

# King County House Price Prediction

Data Analysis and Multiple Regression in Python

Presenter: Bonnie Ma

# Agenda

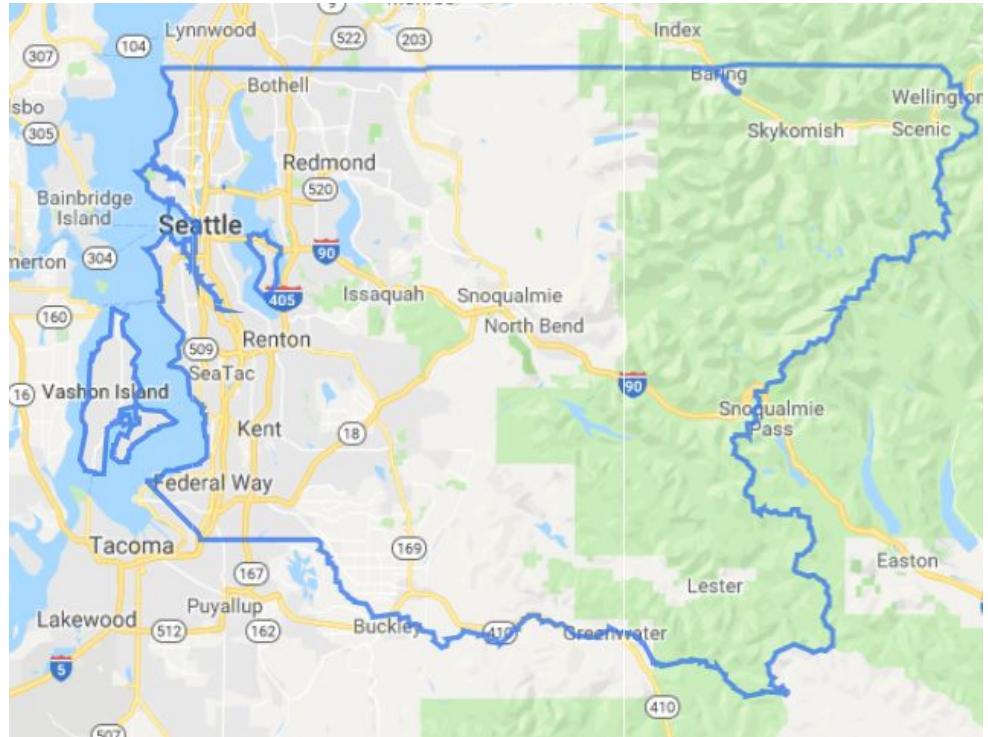
- Business Background
- Analytics Approach
- Data Exploration
- Model Overview
- Model Results
- Insights

# Business Background

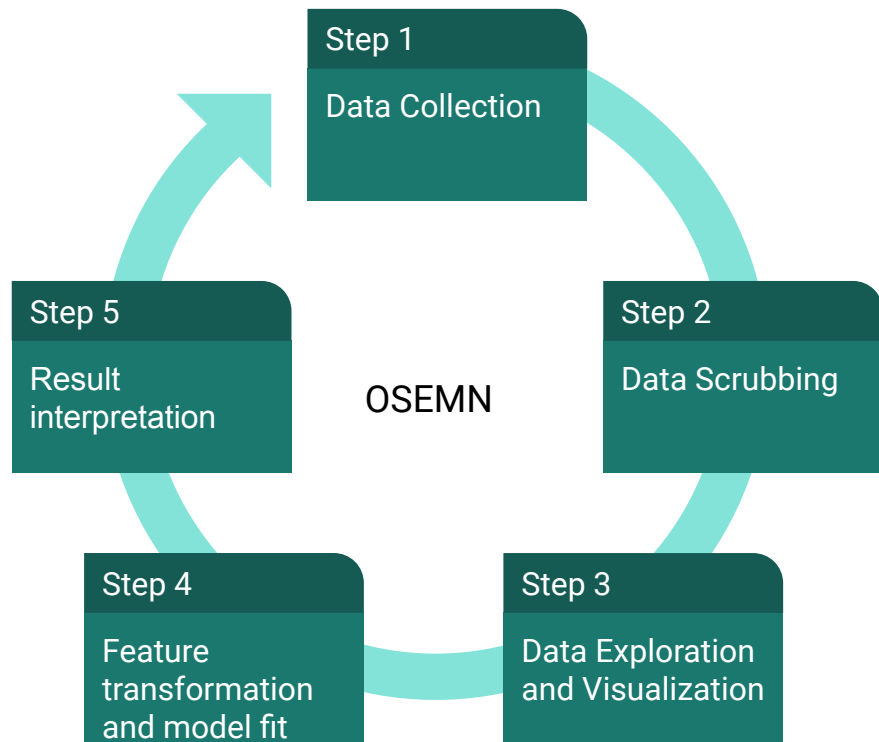
King County is one of the most beautiful and populous area in the United States. It is located in the state of Washington and has the population of 2.3 million as of 2018.

The growth of companies like Amazon and Microsoft has lifted the house price significantly over the past few years.

The goal of this project is to understand the drivers of King County house price in order to predict house price in the future and identify investment opportunities.



# Our Approach – OSEM MN



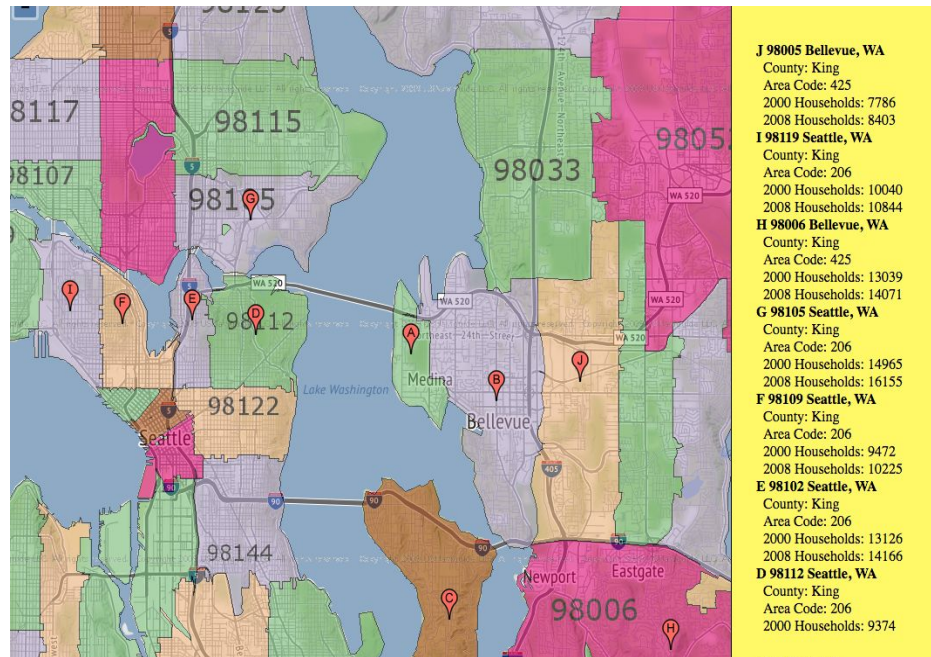
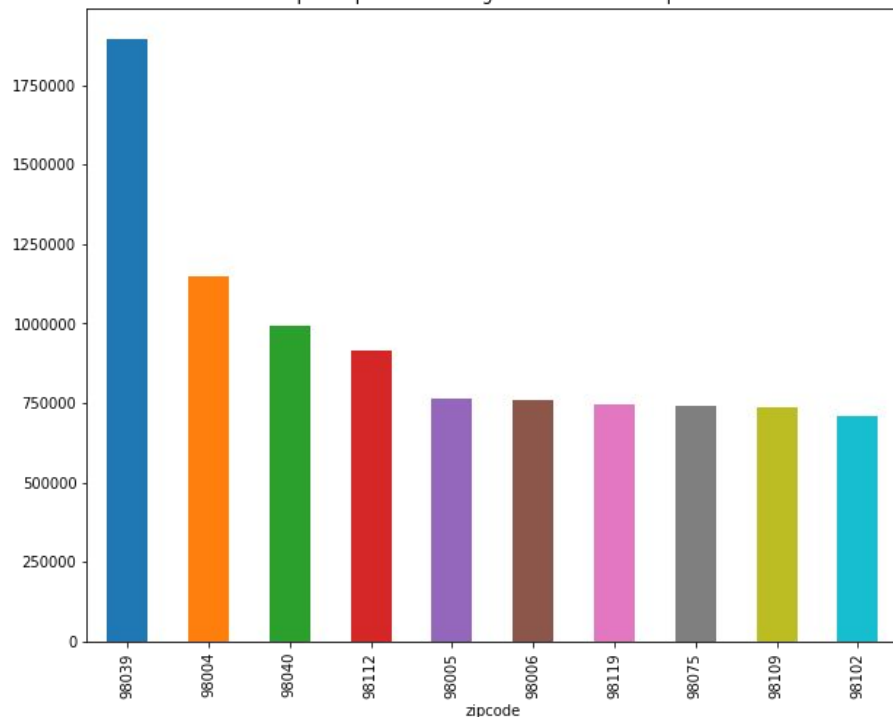
**Goal:** build a robust regression model to predict King County house price as accurate as possible

**KPIs:** Adjusted R-square, Mean Squared Error (MSE)

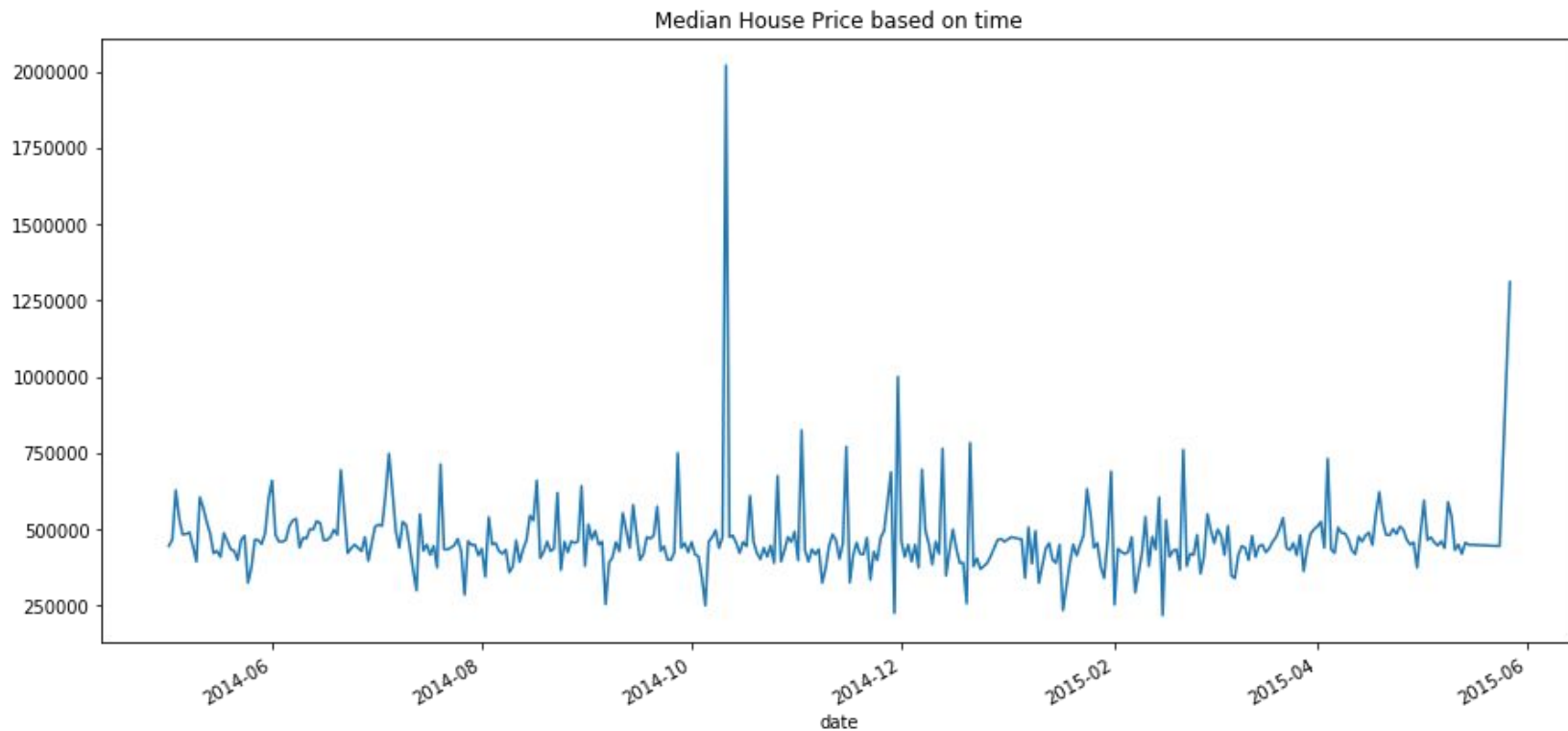
The data set has 21,597 records and 18 variables.

# Top 10 Areas with Highest Median House Price

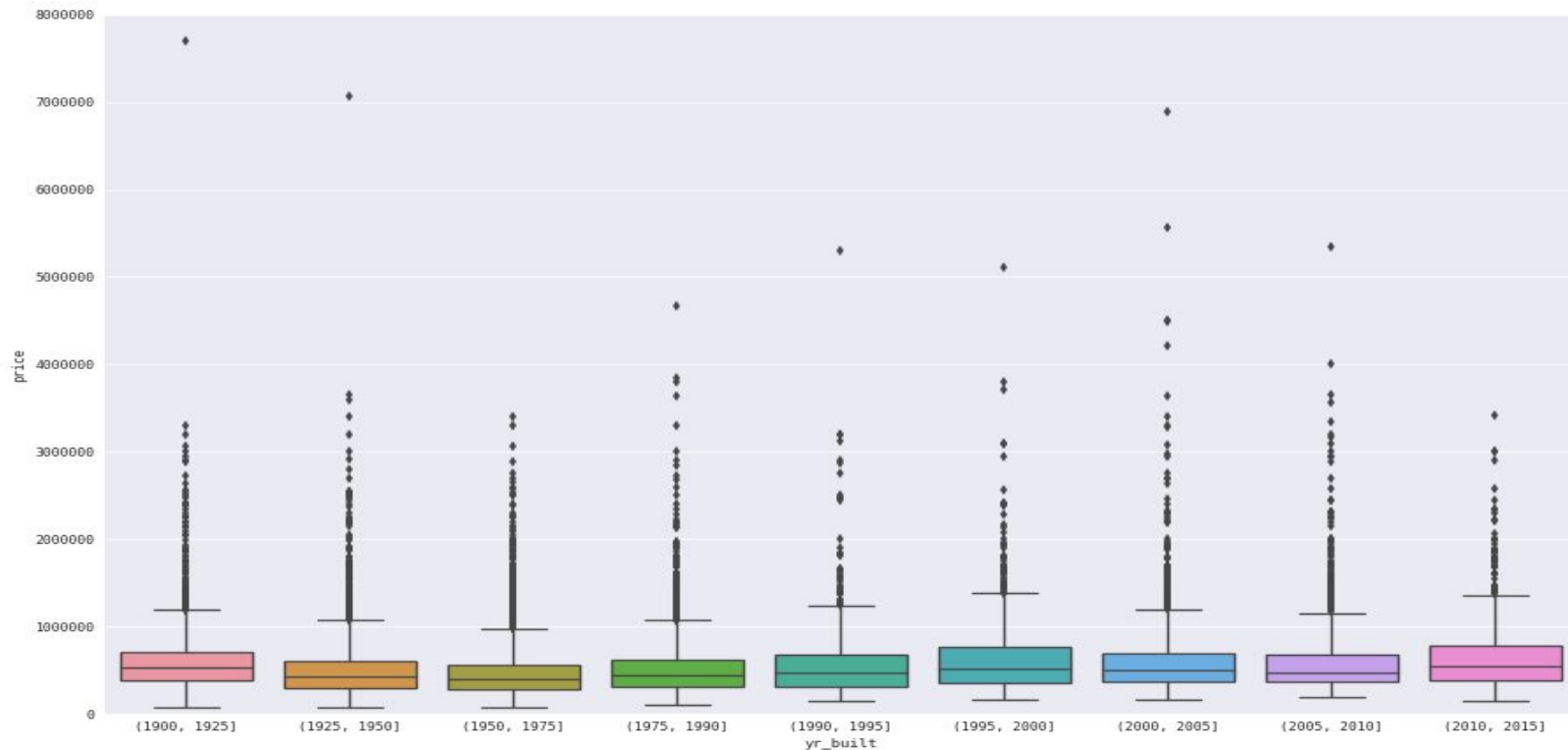
Top 10 zipcodes with highest median house price



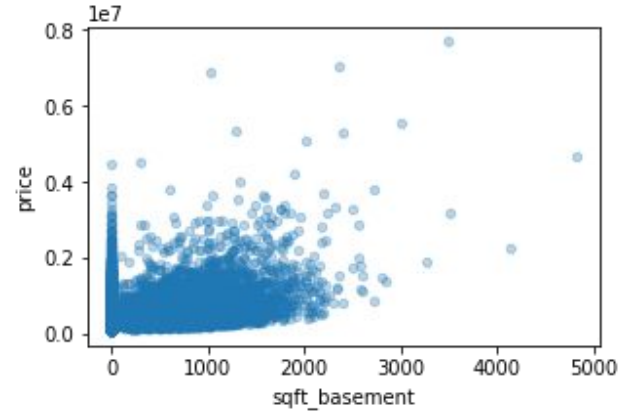
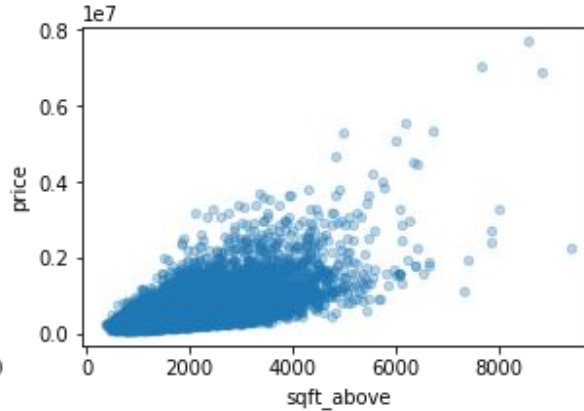
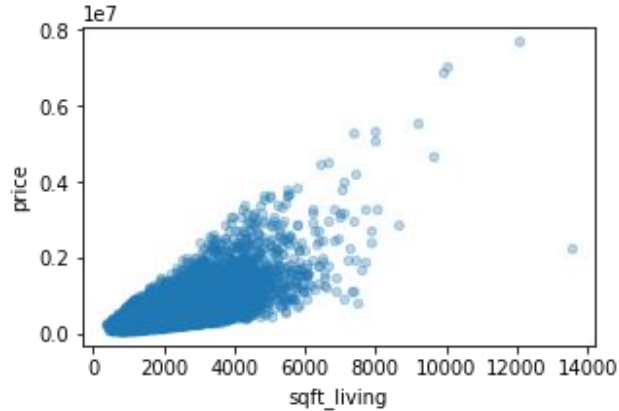
# House Price Change Over Time



# House Price Change Over Time



# House price vs Footage



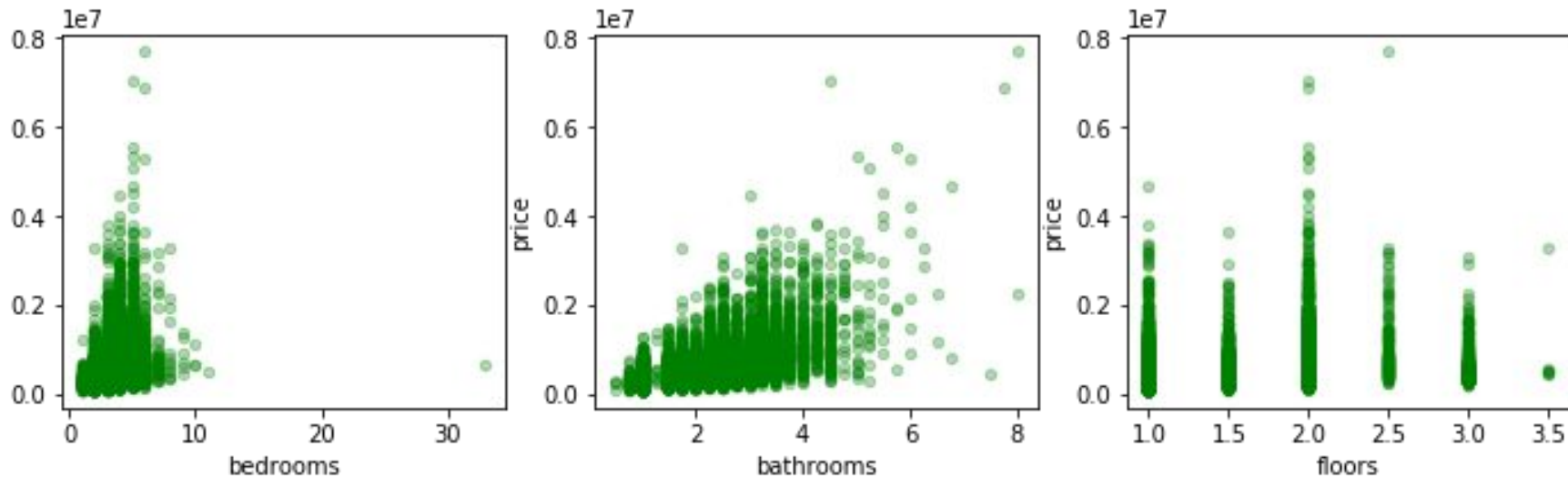
Sqft\_living: Footage of the home

Sqft\_above: footage of house apart from basement

Sqft\_basement: footage of the basement

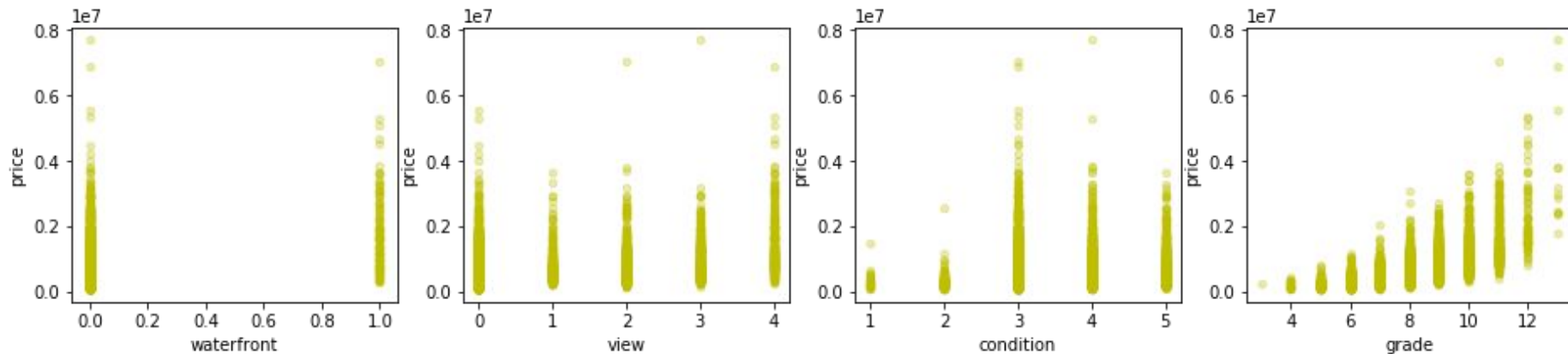


# House price vs Bed/Bathrooms/Floors



No obvious correlation observed between variable and target.

# House price vs other Categorical Variables



Grade has a positive correlation with price.

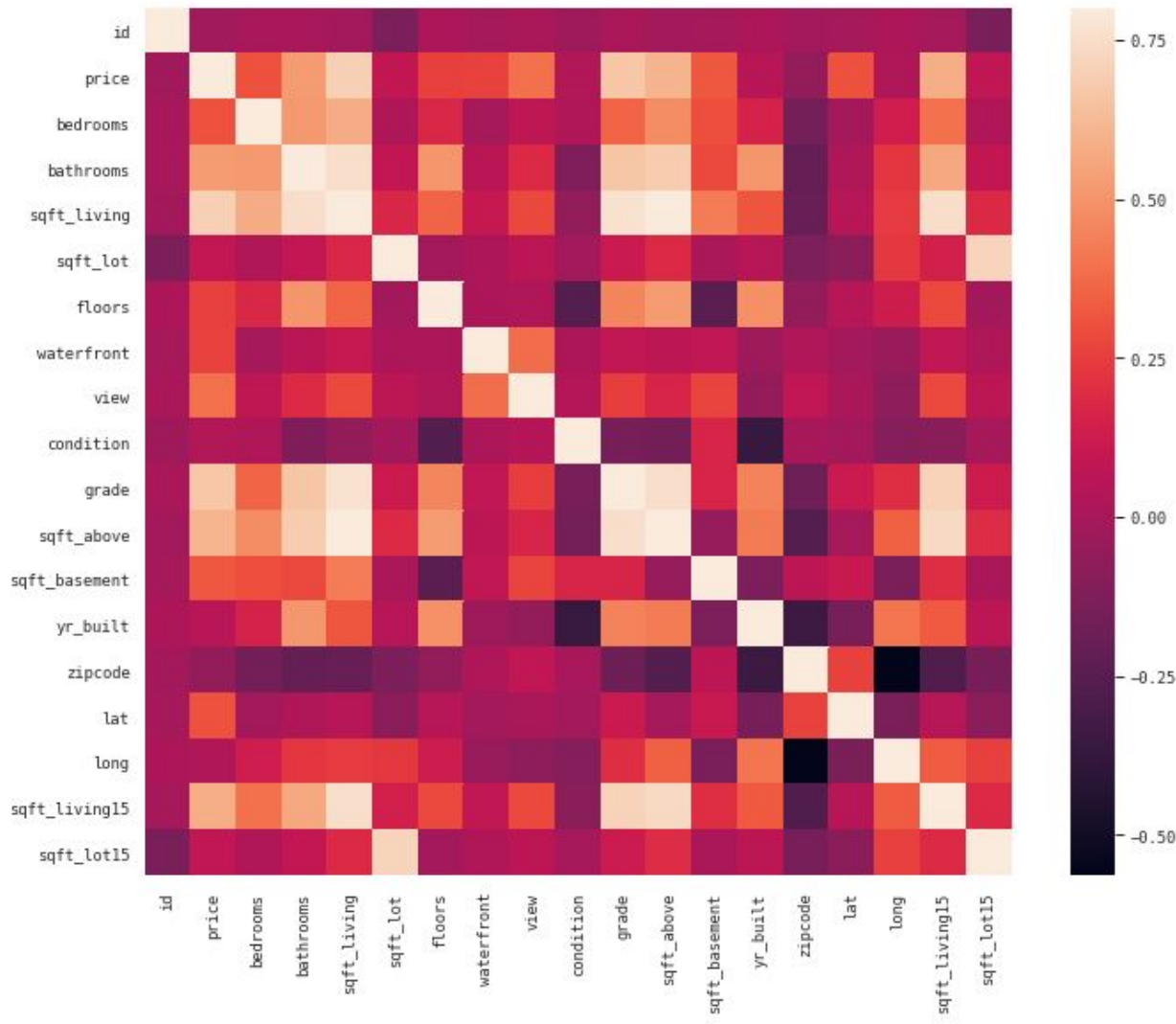
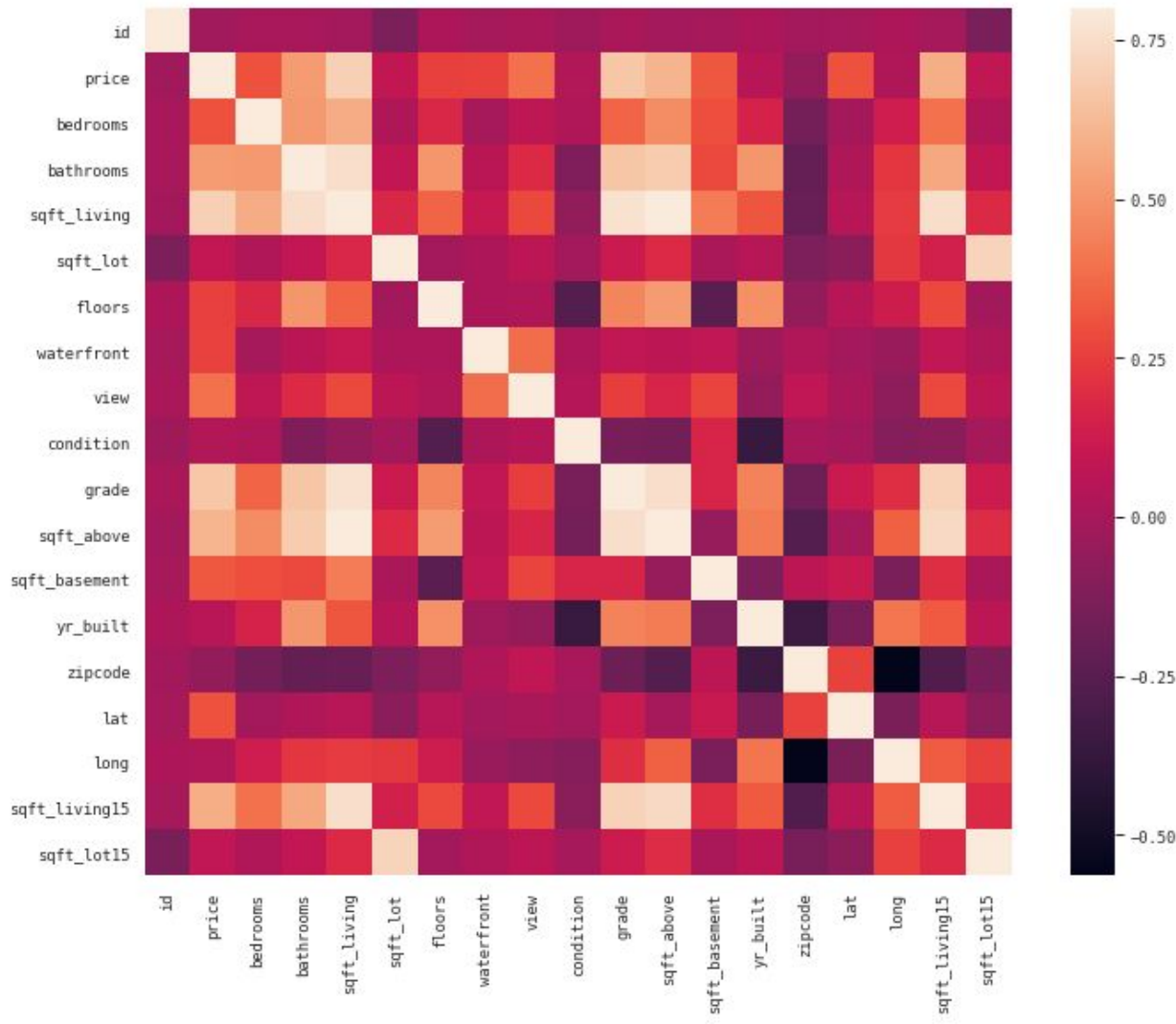
## Correlation between variables

Variables that are highly correlated to price are:

- sqft\_living
- sqft\_above
- grade

Independent variables which are highly correlated with each other:

- sqft\_living with sqft\_living15, bathroom, grade, sqft\_above
- grade with bathroom, sqft\_above, sqft\_living15, sqft\_living



# Model Set up

- Multiple Linear Regression on Price
- Data Transformation:
  - Drop variables with high collinearity or poor quality (sqft\_living15, sqft\_lot15, yr\_renovated)
  - Turn binary and ordinal variables into categorical values (bedrooms, bathrooms, waterfront, view, condition, grade, yr\_built, zip code)
  - Log transformation and standardization on continuous variables
- Package used for model: statsmodels, sklearn
- Model selection method: Forward selection and Stepwise selection
- Model validation: cross-validation

# Model Result – 83 variables

Variables that significant impact to house price:

- Footage of area above basement
- Bathroom: 1-5-2, 2-2.5, 5-8
- Bedroom: 1-2, 3-4, 4-5, 5-6
- Grade: 6, 7, 8, 9, 10, 11, 12, 13
- Zipcode: some of zipcodes
- Year: 1925-1950, 1950-1975, 1975-1990, 1990-1995, 1995-2000

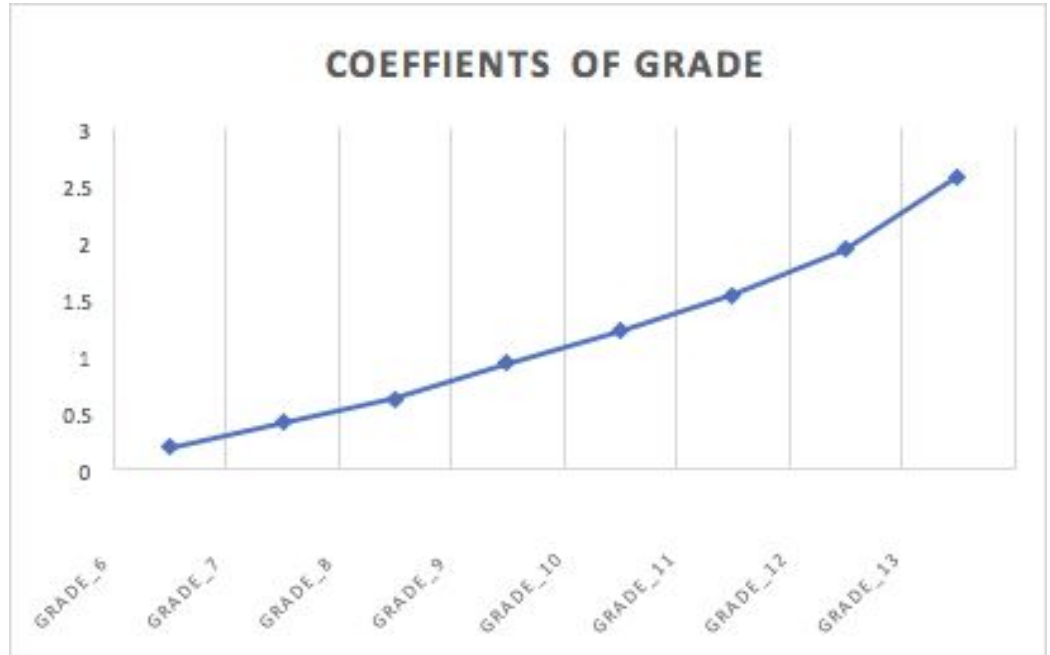
Model KPIs:

- Adjusted R-sqaure: 0.839
- P values: all p values is less than 0.05
- Train MSE: 0.1615
- Test MSE: 0.1607



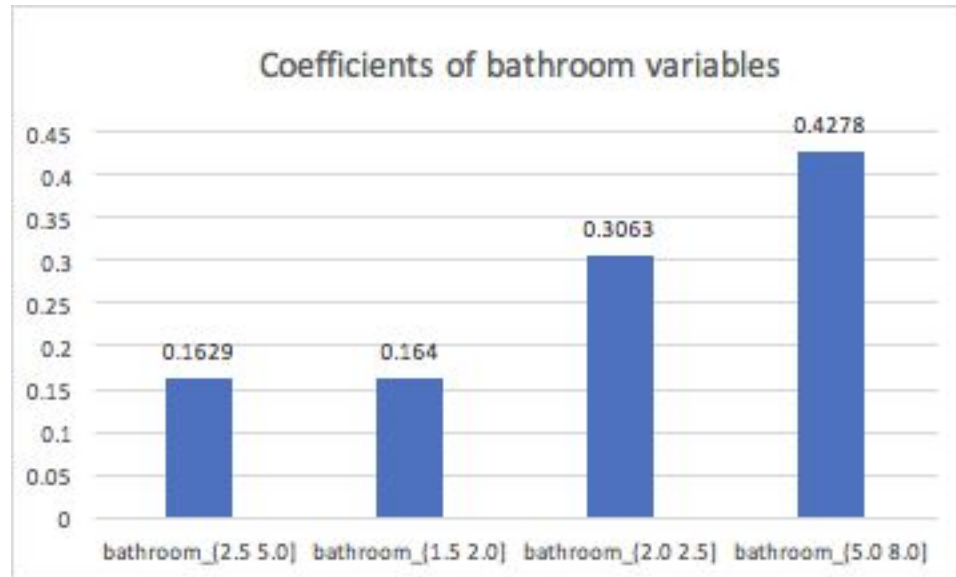
# Insights – grade

Grade begins to have significant impact on the house price till the level reaches 6 and the coefficient accelerates as the grade increases. It is important to pay attention to house grades when making investment decisions.



# Insights – bathrooms

Number of bathrooms per bedroom has different impacts on house price in different bins. When it reaches 5-8 bathrooms per bedroom, the coefficient increased the most which may indicate a different kind of building.



# Insights – zip code

Coefficients of zip code tell which area is more referable by the market. For long term investment, it is better to choose zip code with a positive coefficient.

## Fastest growing areas

98039  
98004  
98112  
98109  
98119

## Fastest declining areas

98022  
98023  
98001  
98003  
98092

