

# HOUSE SALE IN NORTHWESTERN COUNTY





# INTRODUCTION

- Real estate can be defined as land, natural resources and any permanent structures attached to the land including the rights to use and occupy them
- Real estate may be classified into different categories which include ; Residential Real Estate which deals with properties used for living purpose
- The value of the real estate can be influenced by number of factors which include location, size, condition, infrastructure surrounding the real estate, economic conditions, population and government policies



# PROBLEM STATEMENT

- Real estate has a number of stakeholders. These stakeholders include ; buyers, sellers, brokers, developers, financial institutions and the government.
- All these stakeholders are interested in the value of the property and the factors most affecting the value
- They need to easily and with good accuracy, predict the value of the properties they are interested in. This will help the stakeholders arrive on agreement more easily and quickly



# DATA UNDERSTANDING

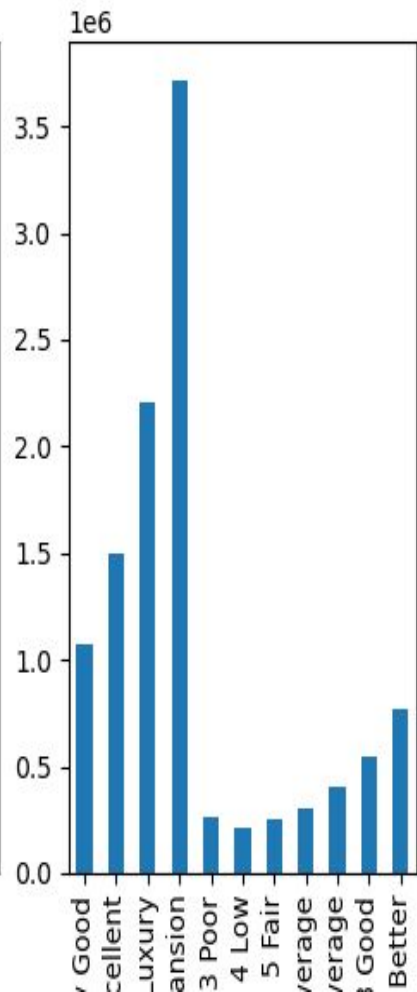
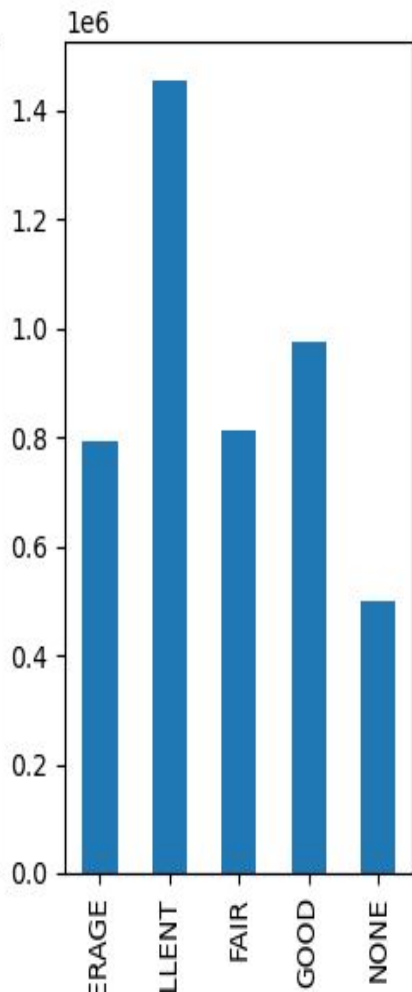
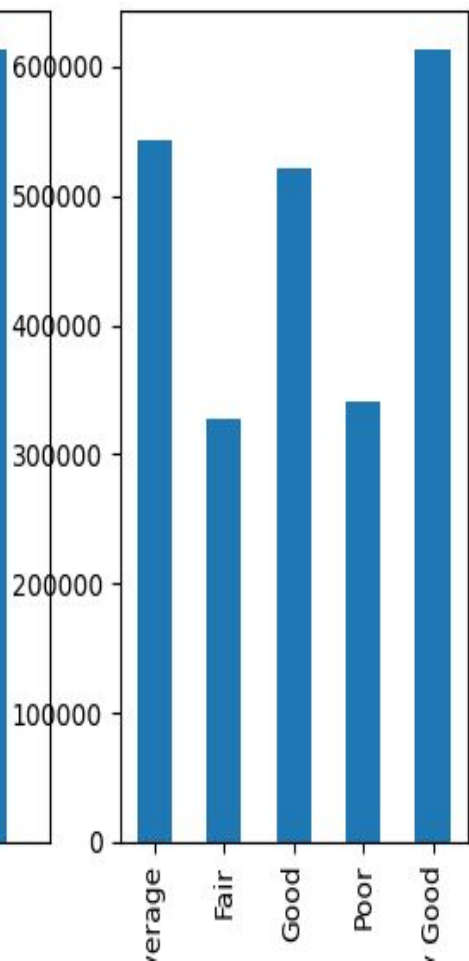
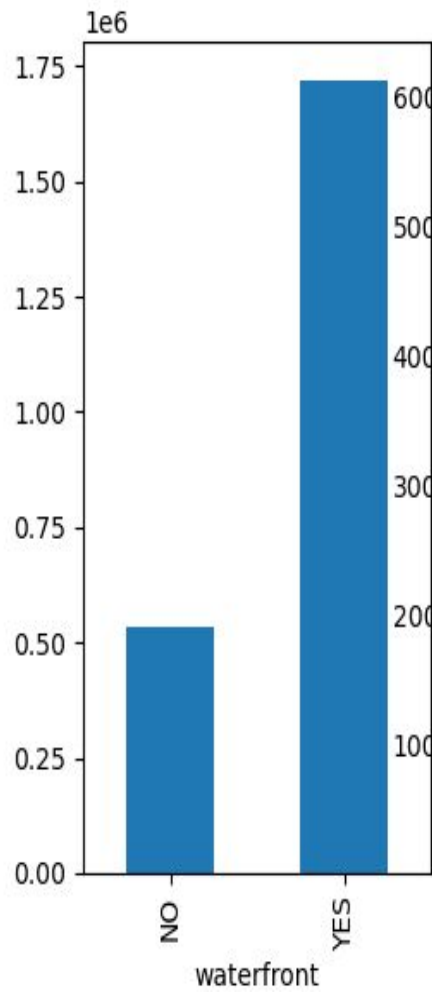
- This project uses the King County House Sales dataset and contains data with 21597 rows and 20 columns
- This data contain information on houses including ; dates, prices,number of rooms and bathrooms,living room , lot and basement areas, waterfronts,the views from the house, house grade, areas of the 15 neighbouring houses.
- All this information from this data set has an impact on the overall price of the houses



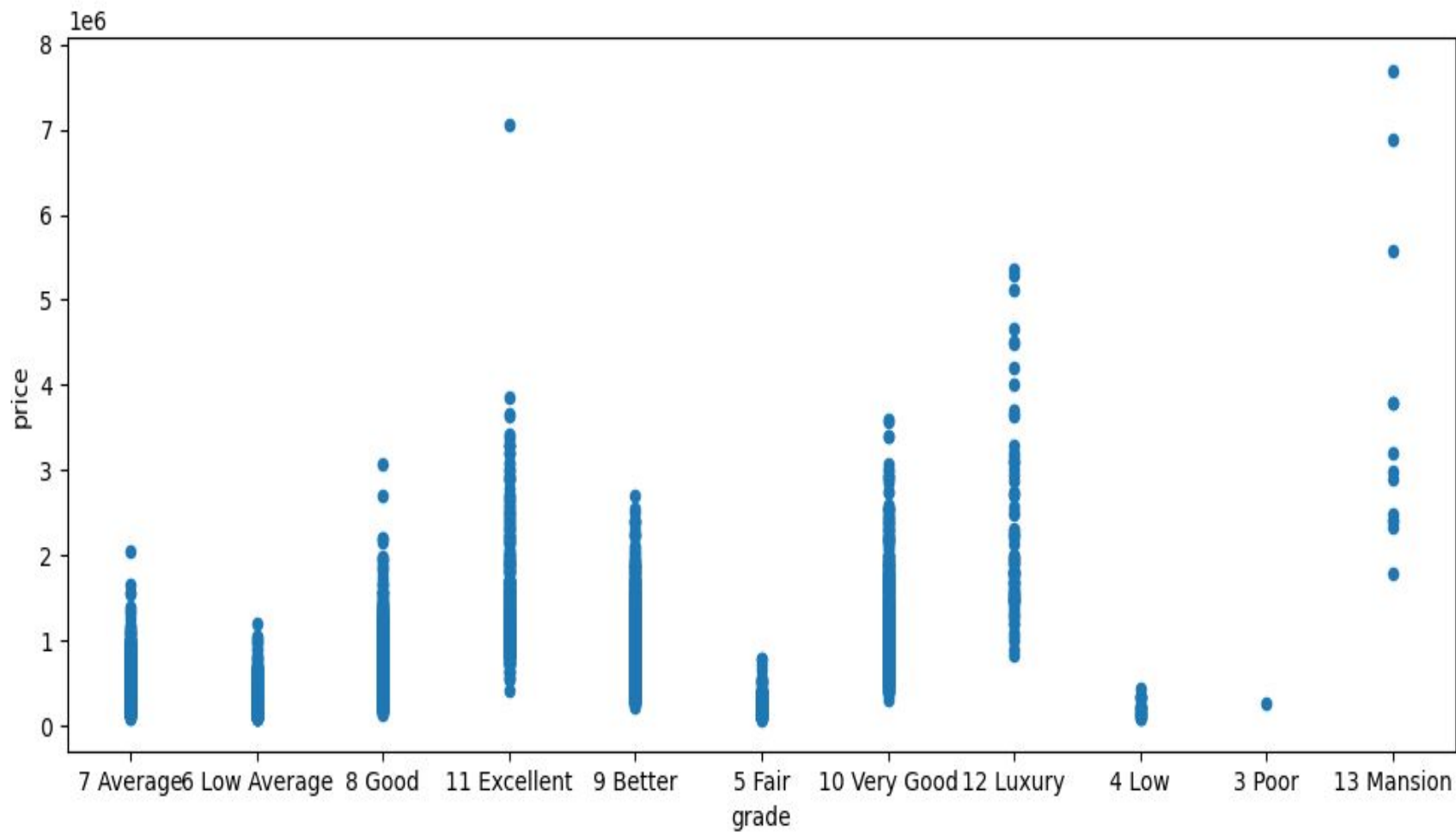
# DATA PREPARATION

- replacing missing data was done on view and waterfront columns.They were replaced with none in each case.
- checking for numeric column most correlated to price and correlation between them was done to find a best predictor and avoid multicollinearity .
- Transforming Categorical Variables .the data had 4 important categorical variables.view,waterfront,condition and grade.
- Checking their impact using graphs;

categorical bar graphs



Scatter Plot grade vs price



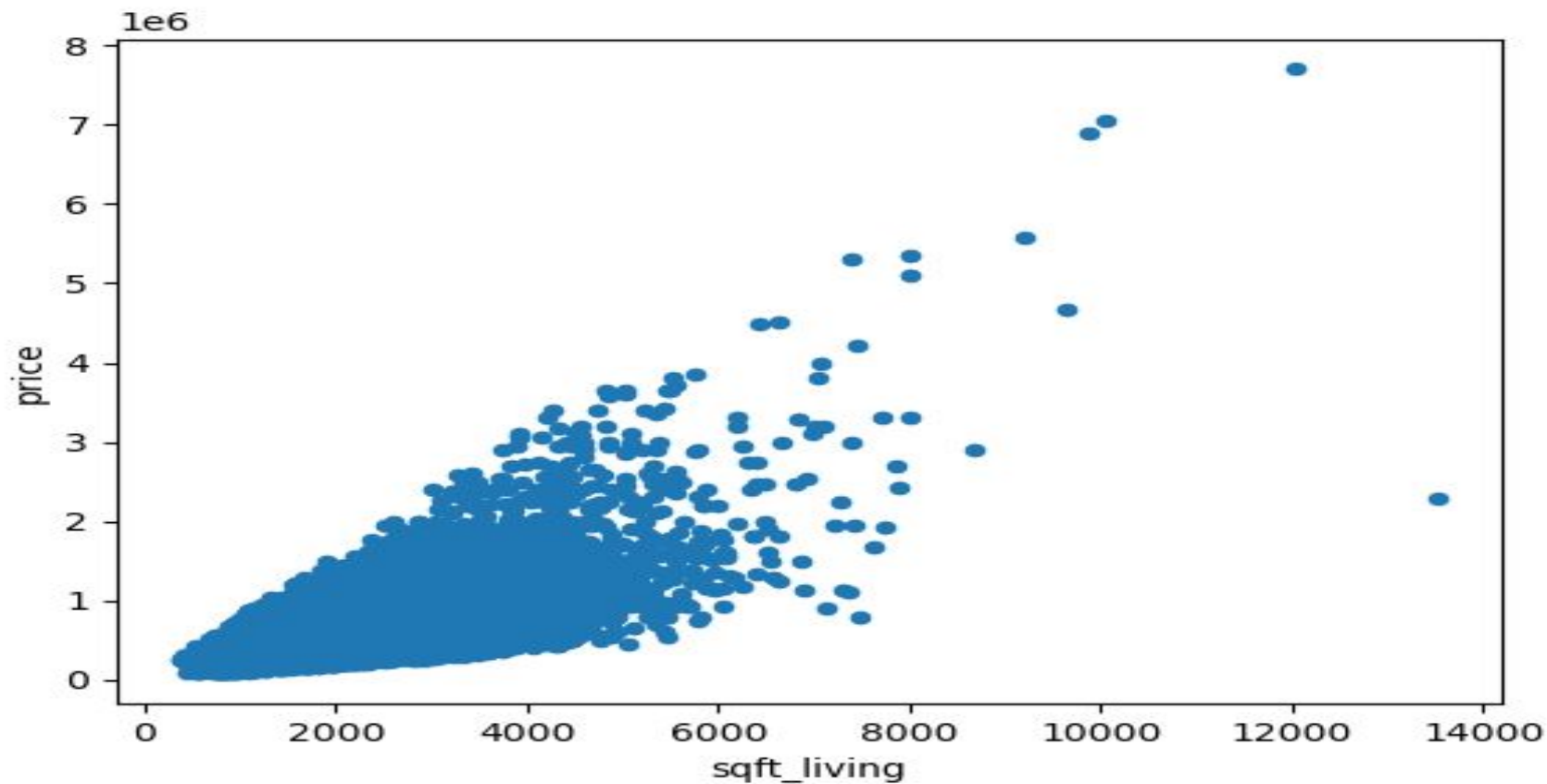


# MODELING

- We will start modeling with a simple linear regression model between the area of the housing living room(sqft\_living) against price.
- This is because sqft\_living column has the highest correlation to price compared to other numeric variables.
- First we will have a scatter plot preview of sqft\_living vs price to get a picture of the correlation between the 2.



scatter plot showing distribution of sqft\_living against price

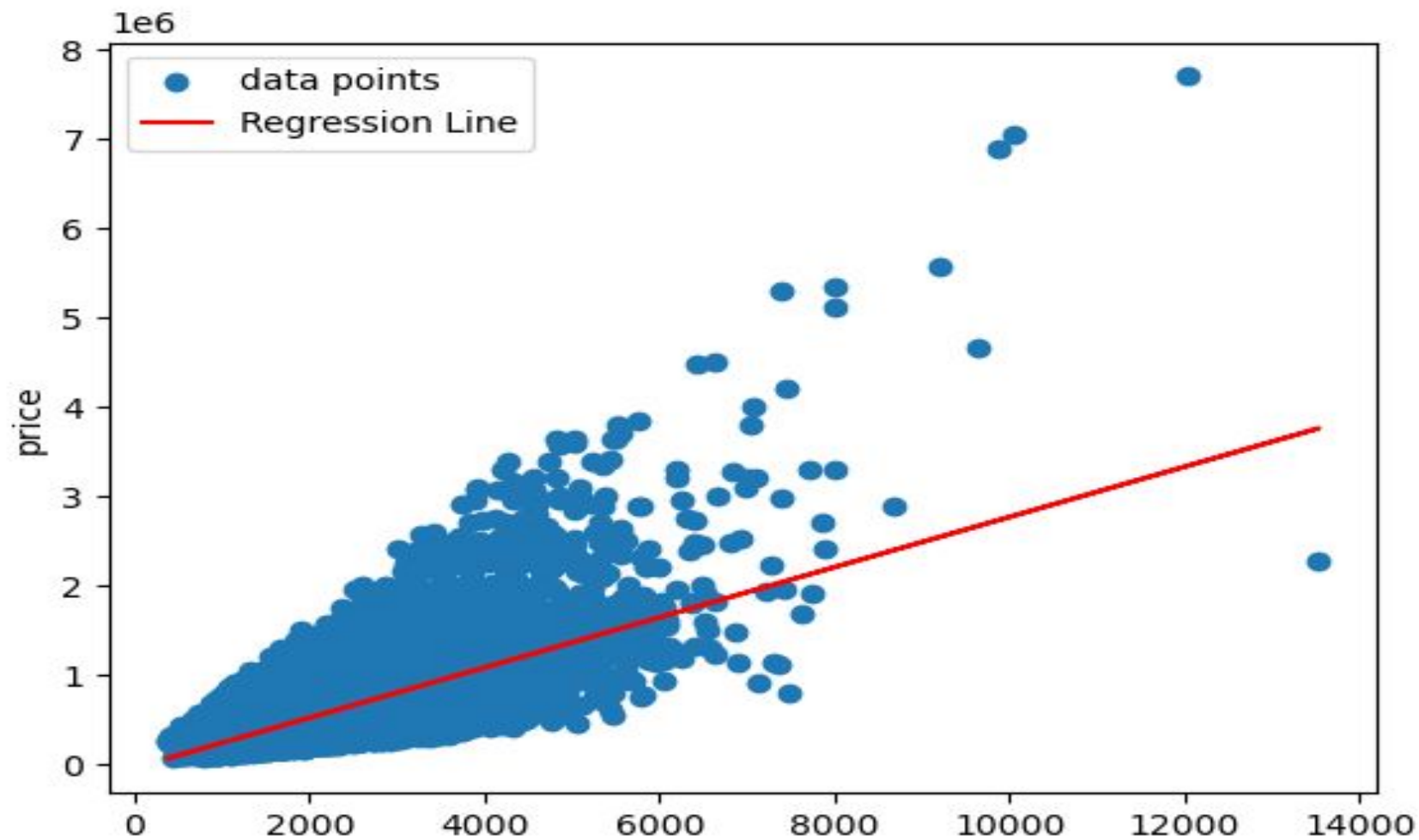




# First simple linear regression model

- **$\text{price} = 281(\text{sqft\_living}) - 43989$**
- It has  $R^2$  of about 0.49. This means 49% of the variation in price is explained. It has an overall p-value of 0 which is below the typical alpha of 0.05. This means that the model is statistically significant.
- The model has a slope of 281 which means an increase of  $\text{sqft\_living}$  by 1 increases the price by 281 units.
- the y-intercept is -43989 which means when the living area is zero the price is -43989. That does not make sense but we will transform our more advanced models for an interpretable y-intercept.

Scatter Plot sqft\_living vs price with regression line





# Improving the model

- In our final model we use `sqft_living_squared` in place of `sqft` as it seems to improve the `r_squared` therefore improving our variation explanation while improving the p values.
- Introduction of categorical variables in our model help improve the `r_squared` further explaining the variation in price.
- Transforming the x axis values to be centred around the mean of `sqft_living_squared` improves our explanation of the y intercept.



# Final model Regression Results

- **R\_squared** and **F\_pvalue** Our final regression model has R\_squared of about 0.653. This means about 65% of the variation in price is explained by our model. Our model has an overall p\_value of 0 which is below the typical alpha of 0.05. This means that our model is statistically significant.
- **Y\_intercept** Our y\_intercept is 418800 ,that is to say the price of an average size cost is about 418800 dollars on price. This is because our data is centered around the mean of sqft\_living



# Regression results continued

- **coefficient of sqft\_living** In the final model sqft\_living is squared therefore 0.0256 is coefficient of sqft\_living squared therefore gradient translates to  $2(0.0256)\text{sqft\_living}$ .which can be simply interpreted as a increase of 1% in Sqft\_living will translate in about 0.0512% increase in price
- **view** there is an increase of 164800 for view fair compared to view none,an increase of 96690 for view average compared to none,an increase of 155800 for view good as compared to none and an increase of 293100 for excellent view compared to none
- **waterfront** There is an increase of 533600 in price for houses with waterfront compared to those that do not have a waterfront.



## Regression results continued

- **condition** we can't say anything about about condition poor and fair as they are not statistically viable. However an increase of 51490 is observed for good condition compared to average condition while an increase of 134600 is observed for houses with very good conditions compared to average
- **Grade** decrease of 122300, 121300, 67130 in prices is observed for houses of grade 4-low, 5-fair and 6-low average respectively compared to houses of grade\_7 Average. increase of 86460, 218900, 386600, 591300, 910900, 1764000 is seen for house of grade 8 Good, 9 Better, 10 Very Good, 11 Excellent, 12 Luxury, 13 Mansion respectively compared to houses of grade\_7 Average.



# Conclusions

- Statistically significant - the final model has a has an overall p\_value of 0 which is below the typical alpha of 0.05.This means that our model is statistically significant.
- Coefficients - an increase in area (living room space which is correlated to other areas) show a improvement of 0.0512% in price for 1% increase in area.grade,condition,view and waterfront seem to greatly improve the prices of the houses.
- limitation - our model provided only about 65% explanation of variation which might be small in some practical cases.





# Recommendations

- All the factors considered in this model should be given much considerations when dealing with house as they have proved to be statistically significant.
- House Grade seems to have the largest impact among categorical data and should be given priority compared to the other categorical data.
- The model should be used in less variation sensitive cases such as price estimation but not price setting as the  $r\_squared$  is not that high.