



# Article House Price Valuation Model Based on Geographically Neural Network-Weighted Regression: The Case Study of Shenzhen, China

Zimo Wang<sup>1</sup>, Yicheng Wang<sup>1</sup>, Sensen Wu<sup>1,2,\*</sup> and Zhenhong Du<sup>1,2</sup>

- <sup>1</sup> School of Earth Sciences, Zhejiang University, Hangzhou 310058, China
- <sup>2</sup> Zhejiang Key Laboratory of Resources and Environmental Information System, Hangzhou 310058, China
- \* Correspondence: wusensengis@zju.edu.cn; Tel.: +86-0571-88273287

Abstract: Confronted with the spatial heterogeneity of the real estate market, some traditional research has utilized geographically weighted regression (GWR) to estimate house prices. However, its predictive power still has some room to improve, and its kernel function is limited in some simple forms. Therefore, we propose a novel house price valuation model, which is combined with geographical neural network-weighted regression (GNNWR) to improve the accuracy of real estate appraisal with the help of neural networks. Based on the Shenzhen house price dataset, this work conspicuously captures the variable spatial regression relationships at different regions of different variables, which GWR has difficulty realizing. Moreover, we focus on the performance of GNNWR, verify its robustness and superiority, and refine the experiment process with 10-fold cross-validation. In contrast with the ordinary least squares (OLS) model, our model achieves an improvement of about 50% on most of the metrics. Compared with the best GWR model, our thorough experiments reveal that our model improves the mean absolute error (MAE) by 13.5% and attains a decrease of the mean absolute percentage error (MAPE) by 13.0% in the evaluation on the validation dataset. It is a practical and powerful way to assess house prices, and we believe our model could be applied to other valuation problems concerning geographical data to promote the prediction accuracy of socioeconomic phenomena.

Keywords: GNNWR; GWR; house price valuation; spatial heterogeneity

# 1. Introduction

Housing prices are closely related to the lives of new urban residents, and they also comprise a vital economic index to which the government needs to pay close attention. Exploring the spatial distribution pattern of housing prices has great practical significance and guiding value for government regulation, individual house purchase, or third-party valuation.

As a country with one of the fastest urbanization processes, China has seen steadily rising housing prices in the past few decades, especially in its major cities. Affected by the COVID-19 pandemic in 2020, the world's major economies entered a liquidity easing cycle, and housing prices in many cities in China rose significantly [1]. On this basis, several Chinese cities, such as Shenzhen, Xi'an, and Chengdu, have implemented second-hand housing transaction reference pricing, which is used to curb house price increases. The reference price provides us with a reasonable valuation for slight housing price bubbles.

In this research, we propose a novel house price valuation model based on the data of Shenzhen, China. In the past, different models have been developed by many scholars to model and estimate house prices. In 1972, Rosen [2] proposed the hedonic model, which aims to measure property prices using a number of environmental factors. Early studies mainly consisted of three components: location traits, structural traits, and neighborhood traits, i.e., housing prices are mainly a function of these three characteristics and are approximately linearly related in an exponentialmanner [3]. A number of subsequent



Citation: Wang, Z.; Wang, Y.; Wu, S.; Du, Z. House Price Valuation Model Based on Geographically Neural Network-Weighted Regression: The Case Study of Shenzhen, China. *ISPRS Int. J. Geo-Inf.* 2022, *11*, 450. https://doi.org/ 10.3390/ijgi11080450

Academic Editors: Maria Antonia Brovelli and Wolfgang Kainz

Received: 1 July 2022 Accepted: 30 July 2022 Published: 18 August 2022

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2022 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (https:// creativecommons.org/licenses/by/ 4.0/). studies have demonstrated the relative validity of this model, and measures of these factors are able to estimate the positive or negative correlation between each independent variable and house price more accurately. For example, MoK et al. [4] modeled the house prices in Hong Kong in 1994, showing that house prices were significantly negatively correlated with the age of the house and the distance from the central business district and significantly positively correlated with the floor. As time progresses, more and more independent variables were taken into account and more statistical indicators were added to test the validity of the model. Further studies have partially incorporated land use planning, as well as accessibility in terms of transportation [5]. In recent years, related residential house price studies have incorporated a variety of external environmental factors such as the natural landscape and neighborhood size to analyze their impact on house prices [6]. However, these models have constantly encountered problems in dealing with spatial heterogeneity, i.e., the same independent variable has different effects on house prices in different regions. Ordinary hedonic models can only model a certain independent variable with constant coefficients, but the real situation is often influenced by spatial factors. For example, in suburban areas, transportation conditions dominate house prices, and the quality of nearby schools does not matter. By contrast, in downtown areas, the quality of schools near homes might be more critical, and nearby transportation conditions are relatively less important. This is something that cannot be analyzed by ordinary hedonic models.

Furthermore, taking into account the spatial heterogeneity of the different influencing factors, geographically weighted regression (GWR) methods are proposed, which allow the coefficients to change at different locations [7,8]. The method can be understood as a local weighted linear regression for each local area, and the coefficients fully take into account the effects of adjacent data points according to the first law of geography proposed by Tobler [9]. In order to build a more satisfying model for the geographically variable regression relationships, Brunsdon and Fotheringham [10] mentioned several key questions that GWR has faced: the selection of the variables, the bandwidth, and the spatial autocorrelation of the error. Many scholars have made attempts to resolve these questions on this basis. For example, Tu et al. [11] used the GWR model to discover the relationship between the spatial variations and the urban public ridership in Shenzhen. In 2011, Geng et al. [12] used the GWR model to model house prices in Shenzhen. Compared with the ordinary least squares (OLS) model, the  $R^2$  improved from 0.56 to 0.79. Zhang et al. [13] used mixed geographically weighted regression to model the rent in Nanjing, i.e., some variables were locally weighted according to the geographic location while some variables were globally weighted, and good results were achieved. Lu et al. [14,15] added non-Euclidean distance to GWR, and for some geographic elements that do not obey the standard linear measure, this model achieved better results on the spatial proximity measurement of London and could have better estimation performance for house prices.

However, the ability of GWR to express spatial relationships is limited. Therefore, many scholars have resorted to artificial intelligence methods, which have developed rapidly in recent years, to model house prices using their superb ability to fit to house price [16,17]. Although the estimation performance of neural network models is usually superior to that of OLS models, the spatial distributions obtained by these models are not entirely reasonable, and the constructed regression relationships are difficult to interpret spatially, because they ignore the spatial properties of housing price regression relationships. Another flaw is that some studies on modeling house prices with the help of neural networks have not introduced a 10-fold validation mechanism. Moreover, as some scholars have suggested, the "black box" approach of neural networks has significantly limited their practical significance in predicting house prices and other socioeconomic problems [18]. Both polynomial regression models and traditional neural network methods depart from this linear structure and have relatively complex expressions, making analysis and prediction much more difficult.

In recent years, based on the idea of geographic weighting of GWR, Du et al. [19] proposed a geographically neural network-weighted regression (GNNWR) model that

combines the OLS and neural network models. Owing to their powerful learning abilities, neural networks can effectively address the potential spatial nonstationarity and complex nonlinear features in regression relations. The concept of GNNWR is similar to that of the transformer model [20], a groundbreaking and popular model proposed in 2017. The principle of the transformer model can be summarized as "attention is all you need". The GNNWR model was proposed in the same period as a simpler approach to evaluate the similarity (attention) between the estimated point and training dataset. Considering the widespread application of attention models in recent years, GNNWR is used in this study to model the ecological environment of nearshore seas [21] and estimate the spatial PM 2.5 concentrations in China [22]. The model performance and explanatory power are noted to be satisfactory.

Unlike ecological phenomena, socioeconomic phenomena such as housing prices are not always continuous in physical space. For example, geographical coordinates are discrete and represent discontinuous location qualities. We speculate that the limitations associated with the complexity of socioeconomic problems can potentially be overcome using valuation strategies based on the GNNWR model. Housing price valuation is a classic problem that involves geographical data. In this study, we consider the data of Shenzhen, which is a representative city in China. On 8 February 2021, the Shenzhen Real Estate and Urban Construction Development Research Center released reference prices for second-hand housing transactions for the city's 3595 residential quarters. This dataset can be used to establish a residential price valuation model considering various factors, such as property endogenous variables, subway, and school district conditions. Consequently, we use the GNNWR framework to build a residential price valuation model to address the spatial heterogeneity and nonstationarity present in the data [19].

In summary, the objective of this study is to introduce the GNNWR model in the socioeconomic field to establish a residential price valuation model based on the reference price data of second-hand housing transactions in Shenzhen, accurately fit the spatial heterogeneity and nonlinear relationships of multiple environmental factors, and obtain a more accurate housing pricing model than the GWR method by considering the spatial distributions of multiple factors and their influence on housing prices. In particular, by improving this pricing model, more patterns can be mined to clarify the importance of each factor and variations in the factor weights.

From the viewpoint of methodological hybridization, the proposed approach provides several novel opportunities for geographical comprehension. As a platform for residential valuation, each hybridization provides a range of new possibilities for driving a paradigm shifts toward sustainability, as indicated by Benessia et al. [23]. The proposed model can provide reference for residential valuation, land auctions, and reference prices of second-hand housing transactions in other cities. Through its integration in a hybrid decision-support system, the proposed method can support the real estate market and sustainable land use development as a diagnostic remedy, especially in the COVID-19 pandemic scenario [24].

#### 2. Study Area, Data Sources and Research Methods

#### 2.1. Shenzhen House Price Profile

The Shenzhen Special Economic Zone was established in 1980. This zone is adjacent to Hong Kong in the south and lies to the west of the Pearl River Estuary in Guangdong Province, China. Owing to its geographical proximity to Hong Kong and policy support, Shenzhen has emerged as the region with the third largest GDP, with nine districts under its jurisdiction. According to the Seventh National Census data, the population of Shenzhen has reached 17.56 million. Even under the impact of COVID-19, Shenzhen's regional GDP was RMB 2767.024 billion in 2020, 3.1% higher than that in 2019.

With the increasing population, the housing prices in Shenzhen are rising owing to its excellent economic conditions and business environment. To suppress the prohibitive increase in housing prices, in February 2021, the Shenzhen Real Estate and Urban Construc-

tion Development Research Center established the reference prices for second-hand housing transactions in 3595 residential quarters, based on the government-recorded transaction prices of second-hand housing and surrounding first-hand housing prices.

In this study, data from Shenzhen are used considering the following factors. First, the reference prices of second-hand housing transactions in this region have been extensively evaluated compared to data in the other regions. The differences across house types and floors are averaged, government-recorded transaction prices and surrounding first-hand housing prices are combined, short-term heat and bubbles are removed, and accurate valuation results for a property can be obtained. Second, Shenzhen's urban development is mostly natural. No important political center, relics, or slums exist that can influence the urban planning. Third, the reference prices are introduced in a uniform batch, with a large amount of data and notable influence. Therefore, the results of modeling the reference prices of second-hand house transactions in Shenzhen can provide reference for other cities to introduce similar measures.

#### 2.2. Experimental Data

In total, 2871 complete and effective initial data records are obtained, covering 2871 residential quarters in Shenzhen as shown as Figure 1. The data source is https://shenzhen. qfang.com (accessed on 20 August 2021). The records include the following three types of data points:

- 1. Latitude and longitude: The latitude and longitude range between 22.484310° N and 22.788011° N and between 113.814605° E and 114.498340° E, respectively. The latitude and longitude coordinates used for GNNWR are converted to the WGS 1984 50N coordinate system after projection conversion.
- 2. Endogenous variables: These variables include the age of the building (AB), number of parking spots (NPS), management fee (MF), green ratio (GR), and plot ratio (PR). AB is calculated as the difference between 2021 and the construction year. If the construction age is a range, the completion time is considered. If MF is a range, it is calculated as the average of the upper and lower bounds. GR and PR are defined as follows:

$$GR = rac{S_{Vegetation}}{S_{Land}}$$
 $PR = rac{S_{Floor}}{S_{Land}}$ 

where  $S_{Vegetation}$  is total green area;  $S_{Land}$  is total land area;  $S_{Floor}$  is total building area of the neighborhood.

3. Environment-related variables: These variables include distance from the sea (SD), quality of available public schools (QAPS), number of subway stations within a radius of 1 km (NSS), and distance to the nearest subway station (DSS). SD is calculated with reference to the location of the nearest coast, and DSS is indicated in units of meters. QAPS is calculated using the following process: We divide schools into four types (provincial key schools, city key schools, district key schools, and ordinary schools) and assign points to each category (1-4 for ordinary, district key, city key, and provincial key junior high schools, respectively; and 1.5, 2.5, 3.5, and 4.5 for ordinary, district key, city key, and provincial key elementary schools, respectively). The points of the best school in a school district are set as the QAPS. The QAPS is designed considering the following aspects. First, according to the real-estate agencies in Shenzhen and Hangzhou, key schools correspond to higher weights than common schools. Second, the real-estate agencies indicate that elementary schools correspond to higher weights than junior high schools because parents are more likely to choose a better school in the early stages of child development. Third, according to comparative analyses, the QAPS is a metric with high statistical significance.



Figure 1. Distributions of Neighborhood Data Points in Shenzhen, China among 10 Districts. (a) 10 Folds; (b) Test Set.

# 2.3. Research Methodology

#### **Geographically Weighted Regression (GWR):**

Based on the first law of geography, some scholars have proposed a geographical weighted regression (GWR) model, trying to change the regression coefficient from global to local, and change the weights of adjacent points according to different distances in the regression framework. GWR model defines spatial nonstationarity as [7,25]:

$$y_i = w_0(u_i, v_i) \times \beta_0 + \sum_{k=1}^p w_k(u_i, v_i) \times \beta_k x_{ik} + \epsilon_i$$

We can denote the coefficient as  $\beta_k(u_i, v_i) = w_k(u_i, v_i) \times \beta_0$ , to substitute the estimated value of ordinary least squares model  $\beta_k$ , the estimated value can be obtained as a uniform linear structure:

$$\widehat{y}_i = \sum_{k=0}^p \widehat{\beta}_k(u_i, v_i) x_{ik}$$

The estimator in matrix form can be expressed as:

$$\widehat{y}_i = \mathbf{x}_i^T (\mathbf{X}^T \mathbf{W}(u_i, v_i) \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}(u_i, v_i) \mathbf{y}$$

The spatial weight matrix  $W(u_i, v_i)$  can be expressed as:

$$\mathbf{W}(u_{i}, v_{i}) \triangleq \begin{bmatrix} w_{1}(u_{i}, v_{i}) & 0 & \cdots & 0 \\ 0 & w_{2}(u_{i}, v_{i}) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & w_{n}(u_{i}, v_{i}) \end{bmatrix}$$

In GWR model, the weight kernels usually use Gaussian, bi-square, tri-cube and exponential functions. These functions can relatively simply express the complex relationship between spatial proximity (e.g., spatial distance) and spatial nonstationarity (i.e., spatial weight).

It should be noted that there are different ways to select the function in the spatial weight matrix, and different selection methods directly affect the final modeling accuracy. Th

$$w_{ii} = e^{-\frac{d_{ij}^3}{b^2}}$$

where  $d_{ij}^s$  is the distance between points *i* and *j*; *b*, the bandwidth, producing a declining effect relative to  $d_{ii}^s$ , has different methods to select: for the fixed Gaussian weight function, the bandwidth is the same at each point and is a constant in the same model; for the adaptive Gaussian weight function, the bandwidth is different at each point, and the point distance closest to the point is often taken as the value of bandwidth. In any case, the Gaussian weight function requires a variable input, that is, the distance range (fixed bandwidth) or the number of major adjacent features (dynamic bandwidth).

The bi-square weighted function can be expressed as:

$$w_{ij} = \begin{cases} [1 - (d_{ij}^s / b_i)^2]^2, & d_{ij}^s < b_i; \\ 0, & \text{the others.} \end{cases}$$

where  $d_{ij}^s$  is the distance between two points;  $b_i$  is the bandwidth. It is also divided into fixed type and adaptive type according to the above method.

This model is built using adaptive functions, i.e., an input variable is needed to select the number of major neighboring elements, and the  $AIC_c$  criterion is used to determine whether it is more preferable [26].

#### Geographically Neural Network-Weighted Regression (GNNWR):

Similarly, based on the nonstationarity in the spatial relationship, GNNWR goes further than GWR, trying to more accurately catch the fluctuation of spatial nonstationarity on the regression relationship at different locations. The key step of GWR is the selection and construction of spatial weight matrix function. On this basis, GNNWR attempts to go further and find an appropriate spatial weight matrix function with the help of neural network.

To accurately fit the complex relationship between spatial distance and spatial weight, GNNWR designs a spatial weighted neural network (SWNN) to achieve the neural network expression of weight kernel function. Specifically, SWNN takes the spatial distance between points as the input layer and the spatial weight matrix as the output layer, and selects the appropriate number of hidden layers according to the modeling needs. The spatial weight calculation of the points corresponds with:

$$\widehat{y}_i = \mathbf{x}_i^T \widehat{\boldsymbol{\beta}}(u_i, v_i) = \mathbf{x}_i^T \mathbf{W}(u_i, v_i) (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

where  $W(u_i, v_i)$  is the spatial weight matrix as:

$$\mathbf{W}(u_{i}, v_{i}) \triangleq \begin{bmatrix} w_{0}(u_{i}, v_{i}) & 0 & \cdots & 0 \\ 0 & w_{1}(u_{i}, v_{i}) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & w_{p}(u_{i}, v_{i}) \end{bmatrix}$$

That is, this matrix is the result of function  $W : R^2 \to R^{(1+p)\times(1+p)}$ . SWNN further considers the existence of an intermediate variable  $[d_{i1}, d_{i2}, d_{i3}, \dots, d_{in}]$  and matrix  $W(u_i, v_i)$  is a function of variable  $[d_{i1}, d_{i2}, d_{i3}, \dots, d_{in}]$ , where  $d_{ij}$  is the distance from point i to sample point j. Thus, the GNNWR-based house price estimation model is shown as Figure 2:



Figure 2. Network Structure for Geographical Neural Network-Weighted Regression Model.

#### 2.4. Indicators of Model Performance

The paper uses the following metrics to evaluate the performance of the model. Firstly, the correction of Akaike information criteria  $(AIC_c)$  [7] is as follows:

$$AIC_{c} = nln(\hat{\sigma}^{2}) + nln(2\pi) + n\frac{n + tr(S)}{n - 2 - tr(S)}$$

The expression of matrix *S* is shown in Appendix B. The method is applicable for both GWR and GNNWR. In practice, the smaller the value, the better the performance of the model [26], and we use  $AIC_c$  to select the appropriate input parameters for GWR model. Other measures of model performance include: coefficient of determination ( $R^2$ ), root mean square error (RMSE), mean absolute error (MAE), and mean absolute percentage error (MAPE). The definitions are as follows:

$$R^{2} = 1 - \frac{\sum_{i=1}^{n} (y_{i} - \widehat{y}_{i})^{2}}{\sum_{i=1}^{n} (y_{i} - \overline{y}_{i})^{2}}$$
$$RMSE = \sqrt{\frac{\sum_{i=1}^{n} (y_{i} - \widehat{y}_{i})^{2}}{n}}$$
$$MAE = \frac{\sum_{i=1}^{n} |y_{i} - \widehat{y}_{i}|}{n}$$
$$MAPE = \frac{1}{n} \sum_{i=1}^{n} |\frac{y_{i} - \widehat{y}_{i}}{y_{i}}| \times 100\%$$

Among them,  $\overline{y}$  is the average of the observed values;  $\hat{\sigma}^2$  is the mean square error of the model and *p* is the effective degree of freedom of the model.

# 2.5. Neural Network Design and Implementation

The GNNWR model uses a classic neural network framework, the process flow of which is illustrated in Figure 2. Additionally, 10-fold cross-validation is performed to ensure the robustness and reliability of the algorithm. All layers of the spatially weighted neural network are fully connected, and the dropout technique proposed by Srivastava et al. [27] is applied to enhance the generalization ability of the model. Each hidden layer is combined with the batch normalization technique, the parameters are initialized using the

method proposed by He et al. [28], and the parametric refined linear unit is used as the activation function.

OLS, GWR, and GNNWR are applied to the max–min normalized housing price data, and 10-fold cross-validation is performed to build the models over the training set, as shown in Figure 3. The results for the validation set are used to calculate the root mean square error (RMSE) as the loss function and evaluate the generalization ability. The result with the highest generalization ability is selected for the three methods, and the predictive ability of the model is tested on the test set. From the 2871 records, 431 (approximately 15%) constitute the test set, and the remaining 2440 records are equally divided into 10 groups (each containing 244 records, approximately 8.5%) for the cross-validation process.





Considering the results of preliminary experiments (Table 1), we use a six-layer neural network structure containing: one input layer; four hidden layers with 512, 128, 64, and 16 neurons; and one output layer. The number of neurons in the input layer is the number of training samples, and the number of neurons in the output layer is the number of parameters in the linear regression model (number of independent variables plus one).

Table 1. Loss for Different Structures.

Structure of Hidden Layers	Validation Loss	Train Loss	Test Loss
(1024, 512, 256, 128, 64, 32)	0.006470	0.002790	0.008867
(512, 128, 64, 16)	0.006427	0.0038040	0.008683
(512, 128, 32)	0.006529	0.0040193	0.008555
(256, 64, 16)	0.006537	0.0043795	0.008555
(256, 32, 8)	0.006527	0.0049904	0.008379
(256, 32)	0.006567	0.0046721	0.008992

After several trials, in the hyper-parameters, the value of learning rate is  $10^{-2.95} \approx 0.00112$ ,  $\beta_1 = 0.8$ ,  $\beta_2 = 0.999$ , batch size = 128. The loss percentage in the dropout layer is 0.9, and the maximum number of epochs is 90,000.

To reflect the optimization in the iterative process, Figure 4 shows the change in the test indicators of one fold during model training. In more than 30,000 epochs, if the loss



for the validation set does not decrease after 9000 epochs, the neural network training is terminated.

**Figure 4.** Performance Variations on the Train (Orange Line) and Validation (Blue Line) Sets in the Train Process of GNNWR Model. (a) The Decrease of  $AIC_c$  on Train Set. (b) The Decrease of Average Absolute Error. (c) The Decrease of Average Relative Error. (d) The Increase of Determination Coefficient. (e) The Decrease of Loss.

To demonstrate the superiority of GNNWR, the results of the OLS, GWR, and GNNWR on the Shenzhen housing price dataset are compared. The GWR method uses the golden search method to identify the most suitable number of neighboring elements according to the Akaike information criterion ( $AIC_c$ ).

As the NSS and QAPS variables are typically small integers, the design matrix used in GWR modeling may exhibit multicollinearity when the number of neighboring elements is small. Therefore, the search lower bound is set as 100 (i.e., at least 100 neighboring elements are involved in the solution of the local regression coefficients).

According to a simple preliminary experiment, the bi-square function significantly outperforms the Gaussian function as the kernel function of GWR over the 10-fold dataset. The parameters are presented in the Table A1 attached in Appendix A. Therefore, in the subsequent analyses, the bi-square function is used as the weight kernel function for the GWR model. In the experiments, we use OLS as a rudimentary contrast. Both OLS and GWR solutions are built on ArcGIS Pro 2.5.2. GNNWR is implemented using the TensorFlow 1.15.0 library under the Python 3.6.13 kernel.

#### 3. Results

#### 3.1. Data Set Analysis and Descriptive Statistics

Table 2 summarizes the results of the correlation analysis and descriptive statistics of Shenzhen housing prices with different variable factors.

The variables can be ranked as follows in decreasing order of the absolute values of the correlation coefficients with housing prices: SD, MF, NSS, DSS, GR, AB, PR, QAPS, and NPS. MF, NSS, GR, PR, QAPS, and NPS are positively correlated with housing prices; and SD, DSS, and AB are negatively correlated with housing prices.

Indicator	Price	AB	NPS	MF	GR	PR	SD	QAPS	NSS	DSS
Mean	62,219.3	17.995	507.196	2.625	0.340	3.113	6586.1	3.704	1.769	930.5
Maximum	132,000	51	5500	36.6	0.990	7.000	24,967.2	4.5	8	25,110.0
Minimum	16,100	1	1	0	0.100	0.100	23.2	0	0	16.8
Std. Dev.	22,986.5	7.138	647.509	1.888	0.130	1.438	4933.7	1.179	1.432	1838.5
Correlation Coefficient	-	-0.118	0.079	0.262	0.216	0.105	-0.504	0.080	0.248	-0.236
Variation Coefficient	2.707	2.521	0.783	1.391	2.620	2.164	0.749	3.143	1.235	0.506
VIF	-	1.622	1.243	1.227	1.122	1.150	1.167	1.204	1.365	1.136
<i>t</i> -test <i>p</i>	-	0	$3.4 imes10^{-8}$	0	0	0.0197	0	0.3208	0	0

Table 2. Exploratory Analysis and Descriptive Statistics of the Experimental dataset.

Basic hypothesis testing is performed for each variable in the global regression equation using R language. The coefficient for each variable is assumed to be zero in the global regression equation, and a test statistic that satisfies the t-distribution when the hypothesis holds is constructed. Correspondingly, the *p*-value can be calculated. However, it can be found that the *p*-value of PR is not significant at the significance level of 0.01 and the *p*-value of QAPS is not significant at the significance level of 0.05 or even 0.1. If a global regression is used, these two independent variables should be excluded. However, two spatial statistical modeling methods, GNNWR and GWR, take them into account in this study, and the significance of each variable in this model can be re-tested with the help of the F2 statistic in the Appendix B. According to the analysis of nonstationarity diagnostics in Appendices B and C, both variables are highly significant, when the coefficients of the linear model are allowed to vary with geographic coordinates. These results demonstrate the superiority of the spatial statistical modeling approach over the existing approaches.

#### 3.2. Comparison of Indicators of House Price Valuation Models

The housing price valuation model is evaluated considering the fitting ability over the training set and prediction ability over the test set. We stochastically divide the 2871 records into the training and validation sets (2,440 records as 10 folds) and a test set (remaining 431 records).

The models are evaluated considering the R<sup>2</sup> value, RMSE, mean absolute error (MAE), mean absolute percentage error (MAPE), AICc, and Pearson correlation coefficient. For the dataset generated after the 10-fold crossover, the following results are obtained after merging the validation sets in Table 3.

Set	Model	R2	RMSE	MAE	MAPE	Mean Err.	Pearson Cor. Coe.
Merged Validation Set	GNNWR GWR OLS	0.840177 0.788728 0.432101	9069.561 10,427.68 17,096.31	6558.630 7581.746 13,003.76	0.111965 0.128538 0.228767	27.88808 -73.9177 -5.60228	0.916637 0.888123 0.657404
Test Set	GWR GNNWR	0.790389 0.817178	11,195.01 10,455.19	7912.005 7108.715	$0.122266 \\ 0.109174$	911.3839 1393.691	0.891319 0.905834

Table 3. Indicators of GNNWR, GWR, and OLS on Merged Validation Set and Test Set.

The results demonstrate the superiority of the GNNWR model. The OLS model achieves the worst prediction, with the lowest R<sup>2</sup> and highest prediction error in terms of the RMSE, MAE, and MAPE. Given the severe spatial nonstationarity, the OLS model cannot detect the intrinsic relations and spatial fluctuations between the housing prices and independent variables. The RMSE and MAE of the GNNWR model are approximately 13.0% and 13.5% lower than those of the GWR model, respectively. The GNNWR model also outperforms the existing models in terms of the R<sup>2</sup> and MAPE values. Additionally,

the mean residual error of the GNNWR model is 62.2% lower than that of the GWR model, which means that its predictions are more unbiased than those of the GWR on this dataset. Overall, the GNNWR model exhibits improved generalization ability compared to the other models.

The performance indicators of the GWR and GNNWR models in the modeling process over the 10 training sets are compared to examine the fitting quality. Table A2 presents the results of these indicators.

The number after the GWR results refers to the number of most suitable neighboring elements selected based on the  $AIC_c$  value. As the training sets are slightly different, the most appropriate number of neighboring elements is re-selected each time the GWR model is rebuilt. Consequently, we use the best GWR model to ensure the fairness of the comparisons.

For all 10 datasets, the GNNWR model considerably outperforms the GWR model for all indicators (AICc, RMSE, R2, and Pearson correlation coefficient). The notable improvement in AICc indicates that the GNNWR model can more accurately predict the housing price and space weight matrix without significant increase in the complexity. In contrast, the GWR model is vulnerable to overfitting, which decreases the correctness of the predictions on the validation sets. In summary, the GNNWR model produces a more capable kernel function than any GWR model, and it can effectively capture the spatial heterogeneity details, estimate the spatial weights, and predict the dependent variables.

Furthermore, the generalization ability is evaluated by predicting the test set. In this analysis, the results of the models with the best generalization ability are compared. Both the GNNWR and GWR models exhibit the highest performance when dataset 4 is used as the validation set, and the other indicators' details are shown in the Appendix A of Table A3.

The GNNWR model outperforms the GWR model in predicting the test dataset: The MAE and MAPE of the GNNWR model are 10.2% and 10.7% lower, respectively; thus, real estate agencies can achieve more accurate estimations than those associated with the GWR model. Moreover, the RMSE of the GNNWR model is decreased by 6.6%, the R2 value and Pearson correlation coefficient are enhanced, while the mean error is increased. In a recent study, experiments on housing price datasets were conducted using the geographically weighted artificial neural network (GWANN) model [29]. The RMSE of the GWANN model was only 3.3% better than the GWR model for the best batch when predicting housing prices. This result also demonstrates the superiority of the proposed framework.

Moreover, we have conducted some hypothesis testing about the spatial nonstationarity in Appendices B and C. The results of experiments and theoretical analyses highlight the importance of establishing a regression model with coefficients variegated among geographical coordinates. We have also extended the analysis to the results of GWR model in Appendix D. The comparison of the two models indicates that the GNNWR model can more accurately detect the spatial heterogeneity in a facile manner.

#### 4. Comparative Analysis and Discussion

#### 4.1. Comparison of Prediction Performances of House Price Valuation Models

The relative error rate for each prediction over the validation and test sets is calculated. The GWR and GNNWR models are compared using the Q–Q plot as Figure 5a,b, derived using MATLAB. It should to be pointed out that no more than 5% of the points are not shown outside this range.



**Figure 5.** Analysis of the Relative Error Rates between GWR and GNNWR Model. (**a**) Q-Q Plot of Relative Error Rates On Merged Validation Set. (**b**) Q-Q Plot of Relative Error Rates On Test Set. (**c**) Histogram of Relative Error Rates On Merged Validation Set. (**d**) Histogram of Relative Error Rates On Test Set. (**e**) The Ratio of the Top N Best/Worst Predictions from 2 Models.

Ordering the relative error rates, it can be found that the relationship between the k<sup>th</sup> value on the validation set is approximately  $\delta_{GWR}^{(k)} = 1.123 \delta_{GNNWR}^{(k)} + 0.0033$ . The relationship between the k<sup>th</sup> value on the test set is approximately  $\delta_{GWR}^{(k)} = 1.160 \delta_{GNNWR}^{(k)} + 0.0003$ . These reference lines that represent the theoretical distribution have a clear deviation with y = x, which enable us to confirm the superiority of GNNWR models.

The histograms of the models are shown in Figure 5c,d. On the validation set, setting the histogram horizontal coordinates between [0,1] and bin width of 0.09, 9 of the 11 bins with error rates less than or equal to 9.9% include more data from the GNNWR model. This trend is also observed in the test set. Similarly, when the histogram horizontal coordinates are set between [0, 1], and the bin width is set as 0.15, five of the seven bins with error rates less than or equal to 10.5% have more data from the GNNWR model.

The prediction data of the validation sets for both models are shown in Figure 5. The numbers of data points in both sets below a certain threshold are calculated, and the ratio of the two numbers is plotted as a blue line on the graph. The ratio of the number of predicted data from GWR to that from GNNWR when the statistical error rate is above a certain value is plotted as the orange line on the graph. For data (two runs of predictions over 2440 records) with a relative error rate of less than 0.203, the number of predictions from GNNWR are 1.34 times larger than those from GWR. In contrast, for data with error rates higher than 0.37, the number of predictions from GWR are 1.62 times larger than those from GWR are 1.62 times larger than those from GWR and GWR predictions account for more high-precision predictions and high-error predictions, respectively.

Comparing with other works [22], it can be found that another study also supports the conclusion that GWR can significantly reduce the prediction error compared to OLS models, indicating that spatial heterogeneity exists. In another study on Shenzhen house prices, the authors used GWR model to increase the R<sup>2</sup> from 0.56 to 0.79 [12]. Some simple artificial intelligence models, such as decision tree models, can even predict worse than OLS if they are not designed properly [30]. In a separate study comparing the OLS model with multiple models, the best stepwise and tuned support vector machine model reduced the RMSE by 25%, the polynomial regression model reduced the RMSE by 8.3%, and even the optimal simple neural network selected from the 1–3 hidden layers increased the RMSE by 66% [31]. Since the 1990s, scholars have been trying to use neural network models to predict house prices and compare them with OLS models. Some studies have demonstrated the superiority of the neural network approach [32–34], but others have found that there is no great need to use neural networks [35,36]. Considering the 47% reduction in RMSE metrics compared to OLS in this study, it is easy to see that simply using complex functions trained by neural networks to approximate the training dataset does not improve the prediction accuracy, and that a GWR-based framework can best capture information on the geographic distribution. These indicate that accurate estimation of spatial heterogeneity is extremely necessary.

According to the literature [22], the prediction error of the GWR is significantly lower than that of the OLS models, indicating the presence of spatial heterogeneity. In a study on Shenzhen housing prices, the GWR model was noted to increase the  $R^2$  from 0.56 to 0.79 [12]. Notably, simple artificial intelligence models, such as decision tree models, may achieve inferior predictions compared to those of the OLS model if designed inappropriately [30]. In another study, the OLS model was compared with multiple models. Compared to the OLS model, the best stepwise and tuned support vector machine model achieved a 25% lower RMSE, the polynomial regression model achieved an 8.3% lower RMSE, and the optimal simple neural network (selected from frameworks with 1-3 hidden layers) achieved a 66% higher RMSE [31]. Since the 1990s, scholars have constantly attempted to use neural network models to predict housing prices and compare them with OLS models. Although several studies have demonstrated the superiority of the neural network approaches [32–34], other studies have highlighted their limitations in certain applications [35,36]. Considering the 47% reduction in RMSE metrics compared to OLS in this study, it is easy to see that simply using complex functions trained by neural networks to approximate the training dataset does not improve the prediction accuracy, and that a GWR-based framework can best capture information on the geographic distribution. These findings highlight the importance of accurately estimating spatial heterogeneity.

#### 4.2. Analysis of Variables Related to House Prices

The GNNWR model is based on the structure of linear regression, in which different coefficients are assigned to different variables based on the location of the property to capture the spatial heterogeneity. For the 10-fold dataset established in this study, the coefficients of different independent variables at each prediction point can be visualized and output after merging the validation sets, as shown in Figure 6. This section describes the analysis of the fluctuations of these coefficients.

As the data are normalized before they are used in GNNWR training, the values for different parameters can be directly compared, as indicated in Table 4. The mean values indicate that the degree of influence of each independent variable on housing prices is different, and the parameters can be ranked as follows in the decreasing order of their absolute values: SD, MF, AB, NPS, DSS, GR, NSS, QAPS, and PR. When spatial heterogeneity is considered, the effect of the NPS and AB on housing prices is more notable than that estimated using the correlation coefficient, and the effect of the NSS is less notable than that estimated using the correlation coefficient. The trends of the positive and negative correlations with the housing prices are the same as those in the previous analysis: MF, NPS, GR, NSS, QAPS, and PR are positively correlated with housing prices; and SD, AB, and DSS are negatively correlated with housing prices. The parameters can be ranked in decreasing order of the standard deviations of the coefficients as follows: DSS, MF, SD, AB, NSS, NPS, GR, PR, and QAPS. Specifically, the public transportation conditions and school district conditions (DSS and QAPS, respectively) exhibit the highest and lowest degrees of spatial heterogeneity, respectively, consistent with the intuition. In a more extensive analysis, the distributions of the coefficients of each variable are considered, as shown in the following figures. The figures are plotted using the natural breakpoint method with inconsistent color ranges for different subplots, and the boundaries near zero are finetuned to highlight the positive and negative correlation features. Owing to the small standard deviations, the data of PR and QAPS are classified into only six levels unlike the eight levels of the other variables. The modeling results based on the 10 training sets are smooth, with few mutations and outliers in the geographic proximity. The predictions for the coefficient distributions are consistent for all sets.

The distributions of the intercept and variables derived from the GNNWR are examined to demonstrate the significance of the proposed model in socioeconomic research.

First, we analyze the endogenous variables, specifically, MF, AB, NPS, GR, and PR, in descending order of their influence on housing prices.





(b)



Figure 6. Cont.



**Figure 6.** Coefficient Weight Distributions of Variables for GNNWR model in Shenzhen Real Estate Market. (a) Intercept; (b) Age of Building (AB); (c) Distance to the Nearest Subway Station (DSS); (d) Green Ratio (GR); (e) Management Fee (MF); (f) Number of Parking Spots (NPS); (g) Number of Subway Stations within 1 km radius (NSS); (h) Plot Ratio (PR); (i) Quality of Available Public Schools (QAPS); (j) Sea Distance (SD).

Coefficients of Variables	AB	NPS	MF	GR	PR	SD	QAPS	NSS	DSS	Intercept
Mean	-0.280	0.170	0.458	0.114	0.002	-0.474	0.021	0.033	-0.160	0.508
Maximum	0.612	1.101	2.675	0.836	0.320	0.383	0.179	0.701	6.763	1.486
Minimum	-1.450	-0.420	-1.322	-0.230	-0.191	-2.035	-0.055	-0.868	-4.627	-0.018
Std. Dev.	0.195	0.179	0.609	0.108	0.057	0.253	0.026	0.187	0.820	0.208

 Table 4. Descriptive Statistics of Coefficients of Variable.

The housing prices are mainly positively correlated with the MF, with representative areas including southwest Nanshan District and the southern coast of Baoan District. The negatively correlated areas include Huanggang in Futian District along the border with Hong Kong. We speculate that marginal districts may exhibit a stronger positive correlation between the MF and housing prices, because the MF may characterize the differences between villas and hotels, as example dwellings.

The increase in AB limits the housing prices in Shenzhen. The negative correlation between the housing prices and AB is the strongest in the coastal Nanshan District with Houhai as the core, central Futian District with Xiangmi Lake's eastern shore as the core, and southern Longhua District with Shenzhen North Station as the core. The strong negative correlation is attributable to the large supply of high-quality new houses near these locations, as old properties are vulnerable to cold markets. At the border of Luohu and Futian districts, the correlation between AB and the housing prices transforms from negative to positive. According to Goodman et al. [37], AB influences housing prices in a nonlinear manner, with a positive effect observed when AB is greater than a certain threshold. The Shenzhen areas explored by the GNNWR model are those with the earliest constructions, and these areas include famous landmarks, such as Dongmen Old Street and Diwang Building.

The NPS and housing prices are positively correlated. The strongest positive correlations are observed in central Nanshan District, Xiangmi Park in western Futian District, and near Caiwuwei in Luohu District. The contribution of NPS to housing prices is expected to be the most significant in middle-class residential areas and wealthy areas. The areas with strong positive correlations, explored by the GNNWR model, coincide with such regions.

The increase in GR increases housing prices, especially in the central Futian District and central Luohu District, which are located in the prosperous part of the city with higher demand for GR. In the suburbs and along the coast, the GR does not tend to increase housing prices, and a weak negative correlation is observed in Longgang District.

The variations in the PR are small, and its influence on housing prices is not significant in terms of the average weight. However, the PR and housing prices are positively correlated in western Luohu District. This finding is contrary to the general perception, potentially because the PR in this region is closely related with the overall appearance of the neighborhood. Western Luohu District is the older urban area of Shenzhen. A low PR is representative of the old and dilapidated characteristic of the neighborhood, whereas a high PR corresponds to new high-rise housing. These characteristics are the potential reason for the positive correlation in this region unlike in the other locations.

Second, we examine the influence of the environmental variables, specifically, SD, DSS, NSS, and QASP, in descending order of their influence on housing prices.

The SD and housing prices are negatively correlated, with typical areas including most of the Nanshan and Luohu Districts. Moreover, a negative correlation is observed in certain inland parcels, such as the southern part of Longgang District and the northern edge of Luohu District. This correlation is attributable to the fact that the SD characterizes the distance from the core urban area.

The correlation between the DSS and housing prices fluctuates considerably. Generally, a larger distance from the subway corresponds to lower housing prices. However, a positive correlation is observed in southeastern Futian District, southern Luohu District,

and southern Nanshan District. Notably, these regions correspond to high employment concentration. Southeastern Futian District houses the Huaqiang North Market, one of the largest cellphone part distribution markets worldwide. Southern Luohu and Nanshan Districts house the Xinxiu Village Industrial Zone and Shekou Industrial Zone, respectively. Consequently, the following inferences can be derived: in residential and suburban areas, a smaller DSS corresponds to more convenient commuting and thus higher housing prices. In contrast, in areas with dense subway entrances, such as central business districts or industrial areas, a larger DSS may drive down housing prices. People who buy houses in this neighborhood are already close to their workplace; thus, the need to commute through subway channels is insignificant. Proximity to the subway entrance may lead to aggravated noise and congestion.

The NSS and housing price are positively correlated. For the Nanshan, Futian, and Luohu Districts, the distribution of positive and negative correlations is the opposite to that of DSS. This result confirms the abovementioned conjecture.

The QASP and housing prices are weakly correlated. Except for western Luohu District, in which the housing prices are strongly positively correlated with the QASP, the effect of the QASP is not significant in the other districts in Shenzhen. Western Luohu District represents the old city of Shenzhen, and the old residential areas are desirable owing to the mature school districts, resulting in a strong positive correlation. In contrast, the high-quality new housing in the new district does not have an established school district, and the housing price is dominated by other factors. Consequently, the influence of the school district is weak.

The distributions of the intercepts in the regression model can be effectively explained. The intercept represents the inherent premium of the house after considering all of the effects of the independent variables. Figure 6 shows that according to the reference prices for the second-hand housing transactions introduced by the government, the highest inherent premium pertains to the coast of Nanshan District with Houhai as the core and middle of Futian District with the east shore of Xiangmi Lake as the core. The market frantically attempts to exploit the scarcity of premium locations and assign higher premiums to the abovementioned regions. In 2020, the highest residential transaction prices for these two sites (USD 50,000–70,000 per square meter) set a new record for housing prices in Shenzhen, and the prices for marginal residences were USD 10,000–20,000. In this context, the reference prices defined by the government helped reflect the inherent premium distribution and narrow the gap between the inherent premiums. Similarly, the GNNWR model can accurately estimate this premium based on the reference price.

The regression functions of the GNNWR and GWR models are similar, which validates the proposed model. However, the GWR model is associated with higher errors. These analyses highlight the value of the proposed valuation model in socioeconomic research.

#### 5. Conclusions

Based on Shenzhen housing price data, we demonstrate the superiority of the GNNWR model over the OLS and GWR models in price valuation. We introduce a 10-fold validation approach and obtain the following results by performing predictions for 1 10th of the data in each fold. The RMSE of the GNNWR model is 13% and 47% lower than those of the GWR and OLS models, respectively. Additionally, the GNNWR outperforms the GWR and OLS strategies in terms of the other metrics. The robustness of the GNNWR model is experimentally evaluated. The 10-fold validation mechanism avoids stochastic interference, and the results of testing over the test set demonstrate the validity of the proposed model. Hypothesis testing is performed to analyze the significance of the spatial heterogeneity. The AICc metric for the training set indicates that the increase in the fitting accuracy of the GNNWR model compared with that of the GWR model (which also models the spatial heterogeneity) is considerably larger than the increase in the complexity of the kernel function. Therefore, the GNNWR model is highly robust. The GNNWR model is noted to

exhibit excellent information mining ability, in the context of the spatial heterogeneity of Shenzhen.

The contributions of this study can be summarized as follows. First, the spatial weight matrix can adequately reflect the dataset characteristics. In contrast, GWR, as a relatively traditional modeling method for spatial analysis, commonly uses kernel functions, the choices of which are limited, for example, to Gaussian and bi-square functions. Therefore, the performance of GWR frameworks is limited. Second, the proposed framework overcomes similar limitations of the existing methods by refining the experiment process and diminishing the black-box property. In particular, the GNNWR model uses the neural network only to fit the kernel function more accurately than that achievable using GWR and reserves the reasonable framework of GWR (i.e., spatial weighted regression). Therefore, the kernel function in the GNNWR model is not a simple linear function or an a priori exponential function similar to that used by the GWR. Instead, the kernel function clarifies the general correlation between the spatial distances and weights and is learned from the data through artificial intelligence. Third, the proposed valuation model can more effectively interpret the geographical coordinates compared to other neural network prediction methods that take into account less geographical location information and thus exhibit unstable performances. Although the neural network predictions are noted to be more accurate than those of the OLS models [32–34] in this study, several studies have highlighted that neural network models often do not outperform the OLS model and its improved variants (e.g., hedonic models that correct the dependent variable through log and polynomial regression models) [35,36]. Nevertheless, the RMSE of a few neural network models is more than 30% lower than those of the OLS models.

The GNNWR model has not been applied in the socioeconomic domain since its invention, and we fill this gap. To address the housing price prediction problem, state-of-the-art methods typically modify the GWR model in a coarse manner to better understand the geospatial information. The resolution is not as fine as that associated with the GNNWR model. In common housing price prediction tasks, for instance, in Kaggle competitions, the input data lack geographical coordinates to prevent the data from confusing neural network and decision tree models [38,39]. Only a few neural network models [29] exhibit a high performance when dealing with geographical coordinates. The proposed model outperforms the existing models because we combine the GWR and neural network models. Specifically, we exploit the GWR framework to comprehend the geospatial coordinates and overcome the limitation associated with the kernel function by using a neural network. Future research can be focused on the following aspects. First, the linearity, homoscedasticity, independence, and normality properties of the error term in the linear model are examined. If these aspects are not satisfied, the dependent variable can be preprocessed using the Box–Cox method. Second, more independent variables can be considered to expand the choice of independent variables. Third, comparable tests can be performed on other datasets, or data from multiple cities can be acquired to build a housing price prediction benchmark.

Author Contributions: Conceptualization, Zimo Wang and Yicheng Wang; methodology, Sensen Wu; software, Zimo Wang and Yicheng Wang; validation, Zimo Wang; formal analysis, Zimo Wang and Yicheng Wang; resources, Sensen Wu; data curation, Zimo Wang; writing—original draft preparation, Zimo Wang and Yicheng Wang; writing—review and editing, Zimo Wang and Yicheng Wang; visualization, Yicheng Wang; supervision, Sensen Wu and Zhenhong Du; project administration, Sensen Wu; funding acquisition, Zhenhong Du. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research was funded by FUNDER National Key R&D Program of China grant number 2021YFB3900902, National Natural Science Foundation of China grant number 42001323, 41871287, Provincial Key R&D Program of Zhejiang grant number 2021C01031.

Institutional Review Board Statement: Not applicable.

Informed Consent Statement: Not applicable.

**Data Availability Statement:** You can find the Chinese document "Shenzhen 2020 National Economic and Social Development Statistical Communiqué" at http://tjj.sz.gov.cn/zwgk/zfxxgkml/tjsj/tjgb/content/post\_8717370.html (accessed on 1 September 2021). The specific price of each neighborhood can be obtained from the Chinese document "Notice of Shenzhen Municipal Housing and Construction Bureau on the Establishment of a Reference Price Release Mechanism for Secondhand Housing Transactions" at http://www.sz.gov.cn/cn/xxgk/zfxxgj/tzgg/content/post\_8545768.html (accessed on 1 September 2021). More details about each neighborhood are collected at https://shenzhen.qfang.com (accessed on 1 September 2021).

Conflicts of Interest: The authors declare no conflict of interest.

# Abbreviations

The following abbreviations are used in this manuscript:

OLSOrdinary Least SquaresGWRGeographical Weighted RegressionGNNWRGeographical Neural Network Weighted Regression

#### Appendix A. Tables Related to Experiment Results

 Table A1. GWR Model Performance for Different Kernel Types.

			Test					
GWR Kernel Type <sup>1</sup>	R2	RMSE	MAE	MAPE	AICc	Correlation Coefficient	R2	Correlation Coefficient
Bi-square (105)	0.8861	7655.203	5623.454	0.094994	51,842.0	0.941789	0.7935	0.892818
Gaussian (101)	0.7471	11,408.08	8372.857	0.141284	52,790.3	0.865382	0.6120	0.783185

<sup>1</sup> The number in parentheses indicates how many neighbors have been used to build the spatial weight matrix

Table A2. Indicators of GWR and GNNWR over Train Sets
---

Train Set <sup>1</sup>	Model <sup>2</sup>	R2	RMSE	MAE	MAPE	Pearson Cor. Coe.	AICc
0	GNNWR	0.9130	6665.74	4907.16	0.084455	0.955890	44,935.93
0	GWR (108)	0.8806	7810.28	5722.90	0.096476	0.938932	46,711.89
1	GNNWR	0.9145	6698.92	4945.16	0.084418	0.956290	44,923.25
1	GWR (101)	0.8881	7662.40	5624.61	0.095042	0.942811	46,714.44
r	GNNWR	0.9168	6516.65	4849.38	0.084081	0.957767	44,799.69
2	GWR (101)	0.8835	7713.28	5663.39	0.095898	0.940476	46,751.67
2	GNNWR	0.9068	6923.85	5118.84	0.087057	0.952315	45,066.73
3	GWR (101)	0.8887	7566.31	5594.39	0.094180	0.943155	46,666.59
4	GNNWR	0.9180	6489.22	4784.29	0.080781	0.958206	44,776.13
4	GWR (101)	0.8842	7711.99	5656.98	0.095410	0.940830	46,748.33
Б	GNNWR	0.9109	6777.58	4990.82	0.086299	0.954517	44,972.38
5	GWR (101)	0.8849	7702.18	5664.53	0.095798	0.941220	46,744.85
6	GNNWR	0.9175	6538.53	4796.90	0.082787	0.958058	44,850.11
0	GWR (106)	0.8814	7839.84	5715.71	0.096394	0.939293	46,751.99
7	GNNWR	0.9183	6529.28	4827.74	0.081521	0.958488	44,795.06
7	GWR (101)	0.8871	7675.50	5658.68	0.095299	0.942348	46,730.64
8	GNNWR	0.9077	6868.49	5104.66	0.087495	0.953352	45,038.13
0	GWR (104)	0.8825	7749.99	5713.72	0.096287	0.939906	46,735.51
9	GNNWR	0.9082	6814.82	4998.50	0.086338	0.953331	45,025.80
)	GWR (107)	0.8796	7802.49	5714.57	0.096653	0.938439	46,719.31

<sup>1</sup> The Train set of 0 means that dataset 0 is excluded and the 1, 2, ..., 9 datasets are selected, and so on. <sup>2</sup> The number in parentheses indicates how many neighbors have been used to build the spatial weight matrix.

Validation Set	R2	RMSE	MAE	MAPE	Pearson Cor. Coe.
0	0.836963	9454.353	6876.564	0.115830	0.915256
1	0.821430	8697.789	6277.030	0.105508	0.907206
2	0.855282	8930.441	6546.944	0.109760	0.924835
3	0.812268	9812.707	6643.969	0.112837	0.901745
4	0.871563	8189.907	5944.871	0.104658	0.933591
5	0.837077	9098.732	6577.607	0.111918	0.915711
6	0.840490	8783.315	6620.163	0.113857	0.917376
7	0.813833	9034.993	6549.102	0.115053	0.902763
8	0.840247	9332.620	6836.936	0.117101	0.916725
9	0.854961	9260.348	6713.113	0.113130	0.925135

Table A3. Prediction Performance of 10 GNNWR Models for Each Validation Set.

# Appendix B. Indicators of Significance Test Statistics for Spatial Nonstationarity

To test whether a relationship exhibits significant spatial nonstationarity, we perform significance tests of GNNWR and GWR modeling results by using the residual sum of squares and its approximated distribution deduced by Leung et al. [40] and Du et al. [19]. We present the formulation for the GNNWR model as an example. The expression for the GWR can be similarly derived.

First, the hat matrix of GNNWR is expressed as:

$$S_{GNNWR} \triangleq \begin{bmatrix} x_1^T W(u_1, v_1) (X^T X)^{-1} X^T \\ x_2^T W(u_2, v_2) (X^T X)^{-1} X^T \\ \vdots \\ x_n^T W(u_n, v_n) (X^T X)^{-1} X^T \end{bmatrix}$$
$$\delta_i \triangleq tr\{[(I-S)^T (I-S)]^i\}, i = 1, 2 \cdots$$

The statistical quantities  $F_1$  is obtained as:

$$F_1 = \frac{RSS_{GNNWR}/\delta_1}{RSS_{OLS}/(n-p-1)}$$

The distribution of *F*1 can be approximated as F distribution, where  $\frac{\delta_1^2}{\delta_2}$  is the degree of freedom of the numerator and n - p - 1 is the degree of freedom of the denominator. That is, given a significance level  $\alpha$ , if the inequality  $F_1 < F_{1-\alpha}(\frac{\delta_1^2}{\delta_2}, n - p - 1)$  holds, it can be determined that the regression relationship has significant spatial non-smoothness, otherwise the spatial non-smoothness is not significant.

Second, the significance of the spatial nonstationarity can also be verified for each independent variable. The null hypothesis is that the coefficients of this independent variable are at all points in space. The alternative hypothesis is that the coefficients of this independent variable differ in at least one point in each part of the space. First, we define the variance of the coefficients of the *k*th independent variable over n data points.

$$V_k^2 \triangleq \frac{1}{n} \sum_{i=1}^n (\widehat{\beta}_{ik} - \frac{1}{n} \sum_{i=1}^n \widehat{\beta}_{ik})^2$$

Moreover, we define  $e_k$  as an n-rank vector with the (k + 1)th element set as 1 and other set as 0. A square matrix of order n is established, with each element set as 1.

$$B_{k} \triangleq \begin{bmatrix} e_{k}^{T}W(u_{1}, v_{1})(X^{T}X)^{-1}X^{T} \\ e_{k}^{T}W(u_{2}, v_{2})(X^{T}X)^{-1}X^{T} \\ \vdots \\ e_{k}^{T}W(u_{n}, v_{n})(X^{T}X)^{-1}X^{T} \end{bmatrix}$$

$$\gamma_{ik} \triangleq tr\{[\frac{1}{n}B_k^T(I-\frac{1}{n}J)B_k]^i\}, i = 1, 2\cdots$$

The statistical quantities  $F_2$  is obtained as:

$$F_2(k) = \frac{V_k^2 / \gamma_{1k}}{\widehat{\sigma}^2}$$

The distribution of  $F_2(k)$  can be approximated as F distribution, where  $\hat{\sigma}^2$  is the mean square error,  $\frac{\gamma_{1k}^2}{\gamma_{2k}}$  is the degree of freedom of the numerator and  $\frac{\delta_1^2}{\delta_2}$  is the degree of freedom of the denominator. That is, given a significance level  $\alpha$ , if the inequality  $F_2(k) > F_{\alpha}(\frac{\gamma_{1k}^2}{\gamma_{2k}}, \frac{\delta_1^2}{\delta_2})$  holds, the null hypothesis can be rejected and the variable k is determined to have significant spatial nonstationarity, otherwise the spatial nonstationarity is not significant.

# Appendix C. Spatial Nonstationarity Diagnosis of House Price Regression Relationship

Based on the spatial heterogeneity diagnostic indicators, the GNNWR results can be analyzed considering two aspects.

First, we examine if the model results have a significant spatial nonsmoothness. For the 10-fold data, the prediction effect parameters of each GNNWR model in the validation set are presented in Table A3. Using RMSE as the index, the best fitting model (model 4) and the worst fitting model (model 3) were selected for hypothesis testing. The hypothesis testing parameters are determined from the previous derivation as the following Table A4.

Table A4. F1 Hypothesis Testing.

F1 Hypothesis Test	F1	$\sigma_1$	$\sigma_2$	Distribution	Significant Level
Best Fitting Model	0.071602	4439.122	4,871,141	F(4.0454, 2186)	$1  imes 10^{-2}$
Worst Fitting Model	0.117877	3118.766	804,246.6	F(12.094, 2186)	$1 \times 10^{-4}$

After determining the  $F_1$  value and F distribution, the *p* value for the hypothesis can be calculated. The results indicate that the hypothesis is rejected, and severe spatial nonstationarity exists in modeling the Shenzhen house price.

Next, we analyze the significance for each independent variable. The null hypothesis is that the coefficient of each independent variable is a constant. This hypothesis includes another hypothesis that the coefficient of this variable is zero everywhere. In this context, the p value of  $F_2$  can reject both the hypotheses if it is adequately small. The results are presented in Table A5.

Model	Variable	Intercept	AB	NPS	MF	SD	GR	PR	QAPS	NSS	DSS
Best Fitting Model	F Value γ <sub>1</sub> γ <sub>2</sub> Significant Level	$\begin{array}{c} 614.58 \\ 0.0198 \\ 0.0004 \\ 1 \times 10^{-10} \end{array}$	$\begin{array}{c} 234.48 \\ 0.0893 \\ 0.0057 \\ 1 \times 10^{-8} \end{array}$	$\begin{array}{c} 381.14 \\ 0.4068 \\ 0.1095 \\ 1 \times 10^{-9} \end{array}$	$562.52 \\ 0.4108 \\ 0.1085 \\ 1 \times 10^{-10}$	$537.77 \\ 0.4200 \\ 0.1108 \\ 1 \times 10^{-10}$	$503.31 \\ 0.4560 \\ 0.1120 \\ 1 \times 10^{-10}$	$\begin{array}{c} 385.10 \\ 0.5951 \\ 0.1629 \\ 1 \times 10^{-10} \end{array}$	$\begin{array}{c} 418.17\\ 0.6102\\ 0.1595\\ 1\times10^{-11}\end{array}$	$\begin{array}{c} 646.79\\ 0.6753\\ 0.1789\\ 1\times 10^{-12} \end{array}$	$502.37 \\ 0.6692 \\ 0.1824 \\ 1 \times 10^{-11}$
Worst Fitting Model	F2 γ <sub>1</sub> γ <sub>2</sub> Significant Level	1017.10 0.0289 0.0008 0	$\begin{array}{c} 471.84\\ 0.1302\\ 0.0115\\ 1\times 10^{-4}\end{array}$	$573.85 \\ 0.5264 \\ 0.1760 \\ 1 \times 10^{-4}$	872.23 0.5279 0.1741 $1 \times 10^{-5}$	845.57 0.5338 0.1758 $1  imes 10^{-5}$	$774.72 \\ 0.5990 \\ 0.1807 \\ 1 \times 10^{-5}$	344.85 1.3826 0.9202 $1 \times 10^{-4}$	367.06 1.3973 0.9138 $1 \times 10^{-4}$	$\begin{array}{c} 560.84 \\ 1.5010 \\ 0.9341 \\ 1 \times 10^{-4} \end{array}$	$\begin{array}{c} 432.05 \\ 1.4639 \\ 0.8760 \\ 1 \times 10^{-4} \end{array}$

Table A5. F2 Hypothesis Testing.

Every independent variable considerably influences housing prices, but the degree of influence varies across regions, which is suggestive of significant spatial nonstationarity.

This simple comparison highlights that a better model may require a higher spatial nonstationarity estimation for variables and lower spatial nonstationarity estimation for the intercept.

#### Appendix D. Comparison of GNNWR and GWR Coefficients

To compare the GNNWR and GWR models, we visualize the results of GWR, implemented using ArcGIS Pro with a group of parameters that lead to best fitting performance on all data. The coefficients of different independent variables at each prediction point for the GWR model are shown in Figure A1, and the statistical properties of different coefficients are listed in Table A6. For most of the variables except QAPS and SD, the weight distributions of GNNWR (Figure 6) are similar to those of GWR (Figure A1). The coefficients of GWR typically fluctuate more abruptly than those of the GNNWR model. For several significant variables such as the SD, intercept, and DSS, the GNNWR model provides a smoother distribution compared to the GWR. When the GWR model is used, the clusters and clumps grow, and the undesirable characteristic of the coefficient distribution lying across zero becomes more significant.

Tables A6 and 4 demonstrate that GNNWR outperforms the GWR model in terms of the statistical characteristics. The differences in the standard deviation of variables such as AB, NPS, MF, GR, PR, and NSS are smaller than 0.02, which means that the degrees of fluctuations of these variables are similar. However, the standard deviations of the other variables for GWR are 2–4 times larger than those for GNNWR. This finding highlights that the GNNWR models the spatial heterogeneity in a more accurate and more facile manner and is less prone to overfitting compared with the GWR model.

Moreover, the difference in the distribution of the SD coefficients reflects the limitations of the GWR model. The kernel function of GWR requires priori conditions, such as bisquare or exponential correlation waning patterns, owing to which similar coefficients are distributed in clusters. Notably, the variable SD characterizes the distance from the coastline, which means that the weight distribution of SD might have a higher similarity in the direction parallel to the coastline than in the direction perpendicular to the coastline. Theoretically, the coefficient distribution should appear as a strip parallel to the coastline. The actual results are also consistent with our speculation, with the GWR results resembling clumps, and the GNNWR results resembling strips parallel to the coastline. The difference in the DSS coefficients are attributable to a similar cause.

Therefore, the comparison results demonstrate the high robustness of GNNWR.

Coefficients of Variables	AB	NPS	MF	GR	PR	SD	QAPS	NSS	DSS	Intercept
Mean	-0.281	0.191	0.420	0.125	0.000	-0.278	0.030	0.009	-0.240	0.460
Maximum	0.189	0.872	2.369	0.676	0.194	4.920	0.396	0.560	6.202	1.534
Minimum	-0.926	-0.275	-0.656	-0.073	-0.236	-4.456	-0.240	-0.819	-7.993	-0.731
Std. Dev.	0.193	0.176	0.589	0.104	0.068	1.090	0.094	0.190	1.814	0.356

Table A6. Descriptive Statistics of Coefficients of Variable from GWR









Å





(**d**)





(e)







**Figure A1.** Coefficient Weight Distributions of Variables for GWR model. (a) Intercept; (b) Age of Building (AB); (c) Distance to the Nearest Subway Station (DSS); (d) Green Ratio (GR); (e) Management Fee (MF); (f) Number of Parking Spots (NPS); (g) Number of Subway Stations within 1 km radius (NSS); (h) Plot Ratio (PR); (i) Quality of Available Public Schools (QAPS); (j) Sea Distance (SD).

# References

- 1. Second-hand residential sales price index for 70 large and medium-sized cities in May 2021. China Real Estate 2021, 80. (In Chinese)
- 2. Rosen, S. Hedonic prices and implicit markets: Product differentiation in pure competition. J. Political Econ. 1974, 82, 34–55.
- 3. Butler, R.V. The specification of hedonic indexes for urban housing. *Land Econ.* **1982**, *58*, 96–108.
- 4. Mok, H.M.; Chan, P.P.; Cho, Y.S. A hedonic price model for private properties in Hong Kong. *J. Real Estate Financ. Econ.* **1995**, 10, 37–48.
- 5. Basu, S.; Thibodeau, T.G. Analysis of spatial autocorrelation in house prices. J. Real Estate Financ. Econ. 1998, 17, 61–85.
- 6. Glumac, B.; Herrera-Gomez, M.; Licheron, J. A hedonic urban land price index. Land Use Policy 2019, 81, 802–812.
- 7. Fotheringham, A.S.; Brunsdon, C.; Charlton, M. *Geographically Weighted Regression: The Analysis of Spatially Varying Relationships*; John Wiley & Sons: Hoboken, NJ, USA, 2003.
- 8. Brunsdon, C.; Fotheringham, A.S.; Charlton, M.E. Geographically weighted regression: A method for exploring spatial nonstationarity. *Geogr. Anal.* **1996**, *28*, 281–298.
- 9. Tobler, W.R. A computer movie simulating urban growth in the Detroit region. Econ. Geogr. 1970, 46, 234–240.
- 10. Brunsdon, C.; Fotheringham, A.S.; Charlton, M. Some notes on parametric significance tests for geographically weighted regression. *J. Reg. Sci.* **1999**, *39*, 497–524.
- 11. Tu, W.; Cao, R.; Yue, Y.; Zhou, B.; Li, Q.; Li, Q. Spatial variations in urban public ridership derived from GPS trajectories and smart card data. *J. Transp. Geogr.* **2018**, *69*, 45–57.
- Geng, J.; Cao, K.; Yu, L.; Tang, Y. Geographically Weighted Regression model (GWR) based spatial analysis of house price in Shenzhen. In Proceedings of the 2011 19th International Conference on Geoinformatics, Shanghai, China, 24–26 June 2011; pp. 1–5.
- 13. Zhang, S.; Wang, L.; Lu, F. Exploring housing rent by mixed geographically weighted regression: A Case study in Nanjing. *ISPRS Int. J. Geo-Inf.* **2019**, *8*, 431.
- 14. Lu, B.; Charlton, M.; Harris, P.; Fotheringham, A.S. Geographically weighted regression with a non-Euclidean distance metric: A case study using hedonic house price data. *Int. J. Geogr. Inf. Sci.* **2014**, *28*, 660–681.
- 15. Lu, B.; Charlton, M.; Fotheringhama, A.S. Geographically weighted regression using a non-Euclidean distance metric with a study on London house price data. *Procedia Environ. Sci.* **2011**, *7*, 92–97.
- 16. Limsombunchai, V. House price prediction: Hedonic price model vs. artificial neural network. In Proceedings of the New Zealand Agricultural and Resource Economics Society Conference, Blenheim, New Zealand, 25–26 June 2004; pp. 25–26.
- 17. Selim, H. Determinants of house prices in Turkey: Hedonic regression versus artificial neural network. *Expert Syst. Appl.* 2009, *36*, 2843–2852.
- 18. McCluskey, W.; Davis, P.; Haran, M.; McCord, M.; McIlhatton, D. The potential of artificial neural networks in mass appraisal: The case revisited. *J. Financ. Manag. Prop. Constr.* **2012**, *3*, 274–292.
- 19. Du, Z.; Wang, Z.; Wu, S.; Zhang, F.; Liu, R. Geographically neural network weighted regression for the accurate estimation of spatial non-stationarity. *Int. J. Geogr. Inf. Sci.* **2020**, *34*, 1353–1377.
- Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A.N.; Kaiser, Ł.; Polosukhin, I. Attention is all you need. In Proceedings of the Advances in Neural Information Processing Systems, Long Beach, CA, USA, 4–9 December 2017; Volume 30.
- 21. Wu, S.; Du, Z.; Wang, Y.; Lin, T.; Zhang, F.; Liu, R. Modeling spatially anisotropic nonstationary processes in coastal environments based on a directional geographically neural network weighted regression. *Sci. Total Environ.* **2020**, *709*, 136097.
- Du, Z.; Wu, S.; Wang, Z.; Wang, Y.; Zhang, F.; Liu, R. Estimating Ground-Level PM2.5 Concentrations Across China Using Geographically Neural Network Weighted Regression. J. Geo-Inf. Sci. 2020, 22, 122–135.

- Benessia, A.; Funtowicz, S.; Bradshaw, G.; Ferri, F.; Ráez-Luna, E.F.; Medina, C.P. Hybridizing sustainability: Towards a new praxis for the present human predicament. *Sustain. Sci.* 2012, 7, 75–89.
- Renigier-Biłozor, M.; Źróbek, S.; Walacik, M.; Janowski, A. Hybridization of valuation procedures as a medicine supporting the real estate market and sustainable land use development during the covid-19 pandemic and afterwards. *Land Use Policy* 2020, 99, 105070.
- 25. Z, D.; Z, W.; S, W. GNNWR: An effective method for analyzing and predicting spatial nonstationarity by combining deep neural networks and ordinary least squares. *IEEE Trans. Neural Netw. Learn. Syst.* **2018**, 185–199.
- 26. Fotheringham, A.S.; Charlton, M.; Brunsdon, C. Measuring spatial variations in relationships with geographically weighted regression. In *Recent Developments in Spatial Analysis*; Springer: Berlin/Heidelberg, Germany, 1997; pp. 60–82.
- 27. Srivastava, N.; Hinton, G.; Krizhevsky, A.; Sutskever, I.; Salakhutdinov, R. Dropout: A simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.* **2014**, *15*, 1929–1958.
- He, K.; Zhang, X.; Ren, S.; Sun, J. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In Proceedings of the IEEE International Conference on Computer Vision, Santiago, Chile, 7–13 December 2015; pp. 1026–1034.
- 29. Hagenauer, J.; Helbich, M. A geographically weighted artificial neural network. Int. J. Geogr. Inf. Sci. 2022, 36, 215–235.
- 30. Thamarai, M.; Malarvizhi, S. House Price Prediction Modeling Using Machine Learning. Int. J. Inf. Eng. Electron. Bus. 2020, 12.
- Phan, T.D. Housing price prediction using machine learning algorithms: The case of Melbourne city, Australia. In Proceedings of the 2018 International Conference on Machine Learning and Data Engineering (iCMLDE), Sydney, Australia, 3–7 December 2018; pp. 35–42.
- 32. Peterson, S.; Flanagan, A. Neural network hedonic pricing models in mass real estate appraisal. *J. Real Estate Res.* 2009, 31, 147–164.
- 33. Nghiep, N.; Al, C. Predicting housing value: A comparison of multiple regression analysis and artificial neural networks. *J. Real Estate Res.* **2001**, *22*, 313–336.
- 34. Lin, C.C.; Mohan, S.B. Effectiveness comparison of the residential property mass appraisal methodologies in the USA. *Int. J. Hous. Mark. Anal.* **2011**, *3*, 224–243.
- McGreal, S.; Adair, A.; McBurney, D.; Patterson, D. Neural networks: The prediction of residential values. *J. Prop. Valuat. Invest.* 1998, 1, 57–70.
- Rossini, P. Application of artificial neural networks to the valuation of residential property. In Proceedings of the Third Annual Pacific-Rim Real Estate Society Conference, Palmerston North, New Zealand, 20–22 January 1997.
- 37. Goodman, A.C.; Thibodeau, T.G. Age-related heteroskedasticity in hedonic house price equations. J. Hous. Res. 1995, Vol.6 Issue 1, 25–42.
- Varma, A.; Sarma, A.; Doshi, S.; Nair, R. House price prediction using machine learning and neural networks. In Proceedings of the 2018 Second International Conference on Inventive Communication and Computational Technologies (ICICCT), Coimbatore, India, 20–21 April 2018; pp. 1936–1939.
- 39. Wang, J.; Hu, S.; Zhan, X.; Luo, Q.; Yu, Q.; Liu, Z.; Chen, T.P.; Yin, Y.; Hosaka, S.; Liu, Y. Predicting house price with a memristor-based artificial neural network. *IEEE Access* **2018**, *6*, 16523–16528.
- Leung, Y.; Mei, C.L.; Zhang, W.X. Statistical tests for spatial nonstationarity based on the geographically weighted regression model. *Environ. Plan. A* 2000, 32, 9–32.