

Samsung Innovation Campus

Artificial Intelligence Course

King County House Prices Prediction Model

Team Members:

Mina Sameh Mahrous

Kerolos Monier

Galal Mohammed

Supervised by:

Dr. Doaa Mahmoud

Eng/ Abdelrahman Sadory



Objectives

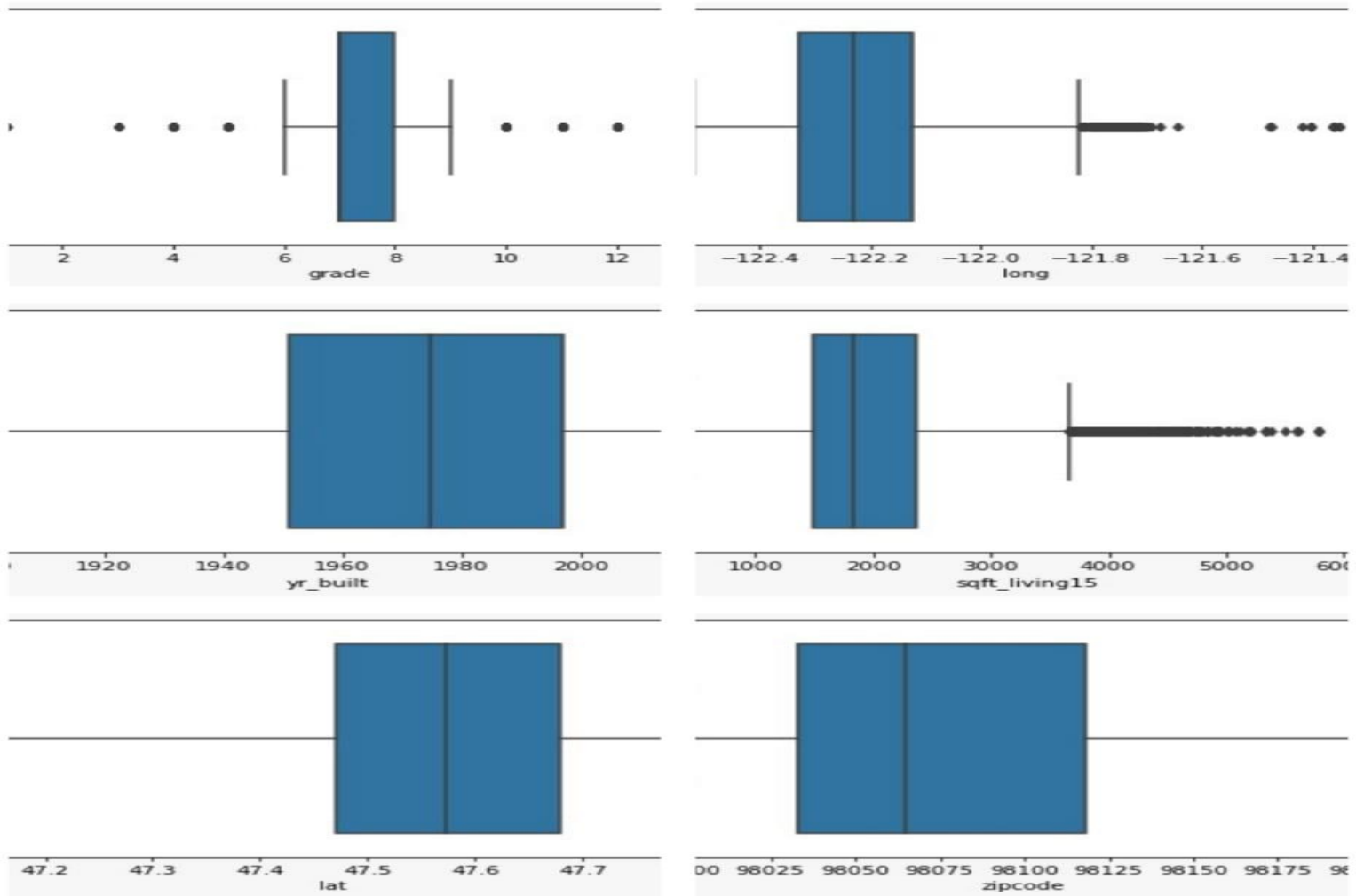
- I. Overview of Data
- II. Data pre-processing
- III. Data visualization and pattern discovery
- IV. Predictive Modeling
- V. Model Implementation
- VI. Plan for future upgrades

I. Overview of Data

Variable	Description
Id	Unique ID for each home sold
Date	Date of the home sale
Price	Price of each home sold
Bedrooms	Number of bedrooms
Bathrooms	Number of bathrooms, where .5 accounts for a room with a toilet but no shower
Sqft_living	Square footage of the apartments interior living space
Sqft_lot	Square footage of the land space
Floors	Number of floors
Waterfront	A dummy variable for whether the apartment was overlooking the waterfront or not
View	An index from 0 to 4 of how good the view of the property was
Condition	An index from 1 to 5 on the condition of the apartment,
Grade	An index from 1 to 13, where 1-3 falls short of building construction and design, 7 has an average level of construction and design, and 11-13 have a high quality level of construction and design
Sqft_above	The square footage of the interior housing space that is above ground level
Sqft_basement	The square footage of the interior housing space that is below ground level
Yr_built	The year the house was initially built
Yr_renovated	The year of the house's last renovation
Zipcode	What zipcode area the house is in
Lat	Latitude
Long	Longitude
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

Outlier Detection: Outliers were detected and analyzed using the Outlier Boxplots.

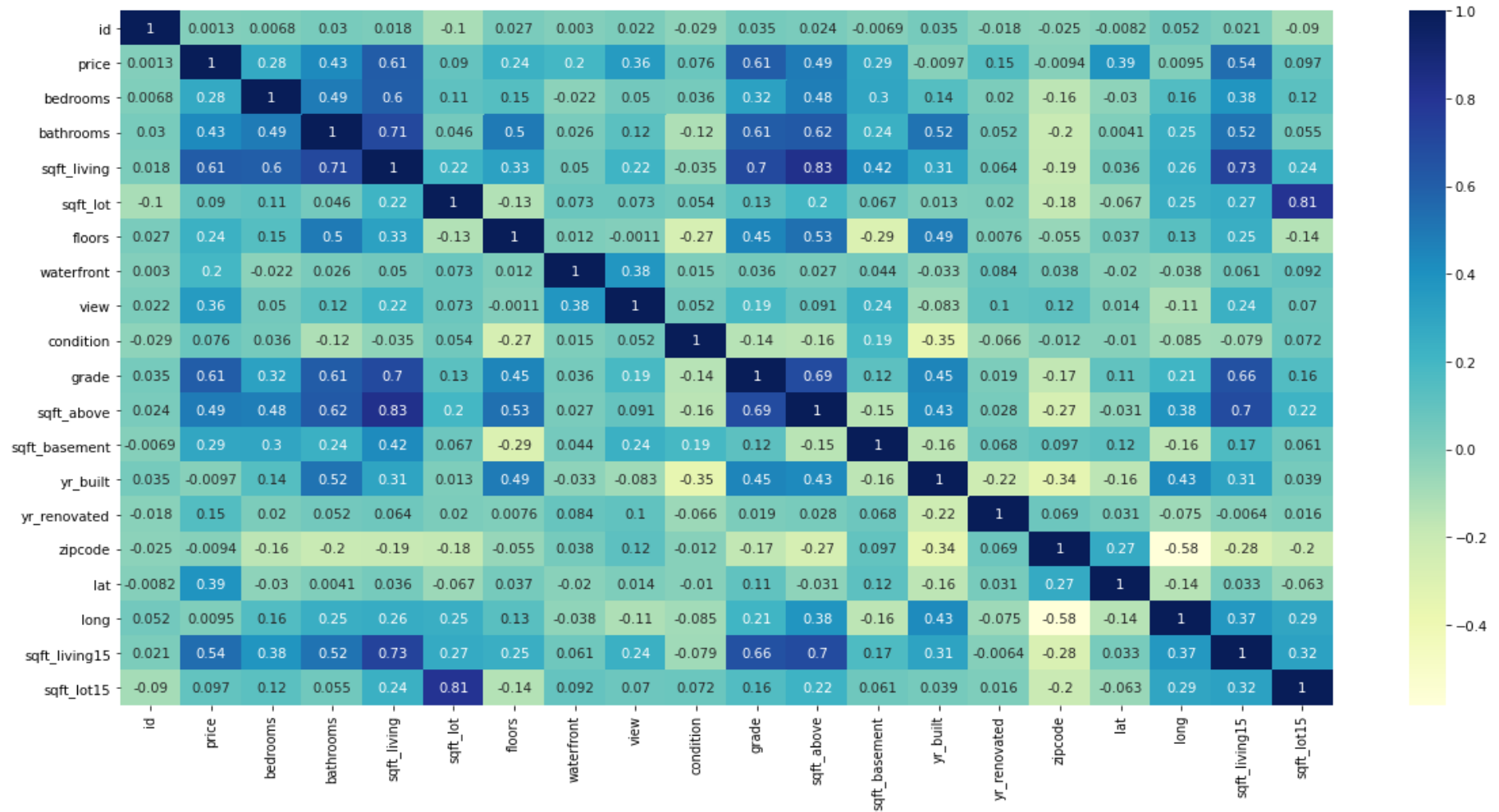
Missing Values Detection: Missing data pattern was used to identify the missing data in the dataset. From the table below it can be observed that the data does not consist of any missing data for any of the variables.



III. Data Visualization and Pattern Discovery

The objective of data visualization and pattern discovery was to reveal relationships between the house features and the response variable, price. We wanted to identify house features that affect price variable and could be potential predictors. Through visualization, we gathered the following information about the data.

Correlation Table: The below correlation table provides a summary of correlation between the continuous variables in the data. The objective behind analyzing correlation between the continuous variables in the data was to identify variables that have significant linear relationship with price and those who don't. Further, the table helps to identify relationship between potential predictors. If two predictors are highly correlated with each other they may explain the same variation in the price variable, leading to over fitting.

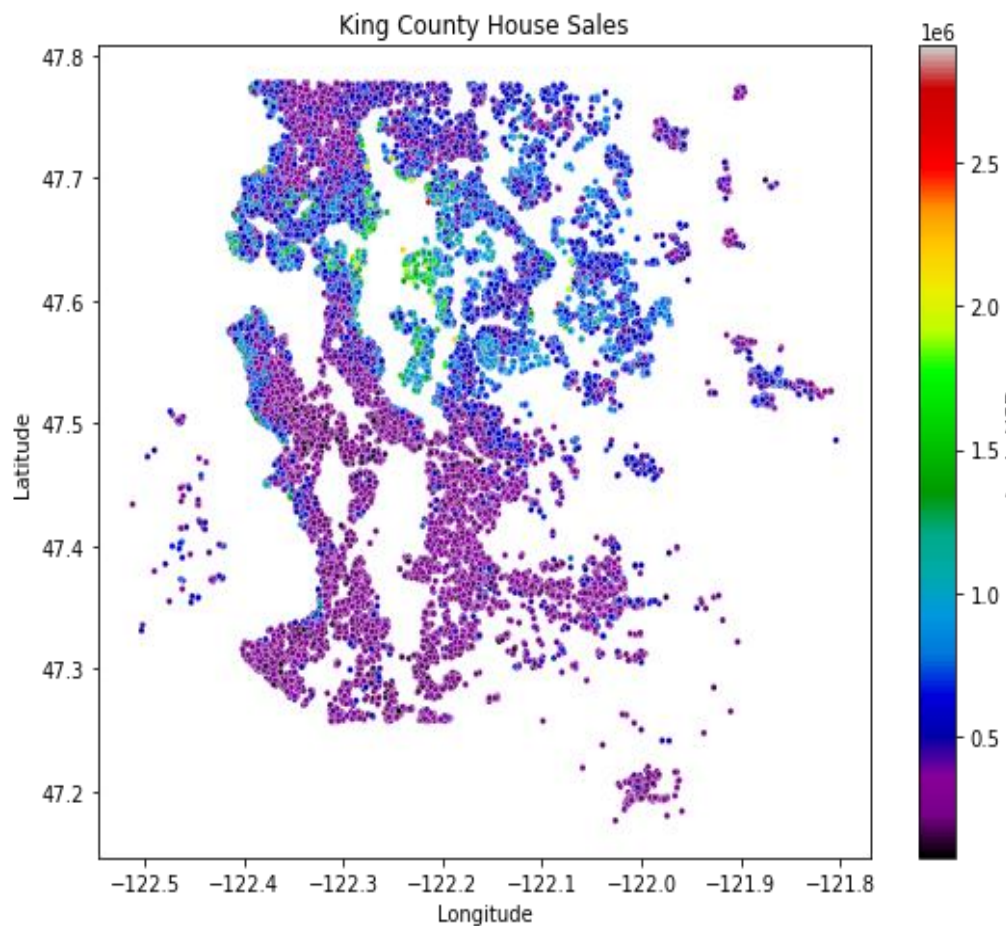


Plots for understanding or Analysis

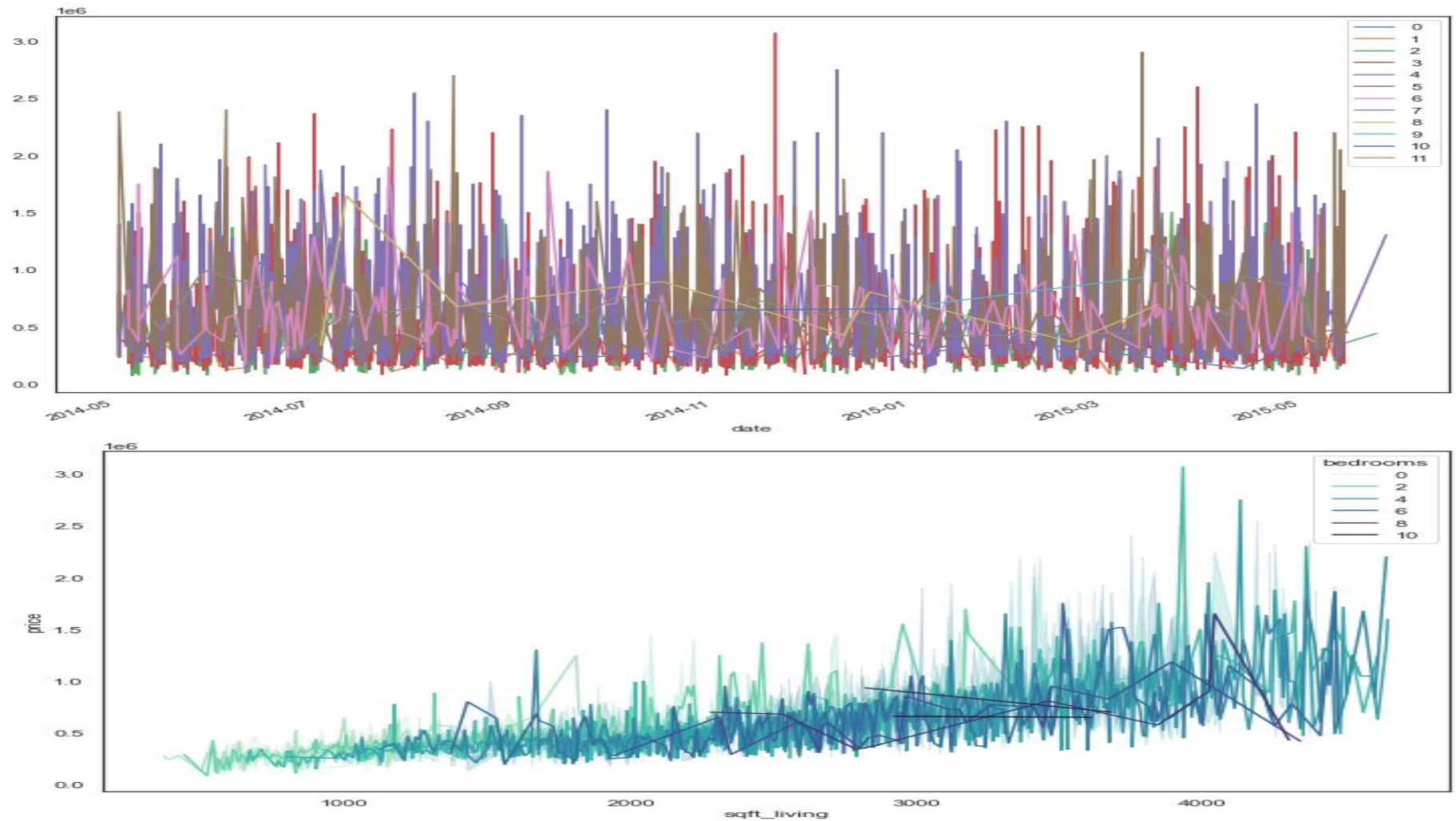
A plot is a graphical technique for representing a data set, usually as a graph showing the relationship between two or more variables. Graphs are a visual representation of the relationship between variables, which are very useful for humans who can then quickly derive an understanding which may not have come from lists of values.

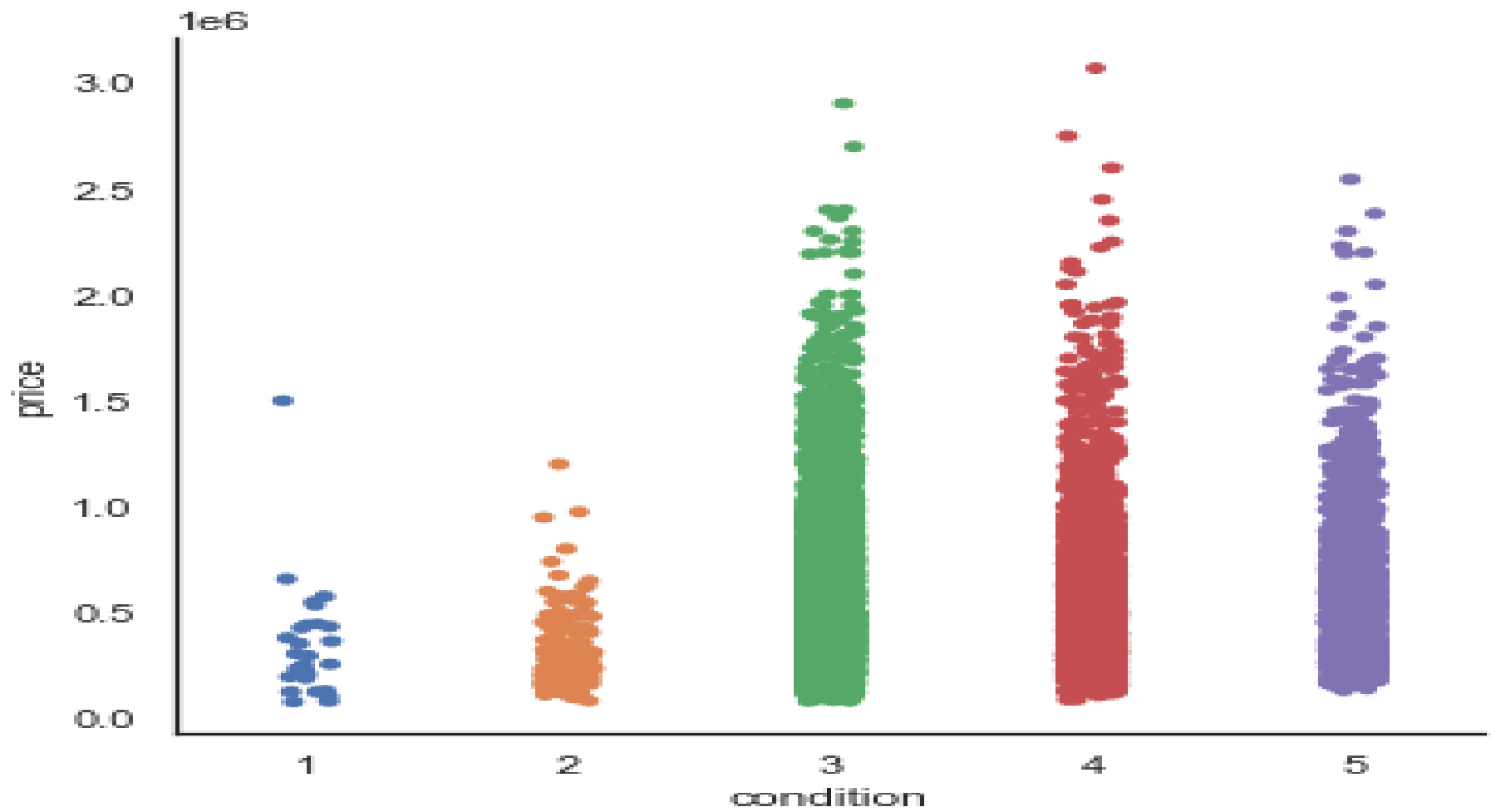
Location

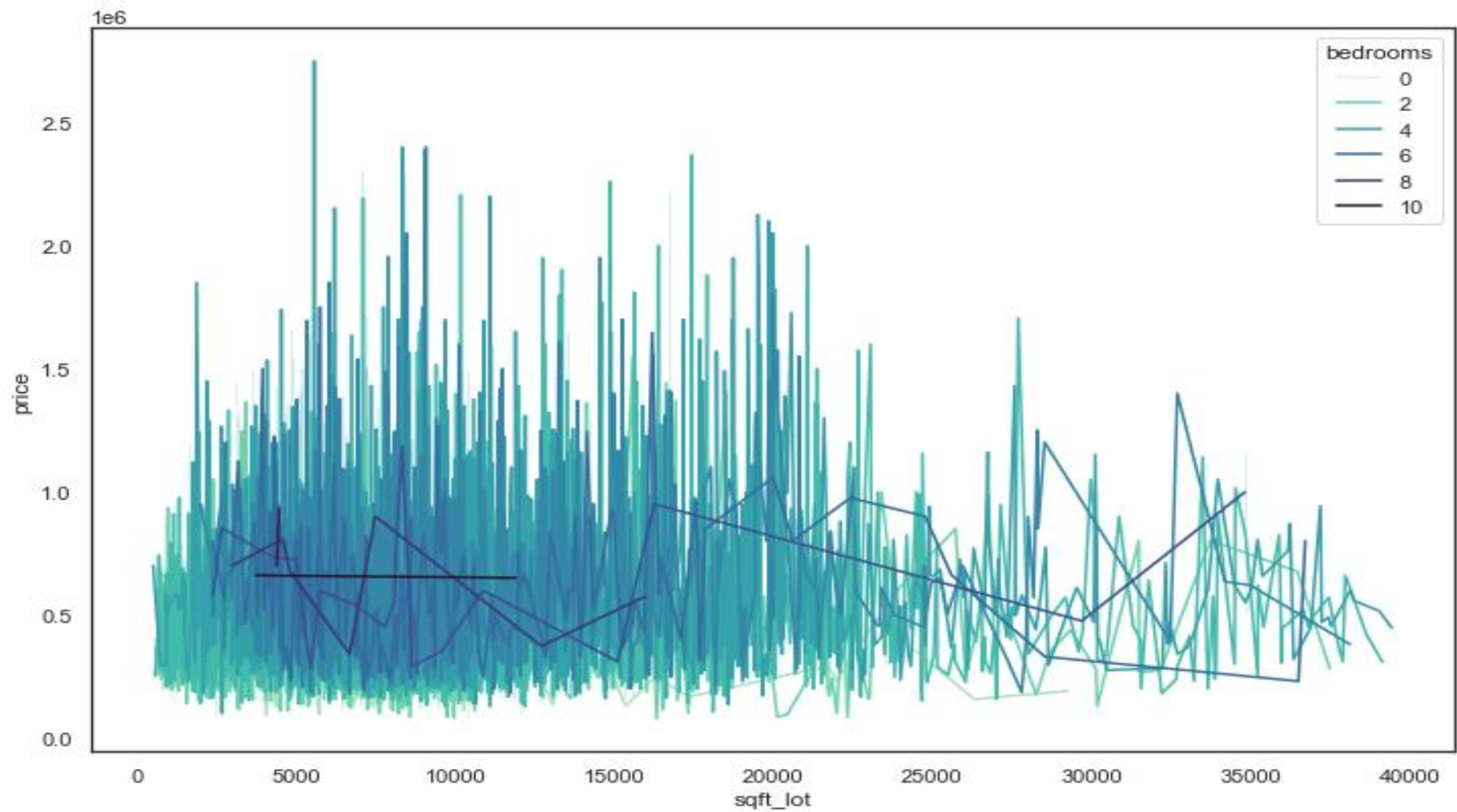
Location is key when it comes to real estate. Our first question seeks to understand the geographical distribution of the homes in our dataset and determine where the highest house sales were recorded. As a starting point, let us create a scatterplot using latitude and longitude features.

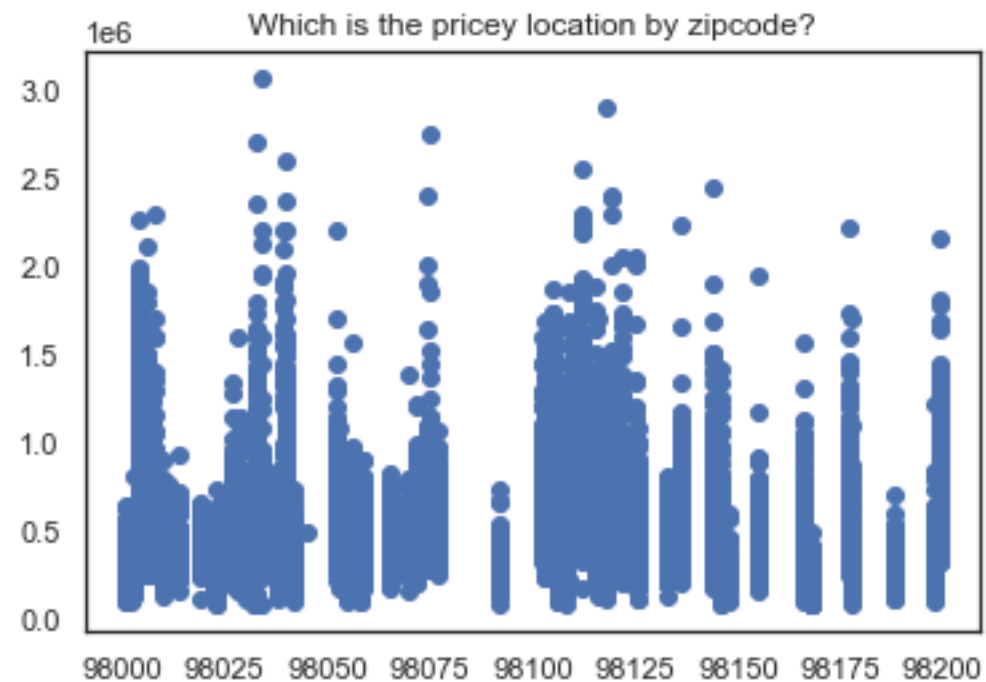
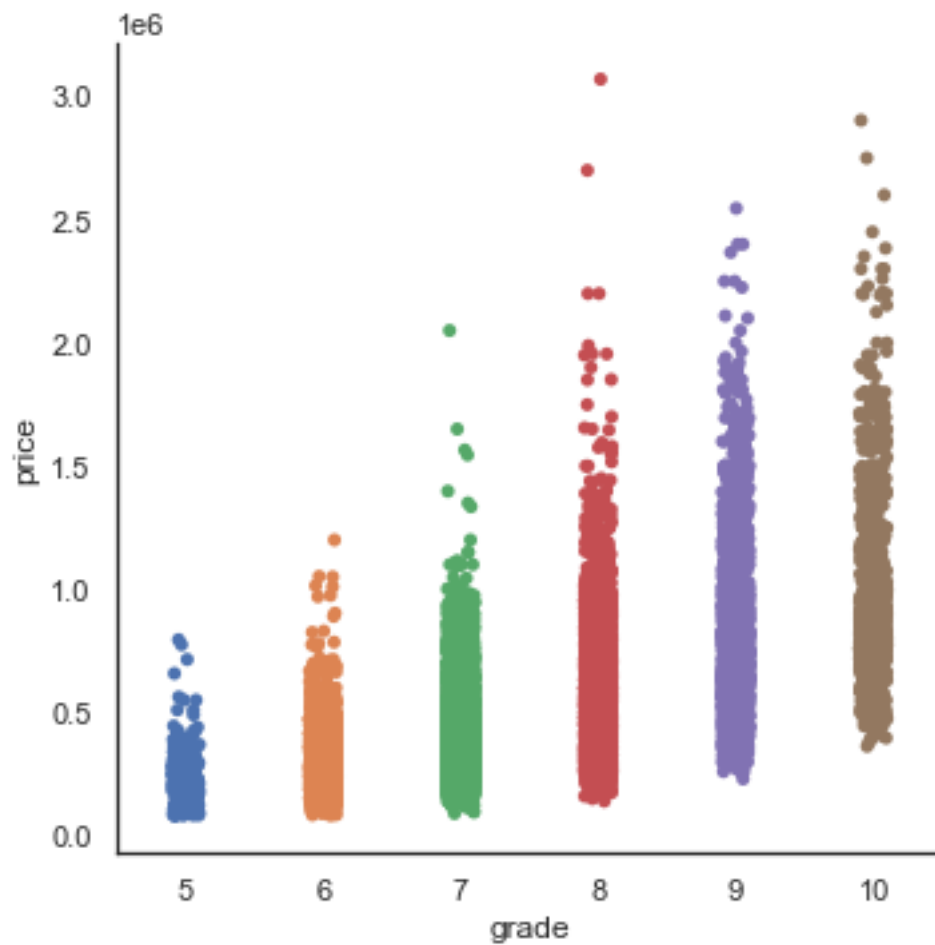


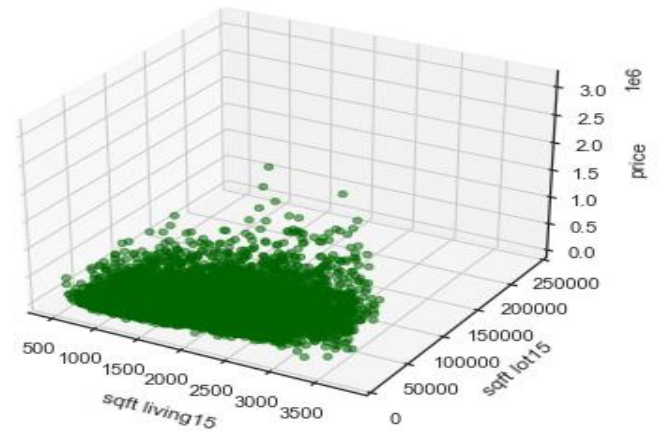
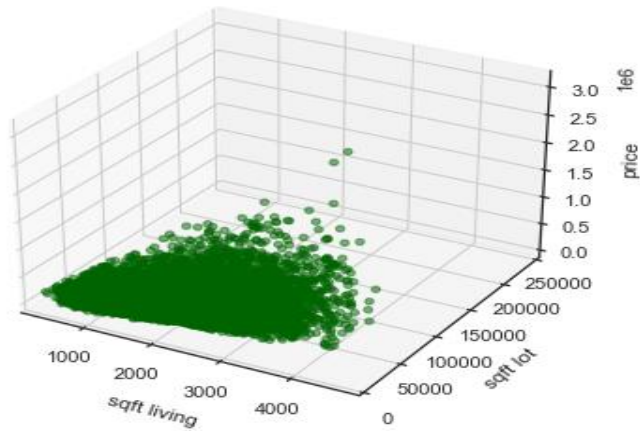
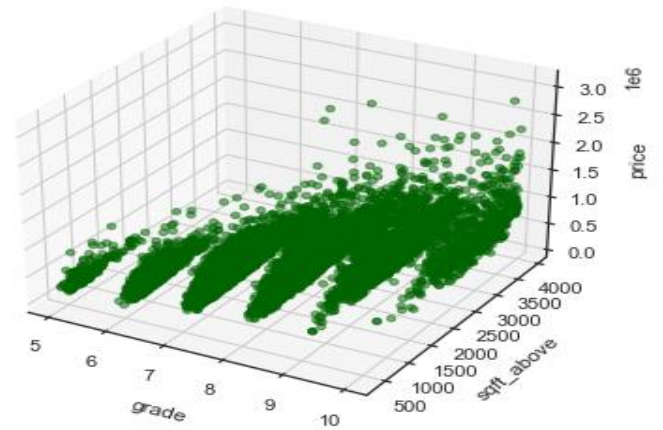
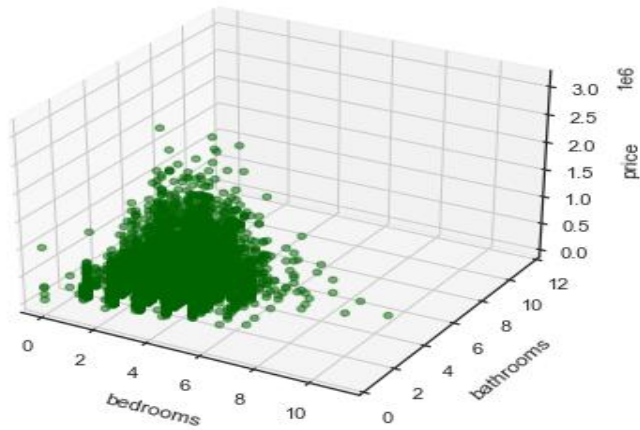
From this visualization we can already draw initial insights based on the houses' geographical locations. The highest house prices are concentrated in the area with latitude around 47.6 and longitude around -122.25. There is a disparity with southern locations achieving lower house prices.











Preliminary Observation

- (1) The frequency of no of bedrooms=3,4 is higher than any other bedrooms. The price of these are mostly similar, but some have higher price than usual because of the other features like bathrooms, location, etc.
- (2) The heatmap identifies the correlations between the features which help us in identifying how the features are dependent on each other which cannot be known by seeing the data. (Example: sqft living is dependent on grade of the house)
- (3) The highest priced houses are sold in months: 9th to 11th. This shows people tend to spend more money on houses which are having more comforts in winter.
- (4) Most of the houses have sqft living in between 500 to 6000 irrespective of no of bedrooms. The higher the living space, the higher is the cost.
- (5) Price of the house is also dependent on sqft of lot (parking) as most people own their own car.
- (6) People are tending to pay less if the condition of the house is bad. They are spending more if the house is in good condition.
- (7) The 3d plot gives relationship between multiple features.

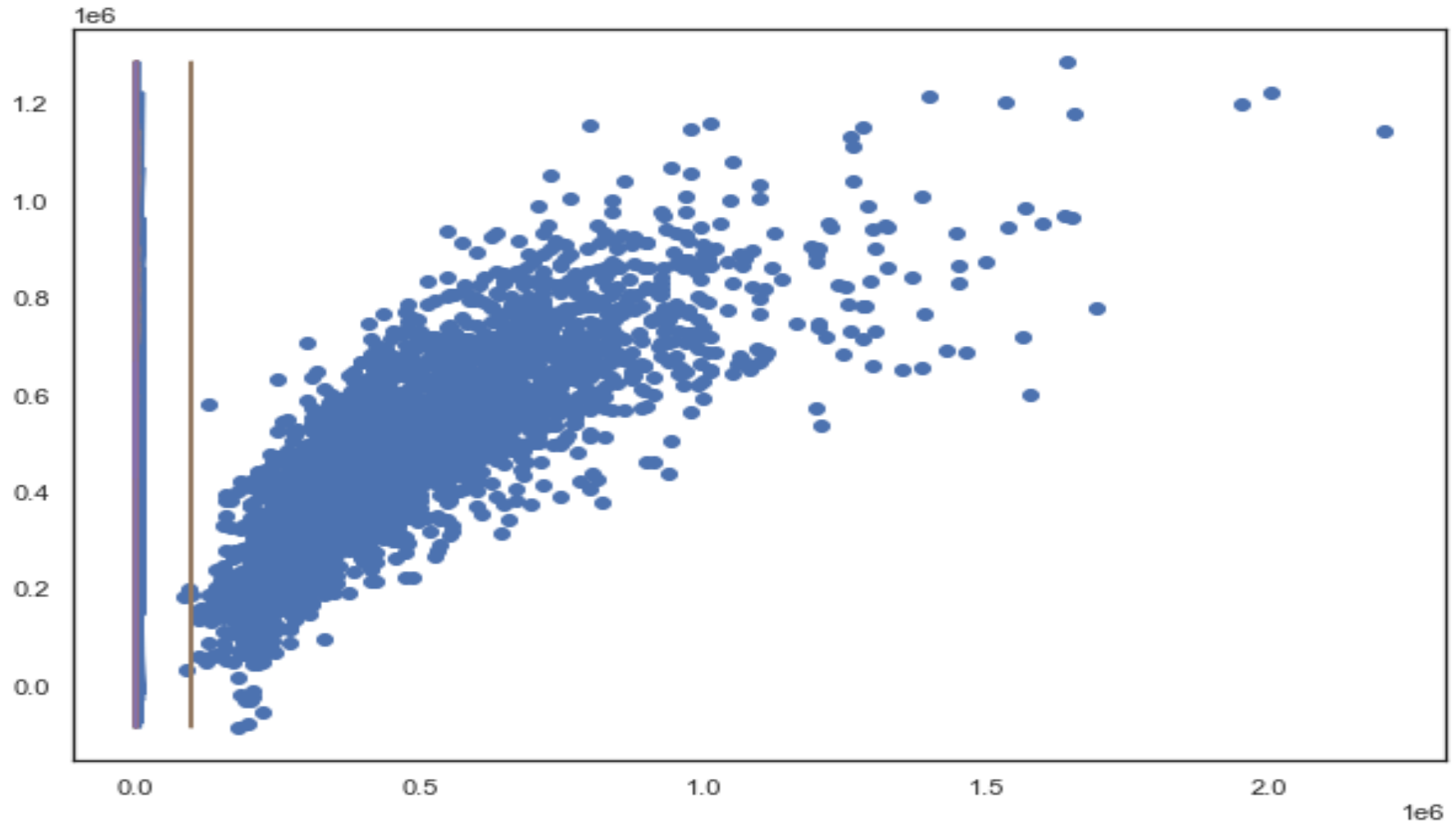
IV. Predictive Modeling

1. Models

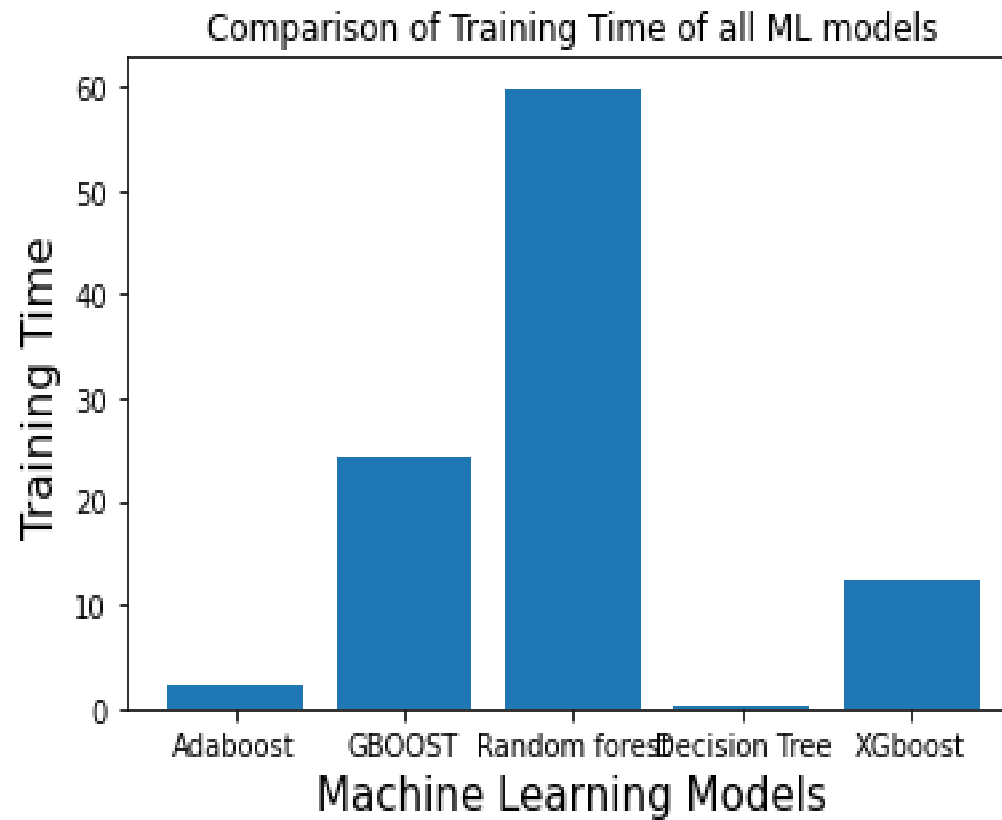
We developed the following models:

- 1. Linear regression**
- 2. Gradient Boosting Recursive Partitioning Model**
- 3. AdaBoost**
- 4. Random Forest**
- 5. 'Decision Tree**
- 6. XGBoost**

Linear Regression



Model	Score	Variance Score	R2 Score	Mean Squared Error	
0	Gradient Boosting	0.872013	0.872022	0.872013	7.623835e+09
4	XGBoost	0.870216	0.870242	0.870216	7.730887e+09
2	Random Forest	0.854403	0.854413	0.854403	8.672788e+09
3	Decision Tree	0.714179	0.714331	0.714179	1.702558e+10
1	AdaBoost	0.675359	0.696531	0.675359	1.933796e+10



A person with short brown hair, wearing a grey and white striped sweater, is seen from behind, looking at a wall covered in numerous sticky notes, diagrams, and sketches. The wall appears to be a brainstorming or project management board. The text is overlaid on the left side of the image.

CONCLUSION

Main point we see from 0 to 1.6 millions is the most frequent range so prices after this need more advertising

Gradient Boosting is best algorithm

Together for Tomorrow! **Enabling People**

Education for Future Generations

©2020 SAMSUNG. All rights reserved.

Samsung Electronics Corporate Citizenship Office holds the copyright of book.

This book is a literary property protected by copyright law so reprint and reproduction without permission are prohibited.

To use this book other than the curriculum of Samsung innovation Campus or to use the entire or part of this book, you must receive written consent from copyright holder.