

Ayoub Amghar, Maroun Boussif, and Supawadee Phakdee

Professor: Alan Burstein

CS521 Information Structures with Python

5/4/2022

## **Predicting Housing Prices**

### **Introduction**

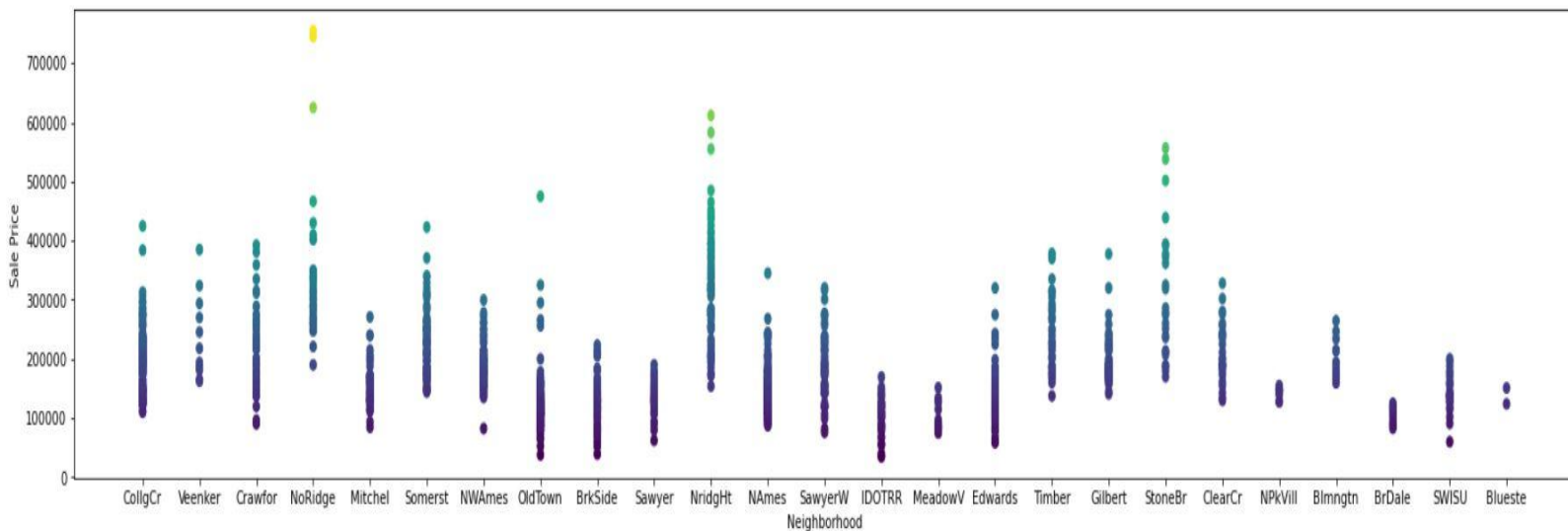
In this project housing prices were chosen to be analyzed. The analysis consisted of implementing a linear regression, then comparing different regression models, followed by forecasting those models, visualizing the results, and finally comparing the results to Zillow. Housing prices were chosen to be predicted because many claim to know how the market works and which way it will go. By using statistics, it can be shown which variables truly matter in the market, and can predict how housing prices will trend in the future.

### **Downloading the Dataset**

### **Model and Data Analysis**

The following scatter plot (figure 1) reflects the relationship between the neighborhood of the house and the sale price variable. From the plot, NoRidge neighborhood has the most expensive houses of any neighborhood in our database, which is over \$700,000. The neighborhood with the lowest home prices is 34,900, which is Idotrr. Because home prices are very much location dependent, there are many neighborhoods with near minimum values. NoRidge neighborhood has the largest range of values. Home prices having a large range of values typically mean a diversity of home types available. Blueste neighborhood had the least

amount of observations. This is mostly likely due to a lack of homes on the market for that particular neighborhood.



**Figure 1: Scatter plot of the neighborhood and housing sale price model**

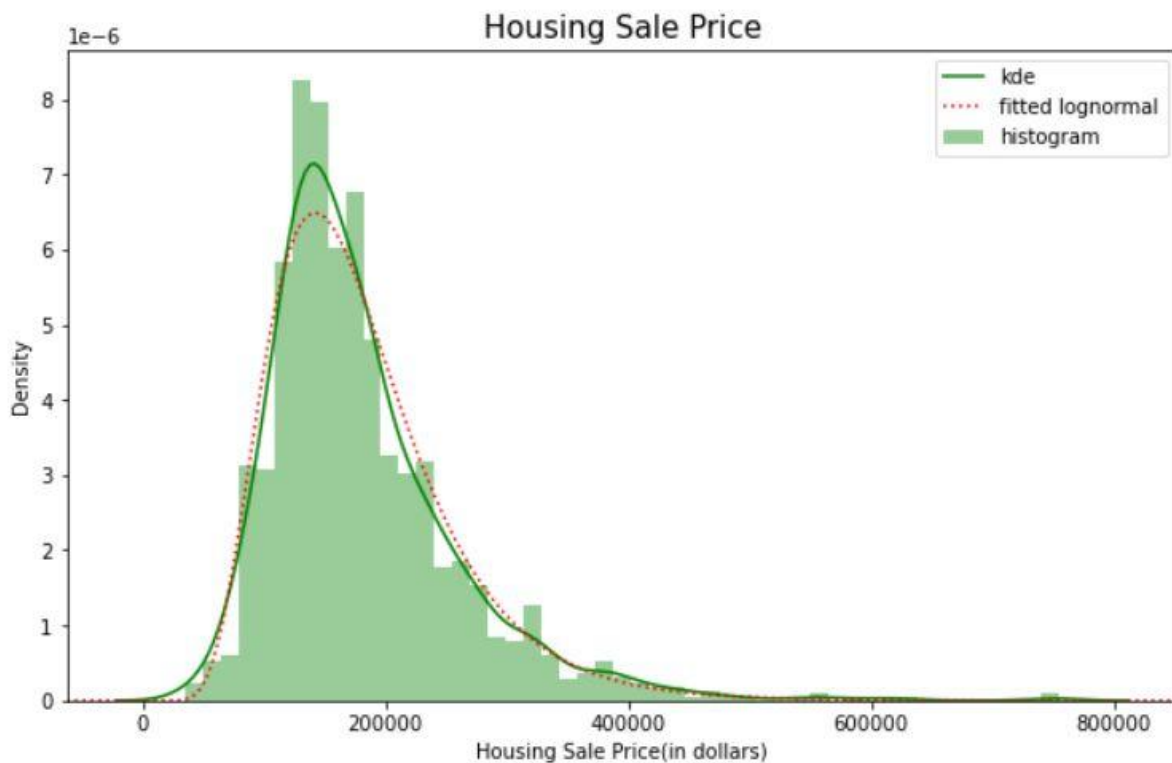
The Histogram represents the sale prices density for houses, as shown in the graph that most of the houses density prices range between 100,000 and 400,000, as the bins describe that in the graph, because the kernel curves follow the flow of the data given in our data.

The peak of the density where we have most houses at the range price between 100,000 and 200,000, as the data distribution were showing in the graph. After the 200,000 we can see that the density started falling down because we have less houses to sell at the range price between 200,000 and 400,000 and the data getting lesse until 600,000.

The blue bins represent the number of houses in the specific price range according to the y axis which present the density.

The green line in the Histogram represents the kernel density estimation, The kernel Density estimation is mathematical by plotting out data to show the changes of the density, the curve

calculated by measuring the distance of all the points of the bins along the distribution. The kernel density calculates the density of the bins points to give us a clear idea about the peak and the lower data concentration of housing prices. The Kernel used to give a clear visualization about the highest peak point of the bins, and also give you a whole idea about the quantity of houses in a specific price range.

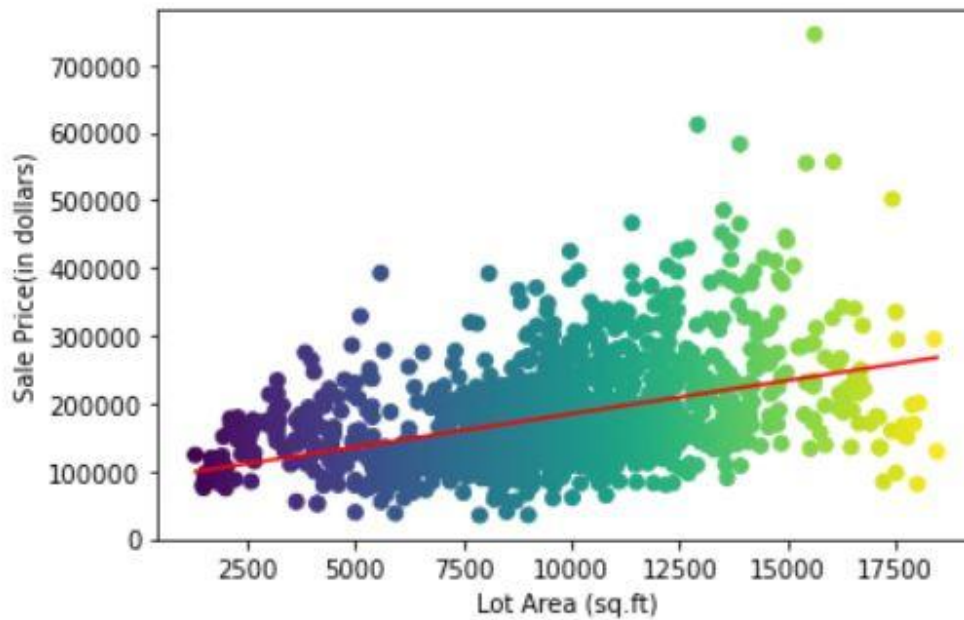


**Figure 2: Histogram and density plot of the housing sale price**

### **Extension 1: Linear Regression**

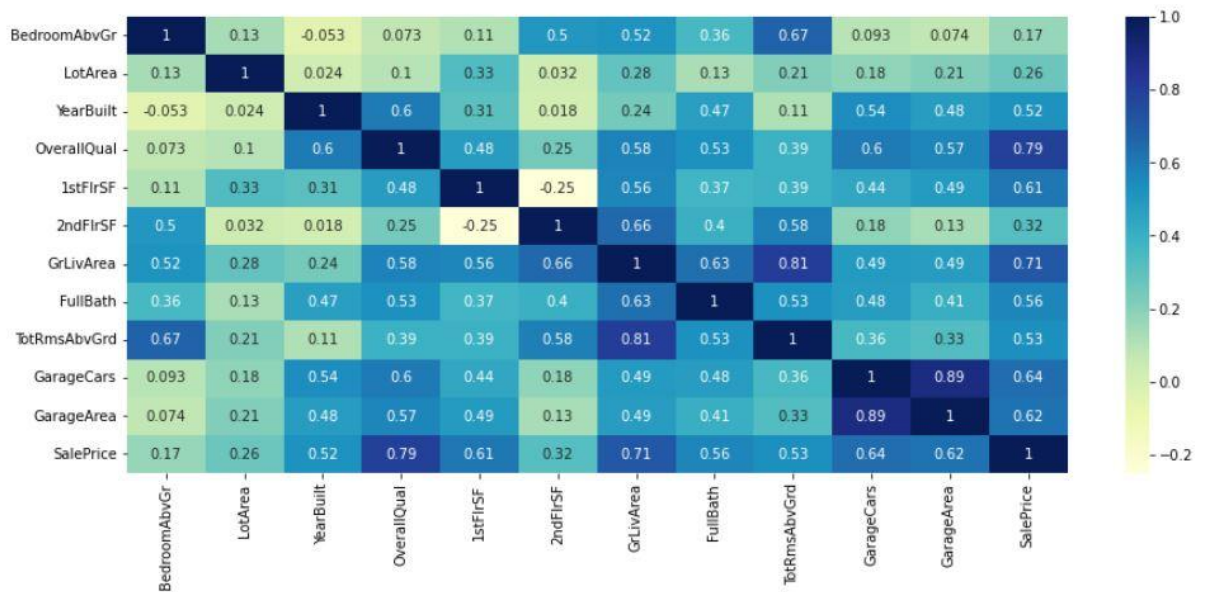
The following linear regression shows the relationship between the sale price of a home and the size of the property. There are a significantly higher number of smaller sized lots or properties, and also a larger number of “low cost” houses. The linear regression shows a strong relationship between the variables for every year, 2006 through 2010. It is obvious that a larger

property will have a higher cost. But as the properties get larger, so does the range of prices for these properties. The majority of the observations are below \$400,000, but once properties start to go over 12,000 sq.ft the variation increases dramatically.



**Figure 3: Linear regression model**

## **Extension 2: Correlation Matrix**

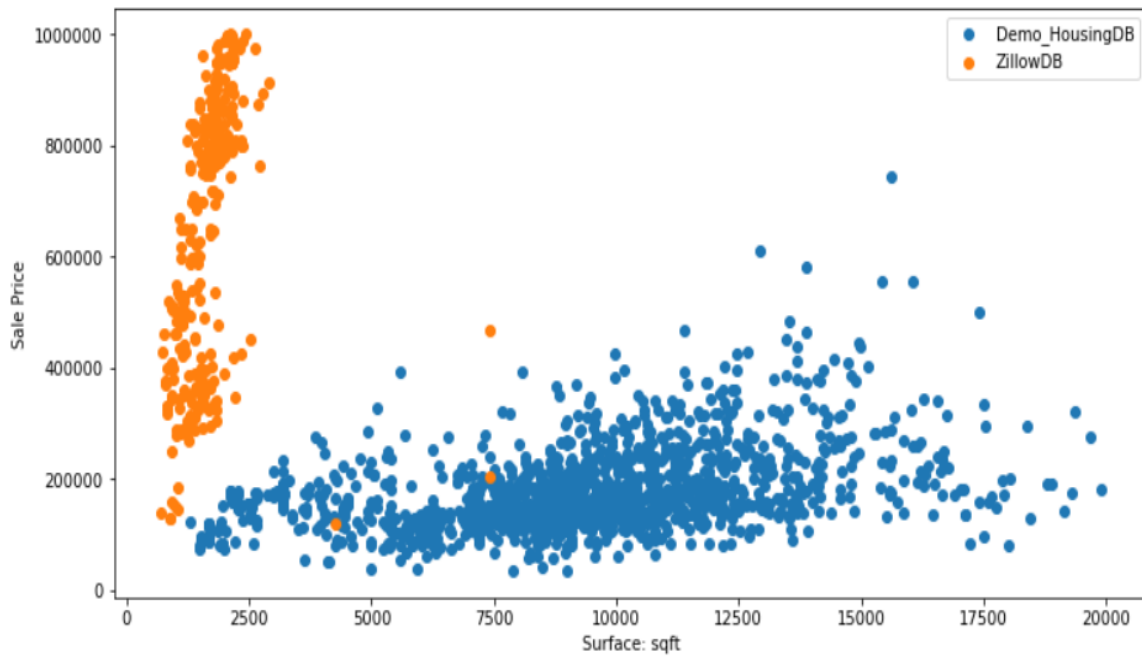


**Figure 4: Correlation**

### Extension 3: Visually and descriptively compare results to Zillow.

The figure 5 draw scatter plot comparing our housing Db of the sale price under 10000000 for the the lot area under 20000 to compare them the zillow DB for the sale price less than 1000000 for the all the lot area less than 20000, the orange plots presents our housing Db and the blue plots presents the zillow DB, as we clear we can see that the x axis presents the surface Sqft and y axis presents the sale price, there is clear correlation between the two data, that more the plot above the ground level, the higher the price of lot is. From the graph we can determine that there are more areas for our demo housing DB at the range sqft was specified than the zillow DB. Also we see that the prices for demo housing DB are cheaper than the zillow even if they get bigger in sqft, which was presented with blue color. On the other hand, the scatter plots for the lot area zillow DB were smaller in sqft size but more expensive which reflected the

location of that lot( city, near transportation...). Also most of lot of areas of the zillow db were between the range of 0 and 2500 sqft, at the opposite in our data we see that lot areas were distributed almost all on the y axis.



**Figure**

## **Conclusion**