John Kundycki

# Predicting Real Estate Prices

King County Housing Dataset

# Data Science Life Cycle: OSEMN

**OSEMN:**

Obtain - Gather Data from relevant resources

Scrub - Clean data to formats that machine understands

Explore - Find significant patterns and trends using statistical methods

Model - Construct models to predict and forecast

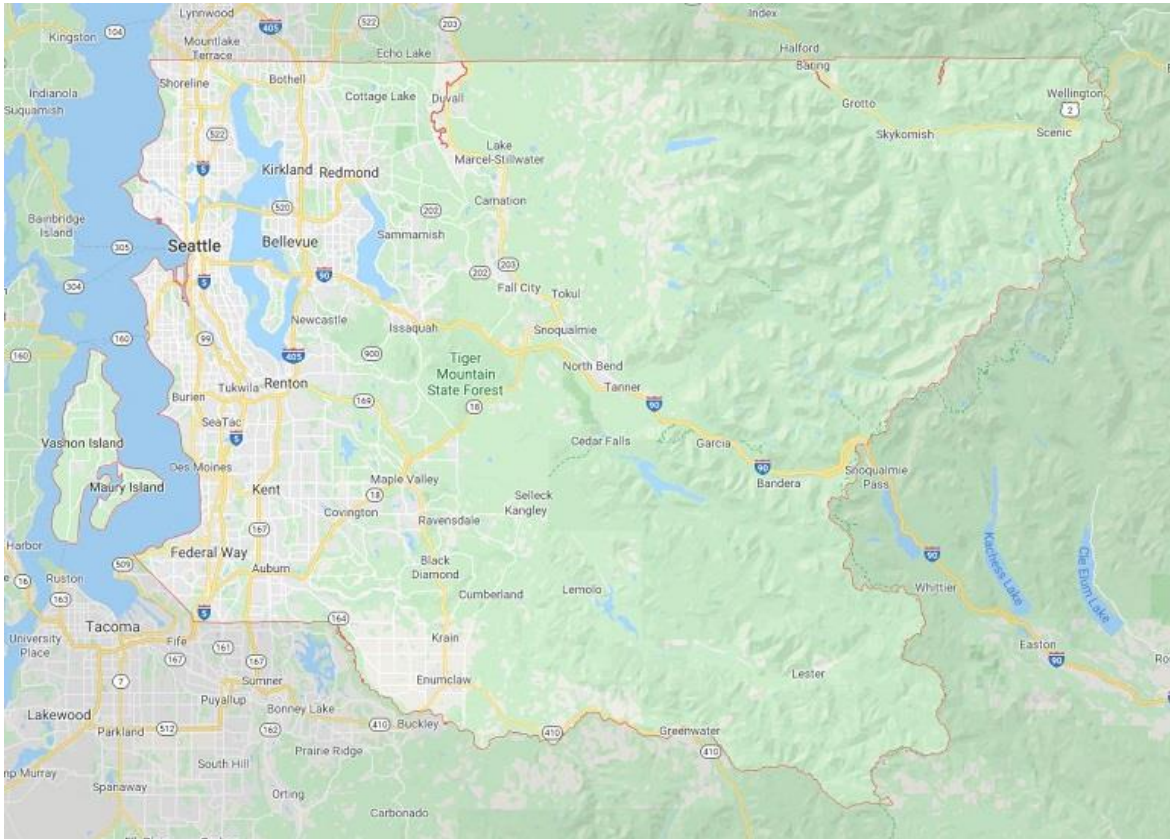Interpret - Put the results into good use

# OBTAIN

Provided data on housing prices and different features for houses in King County, Seattle, WA

```
#First look at the data
data.head()
```

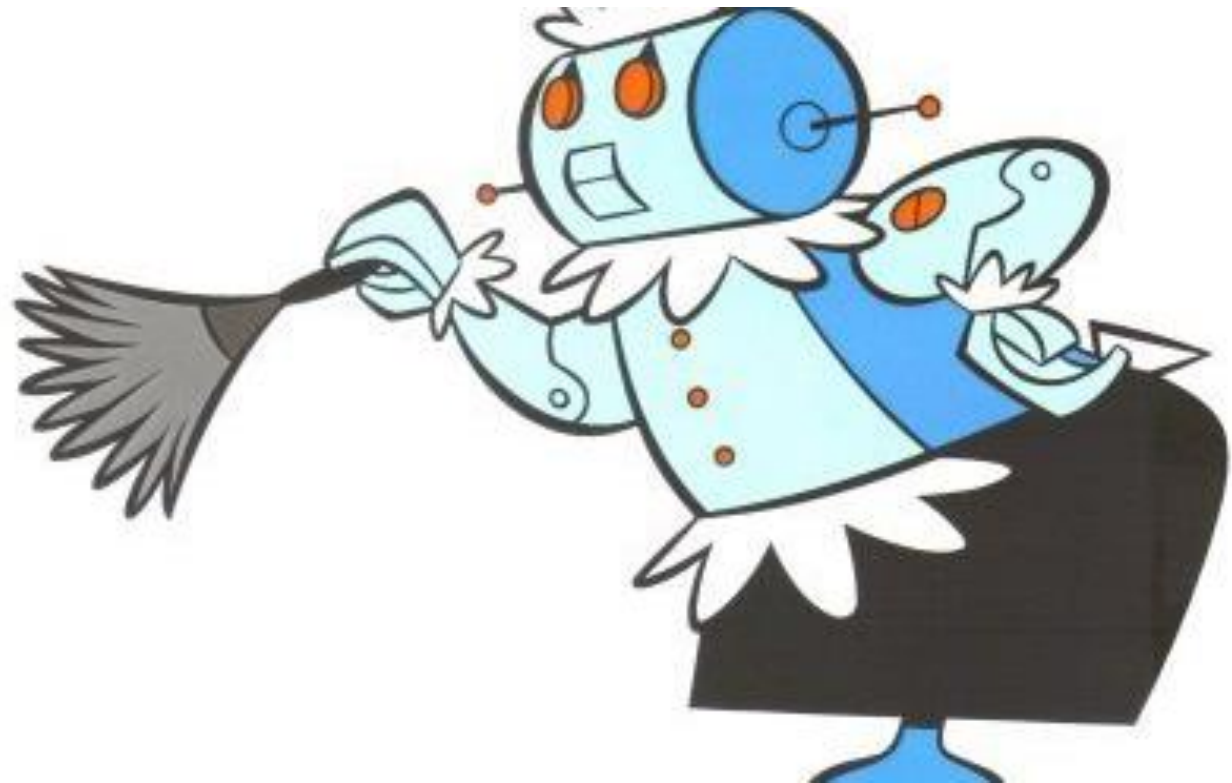| | id | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | view | condition | grade | sqft_above | sqft_basement | yr_bui |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 7129300520 | 10/13/2014 | 221900.0 | 3 | 1.00 | 1180 | 5650 | 1.0 | NaN | 0.0 | 3 | 7 | 1180 | 0.0 | 195 |
| 1 | 6414100192 | 12/9/2014 | 538000.0 | 3 | 2.25 | 2570 | 7242 | 2.0 | 0.0 | 0.0 | 3 | 7 | 2170 | 400.0 | 195 |
| 2 | 5631500400 | 2/25/2015 | 180000.0 | 2 | 1.00 | 770 | 10000 | 1.0 | 0.0 | 0.0 | 3 | 6 | 770 | 0.0 | 193 |
| 3 | 2487200875 | 12/9/2014 | 604000.0 | 4 | 3.00 | 1960 | 5000 | 1.0 | 0.0 | 0.0 | 5 | 7 | 1050 | 910.0 | 196 |
| 4 | 1954400510 | 2/18/2015 | 510000.0 | 3 | 2.00 | 1680 | 8080 | 1.0 | 0.0 | 0.0 | 3 | 8 | 1680 | 0.0 | 198 |

# OBTAIN

More information about the dataset



## Column Names and descriptions for King County Data Set

- **id** - unique identified for a house
- **dateDate** - house was sold
- **pricePrice** - is prediction target
- **bedroomsNumber** - of Bedrooms/House
- **bathroomsNumber** - of bathrooms/bedrooms
- **sqft_livingsquare** - footage of the home
- **sqft_lotsquare** - footage of the lot
- **floorsTotal** - floors (levels) in house
- **waterfront** - House which has a view to a waterfront
- **view** - Has been viewed
- **condition** - How good the condition is ( Overall )
- **grade** - overall grade given to the housing unit, based on King County grading system
- **sqft_above** - square footage of house apart from basement
- **sqft_basement** - square footage of the basement
- **yr_built** - Built Year
- **yr_renovated** - Year when house was renovated
- **zipcode** - zip
- **lat** - Latitude coordinate
- **long** - Longitude coordinate
- **sqft_living15** - The square footage of interior housing living space for the nearest 15 neighbors
- **sqft_lot15** - The square footage of the land lots of the nearest 15 neighbors

# SCRUB

• Cleaning Data means fixing problems with the data that are preventing the information from being interpreted by the user and the machine.

# SCRUB



Is this useable data?

# EXPLORE

ASK MEANINGFUL
QUESTIONS ABOUT THE DATA

MULTICOLINEARITY

SEE HOW FEATURES RELATE
TO PRICE

# EXPLORE

- Question 1: Does renovation have a noticeable effect on price?

- Question 2: Is there a difference in price between a house built in a given time period versus a house renovated in that same time period?

- Question 3: Is there a difference in price based on geographical location in King's County? And if so, why?
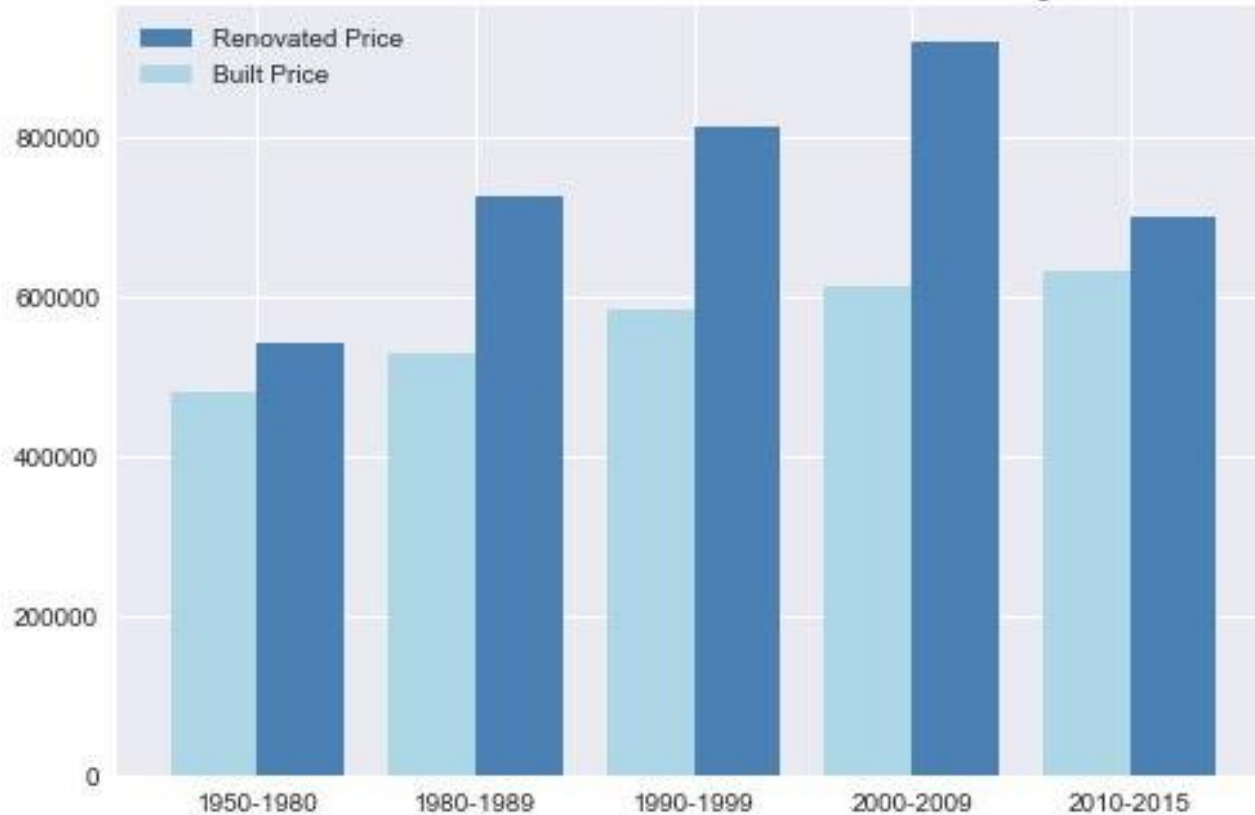
# Question 1: Does renovation have a noticeable effect on price?



## Renovation vs Price

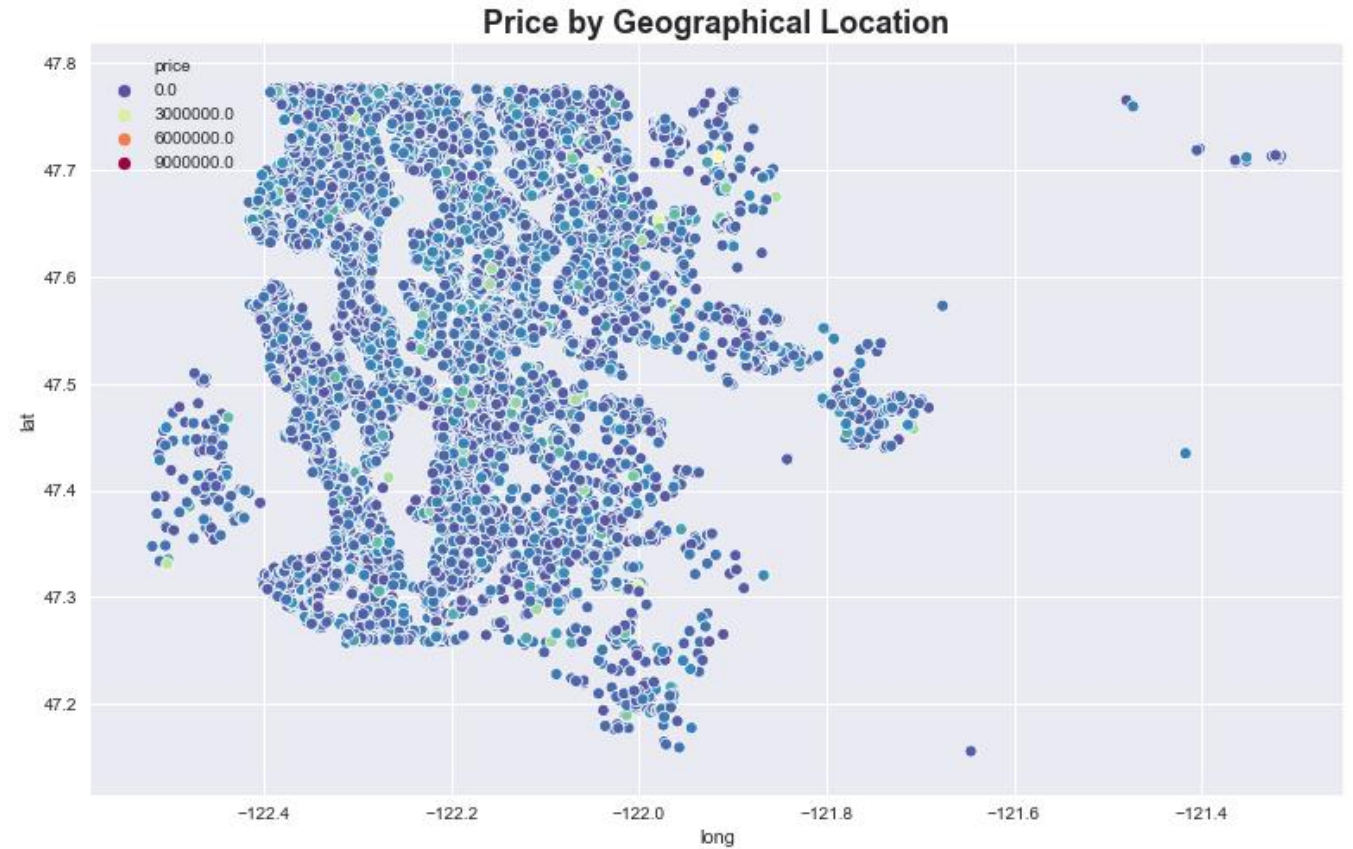- Renovating a house in King County has a mean price increase of $237,423 or 144%

# Question 2: Is there a difference in price between a house built in a given time period versus a house renovated in that same time period?



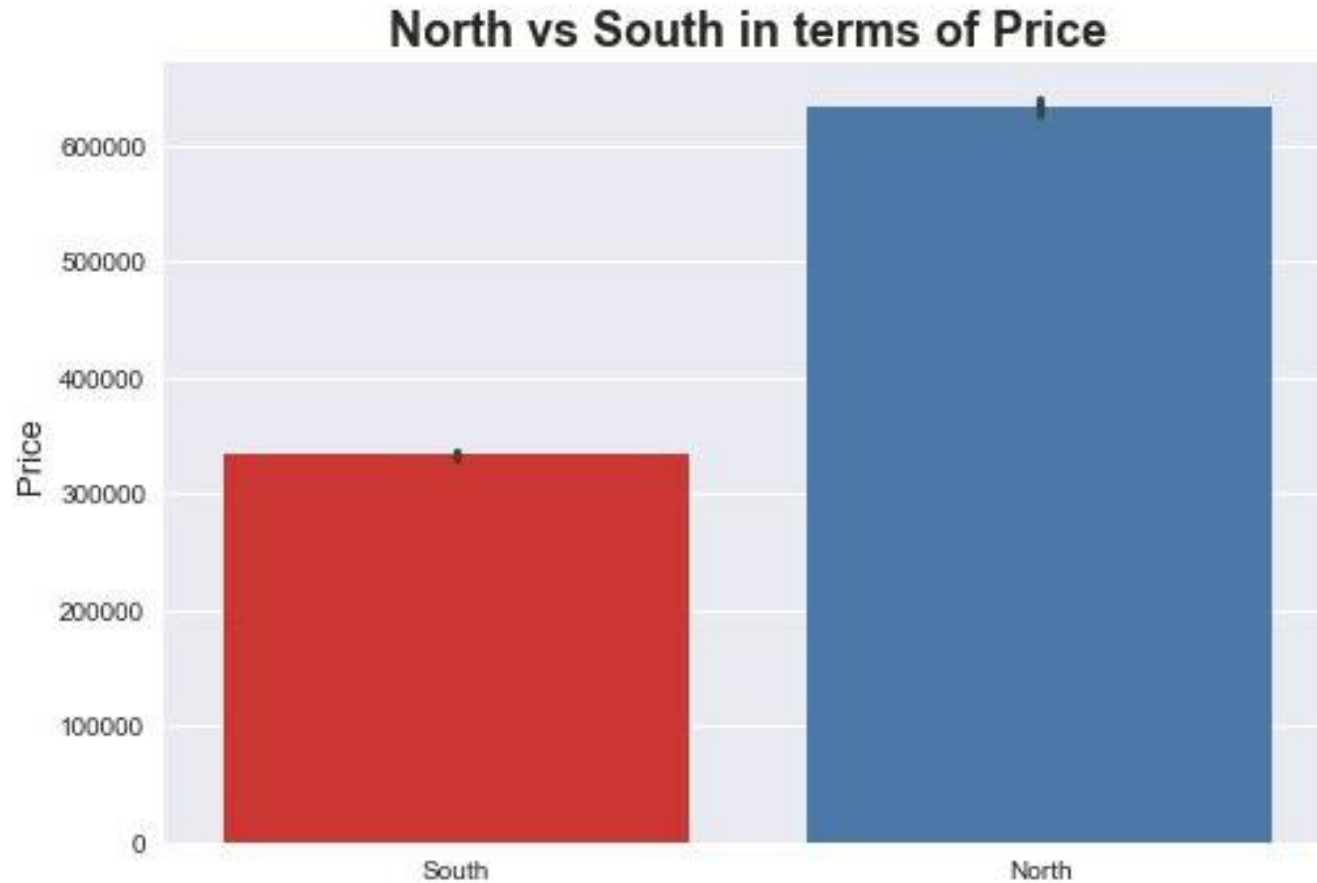Price of a Renovated House vs Built House by Time Period

- Significantly larger price for renovated houses
- Gap between renovated and built houses seems to be growing
- 2010-2015 might need more data points to observe this trend

# Question 3: Is there a difference in price based on geographical location in King's County?

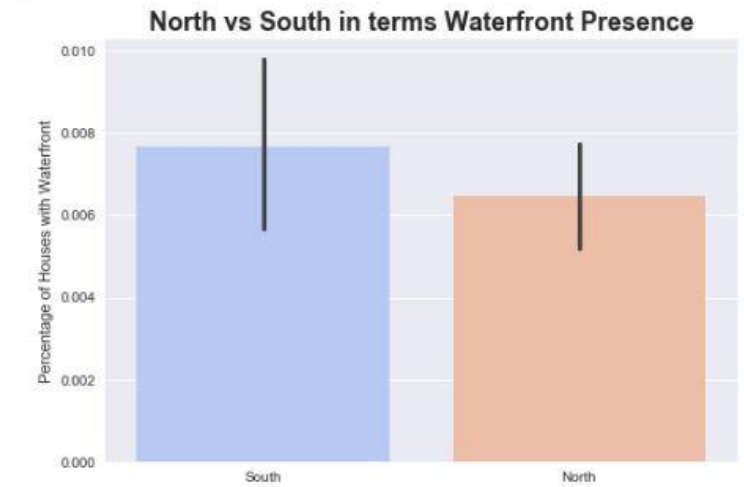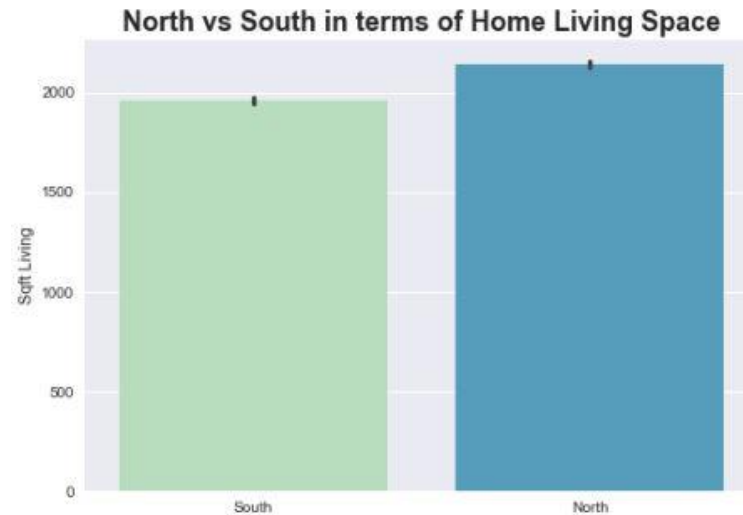# Question 3: Is there a difference in price based on geographical location in King's County?



**North vs South in terms of Price**

- Houses in the Northern half of Seattle are worth 190% those of in the Southern half – Almost double the price

# Question 3 (cont.): Why is there a difference in price between Northern and Southern homes in Seattle?



North vs South in terms of House Grade



North vs South in terms of Home Living Space



North vs South in terms Waterfront Presence

North vs South in terms Waterfront Presence



Price by Geographical Location

# Multicolinearity

# Relationship between house features and Price



Pair Plots of Relationships between Dependent Variables and Price

# MODEL

LINEARITY: The relationship between the independent and dependent variables needs to be linear

NORMALITY: The linear regression analysis requires all variables to be normal.

MULTICOLINEARITY: There is little to no multicolinearity in the data.

HOMOSCEDASTICITY: Residuals must be equal across the regression line.

# MODEL


Price — boxplot of Price in Dollars Amount vs Value Counts

In the linear relationships, there was a skewness due to outliers high in price.

Cut outliers at the upper whisker – $1,120,000

# Linear Relationships before and after removal of Price Outliers



Pair Plots of Relationships between Dependent Variables and Price

Pair Plots of Relationships between Dependent Variables and Price

# Multicolinearity: Round 2

OLS Regression Results

| | | | | |
|---|---|---|---|---|
| Dep. Variable: | price | R-squared (uncentered): | | 0.919 |
| Model: | OLS | Adj. R-squared (uncentered): | | 0.919 |
| Method: | Least Squares | F-statistic: | | 4.619e+04 |
| Date: | Sun, 16 Feb 2020 | Prob (F-statistic): | | 0.00 |
| Time: | 10:44:39 | Log-Likelihood: | | -2.7004e+05 |
| No. Observations: | 20267 | AIC: | | 5.401e+05 |
| Df Residuals: | 20262 | BIC: | | 5.401e+05 |
| Df Model: | 5 | | | |
| Covariance Type: | nonrobust | | | |

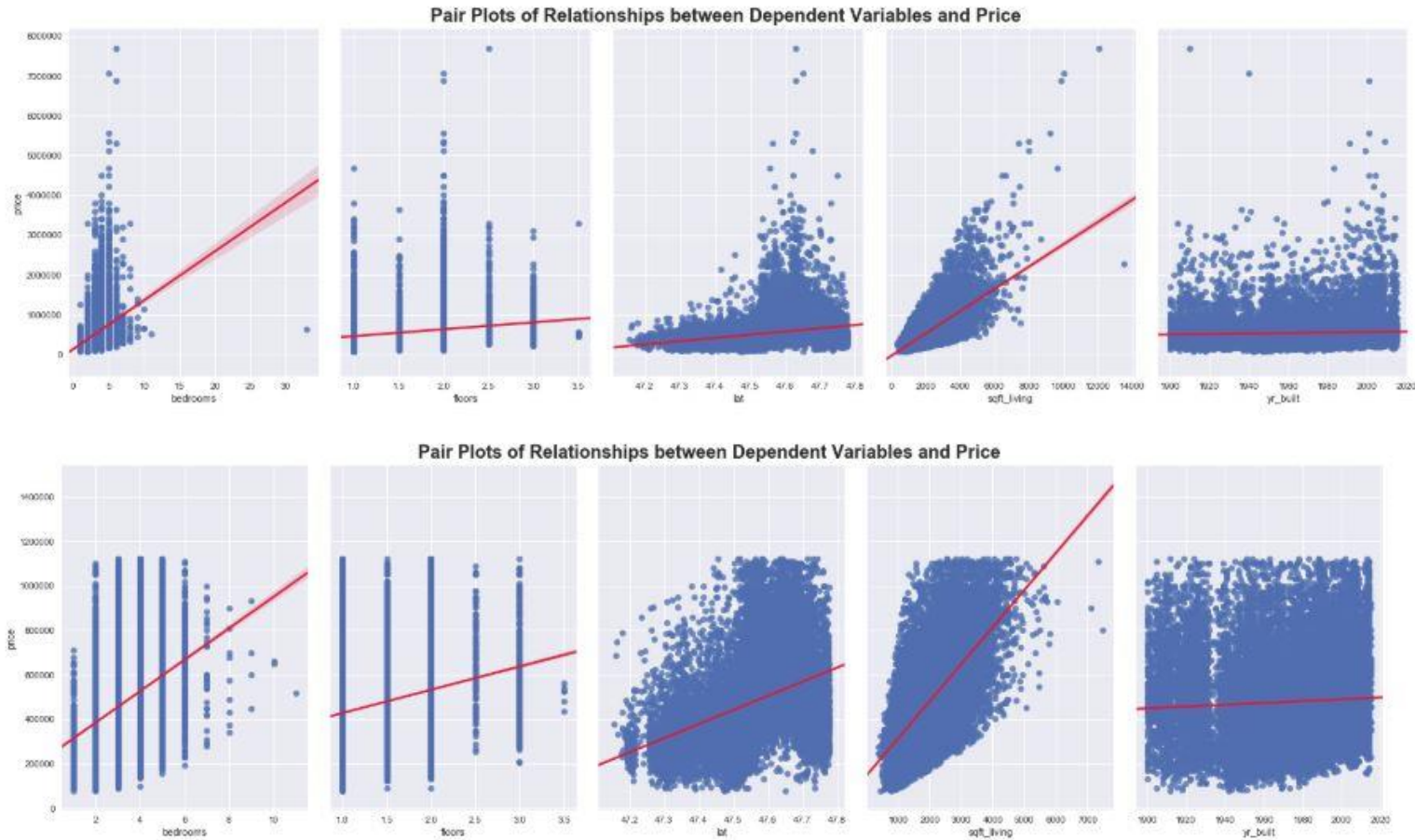| | coef | std err | t | P>\|t\| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| floors | 8e+04 | 2266.511 | 35.297 | 0.000 | 7.56e+04 | 8.44e+04 |
| bedrooms | -2.855e+04 | 1477.412 | -19.327 | 0.000 | -3.14e+04 | -2.57e+04 |
| lat | 1.008e+05 | 1606.914 | 62.701 | 0.000 | 9.76e+04 | 1.04e+05 |
| sqft_living | 198.1701 | 1.798 | 110.220 | 0.000 | 194.646 | 201.694 |
| yr_built | -2399.0821 | 39.682 | -60.457 | 0.000 | -2476.863 | -2321.301 |

| | | | |
|---|---|---|---|
| Omnibus: | 863.525 | Durbin-Watson: | 0.968 |
| Prob(Omnibus): | 0.000 | Jarque-Bera (JB): | 1104.233 |
| Skew: | 0.452 | Prob(JB): | 1.65e-240 |
| Kurtosis: | 3.701 | Cond. No. | 6.58e+03 |

# Multicolinearity: Round 2

OLS Regression Results

| | | | | |
|---|---|---|---|---|
| Dep. Variable: | price | R-squared (uncentered): | | 0.903 |
| Model: | OLS | Adj. R-squared (uncentered): | | 0.903 |
| Method: | Least Squares | F-statistic: | | 9.423e+04 |
| Date: | Sat, 15 Feb 2020 | Prob (F-statistic): | | 0.00 |
| Time: | 21:44:58 | Log-Likelihood: | | -2.7192e+05 |
| No. Observations: | 20267 | AIC: | | 5.438e+05 |
| Df Residuals: | 20265 | BIC: | | 5.439e+05 |
| Df Model: | 2 | | | |
| Covariance Type: | nonrobust | | | |

| | coef | std err | t | P>|t| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| lat | 3179.6410 | 65.790 | 48.330 | 0.000 | 3050.687 | 3308.595 |
| sqft_living | 165.4083 | 1.473 | 112.293 | 0.000 | 162.521 | 168.296 |

| | | | |
|---|---|---|---|
| Omnibus: | 924.899 | Durbin-Watson: | 0.743 |
| Prob(Omnibus): | 0.000 | Jarque-Bera (JB): | 1055.063 |
| Skew: | 0.548 | Prob(JB): | 7.87e-230 |
| Kurtosis: | 3.216 | Cond. No. | 123. |

# MODEL

LINEARITY: THE RELATIONSHIP BETWEEN THE INDEPENDENT AND DEPENDENT VARIABLES NEEDS TO BE LINEAR

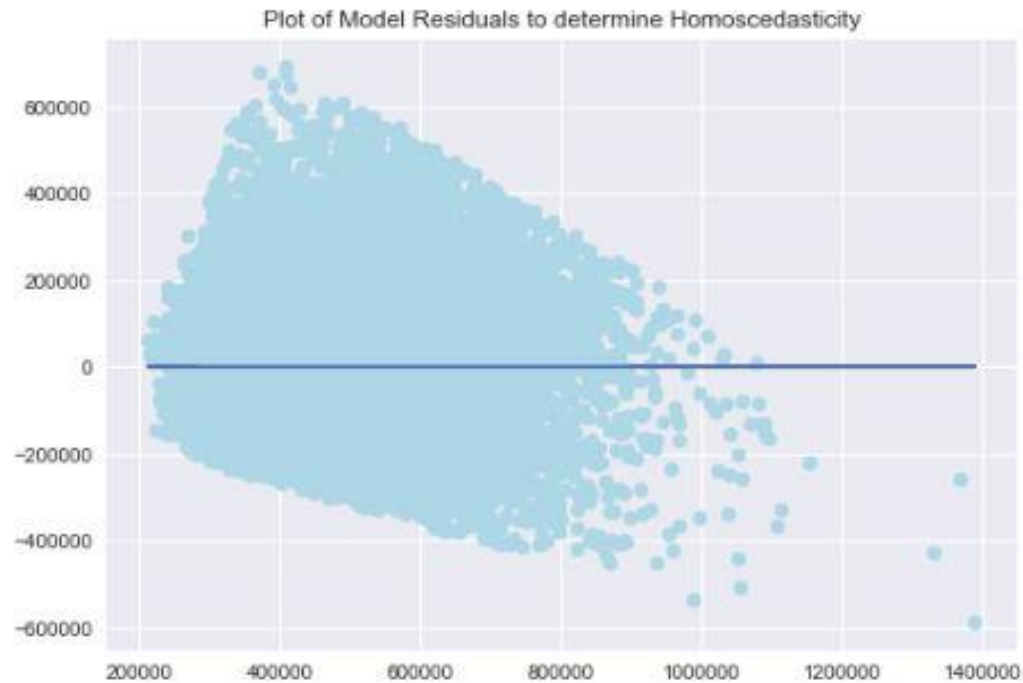NORMALITY: THE LINEAR REGRESSION ANALYSIS REQUIRES ALL VARIABLES TO BE NORMAL.

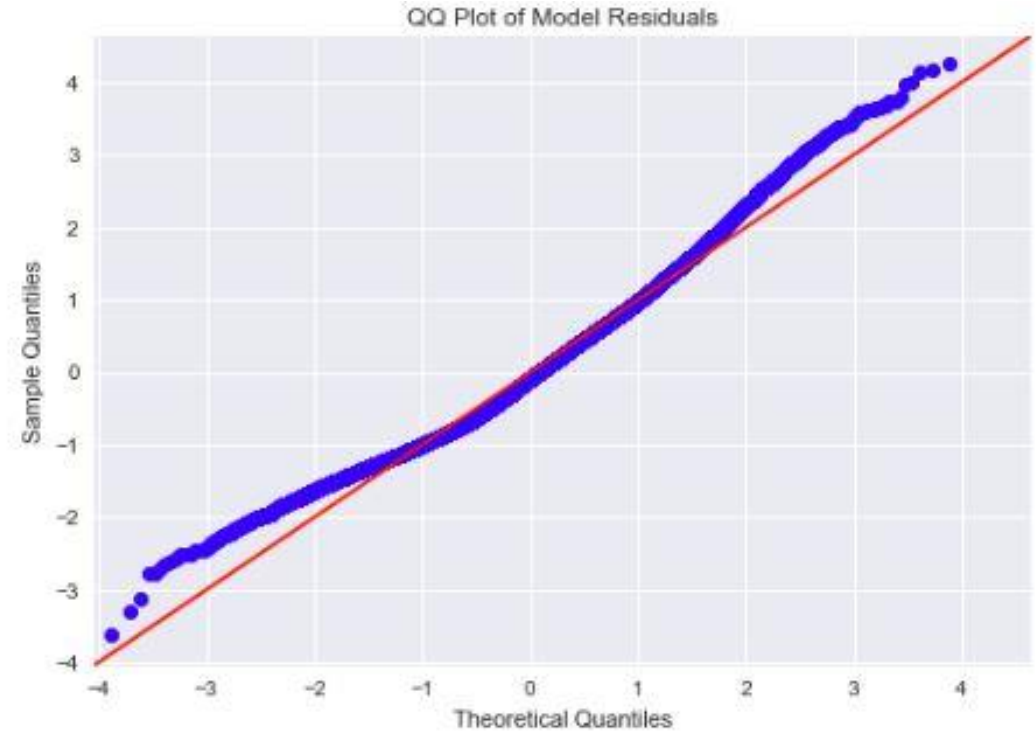MULTICOLINEARITY: THERE IS LITTLE TO NO MULTICOLINEARITY IN THE DATA.

HOMOSCEDASTICITY: RESIDUALS MUST BE EQUAL ACROSS THE REGRESSION LINE.

# Assumptions that can only be tested after the model has been created

**Homoscedasticity**                    **Normality**

# INTERPRET

- **Scores:**

- R-squared: The proportion of the variance in the dependent variable that is predictable from the indenpendent variables. Approximately 90% of the variance in the price of a house can be predicted from the model.

- F-Statistic: The F-statistic is the ratio of the mean regression sum of squares divided by the mean error of sum squares. With a probability of 0, the null hypothesis that the fit of the intercept only model and my model are equal is rejected.

- Coefficients: The coefficients show the relationship between the independent variable and the dependent variable.

- p-Values: The p-values are all 0, meaning the indepedent variables are statistically signficant to the model.

- Skew: The skewness is positive which implies a skewness towards the right. With a value between -0.5 and 0.5, the data is pretty symmetrical.

- Kurtosis: The Kurtosis score is about 3, implying that the data fits a normal distribution.

- Jarque-Bera: The Jarque-Bera test is a goodness-of-fit test. Because the Jarque-Bera score is large, this indicates that errors are not normally distributed.

- Condition Number: Measures multicollinearity, the number is large which implies multicollinearity.

## OLS Regression Results

| Dep. Variable: | price | R-squared (uncentered): | 0.903 |
|---|---|---|---|
| Model: | OLS | Adj. R-squared (uncentered): | 0.903 |
| Method: | Least Squares | F-statistic: | 9.423e+04 |
| Date: | Sat, 15 Feb 2020 | Prob (F-statistic): | 0.00 |
| Time: | 21:44:58 | Log-Likelihood: | -2.7192e+05 |
| No. Observations: | 20267 | AIC: | 5.438e+05 |
| Df Residuals: | 20265 | BIC: | 5.439e+05 |
| Df Model: | 2 | | |
| Covariance Type: | nonrobust | | |

| | coef | std err | t | P>|t| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| lat | 3179.6410 | 65.790 | 48.330 | 0.000 | 3050.687 | 3308.595 |
| sqft_living | 165.4083 | 1.473 | 112.293 | 0.000 | 162.521 | 168.296 |

| | | | |
|---|---|---|---|
| Omnibus: | 924.899 | Durbin-Watson: | 0.743 |
| Prob(Omnibus): | 0.000 | Jarque-Bera (JB): | 1055.063 |
| Skew: | 0.548 | Prob(JB): | 7.87e-230 |
| Kurtosis: | 3.216 | Cond. No. | 123. |

# Questions

John Kundycki

# Predicting Real Estate Prices

King County Housing Dataset