_1} \times \frac{d_k + d_1}{d_k + d_i}, & \text{if } d_k \neq d_1 \\ 1, & \text{if } d_k = d_1 \end{cases} \quad (5)$$

On the basis of Equation (5), the estimated value of a property price is obtained with Equation (6) (Fan et al., 2019):

$$\hat{S}_i = \frac{\sum_{j=1}^k w_{i,y_j} \times S_{y_j}}{\sum_{j=1}^k w_{i,y_j}} \quad (6)$$

2.3 Random Forest

Random Forest (RF) is a classification and regression method based on bagging (Breiman, 2001). RF is a simple yet effective algorithm for improving the accuracy of the decision tree. The

main idea is to reduce the variance to σ^2/n by averaging the set of independent observations of Z_1, \dots, Z_n with variance σ . After learning the group of trees, the predictions are made by combining random forest trees. One specimen is propagated in each tree and that specimen is directed to different T leaves. The v_i vectors associated with T-leaves are described in Equation (7) (Breiman, 2001):

$$v = \frac{1}{T} \sum_{i=1}^T v_i \quad (7)$$

The mean v_i is then the predicted value for the sample.

2.4 Support Vector Machine

Support vector machines (SVMs) were introduced in Boser et al. (1992) and generalized to the non-linear model by Vapnik (1995). For SVM, there are two main categories including support vector regression (SVR) and support vector classification (SVC) (Smola and Schölkopf, 2004). SVR is a three-layer neural network that is able to overcome the sample of small problems and offers better results than the neural network in such problems (Lin and Chen, 2011). A main advantage of SVR is that its computational complexity does not depend on the dimensions of the feature vector. In addition, it has a high generalizability and a predictive accuracy (Awad and Khanna, 2015).

2.5 The Integrated Model

Integration of the models consists of two main parts. In the first part, all basic models are trained using a training data set, while in the second part, the combination rules and decision-making algorithms are included (Ge et al., 2009). In this research, the average output of the models has been used as a final estimate which can be simple averaging (Voting technique) or weighted averaging (Stacking technique).

2.5.1 Voting

Three items must be identified to understand the ensemble learning process in the simple averaging method known as the Voting technique. This includes (Dixit, 2017):

- How to learn: From a dataset, all models are trained;
- How to combine: All models have the same weight;
- Basic models: There should be a variety of models and efficient ones.

In Voting, all the basic models have the same weight in the decision-making process without considering the capabilities of the models.

2.5.2 Stacking

In the Stacking technique, higher weights are given to those basic models that have a better performance, i.e. weights are based upon the uncertainty (Yang and Cao, 2018).

- How to learn:
 - i. The first level of training: from a dataset, all basic models are trained;
 - ii. The second level of training: The second level model is trained to achieve the weight of the models;
- How to combine: Weighted Averaging is undertaken and each model has a different weight in decision making based on its performance (uncertainty);

- Basic models: There should be a variety of models and efficient way of their selection (Dixit, 2017; Zhou, 2019).

2.6 Evaluation

The purpose of performance appraisal of the residential property price estimation models against real data is to test the importance of the parameters and fitness of the models. To do so, the Mean Absolute Error (MAE), the Mean Absolute Error Percentage (MAPE), the Mean Square Error (MSE), the Root Mean Square Error (RMSE) and the Pearson Correlation Coefficient (r) have been employed. The performance evaluation criteria of the models are given in Equations (8-12) (Li and Heap, 2008; Rousseau et al., 2018).

$$MSE = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 \quad (8)$$

$$RMSE = \sqrt{MSE} \quad (9)$$

$$MAE = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i| \quad (10)$$

$$MAPE = \frac{1}{n} \sum_{i=1}^n \left| \frac{\hat{y}_i - y_i}{y_i} \right| \quad (11)$$

$$R = \frac{\sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{n\sigma_{y_i}\sigma_{\hat{y}_i}} \quad (12)$$

where, \hat{y}_i is the estimated price, $\bar{\hat{y}}, \bar{y}_i$ is the mean (\hat{y}_i, y_i), y_i is the actual price, $\sigma_{y_i}, \sigma_{\hat{y}_i}$ is the variance (y_i, \hat{y}_i), and n is the number of properties.

3. METHODOLOGY

In this study, stacking and voting techniques were used to enhance the residential properties' price estimation quality. The voting technique allocates the same weight to all basic models in the decision-making process without paying attention to the advantages of the models. The Stacking technique assigns higher weights to the models that have a better performance (Dixit, 2017; Zhou, 2019). We therefore explored this technique to model uncertainty in property price estimation in this research. For this purpose, the data are initially divided into training, validation, and testing data. Such a data division is undertaken in order to achieve a homogeneous spatial distribution. To model the uncertainty of models, we used two parts. The first part trains the models based on a common data set. In the second part, a linear perceptron model is trained to calculate the weight of each model. In this way, higher weights are allocated to the models with the best performance. Using a training dataset, basic models are trained. In the second stage of training, the output from the models in the first part on the validation data is considered the input for linear perceptrons. A weighted average is used to obtain the final estimate. To investigate the effect of uncertainty modelling, the output of the Stacking technique is compared with the basic models and the Voting technique using the test data. In Figure 1 and algorithm 1, the general structure of uncertainty modelling using the Stacking technique is illustrated.

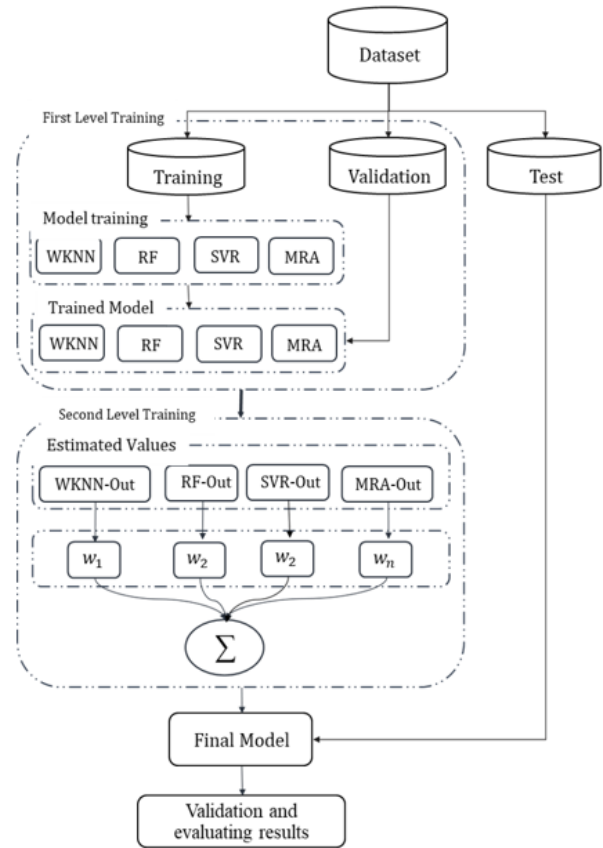


Figure 1. The general structure of the proposed uncertainty modelling using Stacking technique

Algorithm 1. Pseudo code uncertainty modelling using Stacking technique

Inputs: % Train Data
 $D_1 = \{(x_1, y_1), (x_2, y_2), \dots, (x_{n_1}, y_{n_1})\}$
 % Validation Dataset
 $D_2 = \{(x_1, y_1), (x_2, y_2), \dots, (x_{n_2}, y_{n_2})\}$
 % Test Dataset
 $D_3 = \{(x_1, y_1), (x_2, y_2), \dots, (x_{n_3}, y_{n_3})\}$
 First-level learning algorithm L_1, \dots, L_T
 Second-level learning algorithm L

Process:

1. **for** $t = 1, \dots, T$ % First-level learner training by
2. $h_t = L_t(D_1)$; % first-level learning algorithm L_t
3. **end**
4. $D' = \emptyset$; % Generate a new dataset
5. **for** $i = 1, \dots, n_2$
6. **for** $t = 1, \dots, T$
7. $z_{it} = \hat{y}_t(x_i)$;
8. **end**
9. $D' = D' \cup \{(z_{i1}, \dots, z_{iT}), y_i\}$;
10. **end**
11. $h' = L(D')$; % The second-level learner h' by the
 % Second-level learning algorithm L
 % to the new Dataset D'

Output: $Y(x) = \hat{y}'(\hat{y}_1(x), \dots, \hat{y}_T(x))$

4. IMPLEMENTATION

This section describes the study area, the employed data and the process of implementation and validation.

4.1 Study Area

In this research, an experimental study was conducted in the metropolis of Tehran, the capital of Iran. It has a population of approximately 9 million people and an area of 615 km² (*Statistical Report of Tehran Municipality, 2019*). The metropolis includes 22 districts, 123 subdivisions, and 354 neighbourhoods. Figure (2) shows the location of District 5, which is the study area.

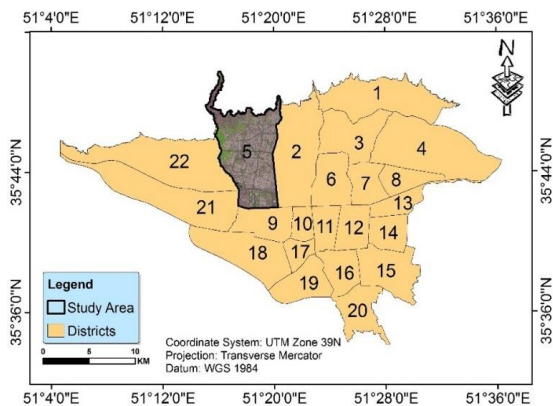


Figure 2. Study Area

The choice of the study area depends upon several factors. Since several spatial parameters have been considered in this study to estimate the residential property valuation, the study area must be large and diverse enough for the spatial phenomena to be measurable. Since District 5 is the second most populous district and the largest district of Tehran, it has been selected as the study area in this research.

4.2 The Employed Data

This study investigates the effect of 31 parameters (9 structural and 21 spatial parameters) on the price of residential properties. These parameters have been selected based upon a literature review and consultation with the concerned experts. Structural parameters are extracted from the residential real estate transaction book. This information includes property price, area, number of rooms, number of parking lots, number of warehouses and the existence of yard and balcony. The collected data are related to residential apartments in the study area from November 6, 2020 until January 16, 2021. It includes 2256 samples of residential property sales data. The dataset is downloaded from the dodota website², which provides free but limited access to information on apartments and property prices. Spatial parameters were extracted from urban maps such as land use maps, road networks, and population density maps produced by the Statistics Center of Iran and the seismic vulnerability map produced in a study undertaken by Sheikhan (2016). Finally, spatial parameters for properties in a geospatial information system (GIS) environment were calculated and extracted using these spatial information. These include:

- The number of educational centers within 200 m from each property and the distance from the nearest one; proximity to schools or universities, due to students' ease of access, which usually has a positive effect on residential property prices.
- The number of health centers (pharmacies, hospitals) within 100 m from each property and the distance

from the nearest one. It is expected that the more sanitary units around the property and the shorter the distance, the lower the prices; because high population density results in poor quality of the environment and the built environment, as well as creation of some parking problems.

- The number of major streets within 200 m from each property and the distance from the nearest one. Proximity to the major street may have a negative impact on property prices due to increased traffic, noise, and environmental pollution, or have a positive impact on property prices due to ease of access to urban facilities and infrastructures.
- The number of services, commercial and leisure centers (such as green spaces, market center, and four squares) within 200 m from each property and the distance from the nearest one. Proximity to service, commercial, and leisure centers are expected to have a positive impact on property prices.

Figure (3) and Table (1) provide distribution and descriptive variables of residential properties.

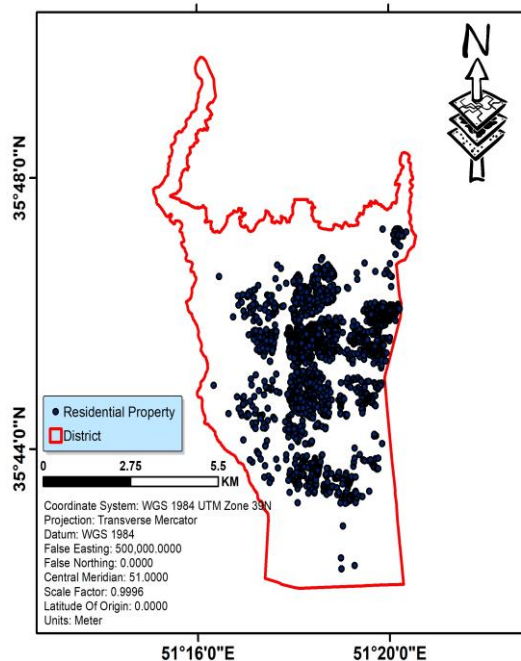


Figure 3. Residential properties distribution.

4.3 Basic Regression Models

To implement the basic regression models, the data were divided into two sub dataset of training and testing according to their spatial distribution. In the next step, the models were trained using 1804 training data samples. The models were then evaluated by considering 452 samples of the test data sets and comparing the estimated price with the real price of residential real estate. In order to improve the generality of the SVR model with the Radial Bias Function (RBF) kernel, the PSO optimization algorithm has been employed to determine the optimal hyper-parameters of this algorithm ($C = 17.2477, \sigma = 5.3198, \epsilon = 0.2758$) (Jafari and Delavar, 2021). Also, the WKNN algorithm for different k may lead to different decisions. Therefore, in order to determine the appropriate number of neighbours, the accuracy of this

² <https://dodota.com>

algorithm is examined for the number of neighbours from 1 to 100 and the best number of neighbours is determined (k=8).

Table 1. Descriptive variables of residential properties

#	Variable	Description	#	Variable	Description
1	area	Apartment total area	17	d_health	Prox. to health centres (such as drug stores and clinics)
2	num_room	Number of rooms in the apartment	18	d_hospital	Prox. to hospitals
3	yr_build	Year of construction	19	d_industry	Prox. to industrial centres
4	elevator	Existence of elevator	20	d_marketing	Prox. to business centres
5	num_parkinglot	Number of parking lots	21	population_density	Population density
6	warehouse	Existence of warehouse	22	d_foursquare	Prox. to foursquare (Azadi Tower)
7	balcony	Existence of balcony	23	n_mainroad_100	number of hospitals within 100 meters of property
8	yard	Existence of yard	24	n_hospital_100	number of major roads within 100 meters of property
9	facility	Existence of facilities (such as sports halls, swimming pools)	25	n_marketing_200	number of major roads within 200 meters of property
10	d_street	Prox. to main city roads	26	n_greenspace_200	number of major roads within 200 meters of property
11	seismic_vulnerability	Seismic vulnerability	27	n_industry_200	number of industry within 200 meters of property
12	d_religious	Prox. to main city religious places	28	n_sport_200	number of sport centres within 200 meters of property
13	d_sport	Prox. to sport facilities	29	n_religious_200	number of religious centres within 200 meters of property
14	d_cultural	Prox. to cultural centres	30	n_cultural_200	number of cultural centres within 200 meters of property
15	d_education	Prox. to educational centres (such as schools/universities)	31	n_education_200	number of education centres within 200 meters of property
16	d_green	Prox. to green space, forest			

The results of the test phase and evaluation criteria are presented in Table (2).

Table 2. Evaluation criteria for the basic models

Model	RMSE	MSE	MAE	MAPE	<i>r</i>
MRA	0.3298	0.1088	0.2699	10.44	0.8815
WKNN	0.4808	0.2312	0.3752	14.83	0.7319
RF	0.3373	0.1138	0.2637	10.10	0.8795
SVR	0.3231	0.1044	0.2619	10.12	0.8889

The results in Table (2) verify that the SVR and WKNN methods have the highest and lowest correlations between residential properties' real and estimated prices, respectively. The MAPE for SVR is 10.12% and that of the WKNN is 14.83%. Therefore, it is proved that regression models for estimating the price of residential properties with the same training data give different results, so these models are not complete. This means that each model has its own strengths and weaknesses, so integrating basic regression models that do not have the same weaknesses can improve the accuracy of the properties price estimation.

This uncertainty in the property valuation estimation can be created in different stages (for example, selection of training samples and estimation of parameters). Therefore, for the final estimation of property prices with uncertainty modelling, by integrating the results of the basic models based on uncertainty analysis, it is expected that property prices that have not been properly met in basic regression models will be improved considering this integration.

4.4 Uncertainty Modelling

For the detailed uncertainty modelling, different integrations of basic regression models were investigated. These integrations include binary, triple, and quadruple basic regression models Table (3) shows the weight of the basic models in the quadruple composition. The SVR model with $w=0.6221$ and the WKNN model with $w=0.0918$ has the highest and lowest weights, respectively. Therefore, the lower the uncertainty of the basic model, the greater would be its weight in the uncertainty modelling process.

Table 3. Weights for the basic models

Model:	WKNN	RF	SVR	MRA
Weight	0.0918	0.2577	0.6221	0.1350

The results presented in Table (4) show that integrating the output of the basic regression models based on their uncertainty has increased the accuracy of estimating the price of residential properties. The estimation accuracy of the binary compounds has been slightly improved compared to the basic models. In the double combination, RF+SVR with MAPE=9.92% has the least error. These binary combinations showed that the stronger the basic models, the better the results. In ternary compounds, an improvement in accuracy was observed in all the cases. By integrating WKNN+RF+MRA with MAPE = 9.80%, the lowest error was observed. A quadruple combination of WKNN+RF+SVR+MRA with MAPE = 9.79, the least errors was observed among these integrations. According to the results, it can be concluded that the more and stronger the number of basic models with the less uncertainty, the better the accuracy. Therefore, a quadruple combination was selected to improve the accuracy of residential property price estimates.

The results of the uncertainty modelling are presented in Table (4).

Table 4. Evaluation criteria of the uncertainty modelling based upon the Stacking technique

#	Combinations	RMSE	MSE	MAE	MAPE	R
1	WKNN+RF	0.3265	0.1066	0.2591	10.05	0.8863
2	WKNN+SVR	0.3165	0.1002	0.2577	10.02	0.8925
3	WKNN+MRA	0.3244	0.1052	0.2647	10.31	0.8865
4	RF+SVR	0.3141	0.0987	0.2564	9.92	0.8940
5	RF+MRA	0.3209	0.1030	0.2634	10.19	0.8888
6	SVR+MRA	0.3186	0.1015	0.2612	10.15	0.8909
7	WKNN+RF+SVR	0.3119	0.0973	0.2530	9.80	0.8958
8	WKNN+RF+MRA	0.3155	0.0996	0.2572	9.97	0.8931
9	WKNN+SVR+MRA	0.3162	0.1000	0.2581	10.04	0.8927
10	RF+SVR+MRA	0.3142	0.0987	0.2569	9.93	0.8939
11	WKNN+RF+SVR+MRA	0.3114	0.0970	0.2532	9.79	0.8960

4.5 Implement the Voting technique to compare results

In this section, the Voting technique combines the basic models. Unlike the Stacking technique, the Voting technique gives all models the same weight in the decision-making process. The result of this technique is shown in Table (5).

The results of Table (5) show that the Voting method has also increased the accuracy of estimating the price of residential properties. Therefore, combining the basic models can be considered as a way to improve the estimate of residential property prices.

Table 5. Evaluation criteria of combination the basic models based upon the Voting technique

#	Combinations	RMSE	MSE	MAE	MAPE	R
1	WKNN+RF	0.3720	0.1384	0.2931	11.52	0.8685
2	WKNN+SVR	0.3617	0.1308	0.2867	11.26	0.8707
3	WKNN+MRA	0.3621	0.1311	0.2868	11.26	0.8705
4	RF+SVR	0.3187	0.1016	0.2562	9.83	0.8926
5	RF+MRA	0.3223	0.1039	0.2600	9.99	0.8895
6	SVR+MRA	0.3204	0.1027	0.2633	10.16	0.8889
7	WKNN+RF+SVR	0.3394	0.1152	0.2690	10.51	0.8891
8	WKNN+RF+MRA	0.3398	0.1155	0.2696	10.54	0.8889
9	WKNN+SVR+MRA	0.3362	0.1130	0.2682	10.48	0.8870
10	RF+SVR+MRA	0.3171	0.1006	0.2574	9.91	0.8929
11	WKNN+RF+SVR+MRA	0.3280	0.1076	0.2612	10.18	0.8937

5. CONCLUSION

This study aimed to model the uncertainty of the basic models in different model integrations. According to the results, we concluded that the output combination of the basic models based upon their uncertainty improves the accuracy of estimating the price of residential properties. Voting and Stacking techniques were used to combine the output of the models. When using the Stacking technique due to weighting, the models based on their uncertainty associated with their modelling performance process, have performed better than that of the Voting technique, which has given the same weight to all models. Therefore, in modelling the uncertainty of the stacking technique, the models that had good performance were given more weight. Integrating efficient basic models have led to better results. A simple linear perceptron model was used to analyse the uncertainty in this study, so it is suggested that other compositional rules theories such as Dempster-Shafer theory be employed in future research. Other more robust basic models could be also used such as Geography Weighted Regression, Spatial Autoregressive Model, and Spatial Error Model which model spatial dependence and spatial heterogeneity.

6. REFERENCES

Awad, M., Khanna, R., 2015. Support Vector Regression, in: Efficient Learning Machines. Apress, Berkeley, CA,

pp. 67–80. https://doi.org/10.1007/978-1-4302-5990-9_4

Boser, B.E., Guyon, I.M., Vapnik, V.N., 1992. A training algorithm for optimal margin classifiers, in: Proceedings of the Fifth Annual Workshop on Computational Learning Theory. pp. 144–152.

Breiman, L., 2001. Random forests. *Machine learning* 45, 5–32.

Čeh, M., Kilibarda, M., Liseč, A., Bajat, B., 2018. Estimating the performance of random forest versus multiple regression for predicting prices of the apartments. *ISPRS international journal of geo-information* 7, 168.

Dixit, A., 2017. Ensemble Machine Learning: A beginner's guide that combines powerful machine learning algorithms to build optimized models. Packt Publishing Ltd.

Dong, S., Wang, Y., Gu, Y., Shao, S., Liu, H., Wu, S., Li, M., 2020. Predicting the turning points of housing prices by combining the financial model with genetic algorithm. *PLoS ONE* 15, e0232478. <https://doi.org/10.1371/journal.pone.0232478>

Fan, G.-F., Guo, Y.-H., Zheng, J.-M., Hong, W.-C., 2019. Application of the Weighted K-Nearest Neighbor Algorithm for Short-Term Load Forecasting. *Energies* 12, 916. <https://doi.org/10.3390/en12050916>

Ge, Y., Li, S., Lakhani, V.C., Lucieer, A., 2009. Exploring uncertainty in remotely sensed data with parallel coordinate plots. *International Journal of Applied Earth Observation and Geoinformation* 11, 413–422.

- Gou, J., Du, L., Zhang, Y., Xiong, T., 2012. A new distance-weighted k-nearest neighbor classifier. *J. Inf. Comput. Sci* 9, 1429–1436.
- Jafari, A. and M. R. Delavar, 2021. Estimating the price of residential properties based on the optimal support vector machine. Presented at the International Conference of GIScience: Basis and Trans/Interdisciplinary Applications, Mashhad, Iran.
- Koohpayma, J., Argany, M., 2020. Estimating the price of apartments in Tehran using extracted compound variables. *International Journal of Housing Markets and Analysis*.
- Li, J., Heap, A.D., 2008. A review of spatial interpolation methods for environmental scientists.
- Lin, H., Chen, K., 2011. Predicting price of Taiwan real estates by neural networks and support vector regression, in: *Proc. of the 15th WSEAS Int. Conf. on Syst.* pp. 220–225.
- LIU, X., Zhe, D., WANG, T., 2011. Real estate appraisal system based on GIS and BP neural network. *Transactions of Nonferrous Metals Society of China* 21, s626–s630.
- Nur, A., Ema, R., Taufiq, H., Firdaus, W., 2017. Modeling House Price Prediction using Regression Analysis and Particle Swarm Optimization Case Study : Malang, East Java, Indonesia. *ijacsa* 8. <https://doi.org/10.14569/IJACSA.2017.081042>
- Pai, P.-F., Wang, W.-C., 2020. Using machine learning models and actual transaction data for predicting real estate prices. *Applied Sciences* 10, 5832.
- Rousseau, R., Egghe, L., Guns, R., 2018. *Becoming metric-wise: a bibliometric guide for researchers*. Chandos Publishing is an imprint of Elsevier, Cambridge, MA.
- Sheikhian, H., 2016. *Tehran Seismic vulnerability assessment using integration of granular computing and artificial neural networks (MSc Thesis)*. University of Tehran, Tehran, Iran.
- Smola, A.J., Schölkopf, B., 2004. A tutorial on support vector regression. *Statistics and computing* 14, 199–222.
- Statistical report of Tehran Municipality, 2019. . Tehran Municipality Information and Communication Technology Organization, Tehran, Iran.
- Vapnik, V.N., 1995. *The nature of statistical learning. Theory*.
- Weisberg, S., 2005. *Applied linear regression*. John Wiley & Sons.
- Wiles, R.C., 1974. Mercantilism and the Idea of Progress. *Eighteenth-Century Studies* 8, 56. <https://doi.org/10.2307/2737891>
- Yang, B., Cao, B., 2018, Research on Ensemble Learning-based Housing Price Prediction Model, School of Software Engineering, Tongji University, Shanghai, China, *BGDDS* 1, 1–8. <https://doi.org/10.23977/bgdds.2018.11001>
- Zhou, Z.-H., 2019. *Ensemble methods: foundations and algorithms*. Chapman and Hall/CRC.