

Linear Regression Analysis for the Kings County's (Seattle, WA) House Market.

by Y. Kostrov

Contents

- 1 Overview
- 2 Business Problem
- 3 Data Description
- 4 My Python Package
- 5 Modeling
- 6 Linear Model
 - Building the Linear Model
 - Checking Statistical Hypotheses
- 7 Multiple Linear Model
 - Explanation of the Model
 - Model
 - Explanation of the Model
- 8 Conclusions
 - Data Modeling
 - Modeling
 - Ways to improve the analysis

Overview

The purpose of this project is to analyze a data set containing data about houses sold in Kings County (Seattle, WA).

Overview

During the analysis:

Overview

During the analysis:

- 1 I will perform necessary data wrangling first.

Overview

During the analysis:

- 1 I will perform necessary data wrangling first.
- 2 I will build a Linear Regression Model with one explanatory variable.:

Overview

During the analysis:

- ① I will perform necessary data wrangling first.
- ② I will build a Linear Regression Model with one explanatory variable.:
 - I will check statistical assumptions for the linear regression model
 - I will explain the model, including intercept, coefficient for the explanatory variable, R^2 , and ANOVA

0

- 1 I will perform necessary data wrangling first.
- 2 I will build a Linear Regression Model with one explanatory variable.:
 - I will check statistical assumptions for the linear regression model
 - I will explain the model, including intercept, coefficient for the explanatory variable, R^2 , and ANOVA
- 3 I will build a Multiple Linear Regression Model with many explanatory variables.

Overview

During the analysis:

- ① I will perform necessary data wrangling first.
- ② I will build a Linear Regression Model with one explanatory variable.:
 - I will check statistical assumptions for the linear regression model
 - I will explain the model, including intercept, coefficient for the explanatory variable, R^2 , and ANOVA
- ③ I will build a Multiple Linear Regression Model with many explanatory variables.
 - I will check statistical assumptions for the multiple linear regression model
 - I will explain the model, including intercept, coefficients for the explanatory variables, R^2 , and ANOVA

Business Problem

- The fair price of the house is a hard quantity to assess.

Business Problem

- The fair price of the house is a hard quantity to assess.
- Both sellers and buyers would like to know the best price for the house.

Business Problem

- The fair price of the house is a hard quantity to assess.
- Both sellers and buyers would like to know the best price for the house.
- Which features of the property would be the best predictors of the value?

Business Problem

- The fair price of the house is a hard quantity to assess.
- Both sellers and buyers would like to know the best price for the house.
- Which features of the property would be the best predictors of the value?
- I will build a regression model that helps predict the value of the house.

Business Problem

- The fair price of the house is a hard quantity to assess.
- Both sellers and buyers would like to know the best price for the house.
- Which features of the property would be the best predictors of the value?
- I will build a regression model that helps predict the value of the house.
- I will, also, check the necessary statistical assumptions for the regression model and explain the model's parameters.

- The file called “kc_house_data.csv” in the data folder of the project holds the data for this project.

Data Description

- The file called “kc_house_data.csv” in the data folder of the project holds the data for this project.
- This project will use this data about Kings County's(Seattle, WA) housing market to create Linear Regression Model.

Data Description

- The file called “kc_house_data.csv” in the data folder of the project holds the data for this project.
- This project will use this data about Kings County’s(Seattle, WA) housing market to create Linear Regression Model.
- The data file contains numerous columns with information about properties sold such as price, size of the living area, size of the basement, number of bedrooms, etc.

My Python Package

- While working on this project, I have created my own Python package with helping functions.

My Python Package

- While working on this project, I have created my own Python package with helping functions.
- The most important function in this package is “evaluate_model.py” (in the “src” folder). This function:
 - creates the model from the data frame.

My Python Package

- While working on this project, I have created my own Python package with helping functions.
- The most important function in this package is “evaluate_model.py” (in the “src” folder). This function:
 - creates the model from the data frame.
 - prints out the model summary of Linear Regression.

My Python Package

- While working on this project, I have created my own Python package with helping functions.
- The most important function in this package is “evaluate_model.py” (in the “src” folder). This function:
 - creates the model from the data frame.
 - prints out the model summary of Linear Regression.
 - performs the checks for the statistical assumptions of the Linear Regression.

My Python Package

- While working on this project, I have created my own Python package with helping functions.
- The most important function in this package is “evaluate_model.py” (in the “src” folder). This function:
 - creates the model from the data frame.
 - prints out the model summary of Linear Regression.
 - performs the checks for the statistical assumptions of the Linear Regression.
 - performs a lot of different visualizations.

Modeling

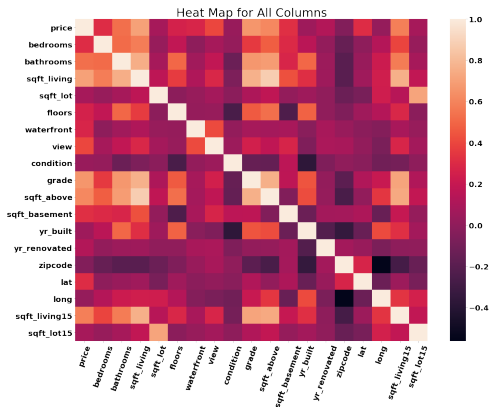
- My first goal was to create linear regression model with one independent variable.

Modeling

- My first goal was to create linear regression model with one independent variable.
- I created the correlation matrix and heat map for visualization purpose.

Modeling

- My first goal was to create linear regression model with one independent variable.
- I created the correlation matrix and heat map for visualization purpose.



- “sqft_living” has the highest correlation of 0.71 with the “price”.

- “sqft_living” has the highest correlation of 0.71 with the “price”.
- I build a regression model for the “price” to be predicted by “sqft_living”.

Downloaded from <http://ajph.org/> on November 10, 2014

- “sqft_living” has the highest correlation of 0.71 with the “price”.
- I build a regression model for the “price” to be predicted by “sqft_living”.
- The model is $\ln(\text{price}) = 11.9524 + 0.0029 \cdot \text{sqft_living}^{0.78}$

Checking Statistical Hypotheses: Linearity

The Null Hypothesis:

The model is linearly predicted by the feature,

Checking Statistical Hypotheses: Linearity

The Null Hypothesis:

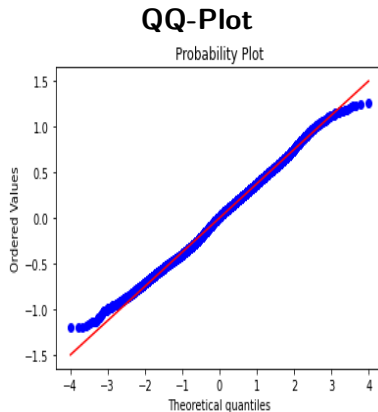
The model is linearly predicted by the feature,

The Alternative Hypothesis:

The model is not linearly predicted by the feature.

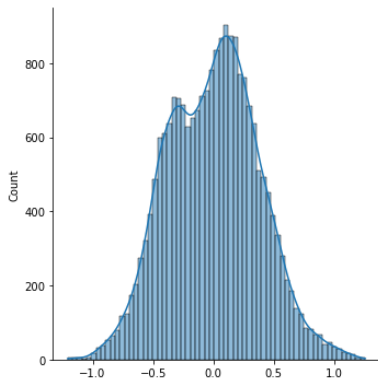
Checking Statistical Hypotheses: Normality 1

To check Normality, I used the following checks:



Checking Statistical Hypotheses: Normality 2

DISTRIBUTIONS PLOT OF RESIDUALS



Checking Statistical Hypotheses: Normality 3

I, also, used D'Agostino Test for Normality:

[illegible][illegible][illegible][illegible]

Checking Statistical Hypotheses: Normality 3

I, also, used D'Agostino Test for Normality:

The Null Hypothesis:

The Residuals are normally distributed

The Alternative Hypothesis:

The Residuals are not normally distributed.

Checking Statistical Hypotheses: Normality 3

I, also, used D'Agostino Test for Normality:

The Null Hypothesis:

The Residuals are normally distributed

The Alternative Hypothesis:

The Residuals are not normally distributed.

- Our p-value for this model is $p = 0.000 < 0.05 = \alpha$.
- We have enough evidence to reject the Null Hypothesis
- We conclude that D'Agostino Test tells us, that residuals are not normally distributed.

Checking Statistical Hypotheses: Normality 3

Conclusion for Normality:

- Based on QQ-Plot, Distributions Plot, and D'Agostino Test, I conclude that the Distribution of Errors is not far away from Normal.

Checking Statistical Hypotheses: Normality 3

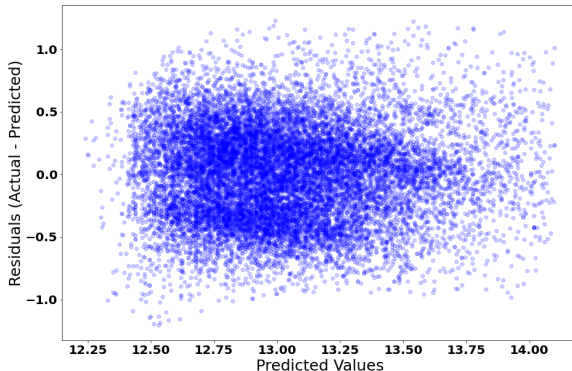
Conclusion for Normality:

- Based on QQ-Plot, Distributions Plot, and D'Agostino Test, I conclude that the Distribution of Errors is not far away from Normal.
- Also, since we have a lot of observations Normality Assumption doesn't play a critical role, since Central Limit Theorem will apply in this case.

Constant Error Variance 1

To check if heteroscedasticity is present in the model,

- I will use Residual-vs-Predicted values plot and Breusch-Pagan test.
- I look at at the Residual-vs-Predicted values plot first.



Constant Error Variance 2

I will use Breusch-Pagan Test:

Constant Error Variance 2

I will use Breusch-Pagan Test:

The Null Hypothesis:

Homoscedasticity is present.

The Alternative Hypothesis:

Constant Error Variance 2

I will use Breusch-Pagan Test:

The Null Hypothesis:

Homoscedasticity is present.

The Alternative Hypothesis:

Homoscedasticity is not present (i.e. heteroscedasticity exists).

- Our p-value for this model is $p = 0.12609 \geq 0.05 = \alpha$.
- Thus, we don't have enough evidence to reject the Null Hypothesis and
- We conclude from Breusch-Pagan Test , that we have don't have heteroscedasticity.

Constant Error Variance 3

Conclusion:

From the Residual-vs-Predicted values plot and Breusch-Pagan Test, I conclude that we don't have Heteroscedasticity in our model.

Overall Conclusion:

I conclude that our model satisfies statistical assumptions for the regression model.

Intercept and slope 1

Our model is

$$\ln(\text{price}) = 11.9524 + 0.0029 \cdot \text{sqft_living}^{0.78}$$

- The model has the sample intercept of 11.9524 and the slope of 0.0029.
- To interpret the slope, we have to transform \hat{x} and \hat{y} towards original *sqft_living* and *price*.
- Let $x = \text{sqft_living}$ and $y = \text{price}$ for the derivation.
- We have $\hat{x} = x^{0.78}$ and $\hat{y} = \ln(y) \implies y = e^{\hat{y}}$.

Intercept and slope 2

Suppose the original sqft_living is x_1 and it moved up to x_2 , then we have the following :

$$y_1 = e^{\hat{y}_1} = e^{11.9524 + 0.0029 \cdot x_1^{0.78}} = e^{11.9524} \cdot e^{0.0029 \cdot x_1^{0.78}}$$

$$y_2 = e^{\hat{y}_2} = e^{11.9524 + 0.0029 \cdot x_2^{0.78}} = e^{11.9524} \cdot e^{0.0029 \cdot x_2^{0.78}}$$

Intercept and slope 2

Suppose the original `sqft_living` is x_1 and it moved up to x_2 , then we have the following :

$$y_1 = e^{\hat{y}_1} = e^{11.9524 + 0.0029 \cdot x_1^{0.78}} = e^{11.9524} \cdot e^{0.0029 \cdot x_1^{0.78}}$$

$$y_2 = e^{\hat{y}_2} = e^{11.9524 + 0.0029 \cdot x_2^{0.78}} = e^{11.9524} \cdot e^{0.0029 \cdot x_2^{0.78}}$$

To understand the change in *price* in percents we will use the following formula:

$$\begin{aligned} 100 \times \left[\frac{y_2 - y_1}{y_1} \right] &= 100 \times \left[\frac{y_2}{y_1} - 1 \right] = 100 \times \left[\frac{\cancel{e^{11.9524}} \cdot e^{0.0029 \cdot x_2^{0.78}}}{\cancel{e^{11.9524}} \cdot e^{0.0029 \cdot x_1^{0.78}}} - 1 \right] \\ &= 100 \times \left[e^{0.0029(x_2^{0.78} - x_1^{0.78})} - 1 \right] \end{aligned}$$

Example

For example, if the *sqft_living* is 1000*ft* and we increase it to 1100*ft*, we will get the change in price of

$$100 \times \left[e^{0.0029(1100^{0.78} - 1000^{0.78})} - 1 \right] \approx 5.02\%.$$

Example

For example, if the *sqft_living* is 1000*ft* and we increase it to 1100*ft*, we will get the change in price of

$$100 \times \left[e^{0.0029(1100^{0.78} - 1000^{0.78})} - 1 \right] \approx 5.02\%.$$

In this particular example, 10% change in *sqft_living* starting from $x_1 = 1000\text{ft}$ forces 5.02% change in price.

R^2

- The model has $R^2 \approx 0.45$.

R^2

- The model has $R^2 \approx 0.45$.
- This means that our model explains about 45% of the variation by using *sqft_living* as independent variable.

ANOVA 1

Is our model with one explanatory variable better than the model with zero explanatory variables?

Our model has $F - statistic = 1.737 \times 10^4$ and $Prob > F$ is 0.000.

The Null Hypothesis:

The slope= 0

The Alternative Hypothesis:

The slope $\neq 0$

Multiple Linear Model

Now, I will build a Multiple Regression Model.

The goals for Multiple Linear Model:

- I want to improve R^2 .
- I want to use more than one explanatory variable.

Choice of Explanatory Variables

I will use the highly correlated with the price features from the correlation matrix (see heat map above, in the beginning of the presentation) for the Multiple Linear Model model.

Results for Multiple Linear Regression Model

At the end of my analysis, I came up with the following formula:

Results for Multiple Linear Regression Model

At the end of my analysis, I came up with the following formula:

$$\begin{aligned}\ln(\text{price}) = & 10.2082 + 0.3618 \cdot \text{waterfront} - 0.0160 \cdot \text{bedrooms} \\ & - 0.0153 \cdot \text{bathrooms} + 0.1400 \cdot \text{sqft_living}^{0.3} + 0.0088 \cdot \text{floors} \\ & + 0.1494 \cdot \text{view}^{0.5} + 0.0105 \cdot \text{grade}^2 + 0.1187 \cdot \ln(\text{sqft_living}^{15})\end{aligned}$$

Linearity

The Null Hypothesis:

The model is linearly predicted by the feature,

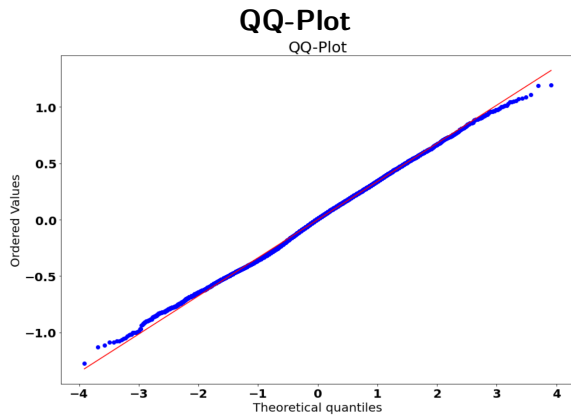
The Alternative Hypothesis:

The model is not linearly predicted by the feature.

- ① Our p-value for this model is $p = 0.933 > 0.05 = \alpha$.
- ② We don't have enough evidence to reject **The Null Hypothesis**
- ③ We conclude that our model satisfies Linearity Assumption.

Normality Assumption for Errors 1

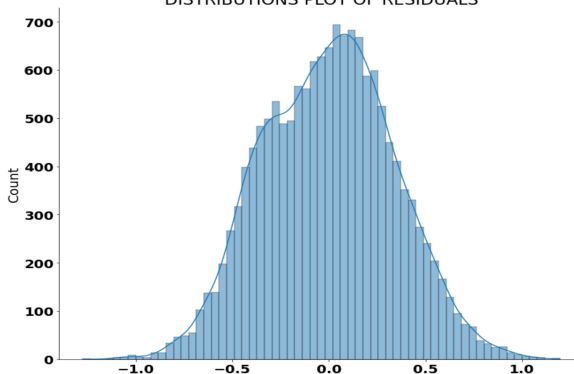
To check Normality, I used the following checks:



Normality Assumption for Errors 2

DISTRIBUTIONS PLOT OF RESIDUALS

DISTRIBUTIONS PLOT OF RESIDUALS



Normality Assumption for Errors 3

I also, used D'Agostino Test for Normality:

The Null Hypothesis:

The Residuals are normally distributed,

The Alternative Hypothesis:

The Residuals are not normally distributed.

- Our p-value for this model is $p = 0.000 < 0.05 = \alpha$.
- We have enough evidence to reject the Null Hypothesis
- We conclude that D'Agostino Test tells us, that residuals are not normally distributed.

Normality Assumption for Errors 4

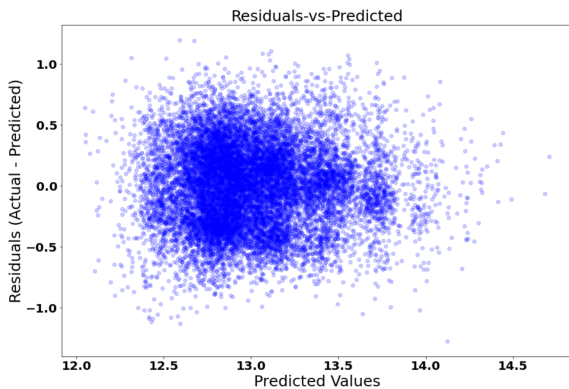
Conclusion for Normality:

9

- Based on QQ-Plot, Distributions Plot, and D'Agostino Test, I conclude that the Distribution of Errors is not far away from Normal.
- Also, since we have a lot of observations Normality Assumption doesn't play a critical role, since Central Limit Theorem will apply in this case.

Constant Error Variance 1

To check if heteroscedasticity is present in the model, I used Residual-vs-Predicted values plot and Breusch-Pagan test. I look at at the Residual-vs-Predicted values plot first.



Constant Error Variance 2

I used Breusch-Pagan Test:

The Null Hypothesis:

Homoscedasticity is present.

The Alternative Hypothesis: Homoscedasticity is not present
(i.e. heteroscedasticity exists).

Constant Error Variance 2

I used Breusch-Pagan Test:

The Null Hypothesis:

Homoscedasticity is present.

The Alternative Hypothesis: Homoscedasticity is not present (i.e. heteroscedasticity exists).

- Our p-value for this model is $p = 0.000 < 05 = \alpha$.
- We have enough evidence to reject the Null Hypothesis.
- We conclude from Breusch-Pagan Test, that we have don't have heteroscedasticity.

Constant Error Variance : Conclusion

From the Residual-vs-Predicted values plot and Breusch-Pagan Test, I conclude that we have some Heteroscedasticity in our model, but it is not very bad.

Overall Conclusion:

I conclude that our model almost satisfies statistical assumptions for the regression model.

Model

$$\begin{aligned}\ln(\text{price}) = & 10.2082 + 0.3618 \cdot \text{waterfront} - 0.0160 \cdot \text{bedrooms} \\ & - 0.0153 \cdot \text{bathrooms} + 0.1400 \cdot \text{sqft_living}^{0.3} + 0.0088 \cdot \text{floors} \\ & + 0.1494 \cdot \text{view}^{0.5} + 0.0105 \cdot \text{grade}^2 + 0.1187 \cdot \ln(\text{sqft_living}^{15})\end{aligned}$$

Explanation of the Model

Before I begin explain the coefficients, I notice that $P > |t|$ for the *floors* variable is 0.155, which makes *floors* insignificant for the analysis.

Intercept

- The model has the sample intercept of 10.2082.

Intercept

- The model has the sample intercept of 10.2082.
- If we assume that all explanatory variables are zeros, this would mean that the price would be $e^{10.2082} \approx 27,124$

0.3618 is the coefficient for *waterfront*

- *Waterfront* is a categorical variable coded as 0 or 1, a one unit difference represents switching from one category to the other. 10.2082 is then the average difference in *price* between the category for which *waterfront* = 0 (no waterfront) and the category for which *waterfront* = 1 (the house has a waterfront).

0.3618 is the coefficient for *waterfront*

- *Waterfront* is a categorical variable coded as 0 or 1, a one unit difference represents switching from one category to the other. 10.2082 is then the average difference in *price* between the category for which *waterfront* = 0 (no waterfront) and the category for which *waterfront* = 1 (the house has a waterfront).

So compared to $\ln(\text{price})$ of the house with no waterfront, we would expect the $\ln(\text{price})$ for the house with waterfront to be 0.3618 higher, on average, if we fix all other explanatory variables.

0.3618 is the coefficient for *waterfront*

Let y_2 = price of the house with waterfront and
 y_1 = price of the house with no waterfront, then

$$\ln y_2 - \ln y_1 = 0.3618 \implies \ln \frac{y_2}{y_1} = 0.3618 \implies \frac{y_2}{y_1} = e^{0.3618}$$

To understand the change in *price* in percents, if we switch from the house with no waterfront to the house with waterfront while keeping all other variables the same, will use the following formula:

$$100 \times \left[\frac{y_2 - y_1}{y_1} \right] = 100 \times \left[\frac{y_2}{y_1} - 1 \right] = 100 \times [e^{0.3618} - 1] \approx 43.59\%$$

0.3618 is the coefficient for *waterfront*

Switching from the house with no waterfront to the house with waterfront while keeping all other explanatory variables fixed, will increase the price by 43.59%.

–0.0160 is the coefficient for number of *bedrooms*

Let y_2 = price for the house with x_2 number of bedrooms and y_1 = price for the house with x_1 number of bedrooms, then

$$\begin{aligned}\ln y_2 - \ln y_1 &= -0.016x_2 - (-0.016x_1) \\ &= -0.016(x_2 - x_1) \implies \ln \frac{y_2}{y_1} = -0.016(x_2 - x_1) \\ \implies \frac{y_2}{y_1} &= e^{-0.016(x_2 - x_1)}\end{aligned}$$

—0.0160 is the coefficient for number of *bedrooms*

If we increase the number of bedrooms by 1 while keeping the other variables fixed and use percents, we will get the following

$$100 \times \left[\frac{y_2}{y_1} - 1 \right] = 100 \times [e^{-0.016} - 1] \approx -1.58\%$$

- Thus, increasing the number of bedrooms by 1 while keeping the other variables fixed, will decrease the price of the house by 1.58%.

—0.0153 is the coefficient for number of *bathrooms*

- 1 The same explanation we have for -0.0153 which a coefficient for number of bathrooms.

—0.0153 is the coefficient for number of *bathrooms*

- 1 The same explanation we have for -0.0153 which a coefficient for number of bathrooms.
- 2 If we increase the number of bathrooms by 1 while keeping the other variables fixed, will decrease the price of the house by 1.51%.

0.1400 is the coefficient for *sqft_living*

Let y_2 = price for the house with x_2 *sqft_living* and y_1 = price for the house with x_1 = *sqft_living*, then

$$\ln y_2 - \ln y_1 = 0.1400x_2^{0.3} - 0.1400x_1^{0.3} \implies \ln \frac{y_2}{y_1} = 0.1400 (x_2^{0.3} - x_1^{0.3})$$

$$\implies \frac{y_2}{y_1} = e^{0.1400(x_2^{0.3} - x_1^{0.3})}$$

If we increase the *sqft_living* from 1000ft to 1100ft while keeping the other variables fixed, we will get the following change in price in percents:

$$100 \times \left[\frac{y_2}{y_1} - 1 \right] = 100 \times \left[e^{0.1400(1100^{0.3} - 1000^{0.3})} - 1 \right] \approx 3.28\%$$

0.1400 is the coefficient for *sqft_living*

- If the *sqft_living* is 1000 ft and we increase it to 1100 ft while keeping the other variables fixed, we will get the change in price of 3.28%.

0.1400 is the coefficient for *sqft_living*

- If the *sqft_living* is 1000*ft* and we increase it to 1100*ft* while keeping the other variables fixed, we will get the change in price of 3.28%.
- In this particular example 10% change in *sqft_living* starting from $x_1 = 1000\text{ft}$ forces 3.28% change in price.

0.1494 is the coefficient for the *view*

The coefficient for the *view* has the same explanation as the *sqft_living*.

0.1494 is the coefficient for the *view*

The coefficient for the *view* has the same explanation as the *sqft_living*.

If we increase the *view* by 1 unit from 2 to 3 while keeping the other variables fixed, we will get the following:

$$100 \times \left[\frac{y_2}{y_1} - 1 \right] = 100 \times \left[e^{0.1494(3^{0.5} - 2^{0.5})} - 1 \right] \approx 4.86\%$$

In this particular example if the *view* will increase from 2 to 3, the price will increase by 4.86%.

0.0105 is the coefficient for the *grade*

The coefficient for the *grade* has the same explanation as the *sqft_living*.

0.0105 is the coefficient for the *grade*

The coefficient for the *grade* has the same explanation as the *sqft_living*.

If we increase the *grade* by 1 unit from 2 to 3 while keeping the other variables fixed, we will get the following:

$$100 \times \left[\frac{y_2}{y_1} - 1 \right] = 100 \times \left[e^{0.0105(3^2 - 2^2)} - 1 \right] \approx 5.39\%$$

In this particular example if the *grade* will increase from 2 to 3 while other variables stay the same, the price will increase by 5.39%.

0.1187 is the coefficient for the *sqft_living15*

Let y_2 = price for the house with x_2 *sqft_living15* and y_1 = price for the house with x_1 *sqft_living15*, then

$$\ln y_2 - \ln y_1 = 0.1187 \ln x_2 - 0.1187 \ln x_1 \implies \ln \frac{y_2}{y_1} = 0.1187 \ln \frac{x_2}{x_1}$$

$$\implies \frac{y_2}{y_1} = \left(\frac{x_2}{x_1} \right)^{0.1187}$$

The change in price in percent will be:

$$100 \times \left[\frac{y_2}{y_1} - 1 \right] = 100 \times \left[\left(\frac{x_2}{x_1} \right)^{0.1187} - 1 \right]$$

0.1187 is the coefficient for the *sqft_living15*

In this particular example, if we increase the *sqft_living15* by 100 units from 1000 to 1100 while keeping the other variables fixed, we will get the following:

$$100 \times \left[e^{0.1187} \times \frac{1100}{1000} - 1 \right] \approx 1.14\%$$

Thus, 10% increase in *sqft_living15* will lead to 1.14% increase in *price*.

R^2

- The model has $R^2 \approx 0.5$.

R^2

- The model has $R^2 \approx 0.5$.
- This means that our model explains about 50% of the variation by using *sqft_living* as independent variable.

ANOVA

Is our model with many explanatory variable better than the model with zero explanatory variables?

Our model has $F - statistic = 1.737 \times 10^4$ and $Prob > F$ is 0.000.

Our model has $F - statistic = 1.737 \times 10^4$ and $Prob > F$ is 0.000.

The Null Hypothesis: The slope = 0

ANOVA

Is our model with many explanatory variable better than the model with zero explanatory variables?

Our model has $F - statistic = 1.737 \times 10^4$ and $Prob > F$ is 0.000.

The Null Hypothesis: The slope= 0

The Alternative Hypothesis: The slope \neq 0

ANOVA

- Our p-value for this model is $p = 0.000 < 0.05 = \alpha$.
- We have enough evidence to reject the Null Hypothesis at 5% level of significance
- We conclude that the Test tells us, that at least one of the coefficients is not 0.
- Since our p-value is 0, there is a 0% probability that the improvements that we are seeing with our independent variables model are due to random chance alone.

Conclusions: Data Modeling

I used the following steps during data modeling:

- Dropped 2376 rows where *waterfront* is NaN.
- Dropped 63 rows where *view* is NaN.
- Dropped 3842 rows where *yr_renovated* is NaN.
- I converted *sqft_basement* string format into float.
- During modeling dropped outliers when necessary.

Conclusions: Modeling

- I built two models:
 - Linear Regression Model
 - Multiple Linear Regression Model
- I checked whether the models satisfy statistical assumptions of Linear Regression
- I explained the models.
- Models can be used for interpolation given the data about a particular property.

Conclusions: Ways to Improve the Analysis

- More data wrangling is need to remove *heteroscedasticity* from the Multiple Linear Regression Model.
- Include more explanatory variables.
- Scrape webpages for more data such as school grade, crime rate, etc. for properties.

THE END
THANK YOU!