



The Deep Learning Approach to Property Valuation – An Application of a Multilayer Neural Net Model for Estimating House Prices

딥러닝 방식에 기초한 부동산 가격평가 : 다층신경망 모형을 활용한 주택 가격 추정

저자 (Authors)	Lee, Changro, Kim, Se Hyong
출처 (Source)	한국지역개발학회지 30(4) , 2018.11, 179–201 (23 pages) Journal of The Korean Regional Development Association 30(4) , 2018.11, 179–201 (23 pages)
발행처 (Publisher)	한국지역개발학회 The Korean Regional Development Association
URL	http://www.dbpia.co.kr/Article/NODE07566773
APA Style	Lee, Changro, Kim, Se Hyong (2018). The Deep Learning Approach to Property Valuation. 한국지역 개발학회지, 30(4), 179–201.
이용정보 (Accessed)	송실대학교 222.107.238.*** 2019/01/03 19:54 (KST)

저작권 안내

DBpia에서 제공되는 모든 저작물의 저작권은 원저작자에게 있으며, 누리미디어는 각 저작물의 내용을 보증하거나 책임을 지지 않습니다. 그리고 DBpia에서 제공되는 저작물은 DBpia와 구독계약을 체결한 기관소속 이용자 혹은 해당 저작물의 개별 구매자가 비영리적으로만 이용할 수 있습니다. 그러므로 이에 위반하여 DBpia에서 제공되는 저작물을 복제, 전송 등의 방법으로 무단 이용하는 경우 관련 법령에 따라 민, 형사상의 책임을 질 수 있습니다.

Copyright Information

Copyright of all literary works provided by DBpia belongs to the copyright holder(s) and Nurimedia does not guarantee contents of the literary work or assume responsibility for the same. In addition, the literary works provided by DBpia may only be used by the users affiliated to the institutions which executed a subscription agreement with DBpia or the individual purchasers of the literary work(s) for non-commercial purposes. Therefore, any person who illegally uses the literary works provided by DBpia by means of reproduction or transmission shall assume civil and criminal responsibility according to applicable laws and regulations.

The Deep Learning Approach to Property Valuation: An Application of a Multilayer Neural Net Model for Estimating House Prices

Lee, Changro* · Kim, Se Hyong**

딥러닝 방식에 기초한 부동산 가격평가:
다층신경망 모형을 활용한 주택 가격 추정

이창로 · 김세형

한국지방세연구원 부연구위원 · 한국감정원 기획조정실장

국문요약: 부동산 가격평가를 위한 전통적 접근은 선형회귀모형을 활용하는 것으로서, 이러한 접근은 설명변수와 종속변수(주택 가격)간 선형의 관계를 가정한다는 한계가 있다. 본 연구는 이러한 한계를 극복하고자 딥러닝 방식, 즉 다층신경망 모형(multilayer neural net model)을 활용하였다. 다층신경망 모형이 전통적 회귀모형보다 주택 가격 예측 성능이 일관되게 탁월함을 4개 사례지역(서울 강남구, 경남 김해시, 전주 덕진구, 전남 해남군) 모두에서 확인하였다. 이러한 성능 향상의 주원인은 다층신경망 모형이 설명변수와 종속변수 간 비선형 관계를 효율적으로 포착할 수 있기 때문인 것으로 해석된다. 설명변수 중 하나인 지리좌표(X, Y)값과 주택 가격 간의 관계를 효율적으로 모델링한 것을 비선형 관계 포착의 예로 들 수 있다. 또 다른 특이점은 사례지역의 이질성이 강해질수록 두 모형 간 성능 격차가 커지는 것을 발견할 수 있었다. 즉, 사례지역에 소재하는 주택의 이질성이 심한 경우, 전통적 모형은 가격 예측의 정확성이 상당히 감소하였으나 다층신경망 모형은 상대적으로 강건한 예측 성능을 보여주었다.

주제어: neural net model(신경망 모형), hidden layer(은닉층), non-linear modeling(비선형 모델링), house price(주택가격), property valuation(부동산 평가)

* Associate Research Fellow, Korea Institute of Local Finance

** Director of Planning & Coordination Office, Korea Appraisal Board

1. Introduction

Property valuation plays a vital role in real estate industries spanning from investment, development, management, tax assessment to other businesses, such as mortgage lending and compulsory acquisition. Numerous studies have attempted to predict the property price as accurate as possible, and the hedonic pricing models have been commonly employed in the studies. The hedonic pricing models are usually specified and fitted by regression analysis, and it became a well-known old criticism that the appropriate functional form for the models cannot be guided on theoretical grounds (Halvorsen & Pollakowski, 1981). Since there exists little support for describing the optimal functional form, that is, the relationship between the predictors and response variables, a linear relationship between them is commonly assumed in specifying the hedonic pricing models. This practice has led to inaccuracy and inefficiency in price estimation.

We try to overcome the limitation mentioned above by using a deep learning model, that is, a multilayer artificial neural net model. We investigate the predictive accuracy of the deep learning model in comparison with the traditional linear model, and analyze the reasons for the difference in predictive performance. To be more specific, we choose a neural net model and apply it to the transaction data of single-family houses in the four study areas. Then its predictive accuracy of the price is compared with that of a traditional linear regression model. The difference in model performance, if any, will be explored and its implications for practice of property valuation will be discussed.

Our application of the neural net model to the house price alleviated the limitation imposed by assuming a linear relationship between variables. The core advantage of the neural net model is its ability to learn non-linear relationships directly from data, and this advantage was proven empirically in the study. By extension, we also examined the difference in performance between the neural net model and the linear regression model, analyzed the reasons, and provided insights to property valuation. The investigation of contextual effects related to regional features of the four study areas is relatively novel and could be found out rarely in past studies.

This paper proceeds from a review of previous studies on property valuation methods, and the third section explains the data used and model specification. Model fit results

and their contextual meaning are discussed in the fourth section, and finally the conclusion summarizes our study results and gives suggestions for future research.

2. Literature review

Machine learning is an area of study on computer sciences, and can be characterized as applying algorithms on data to acquire knowledge. There are ways to implement machine learning, for example random forests, support vector machines, and boosting. Applications of machine-learning tools have been reported in many diverse fields such as pattern recognition, medical diagnoses, finance, marketing, etc. There has also been an emerging trend to apply machine learning to property valuation; Antipov & Pokryshevskaya (2012) applied the random forest to apartments in Saint-Petersburg, Russia. Lasota et al. (2011) reported that the random forest outperformed other tested methods when they tried to predict apartment prices in Poland. Similar studies using the random forest or support vector machines can be found out in the field of property valuations in S. Korea (Kim, 2016; Won et al., 2017).

A neural net model is another approach to machine learning, along with random forest or support vector machines. The neural net model is an approach that models data using artificial neurons that mimics how a neuron in the human brain works. It is a computing device inspired by the neurology of the brain (Bishop, 1995). The neural net model delivered state-of-the art performance for image recognition, speech recognition, property valuation, and other applications (Tsai et al., 1995; Trentin et al., 1998; Yang et al., 2000; Selim, 2009). The neural net models in the aforementioned literature had been trained with the back-propagation algorithm. With back-propagation, the input data is repeatedly fed to the neural network, and the output of the neural network is compared to the desired output and an error is calculated. This error is then fed back (back-propagated) to the neural network and used to adjust the weights such that the error becomes smaller with each iteration and it becomes closer and closer to yielding the desired output.

In property valuation, most of neural net models utilized in early years consisted of

three layers: the input layer, the hidden layer and the output layer. Studies conducted in those times frequently employed the neural net model with a single hidden layer arguing that the one hidden layer is sufficient for modelling property prices (McCluskey et al., 2012). However, studies using a single hidden layer reported contradicting results that the neural net model is superior to the traditional model such as a linear regression model or vice versa (Worzala et al., 1995; Lenk et al., 1997; McGreal et al., 1998; McCluskey et al., 2013; Grover, 2016).

These contradicting results dampened popularity of neural net models, and led to a drastic decrease in interest in neural nets. However, the use of multiple hidden layers and modern computational capabilities have made the second prime day for neural net models since 2010s. In terms of algorithms, employing the multiple hidden layers, thus the name of 'deep learning', have played a key role in overcoming the limitations noticed in early days: the result is unstable and it is only efficient in solving the simple problems, etc. The reason why the multiple hidden layers are efficient is that they can find features inherent in the data and allow subsequent layers to operate on those features, rather than the noisy raw data (Haykin, 2003).

In property valuation, there exist a few studies that employed neural net models with a single hidden layer (Lomsombunchai et al., 2004; Peterson & Flanagan, 2009; Morano, et al., 2015). However, studies with multi-hidden layers are rare, and the exceptional case is the study by Khalafallah (2008). Several network structures were tested by varying the number of hidden layers in order to predict the behavior of the housing market in the study. Specifically, he used the ratio between a house's selling and asking prices as the output variable, and tested a single hidden layer and two hidden layer models. Although the study of Khalafallah (2008) explicitly took into account the number of hidden layers, he concluded that the best result was obtained from the network with a single hidden layer.

Therefore, empirical studies of a multi-hidden layer neural net model are hardly found in literature, and we employ the multi-hidden layer model explicitly in order to predict house prices. In addition to comparing model performances, we take one step further and investigate the relationships between model performances and data characteristics related to the study areas. That is, we discuss the contextual effects of the regions on the model performance, which has not been attempted in any past studies.

3. Data and model specification

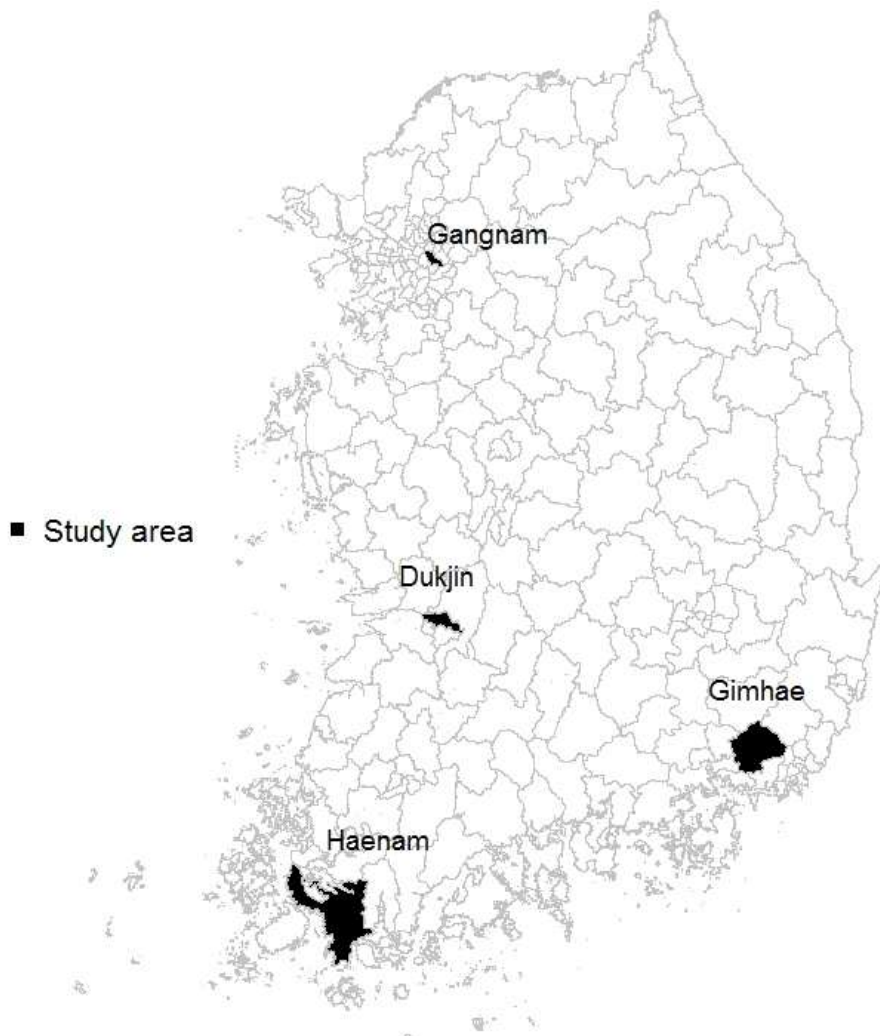
3.1. Study area

Four cities are chosen as the study areas: Gangnam, Gimhae, Dukjin and Haenam. Gangnam is one of the 25 local government districts in Seoul, and is well known for its high standard of living and heavily concentrated wealth. Most importantly, Gangnam is widely known for extremely expensive property prices. Gimhae is a city in the southeastern part of South Korea, whereas Dukjin is a city in the southwestern part of it. Both cities are local centers famous for tourist attractions and innovative festivals. They are similar to each other in that they are both urban and rural due the closeness of forest and farm land area. Haenam is a county at the southern end of the Korean Peninsula, and a typical farming area with rice and radish being the two most representative crops. Table 1 provides city profiles for each study area, and Figure 1 shows their locations. As seen in the table, Gangnam shows the highest population density (14,344 per km²), followed by Dukjin and Gimhae, with Haenam being the last one far behind these three cities. Thus it can be summarized that Gangnam was chosen to represent a large global city, Dukjin and Gimhae was selected as middle sized local cities, and Haenam to represent a typical farming area.

〈Table 1〉 City profiles (2017)

City	Area (km ²)	Population	Pop. Density (per km ²)	Single-family houses (units)
Gangnam	39.5	566,590	14,344	7,470
Gimhae	463.4	529,756	1,143	27,203
Dukjin	110.8	293,142	2,646	20,524
Haenam	1013.3	74,701	74	24,497

Source: Ministry of the Interior and Safety, Yearbook 2017



〈Figure 1〉 Location of study areas

3.2. Data and characteristics

The data were obtained from the Real estate Transaction Management System (RTMS). The RTMS is operated by the government, and discloses the sales prices of houses on a monthly basis. We collected the single-family house data that were traded between 2012 and 2015, and the basic characteristics of the samples are presented in Table 2. The

average prices of houses sold are 2,306 million KRW in Gangnam, 300 million KRW in Gimhae, 185 million KRW in Dukjin, and 66 million KRW in Haenam. Sales prices in Gangnam are, on average, 35 times higher than in Haenam. The floor area (Bldg area in Table 2) tends to decrease as it goes from Gangnam to Haenam, and it can be interpreted that the house in Gangnam is usually constructed in a multi-story type, generally four floors or higher, whereas the house in rural area such as Haenam is likely to be built in a single story structure on a relatively large site.

〈Table 2〉 Descriptive statistics

City		Min.	Mean	Median	Max.
Gangnam (n=466)	Sales price (mil. KRW)	292	2306	2007	13000
	Site area(m ²)	52	246	228	908
	Bldg area(m ²)	66	358	315	2088
	Age (year)	1	24	24	42
Gimhae (n=2,131)	Sales price (mil. KRW)	11	300	235	3515
	Site area(m ²)	12	279	229	9477
	Bldg area(m ²)	18	220	172	1742
	Age (year)	1	24	20	113
Dukjin (n=2,174)	Sales price (mil. KRW)	10	185	115	2100
	Site area(m ²)	12	238	194	14918
	Bldg area(m ²)	22	188	114	1401
	Age (year)	1	30	33	114
Haenam (n=678)	Sales price (mil. KRW)	3	66	35	1550
	Site area(m ²)	42	487	368	16748
	Bldg area(m ²)	16	120	87	11384
	Age (year)	2	35	33	109

We emphasize that most past studies that attempted to predict house prices focused on the apartment, since it is a highly standardized commodity among residential properties, and thus relatively easy to model the price. In contrast, single-family houses vary in their locations and building features, thus more challenging to model the prices. Another reason why we chose the single-family house as the target property is to show non-linear effects of geographical locations on the house price. In contrast to the single-family house, the apartment unit is generally nested in the complex, thus the location and quality of the complex itself might play a more important role than locations of individual apartment units.

3.3. Predictors

In the model specification, the sales price (million KRW) was designated as the response variable, while 13 predictors were employed to predict the sales price. Selecting predictors is always a compromise between theory and data availability. The following information has been obtained from the property transaction data: the road width the house abuts on, the age of the property, the site area, the sum of floor areas, the year the transaction is reported in, land zoning, the shape of the site, the bearing of the property, the construction structure of the building, the roof type of the building, the neighborhood characteristics, and geographical coordinates (longitude and latitude). 6 out of 13 predictors are categorical variables, and their possible levels are provided in Table 3. We believe that we include all the basic and fundamental variables affecting the prices of single-family houses under the constraints of data collectability.

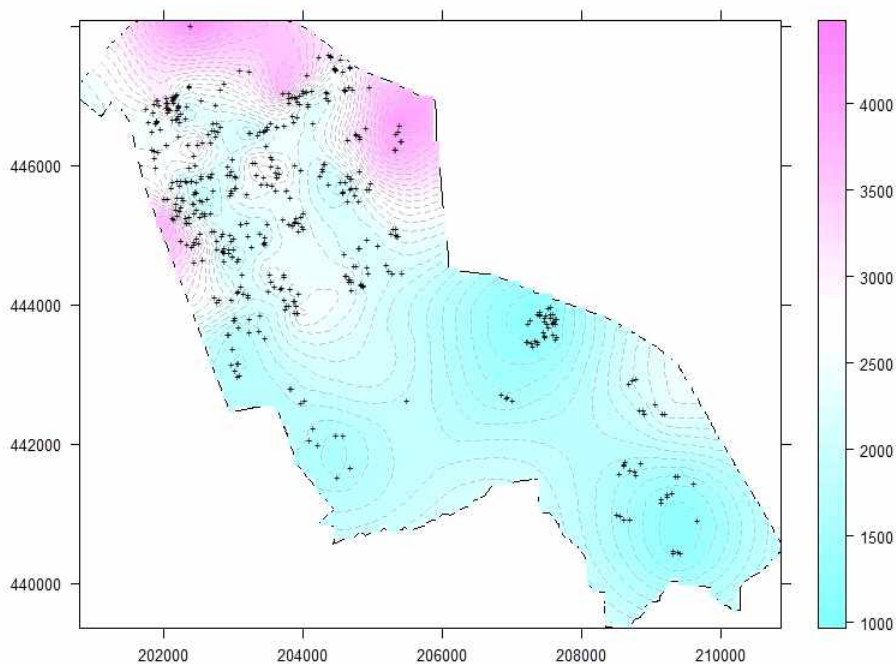
〈Table 3〉 Summary of categorical variables

Variable name	Number of levels	Possible levels
Land zoning	24	Residential, commercial, agricultural, etc.
Neighborhood characteristics	6	Urban, rural, mountainous, etc.
Shape of the site	4	Rectangular, trapezoidal, irregular, etc.
Bearing	4	East, west, north, south
Construction structure	8	Reinforced concrete, brick, wood, etc.
Roof type	3	Slab, shingle, others

The neural net model is generally considered as a toolkit to analyze a bid data, and it is true that the performance of the model improves when more training data are available. But the neural net model is able to predict complex targets within a small dataset framework as well, which is the case of this study.

It is worth noting that we utilized geographical coordinates as one of the predictors while specifying the model. The importance of geographical locations cannot be overemphasized, and it is well known that perfectly identical houses can be traded for vastly different prices depending on location. This spatial effect on the house price has been ignored or minimally considered in the traditional linear regression model. In the studies, if any, where the spatial component was explicitly taken into account, the

typical approach was to incorporate the spatial autocorrelation inherent in the house price data into the error term structure (Kim et al., 2003; Anselin & Lozano-Gracia, 2008; Cohen & Coughlin, 2008; Sander et al., 2010; Lazrak et al., 2014). In contrast, we directly employed geographical coordinates as the predictor in model specification, and this treatment of the coordinates can be justified by the neural network's capability of modeling extremely complex non-linear relationships between predictors and house prices. The pattern of change of house prices with respect to the X-Y geographical coordinates would be incorporated into the neural net model. For example, the downward gradient of house prices along the Y coordinates (North - South direction) would be learned while training the neural net model. However, as shown in Figure 2, assuming a linear relationship between the geographical coordinates and house prices does not seem to be relevant, and thus the neural net model could be a useful tool to deal with these complex non-linearity.



* "+" indicates locations of samples.

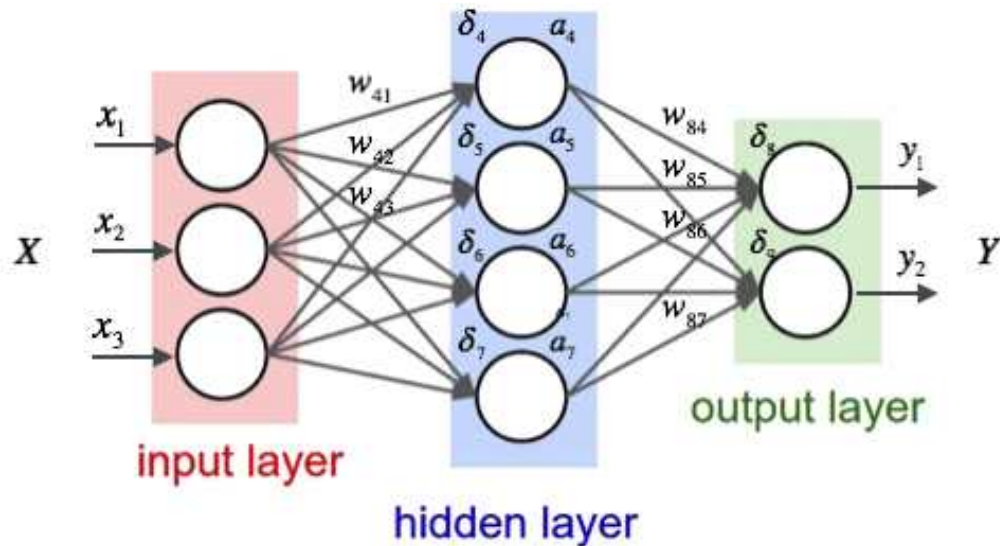
〈Figure 2〉 Geographical distribution of house prices in Gangnam (units: million KRW)

3.4. Specification of the linear regression model

The linear regression model used in the study is not wholly classical, but instead, involves regularization, which is a useful way to avoid overfitting. The regularization can be implemented in various manners, but the typical ones are L1 (lasso regression) and L2 (ridge regression). L1 regularization will set some of regression coefficients to zero, whereas L2 regularization keeps all the coefficients close to zero, but nonzero. When the two regularizations are mixed together, the regularization technique is called an elastic net. Roughly speaking, if data have a lot of predictors, but which ones are important is unknown, L1 will work better. In contrast, if data is dense, meaning all predictors are likely to explain something about the response variable, L2 will be a useful option. The elastic net could prove its worth when data features are situated between the two conditions given above (Cook 2016, p.167). In the study, a lasso regression was adapted in Gangnam and Haenam, and a ridge regression was used in Dukjin. As for the area of Gimhae, an elastic net was chosen. These model choices were made based on the predictive performance on test data of each study area.

3.5. Specification of the neural net model

The deep learning approach employed in the study is the use of a feed-forward multilayer artificial neural net model, and has been trained with the back-propagation algorithm. When applying the algorithm, the input data is provided to the neural net repeatedly, the data is transformed in the hidden layers, and the output of the neural net is compared to the target and an error is computed. This error is fed back to the neural net and used to adjust the weights so that the error decreases with each iteration. The neural net model finally achieves the desired level of predictive accuracy when enough iterations have been done, and this process is represented at Figure 3.



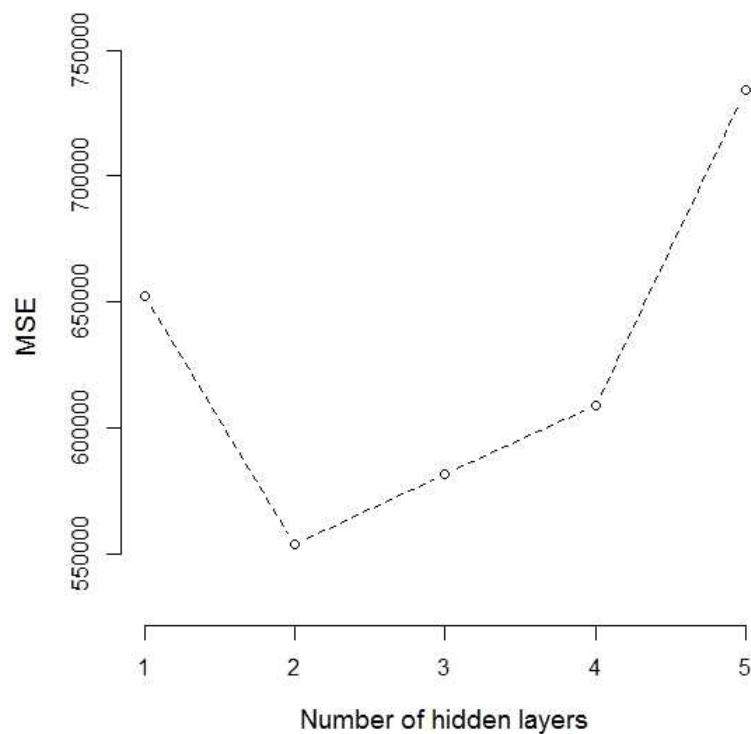
(source: <https://medium.com/@curiously/tensorflow-for-hackers-part-iv-neural-network-from-scratch-1a4f504dfa8>)

〈Figure 3〉 A graphical representation of a neural net model with a single hidden layer

The first and critical component to decide, when using a neural net model, is the shape of the network, that is, the number of layers. When you have enough training data, you can use multiple hidden layers, but just one hidden layer often works best with a small dataset. Recent research in the neural net model shows that many hidden layers can be effective for complex tasks, such as image recognition and learning for text and sequences (Oyedotun & Khashman, 2017). Using several hidden layers enables the neural net model to create several layers of data abstraction necessary to effectively model complex underlying relationships in the input data (Semwal et al., 2017). This layer depth is what enables the model to truly utilize deep learning, one of the most powerful learning tools in the literature (Shwartz-Ziv & Tishby, 2017). However, with each additional hidden layer, the computational complexity of training larger models requires much larger amounts of input data. Thus, proper hidden layers has been empirically chosen as compromise between high-level feature abstraction and limited input data (O'Connell et al., 2018).

While theoretically one hidden layer is enough to represent all the data, whether

involving non-linear relationships or not, it is not such a case in practice. Thus we start with a neural net model with a single hidden layer, and explore the model performance by varying the number of hidden layers. Figure 4 shows average cross-validation mean squared errors (MSE) from Gangnam of which models were fitted by using the 13 predictors mentioned in the previous section. As seen in the figure, three or more hidden layers did not contribute much to improvement in prediction performance, and similar patterns were found out in the other three study areas. Thus we fixed the number of hidden layers at two. 324 neurons were used for each hidden layer in the figure, and the reason for it will be explained in the following paragraph.



〈Figure 4〉 Average cross-validation MSE from Gangnam (324 neurons were used for each hidden layer)

The number of neurons in each hidden layer was determined by optimizing the network structure that best fits the data using the grid search. Table 4 shows the number of neurons in the hidden layer and corresponding MSEs from Gangnam. Employing 200

neurons in the hidden layer is a default setting in the statistical software used in the study.¹⁾ The results of the grid search show that the neural net model tends to benefit from using more than 200 neurons, and we chose 324 neurons as the final number of neurons for each hidden layer after reviewing overall results of the grid search including the other three study areas. The exact number of neurons (50, 200, 324, and 648) in Table 4 does not have any particular meaning, and it only indicates that using more than 200 neurons will improve the model performance. The results from using 324 neurons for each hidden layer were presented in Figure 3 in the same context.

〈Table 4〉 Number of neurons in the hidden layer and corresponding MSEs from Gangnam

Number of neurons	(Second hidden layer) 50	200	324	648
(First hidden layer) 50	770,636	720,890	753,471	717,596
200	755,383	695,990	701,149	687,799
324	729,945	642,996	646,929	770,814
648	789,021	751,305	735,402	773,726

There exist many other tuning parameters in the model, but most of them seem to affect the model performance trivially. With the learning rate, higher value means that the fitting process will be less stable, while lower one implies that it will take longer to converge. The default in the software used is 0.005. With the loss function, the default choice for Gaussian distribution which is the case in our study is quadratic loss function. And those two parameters seem to make little effect to model performance, thus we left them as the default values. As for stopping criteria, MSE, mentioned in the preceding section, is used as stopping metric, and the fitting process is set to stop if relative improvement is not at least 0.0005 (0.05%) over last 5 rounds. With epochs, the default is 10, and we increased it considerably by adjusting the stopping criteria as above.²⁾ A more detailed process of the application of the neural net model can be found in Cook (2016).

Table 5 shows the details of the neural net model developed in the study. With the activation function in the table, three options are available: hyperbolic tangent, max-out, and rectifier. The hyperbolic tangent function showed better performance across all the datasets, but the improvement was trivial.

〈Table 5〉 Details of the neural net model

Parameters	Details
Network architecture	4-layer (input-hidden-hidden-output)
Algorithm	Backpropagation
Activation function	Hyperbolic tangent
Training and testing ratio of data	80:20 (random sampling)
Validation	5-fold cross-validation

In summary, we constructed the neural net model specified in Table 6. The number of neurons in the first input layer increased from the original 13 predictors since each categorical predictor was converted into a binary representation while fitting each model. For example, we have 7 possible levels in the categorical predictor zone in Gangnam, whereas 15 levels in Gimhae, and thus the number of input neurons varies across the study areas ranging from 36 to 52 neurons. The second and third layers (hidden layers) have 324 neurons for each. And the last output layer has 1 neuron that is the house price. To ensure that the neural network do not overfit the training dataset, a dropout ratio of 50% was used on the hidden layers while training, meaning that we employed a neural net model with partially-connected layers.

〈Table 6〉 Network architecture

Layer	Neurons	Dropout ratio
Input	36(Gangnam), 52(Gimhae), 47(Dukjin), 46(Haenam)	0%
Hidden #1	324	50%
Hidden #2	324	50%
Output	1	-

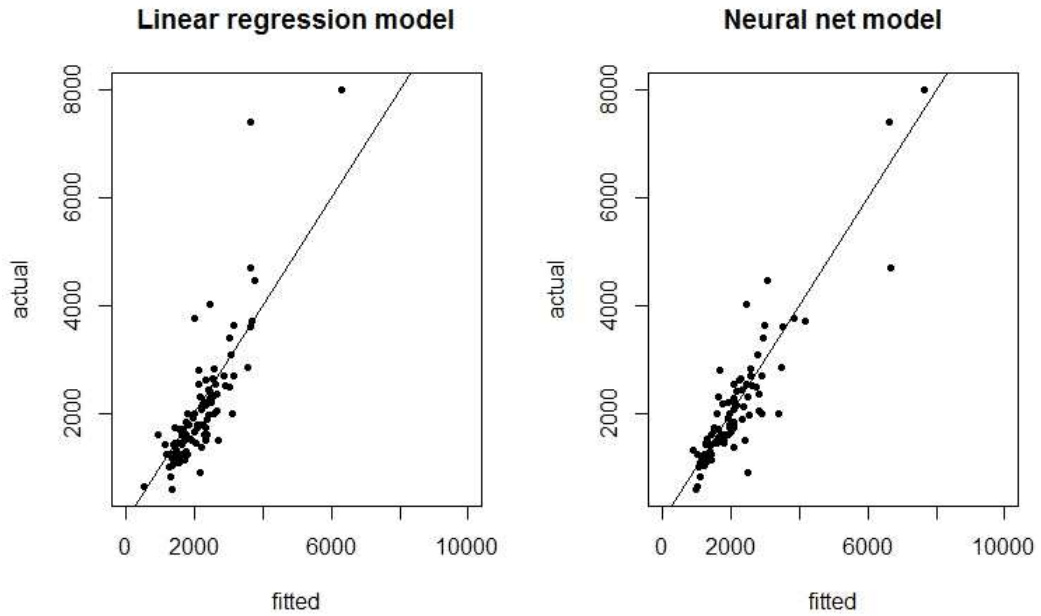
4. Model fit results and discussion

4.1. Model fit results

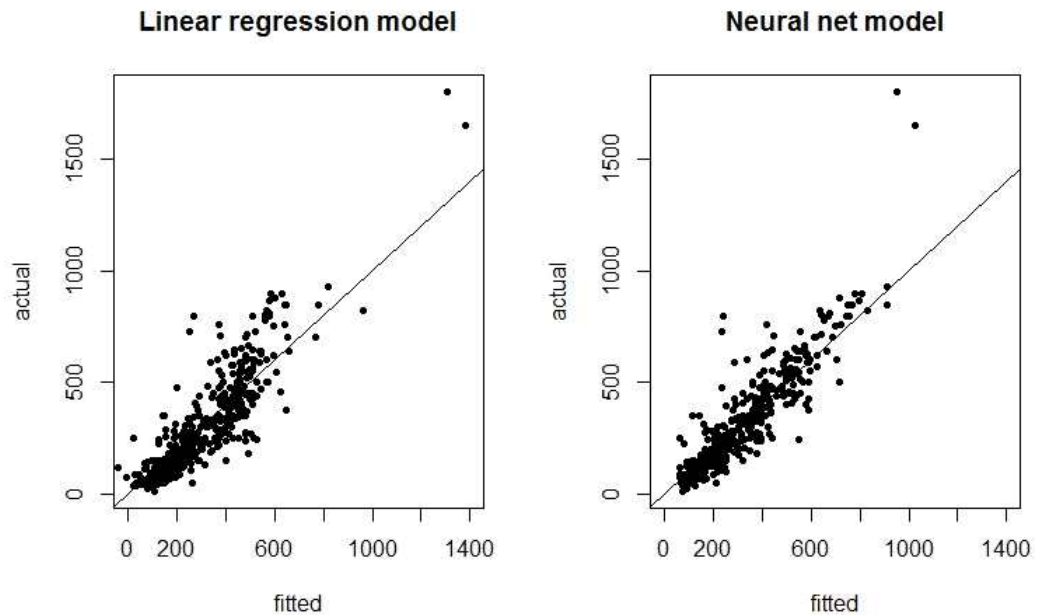
Figure 5 shows the goodness of the fit of models. In order to test the quality of predictions, we divided the data to 80% training sample and 20% test sample in a random sampling manner, and the graphs in the figure were created based on the 20%

holdout sample for each area. Overall, Figure 5 indicates that the goodness of the fit improves as it goes from a linear regression model to a neural net model.

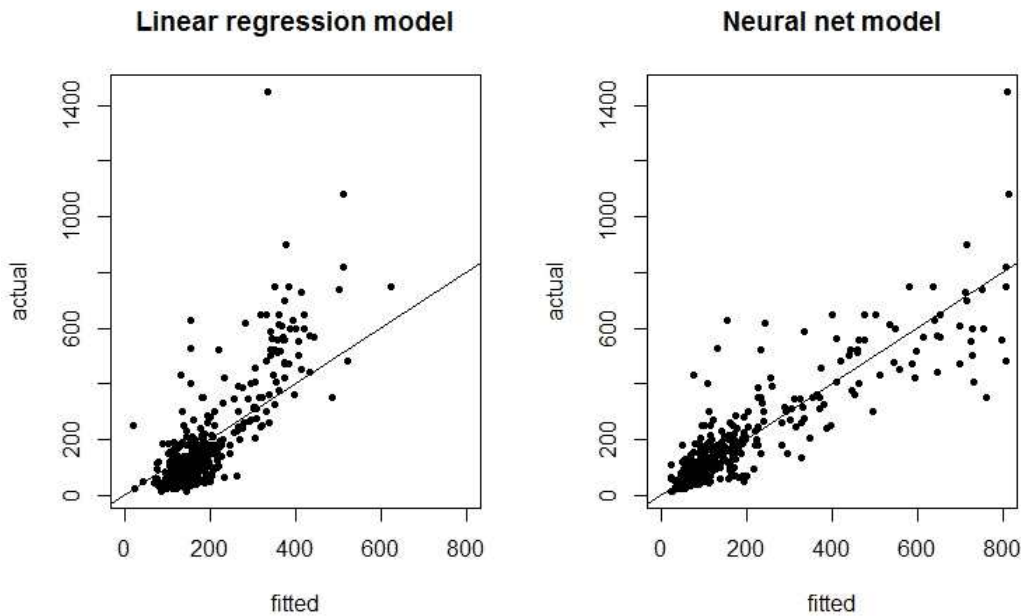
① Gangnam



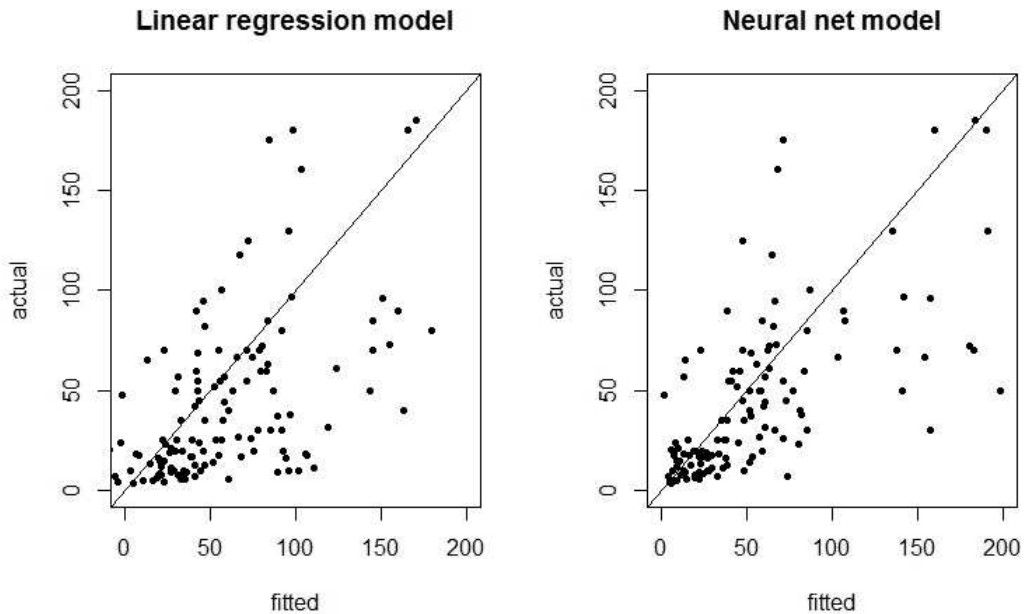
② Gimhae



③ Dukjin



④ Haenam



〈Figure 5〉 Goodness of the fit of models (test data)

The graphs above only show the model performance visually, and it would be more reasonable to compare the model performance by using an established performance measure. There exist a number of model performance metrics in the literature. Coefficient of dispersion (COD) is adopted in the study, and this is to ascertain that our model can satisfy established standards in the appraisal domain. COD is a measure of uniformity of tax appraisal, and frequently used in tax assessment performance (IAAO, 2013). One advantage of it is that it can be used to compare assessment quality across different regions. COD is calculated as following:

First, calculate the ratio of the predicted values over observed ones.

Second, subtract the median from each ratio.

Third, take the absolute value of the calculated differences.

Fourth, sum the absolute differences and divide by the number of ratios to obtain the average absolute deviation.

Finally, divide by the median and multiply by 100.

The formula for COD can be represented as following:

$$COD = \frac{\left(\frac{\sum | \text{each ratio} - \text{median} |}{\text{number of ratios}} \right)}{\text{median}} \times 100.$$

COD values closer to zero are interpreted as to have the better performance. Table 7 shows COD values for each model and study area, and it is confirmed that the neural net model outperforms the linear regression model in all study areas. The far right column in the table presents the Gini coefficient of house prices in each area, which will be explained further in the later section.

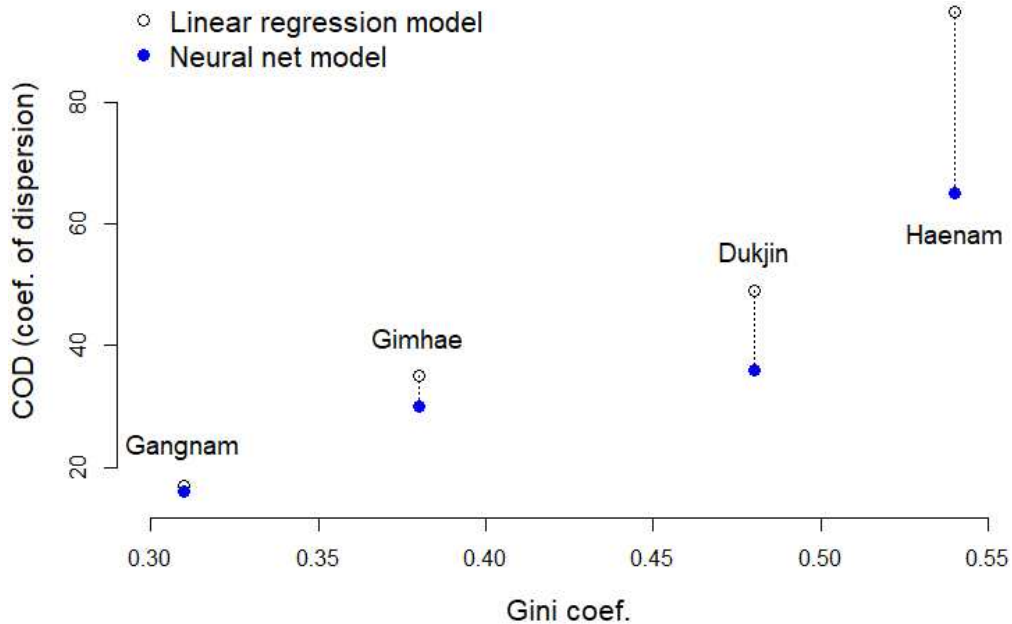
〈Table 7〉 Model performance (test data)

City	COD from linear regression model	COD from neural net model	Gini coef.
Gangnam	17	16	0.31
Gimhae	35	30	0.38
Dukjin	49	36	0.48
Haenam	95	65	0.54

4.2. Discussion

The true power of the neural net model lies in its ability to capture non-linear relationships and learn these relationships directly from the data. For many years, linear modeling has been the norm in capturing functional relationships between predictors and response variables, largely because of its easy-to-explain optimization processes. In the problem where the linear approximation of a relationship was not valid, which happens frequently in empirical research, the models suffered severely. As confirmed in Figure 2, geographical coordinates hardly have a linear relationship with the house prices. Longitude and latitude were directly employed as the predictor to model the non-linear relationship in our study, and we interpret that the better performance of the neural net model could be largely attributed to efficient non-linear modeling of predictors such as geographical coordinates, as well as the property age, the site area, the sum of floor areas, etc.

The Gini coefficient is a measure of dispersion that represents the income distribution of a nation's residents, and it is the most commonly used index of inequality. It can serve to indicate how much property prices vary across the area, and higher value means that property prices show a considerable difference among houses. Lower value implies that most houses within a particular area command same or similar prices. From the valuation perspective, as the Gini coefficient increases, it means that property prices become more and more heterogeneous, and thus difficult to model. Table 7 shows COD and Gini coefficient values for each study area, and their relationships are demonstrated graphically in Figure 6.



〈Figure 6〉 COD and Gini coefficient in each study area

The noticeable pattern in the figure is that the gap of the model performance between the linear regression and the neural net model widens as the Gini coefficient value increases from 0.31 to 0.54. The city with a lower Gini coefficient can be interpreted as the area where homogeneous single-family houses are located, with the example being Gangnam. In Gangnam, most single-family houses are four or five story buildings with reinforced concrete structure. The maximum sales price is only 45 times higher than the minimum one in Gangnam, and this is drastically contrasted to those in the other study areas; 320 times in Gimhae, 210 times in Dukjin, and 517 times in Haenam, which can be confirmed in Table 2. Thus it can be said that the neural net model outperforms the linear regression model in all areas, and the difference would be trivial in the area with relatively similar houses situated. In contrast, the neural net model outperforms the linear regression model significantly in the area with a higher Gini coefficient, that is, the area with heavily heterogeneous houses traded. For example, Haenam is the county where various types of houses exist from a plank house with coarse thatch roof to a luxurious one with extensive balconies and spacious grass. This

heterogeneity of house prices can also be confirmed by the range of sales price in Haenam; the maximum sales price is 517 times higher than the minimum one. In summary, Figure 6 suggests that the neural net model exceeds the linear regression model in the model performance across all study areas, and the degree of excess becomes noticeably large as the area goes from a relatively homogeneous to more complex and heterogeneous one. Therefore, we conclude that the neural net model could perform more efficiently when predicting house prices in the heterogeneous area where a traditional regression model has suffered severely.

5. Conclusion

We pointed out the limitation of linear modeling and tried to overcome it by employing a deep learning approach, that is, a multilayer neural net model. We specified and fitted the neural net model to predict single-family house prices, and its predictive accuracy was compared with that of a linear regression model. We have presented results from four study areas where the neural net model outperformed the linear regression model across all study areas. We interpreted that the better performance of the neural net model might stem from non-linear modeling of predictors with the example being geographical coordinates. The noticeable finding was that the difference of the model performance became large as the study area went more heterogeneous, that is, property characteristics of houses became more various and complicated. We concluded that the neural net model could perform more efficiently in the case that a traditional linear model suffered. It is expected that the neural net model, if specified with due care for its several tuning parameters, could serve as an alternative tool when predicting the house prices of the heterogeneous area which can be characterized by the mix of older and newer properties, small and large sized houses, active and depressed real estate submarkets, etc.

When specifying a model, we employed predictors that were available in the dataset, and thus our final model was able to include variables related to physical attributes such as the property age or the site area. We expect that a future research could expand

predictors to socioeconomic variables such as the income of local residents, the crime frequency, etc. Our neural net model took as the input the vectorized data of 13 columns (predictors), which is relatively simple input shape. The true power of the neural net model lies in the capability of handling the input data of more complex shape, such as image, video and texts. A future research needs to develop a more sophisticated model, for example, utilizing aerial photographs as a predictor to estimate building prices. Finally we hope to see that the neural net model would be used extensively to model more challenging properties such as office buildings, shopping centers and industrial plants, and join the ranks of the established tools in the practice of property valuation.

Note

- 1) Associate Research Fellow, Korea Institute of Local Finance.
- 2) Director of Planning & Coordination Office, Korea Appraisal Board.
- 3) The H2O open-source software was used, and the software was called from the statistical package R environment.
- 4) Epochs are 252 in Gangnam, 156 in Gimhae, 175 in Dukjin, and 168 in Haenam.

References

- Anselin, L., & Lozano-Gracia, N. (2008). Errors in variables and spatial effects in hedonic house price models of ambient air quality. *Empirical economics*, 34(1), 5-34.
- 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.
- Bishop, C. M. (1995). *Neural networks for pattern recognition*. Oxford university press.
- Cohen, J. P., & Coughlin, C. C. (2008). Spatial hedonic models of airport noise, proximity, and housing prices. *Journal of Regional Science*, 48(5), 859-878.
- Cook, D. (2016). Practical machine learning with H2O: powerful, scalable techniques for deep learning and AI. O'Reilly Media.
- Grover, R. (2016). Mass valuations. *Journal of Property Investment & Finance*, 34(2), 191-204.
- Halvorsen, R., & Pollakowski, H. O. (1981). Choice of functional form for hedonic price equations. *Journal of urban economics*, 10(1), 37-49.
- Haykin, S. (2003). *Neural Networks: A Comprehensive Foundation*, second edition, New Jersey:

- Prentice Hall.
- International Association of Assessing Officers (IAAO). (2013). Standard on Ratio Studies.
- Khalafallah, A. (2008). Neural network based model for predicting housing market performance. *Tsinghua Science & Technology*, 13, 325-328.
- Kim, G. M. (2016). Analysis for Factors Determining the Price of Multi-family Housing through Machine Learnings. *Journal of the residential environment institute of Korea*, 14(3), 29-40.
- Kim, C. W., Phipps, T. T., & Anselin, L. (2003). Measuring the benefits of air quality improvement: a spatial hedonic approach. *Journal of environmental economics and management*, 45(1), 24-39.
- Lasota, T., Łuczak, T., & Trawiński, B. (2011). Investigation of random subspace and random forest methods applied to property valuation data. In *International Conference on Computational Collective Intelligence* (pp. 142-151). Springer, Berlin, Heidelberg.
- Lazrak, F., Nijkamp, P., Rietveld, P., & Rouwendal, J. (2014). The market value of cultural heritage in urban areas: an application of spatial hedonic pricing. *Journal of Geographical Systems*, 16(1), 89-114.
- Lenk, M. M., Worzala, E. M., & Silva, A. (1997). High-tech valuation: should artificial neural networks bypass the human valuer?. *Journal of Property Valuation and Investment*, 15(1), 8-26.
- Lomsombunchai, V., Gan, C., & Lee, M. (2004). House price prediction: Hedonic price models vs. artificial neural nets. *American Journal of Applied Statistics*, 1(3), 193-201.
- McCluskey, W., Davis, P., Haran, M., McCord, M., & McIlhatton, D. (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.
- McCluskey, W. J., McCord, M., Davis, P. T., Haran, M., & McIlhatton, D. (2013). Prediction accuracy in mass appraisal: a comparison of modern approaches. *Journal of Property Research*, 30(4), 239-265.
- McGreal, S., Adair, A., McBurney, D., & Patterson, D. (1998). Neural networks: the prediction of residential values. *Journal of Property Valuation and Investment*, 16(1), 57-70.
- Ministry of the Interior and Safety. (2017). Yearbook 2017.
- Morano, P. I. E. R. L. U. I. G. I., Tajani, F. R. A. N. C. E. S. C. O., & Torre, C. M. (2015). Artificial intelligence in property valuations an application of artificial neural networks to housing appraisal. *Advances in Environmental Science and Energy, Planning*, 23-29.
- O'Connell, J., Li, Z., Hanson, J., Heffernan, R., Lyons, J., Paliwal, K., ... & Zhou, Y. (2018). SPIN2: Predicting sequence profiles from protein structures using deep neural networks. *Proteins: Structure, Function, and Bioinformatics*, 86(6), 629-633.
- Oparaji, U., Sheu, R. J., Bankhead, M., Austin, J., & Patelli, E. (2017). Robust artificial neural network for reliability and sensitivity analyses of complex non-linear systems. *Neural Networks*, 96, 80-90.
- Oyedotun, O. K., & Khashman, A. (2017). Deep learning in vision-based static hand gesture recognition. *Neural Computing and Applications*, 28(12), 3941-3951.
- Peterson, S., & Flanagan, A. (2009). Neural network hedonic pricing models in mass real estate appraisal. *Journal of Real Estate Research*, 31(2), 147-164.
- Sander, H., Polasky, S., & Haight, R. G. (2010). The value of urban tree cover: A hedonic property price model in Ramsey and Dakota Counties, Minnesota, USA. *Ecological Economics*, 69(8), 1646-1656.
- Selim, H. (2009). Determinants of house prices in Turkey: Hedonic regression versus artificial neural network. *Expert Systems with Applications*, 36(2), 2843-2852.

- Semwal, V. B., Mondal, K., & Nandi, G. C. (2017). Robust and accurate feature selection for humanoid push recovery and classification: deep learning approach. *Neural Computing and Applications*, 28(3), 565-574.
- Shwartz-Ziv, R., & Tishby, N. (2017). Opening the black box of deep neural networks via information. *arXiv preprint arXiv:1703.00810*.
- Trentin, E., Bengio, Y., Furlanello, C., & De Mori, R. (1998). Neural networks for speech recognition. *Spoken Dialogues with Computers*, 311-361.
- Tsai, I. S., Lin, C. H., & Lin, J. J. (1995). Applying an artificial neural network to pattern recognition in fabric defects. *Textile Research Journal*, 65(3), 123-130.
- Won, S. W., Lee, C. G., Park, J. M. (2017). A Study on the Prediction of Land Price with Machine Learning Technique. *Journal of the Korean Association of Professional Geographers*, 51(4), 347-355.
- Worzala, E., Lenk, M., & Silva, A. (1995). An exploration of neural networks and its application to real estate valuation. *Journal of Real Estate Research*, 10(2), 185-201.
- Yang, C. C., Prasher, S. O., Landry, J. A., & DiTommaso, A. (2000). Application of artificial neural networks in image recognition and classification of crop and weeds. *Canadian agricultural engineering*, 42(3), 147-152.

(논문접수일: 2018. 08. 06 / 논문수정일: 2018. 09. 28 / 게재확정일: 2018. 09. 12)

※ **이창로(李昌魯)**는 서울대학교에서 지리학 박사학위를 취득하고, 현재 한국지방세연구원에서 재직 중이다. 부동산 가치의 추정, 기계학습 모형의 응용, 과세표준 산정에서의 딥러닝 기법의 적용 등이 주요 관심분야이다. 주요 논문으로 Analyzing the rent-to-price ratio for the housing market at the micro-spatial scale (International journal of strategic property management, 2018) 등이 있다(spatialstat@naver.com).

※ **김세형(金世衡)**은 상명대학교에서 박사과정을 수료하고, 현재 한국감정원에서 재직 중이다. 감정평가 기법 연구, 부동산 가치의 추정, 과세표준 산정에 있어 대량평가모형 연구 등이 주요 관심분야이다. 주요 논문으로는 Exploring Spatial Heterogeneity of Real Estate Price Formation Factors(한국 감정평가학회, 감정평가학 논집 13권 1호, 2014) 등이 있다(k23584@hanmail.net).