

FORECASTING HOUSE PRICES WITH REGRESSION MODELS

S.NO.	NAME	EMAIL ID	ROLE
1.	Achyuth Kumar Miryala	AchyuthKumarMiryala@my.unt.edu	Coding and Presentation
2.	Rajkumar Vallepu	RajkumarVallepu@my.unt.edu	Researching and Editing
3.	Bhanu Prakash Gadala	bhanuprakashgadala@my.unt.edu	Training and Testing
4.	Sailendra Chowdary Sabbineni	Sailendrachowdarsabbineni@my.unt.edu	Determining Accuracy and Documentation
5.	Naga Swetha Gollapudi	NagaSwethaGollapudi@my.unt.edu	Data Cleaning and Data Collection

Workflow:

By splitting out the work, we were able to collaborate on the project and complete tasks like data preprocessing, visualization, and forecasting the final result. subsequently talking about the project tasks' updates and progress. This cooperative endeavor has encompassed exchanging study results and cooperating on tasks to create and enhance methods and strategies that enhance performance. We studied missing data handling methodologies. I've included those to the references within this document's section.

Abstract

The Real Estate Agency is interested in building a model that can accurately predict the sale price of a house based on its characteristics. To address this problem, we will be using the Zillow house dataset which contains information on the sale of houses in the USA.

There is an observed increasing demand for accurate and reliable estimates of the sale price of homes. However, predicting the sale price of a house can be a complex and challenging task, as it depends on various factors such as location, size, condition, and amenities. To address this challenge, the Real Estate Agency has decided to build a model that can accurately predict the sale price of a house based on its characteristics.

One of the major challenges in building such a model is identifying which features have the most impact on the sale price of a house. There may also be missing or incorrect data in the dataset, which will need to be handled appropriately. Additionally, some features may be categorical in nature (e.g., zip code) and will need to be converted to numerical values for use in the model.

We propose building a machine learning model using regression techniques to predict the sale price of a house based on its features. We will preprocess the dataset to handle missing or incorrect data and convert any categorical features to numerical values. We will then split the data into training and testing sets and use various regression techniques (such as linear regression, random forest regression, and XG Boost) to train and evaluate the model. We will select the best performing model and use it to predict the sale price of new houses.

By building this model, the Real Estate Agency will be able to provide accurate estimates of house prices to their clients. Additionally, by understanding which features have the most impact on sale price, the Real Estate Agency can provide recommendations to homeowners on how to increase the value of their homes through renovations or improvements. Overall, this project will provide valuable insights into the factors that influence the sale price of a house in the USA.

Data Specification

We'll be using the Zillow dataset.

The dataset we used is a valuable resource for understanding the real estate market in USA. The data can be used to track trends in sale prices, square footage, and other factors over time. It can also be used to compare different neighborhoods and to identify properties that are under- or overpriced.

The data has:

- **21597** rows which are the number of houses sold.
- **21** columns which represent the house's features.

The features are either strings, floats or integers.

The dataset contains:

id - unique identifier for a house

date - date the house was sold.

price - Sale price of the house (in dollars)

bedrooms - number of Bedrooms in a house

bathrooms - number of bathrooms in a house

sqft_living - square footage of the house

sqft_lot - square footage of the lot

floors - total floors (levels) in the house

waterfront - House which has a view to the waterfront.

view - Quality of view from house

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 - Year the house was built

yr_renovated - Year the house was renovated.

zip code - zip code in which the house is located.

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

Problem Statement

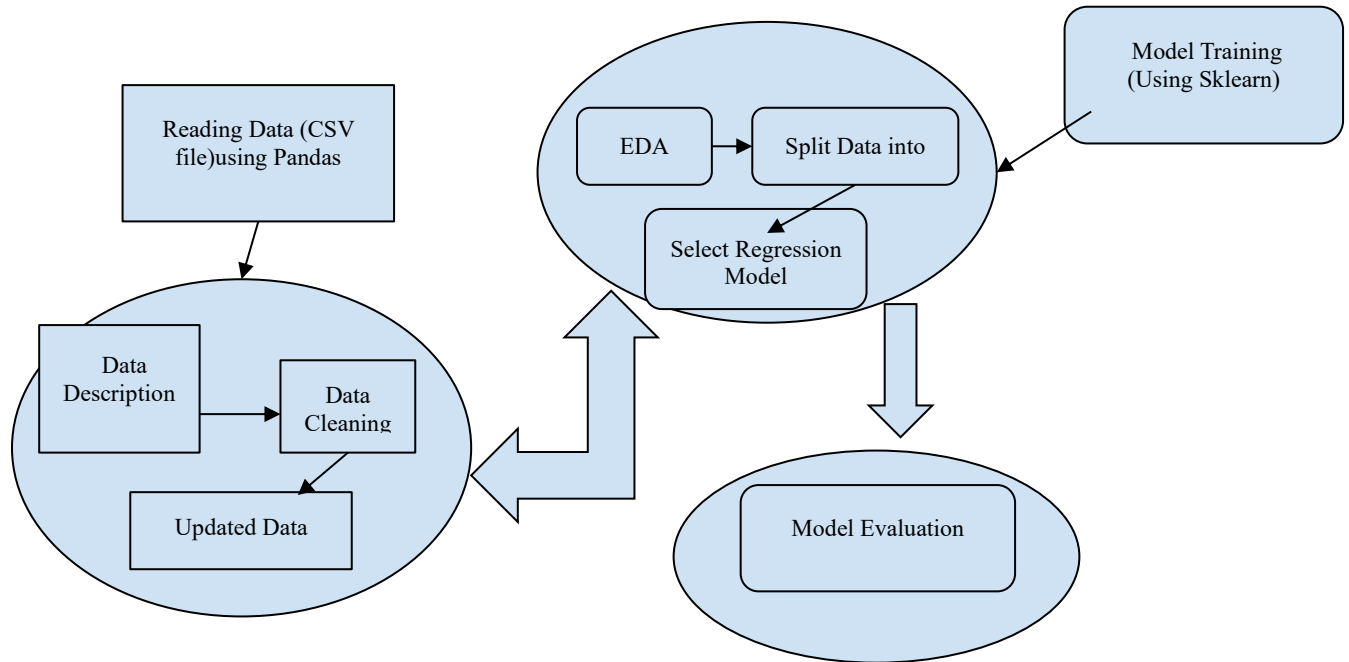
The Real Estate Agency wants to build a model that can accurately predict the sale price of a house based on its characteristics using the Zillow dataset. The agency faces the challenge of identifying which features have the most impact on the sale price of a house, handling missing or incorrect data in the dataset, and converting categorical features to numerical values.

Objectives

- To identify the features that have the most impact on the sale price of a house in USA.
- To preprocess the dataset and handle any missing or incorrect data.
- To convert any categorical features to numerical values for use in the model.
- To build a machine learning model using regression techniques to predict the sale price of a house based on its features.
- To evaluate the performance of various regression techniques and select the best performing model.
- To use the selected model to predict the sale price of new houses with high accuracy.
- To provide insights to the Real Estate Agency on how different features impact the sale price of a house and provide recommendations to homeowners on how to increase the value of their homes through renovations or improvements.

Project Design:

Flowchart of Design



Implementation of data Processing techniques:

Data Preprocessing:

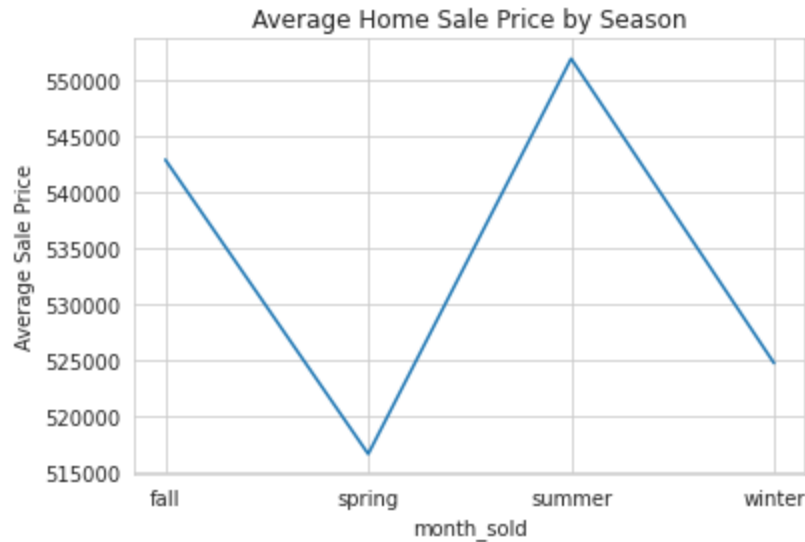
Data Cleaning and EDA

Minimal data cleaning was required, however there was a need to manage missing values, duplicates, and outliers. The exploratory data analysis(EDA) sought answers to the following questions:

How does the price vary with the month sold?

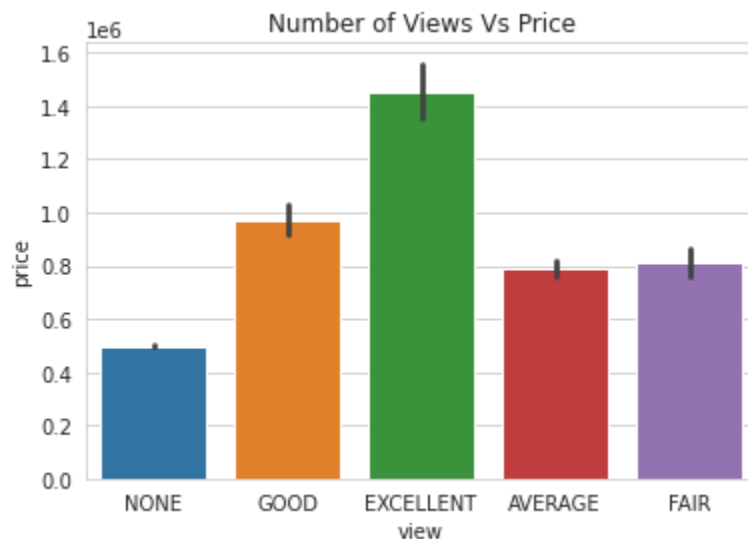
How does the view affect price?

Which factors are most correlated with price?



The above screenshot shows the home sale prices within our first time home buyers criteria by the season the house was sold. The seasons are split between spring, summer, fall and winter.

We see that prices start rising from spring to summer and start dropping during summer to spring.



From the above screenshot houses with better views have higher prices.

Based on the questions we provided, here are some possible answers that the exploratory data analysis (EDA) may have found:

How does the price vary with the month sold?

There is a seasonal effect on the home sale prices, with prices being higher during certain months and lower during others. For example, the analysis may have found that prices tend to be higher in the spring and summer months when the weather is better and people are more likely to be

interested in buying a house. On the other hand, prices may be lower in the winter months when there are fewer buyers in the market.

How does the view affect the price?

The better the view the higher the price.

Which factors are most correlated with price?

Certain factors are more strongly correlated with the sale price of a house than others. For example, the size of the house (measured in square feet) may be strongly correlated with the sale price, with larger houses selling for higher prices. Other factors that may be strongly correlated with price include the location of the house (e.g., proximity to schools or public transportation), the number of bedrooms and bathrooms, and the condition of the house.

Modeling and Evaluation

Screenshots for some of our code

```
XG Boost Regressor

[341]: from xgboost import XGBRegressor
xgb = XGBRegressor(n_estimators=1000, learning_rate=0.01)
xgb.fit(X_train, y_train)
predictions = xgb.predict(X_test)

mae, mse, rmse, r_squared = evaluation(y_test, predictions)
print("MAE:", mae)
print("MSE:", mse)
print("RMSE:", rmse)
print("R2 Score:", r_squared)
print("-"*30)
rmse_cross_val = rmse_cv(xgb)
print("RMSE Cross-Validation:", rmse_cross_val)

new_row = {"Model": "XGBRegressor", "MAE": mae, "MSE": mse, "RMSE": rmse, "R2 Score": r_squared, "RMSE (Cross-Validation)": rmse_cross_val}
models = models.append(new_row, ignore_index=True)

MAE: 153913.37405237267
MSE: 63215152224.8582
RMSE: 251426.23615060182
R2 Score: 0.5145372593211468
-----
RMSE Cross-Validation: 244244.36949240445
```

Random Forest Regressor

```
[ ]: from sklearn.ensemble import RandomForestRegressor
random_forest = RandomForestRegressor(n_estimators=100)
random_forest.fit(X_train, y_train)
predictions = random_forest.predict(X_test)

mae, mse, rmse, r_squared = evaluation(y_test, predictions)
print("MAE:", mae)
print("MSE:", mse)
print("RMSE:", rmse)
print("R2 Score:", r_squared)
print("-"*30)
rmse_cross_val = rmse_cv(random_forest)
print("RMSE Cross-Validation:", rmse_cross_val)

new_row = {"Model": "RandomForestRegressor", "MAE": mae, "MSE": mse, "RMSE": rmse, "R2 Score": r_squared, "RMSE (Cross-Validation)": rmse_cross_val}
models = models.append(new_row, ignore_index=True)
```

MAE: 159514.81383475527

MSE: 62146803995.02415

RMSE: 249292.60718084712

R2 Score: 0.5227416730004782

RMSE Cross-Validation: 247600.34588183803

Ridge Regression

```
[337]: from sklearn.linear_model import Ridge
ridge = Ridge()
ridge.fit(X_train, y_train)
predictions = ridge.predict(X_test)

mae, mse, rmse, r_squared = evaluation(y_test, predictions)
print("MAE:", mae)
print("MSE:", mse)
print("RMSE:", rmse)
print("R2 Score:", r_squared)
print("-"*30)
rmse_cross_val = rmse_cv(ridge)
print("RMSE Cross-Validation:", rmse_cross_val)

new_row = {"Model": "Ridge", "MAE": mae, "MSE": mse, "RMSE": rmse, "R2 Score": r_squared, "RMSE (Cross-Validation)": rmse_cross_val}
models = models.append(new_row, ignore_index=True)
```

MAE: 168826.60557282393

MSE: 65185562261.50544

RMSE: 255314.63385694413

R2 Score: 0.49940543375424207

RMSE Cross-Validation: 257689.51956076143

Lasso Regression

```
338]: from sklearn.linear_model import Lasso
lasso = Lasso()
lasso.fit(X_train, y_train)
predictions = lasso.predict(X_test)

mae, mse, rmse, r_squared = evaluation(y_test, predictions)
print("MAE:", mae)
print("MSE:", mse)
print("RMSE:", rmse)
print("R2 Score:", r_squared)
print("-"*30)
rmse_cross_val = rmse_cv(lasso)
print("RMSE Cross-Validation:", rmse_cross_val)

new_row = {"Model": "Lasso", "MAE": mae, "MSE": mse, "RMSE": rmse, "R2 Score": r_squared, "RMSE (Cross-Validation)": rmse_cross_val}
models = models.append(new_row, ignore_index=True)

MAE: 168826.7392325604
MSE: 65185737151.21127
RMSE: 255314.97635511175
R2 Score: 0.49940409068327063
-----
RMSE Cross-Validation: 257689.54636797262
```

Linear Regression

```
336]: from sklearn.linear_model import LinearRegression
lin_reg = LinearRegression()
lin_reg.fit(X_train, y_train)
predictions = lin_reg.predict(X_test)

mae, mse, rmse, r_squared = evaluation(y_test, predictions)
print("MAE:", mae)
print("MSE:", mse)
print("RMSE:", rmse)
print("R2 Score:", r_squared)
print("-"*30)
rmse_cross_val = rmse_cv(lin_reg)
print("RMSE Cross-Validation:", rmse_cross_val)

new_row = {"Model": "LinearRegression", "MAE": mae, "MSE": mse, "RMSE": rmse, "R2 Score": r_squared, "RMSE (Cross-Validation)": rmse_cross_val}
models = models.append(new_row, ignore_index=True)

MAE: 168826.95509923415
MSE: 65185818516.93465
RMSE: 255315.1356988744
R2 Score: 0.4994034658326505
-----
RMSE Cross-Validation: 257689.58087194673
```

Support Vector Machines

```
[342]: from sklearn.svm import SVR
svr = SVR(C=100000)
svr.fit(X_train, y_train)
predictions = svr.predict(X_test)

mae, mse, rmse, r_squared = evaluation(y_test, predictions)
print("MAE:", mae)
print("MSE:", mse)
print("RMSE:", rmse)
print("R2 Score:", r_squared)
print("-"*30)
rmse_cross_val = rmse_cv(svr)
print("RMSE Cross-Validation:", rmse_cross_val)

new_row = {"Model": "SVR", "MAE": mae, "MSE": mse, "RMSE": rmse, "R2 Score": r_squared, "RMSE (Cross-Validation)": rmse_cross_val}
models = models.append(new_row, ignore_index=True)

MAE: 164818.0298561094
MSE: 78277475481.65025
RMSE: 279781.1206669425
R2 Score: 0.39886567629270786
-----
RMSE Cross-Validation: 285141.669698903
```

Project Milestones

There are some achievements that mark significant progress towards project completion. They are listed below:

Milestones:

Milestone 1: Data Preprocessing

Timeline: 1 week

Tasks:

- Import the Zillow dataset.
- Identify and handle missing data.
- Convert categorical features to numerical values.
- Exploratory data analysis (EDA).

Milestone 2: Feature Selection

Timeline: 1 week

Tasks:

- Perform feature selection techniques to identify the most relevant features for predicting house prices.
- Analyze the relationship between features and house prices.
- Evaluate the impact of feature selection on model performance.

Milestone 3: Model Training and Evaluation

Timeline: 2 weeks

Tasks:

- Train various regression models (e.g., linear regression, random forest regression, XG Boost) using the selected features.
- Evaluate the performance of each model using metrics such as mean squared error (MSE) and R-squared.
- Select the best performing model based on evaluation metrics.

Incremental Features:

Incorporate additional data sources:

- Data sources: Real estate market trends, local economic indicators, demographic data.
- Benefits: Improve the accuracy and comprehensiveness of the model.

We will use Linear Regression, Random Forest Regression, Ridge Regression, Lasso Regression, Support Vector Regression and XG Boost Regression for our modeling and select which best suits our data set.

Based on the provided evaluation results, here's a summary of the performance metrics for each model:

Based on the RMSE Cross Validation metric, the XG Boost Regressor and Random Forest Regression model has the lowest value, followed by the Linear Regression and Ridge Regression models. The Linear Regression and SVR Regression models have relatively higher RMSE cross validation values.

Based on the RSME metric, the XG Boost Regression model and Random Forest Regression model has the lowest value, followed by the Random Forest Regression model. Linear Regression, Ridge Regression, and Lasso Regression models have relatively higher RSME values.

Based on the provided evaluation results, the Random Forest Regression model and XGBoost Regression model generally have lower MSE and MAE values compared to the other models. This suggests that they may have better predictive performance.

Regression Models	MAE	MSE	RSME	R2 Score	RSME (Cross-Validation)
XG Boost	153913.374052	6.321515e+10	251426.236151	0.514537	244244.369492
Random Forest	159772.297337	6.265113e+10	250302.083273	0.518869	247275.352598
Ridge	168826.605573	6.518556e+10	255314.633857	0.499405	257689.519561
Lasso	168826.739233	6.518574e+10	255314.976355	0.499404	257689.546368
Linear	168826.955099	6.518582e+10	255315.135699	0.499403	257689.580872
SVR	164818.029856	7.827748e+10	279781.120667	0.398866	285141.669699

Result and Conclusion:

1. Based on the evaluation results, the Random Forest Regression model and XGBoost Regression model show promising performance in terms of MSE and MAE. Further analysis and comparison of these models, considering additional factors, can help determine the most suitable model for your dataset and prediction task.
2. The development of an accurate house price prediction model can provide substantial benefits to a Real Estate Agency and its clients. It can enhance customer satisfaction by providing more precise price estimates and valuable insights into the factors influencing property values. Additionally, the model can support informed decision-making regarding renovations or improvements, leading to increased sales and improved outcomes for homeowners and the agency alike.

Repository /Archive:

https://github.com/Achyuth123456/Group-7_5502.git

Appendix:

Code

```
3. #import important libraries
4. import pandas as pd
5. import numpy as np
6. import matplotlib.pyplot as plt
7. import seaborn as sns
8. sns.set_style("whitegrid")
9. %matplotlib inline
10.
11. from scipy import stats
12. import statsmodels.api as sm
13. import statsmodels.formula.api as smf
14. from statsmodels.formula.api import ols
15. import statsmodels.stats.api as sms
16. from scipy.stats import norm
17. from sklearn.metrics import mean_absolute_error
18. from sklearn.metrics import mean_squared_error
19. from sklearn.metrics import r2_score
20. from sklearn.model_selection import cross_val_score
21. import warnings
22. warnings.filterwarnings(action='ignore', category=UserWarning)
23. import folium
24. #code to display all the columns without truncation
25. pd.set_option('display.max_columns', None)
26. # Reading the dataset
27. df = pd.read_csv("zillow_house_data.csv")
28. #Display the first few rows in the dataframe
29. df.head()
30. #Display the last few rows in the dataframe
31. df.tail()
32. #Display of the number of rows and columns in the dataframe
33. df.shape
34. #Summary of the dataframe
35. df.info()
36. #plotting histograms
37. def histo(df):
38.     return df.hist(figsize=(20,20))
39. histo(df);
40. #summary statistics of the numerical columns in a dataframe
41. def data_description(df):
42.     return (df.drop('id', axis = 1)).describe()
43. data_description(df)
44. list(df['sqft_basement'][:10])
45. def clean_basement(col):
46.     df[col] = df[col].replace('?', '0.0')
47.     df[col] = df[col].astype(float)
```

```

48. print(list(df['sqft_basement'][:10]))
49. clean_basement('sqft_basement')
50. def renaming_columns(df,cols_old,cols_new):
51.     return df.rename(columns={cols_old: cols_new
52.         }, inplace=True)
53. renaming_columns(df,'date','date_sold')
54. #The following function returns all columns with missing values alongside the quantity and percentage
55.
56. def missing_values(data):
57.     """A simple function to identify data has missing values"""
58.     # identify the total missing values per column
59.     # sort in order
60.     miss = data.isnull().sum().sort_values(ascending = False)
61.
62.     # calculate percentage of the missing values
63.     percentage_miss = (round((data.isnull().sum() / len(data)*100),2)).sort_values(ascending = False)
64.
65.     # store in a dataframe
66.     missing = pd.DataFrame({"Missing Values": miss, "Percentage(%)": percentage_miss})
67.     # remove values that are missing
68.     missing.drop(missing[missing["Percentage(%)"] == 0].index, inplace = True)
69.
70.     return missing
71.
72.
73. missing_data = missing_values(df)
74. missing_data
75. df['yr_renovated'].value_counts()
76. df['yr_renovated'].fillna(value= 0.0, inplace=True)
77. df['view'].fillna(value='NONE', inplace=True)
78. #Count of the unique values in the waterfront column
79. df.waterfront.value_counts(normalize=True)
80. # replacing the missing values using their probability
81. def replace_missing_val(df,cols):
82.     df[cols] = df[cols].fillna(value=(np.random.choice(['YES', 'NO'], p=[0.01, 0.99])))
83.     # checking if the missing values have been replaced
84.     print('There are now', df[cols].isnull().sum(), 'missing values in waterfront column')
85. replace_missing_val(df,'waterfront')
86. #finding total number of duplicates
87. # Duplicated entries
88. def identify_duplicates(df):
89.     """Simple function to identify any duplicates"""
90.     # identify the duplicates (dataframe.duplicated() , can add .sum() to get total count)
91.     # empty list to store Bool results from duplicated
92.     duplicates = []
93.     for i in df.duplicated():
94.         duplicates.append(i)
95.     # identify if there is any duplicates. (If there is any we expect a True value in the list duplicates)
96.     duplicates_set = set(duplicates)

```

```

97.     if (len(duplicates_set) == 1):
98.         print("The Data has no duplicates")
99.     else:
100.         no_true = 0
101.         for val in duplicates:
102.             if (val == True):
103.                 no_true += 1
104.         # percentage of the data represented by duplicates
105.         duplicates_percentage = np.round(((no_true / len(df)) * 100), 3)
106.         print(f'The Data has {no_true} duplicated rows.\nThis constitutes {duplicates_percentage}% of the
            data set.")
107. identify_duplicates(df)
108. # checking for duplicates in unique columns
109. def unique_column_duplicates(df, column):
110.     df.id.duplicated().sum()
111.     df = pd.DataFrame(df)
112. # count number of duplicated values in the 'id' column
113. num_duplicates = df['id'].duplicated().sum()
114. print('Number of duplicates in "id" column:', num_duplicates)
115. # Lets see the rows that have duplicates as per the id column
116. pd.concat( g for _, g in df.groupby("id") if len(g) > 1).head()
117. # summary statistics of the numerical columns in a dataframe
118. df.describe()
119. # we can also use a box plot to identify outliers
120. def create_boxplot(dataframe, x_col):
121.     sns.boxplot(data=dataframe, x=x_col)
122. create_boxplot(df, "bedrooms");
123. # let's count the number of occurrences of each unique value in the 'bedrooms' column
124. df['bedrooms'].value_counts()
125. # Let's only show the row where the value in the 'bedrooms' column is equal to 33.
126. df[df['bedrooms'] == 33]
127. # Let's drop the 33 bedroom row
128. df.drop(index=15856, inplace=True)
129.
130. # And verify that it is gone
131. df['bedrooms'].value_counts()
132. # Checking for outliers in bathrooms
133. df['bathrooms'].value_counts()
134. df.isnull().sum()
135. # Extracting an individual month column
136. df['month_sold'] = pd.DatetimeIndex(df['date_sold']).month
137.
138. # set columns for seasons sold
139. df.loc[(df['month_sold'] >= 1) & (df['month_sold'] <= 3), 'season'] = 'spring'
140. df.loc[(df['month_sold'] >= 3) & (df['month_sold'] <= 6), 'season'] = 'summer'
141. df.loc[(df['month_sold'] >= 6) & (df['month_sold'] <= 9), 'season'] = 'fall'
142. df.loc[(df['month_sold'] >= 9) & (df['month_sold'] == 12), 'season'] = 'winter'
143. # Plotting a line graph
144. plt.plot(df.groupby('season')['price'].mean().round(2))

```

```

145.plt.title('Average Home Sale Price by Season')
146.plt.ylabel('Average Sale Price')
147.plt.xlabel('month_sold');
148.#Plotting a Scatter plot for age and price
149.year_sold = pd.DatetimeIndex(df['date_sold']).year
150.age = year_sold - df['yr_built']
151.df['age'] = age
152.plt.scatter(age, df['price'])
153.plt.xlabel('Age')
154.plt.ylabel('Price')
155.plt.title('Relationship between Age and Price');
156.#Extracting value counts
157.df['view'].value_counts()
158.#Plotting a barplot of 'view' against 'price'
159.sns.barplot(x=df['view'], y=df['price']).set(title='Number of Views Vs Price');
160.df['condition'].value_counts()
161.plt.scatter(df['condition'], df['price'], color='brown')
162.plt.title('Relationship between price and condition')
163.plt.xlabel('Condition rating')
164.plt.ylabel('Price')
165.plt.grid(color='purple',
166.         alpha=0.1,
167.         linestyle='-.',
168.         linewidth=2);
169.# Lets create a heatmap to check for correlations
170.sns.pairplot(data=df, x_vars=['sqft_lot', 'sqft_above', 'sqft_living', 'sqft_living15', 'sqft_lot15'],
171.              y_vars=["price"]);
172.corr = df.corr()
173.plt.figure(figsize=(7,6))
174.mask = np.triu(np.ones_like(corr, dtype=bool))
175.sns.heatmap(df.corr(), cmap= 'coolwarm', mask= mask, linewidth= 1)
176.plt.title("Correlation between features",weight='bold',fontsize=18)
177.plt.show()
178.
179.fig, ax = plt.subplots(figsize=(15,8))
180.ax= sns.heatmap(df.corr(), annot=True, cmap="viridis")
181.plt.title("Correlation between features",weight='bold',fontsize=18);
182.#The following function displays columns that are highly correlated with one other
183.
184.def check_multiCorr(df):
185.    num_cols = df.select_dtypes(include=['float64','int64'])
186.    ncc = list(num_cols.columns)
187.    for col1 in ncc:
188.        for col2 in ncc:
189.            if col1 != col2:
190.                correlation = df[col1].corr(df[col2])
191.                if correlation > 0.75:
192.                    print(col1, " is correlated with ", col2, "by", correlation)

```



```

193.         ncc.remove(col1)
194.check_multiCorr(df)
195.df_num = df.select_dtypes(include=['float64', 'int64'])
196.df_cat = df.select_dtypes(include=['object'])
197.df['view'] = df['view'].replace(regex='GOOD', value="AVERAGE")
198.df['view'] = df['view'].replace(regex='NONE', value="POOR")
199.df['view'] = df['view'].replace(regex='FAIR', value="AVERAGE")
200.df['view'].hist();
201.def bin_column(data, col):
202.    list_years= list(df[col].unique())
203.    for x in list_years:
204.        if x !=0:
205.            df[col] = df[col].replace(x, "YES")
206.        else:
207.            df[col] = df['yr_renovated'].replace(x,"NO")
208.
209.bin_column(df, 'yr_renovated')
210.fig, ax = plt.subplots(figsize=(4,5))
211.ax.hist(df['yr_renovated'], color='purple');
212.ax.set_title("Renovated vs Not Renovated");
213.df_cat.head()
214.X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
215.def rmse_cv(model):
216.    rmse = np.sqrt(-cross_val_score(model, X, y, scoring="neg_mean_squared_error", cv=5)).mean()
217.    return rmse
218.
219.
220.def evaluation(y, predictions):
221.    mae = mean_absolute_error(y, predictions)
222.    mse = mean_squared_error(y, predictions)
223.    rmse = np.sqrt(mean_squared_error(y, predictions))
224.    r_squared = r2_score(y, predictions)
225.    return mae, mse, rmse, r_squared
226.models = pd.DataFrame(columns=["Model", "MAE", "MSE", "RMSE", "R2 Score", "RMSE (Cross-
    Validation)"])
227.from sklearn.linear_model import LinearRegression
228.lin_reg = LinearRegression()
229.lin_reg.fit(X_train, y_train)
230.predictions = lin_reg.predict(X_test)
231.
232.mae, mse, rmse, r_squared = evaluation(y_test, predictions)
233.print("MAE:", mae)
234.print("MSE:", mse)
235.print("RMSE:", rmse)
236.print("R2 Score:", r_squared)
237.print("-"*30)
238.rmse_cross_val = rmse_cv(lin_reg)
239.print("RMSE Cross-Validation:", rmse_cross_val)
240.

```

```

241.new_row = {"Model": "LinearRegression", "MAE": mae, "MSE": mse, "RMSE": rmse, "R2 Score":
    r_squared, "RMSE (Cross-Validation)": rmse_cross_val}
242.models = models.append(new_row, ignore_index=True)
243.from sklearn.linear_model import Ridge
244.ridge = Ridge()
245.ridge.fit(X_train, y_train)
246.predictions = ridge.predict(X_test)
247.
248.mae, mse, rmse, r_squared = evaluation(y_test, predictions)
249.print("MAE:", mae)
250.print("MSE:", mse)
251.print("RMSE:", rmse)
252.print("R2 Score:", r_squared)
253.print("-"*30)
254.rmse_cross_val = rmse_cv(ridge)
255.print("RMSE Cross-Validation:", rmse_cross_val)
256.
257.new_row = {"Model": "Ridge", "MAE": mae, "MSE": mse, "RMSE": rmse, "R2 Score": r_squared, "RMSE
    (Cross-Validation)": rmse_cross_val}
258.models = models.append(new_row, ignore_index=True)
259.from sklearn.linear_model import Lasso
260.lasso = Lasso()
261.lasso.fit(X_train, y_train)
262.predictions = lasso.predict(X_test)
263.
264.mae, mse, rmse, r_squared = evaluation(y_test, predictions)
265.print("MAE:", mae)
266.print("MSE:", mse)
267.print("RMSE:", rmse)
268.print("R2 Score:", r_squared)
269.print("-"*30)
270.rmse_cross_val = rmse_cv(lasso)
271.print("RMSE Cross-Validation:", rmse_cross_val)
272.
273.new_row = {"Model": "Lasso", "MAE": mae, "MSE": mse, "RMSE": rmse, "R2 Score": r_squared, "RMSE
    (Cross-Validation)": rmse_cross_val}
274.models = models.append(new_row, ignore_index=True)
275.from sklearn.ensemble import RandomForestRegressor
276.random_forest = RandomForestRegressor(n_estimators=100)
277.random_forest.fit(X_train, y_train)
278.predictions = random_forest.predict(X_test)
279.
280.mae, mse, rmse, r_squared = evaluation(y_test, predictions)
281.print("MAE:", mae)
282.print("MSE:", mse)
283.print("RMSE:", rmse)
284.print("R2 Score:", r_squared)
285.print("-"*30)
286.rmse_cross_val = rmse_cv(random_forest)

```

```

287.print("RMSE Cross-Validation:", rmse_cross_val)
288.
289.new_row = {"Model": "RandomForestRegressor", "MAE": mae, "MSE": mse, "RMSE": rmse, "R2 Score":
    r_squared, "RMSE (Cross-Validation)": rmse_cross_val}
290.models = models.append(new_row, ignore_index=True)
291.from xgboost import XGBRegressor
292.xgb = XGBRegressor(n_estimators=1000, learning_rate=0.01)
293.xgb.fit(X_train, y_train)
294.predictions = xgb.predict(X_test)
295.
296.mae, mse, rmse, r_squared = evaluation(y_test, predictions)
297.print("MAE:", mae)
298.print("MSE:", mse)
299.print("RMSE:", rmse)
300.print("R2 Score:", r_squared)
301.print("-"*30)
302.rmse_cross_val = rmse_cv(xgb)
303.print("RMSE Cross-Validation:", rmse_cross_val)
304.
305.new_row = {"Model": "XGBRegressor", "MAE": mae, "MSE": mse, "RMSE": rmse, "R2 Score":
    r_squared, "RMSE (Cross-Validation)": rmse_cross_val}
306.models = models.append(new_row, ignore_index=True)
307.from sklearn.svm import SVR
308.svr = SVR(C=1000000)
309.svr.fit(X_train, y_train)
310.predictions = svr.predict(X_test)
311.
312.mae, mse, rmse, r_squared = evaluation(y_test, predictions)
313.print("MAE:", mae)
314.print("MSE:", mse)
315.print("RMSE:", rmse)
316.print("R2 Score:", r_squared)
317.print("-"*30)
318.rmse_cross_val = rmse_cv(svr)
319.print("RMSE Cross-Validation:", rmse_cross_val)
320.
321.new_row = {"Model": "SVR", "MAE": mae, "MSE": mse, "RMSE": rmse, "R2 Score": r_squared, "RMSE
    (Cross-Validation)": rmse_cross_val}
322.models = models.append(new_row, ignore_index=True)
323.models.sort_values(by="RMSE (Cross-Validation)")
324.plt.figure(figsize=(12,8))
325.colors = ['aquamarine', 'yellow', 'red', 'pink', 'blue', 'green']
326.sns.barplot(x=models["Model"], y=models["RMSE (Cross-Validation)"], palette=colors, alpha = 0.8)
327.
328.plt.title("Models' RMSE Scores (Cross-Validated)", size=15)
329.plt.xticks(rotation=30, size=12)
330.plt.show()
331.

```

References

- We took Dataset from 'Zillow'.
<https://www.zillow.com>
- We learned Data Cleaning and Preprocessing from 'medium' website.
<https://medium.com/analytics-vidhya/data-cleaning-and-preprocessing-a4b751f4066f>
- We referenced our project from 'the clever programmer' website.
<https://thecleverprogrammer.com/2022/08/15/house-rent-prediction-with-machine-learning/>
- We referenced the Numpy and Pandas syntax from 'codecademy' website for our project.
<https://www.codecademy.com/article/introduction-to-numpy-and-pandas>

Future Work

- Improve accuracy by reducing noise in the data to improve the accuracy of our predictions.
- Look at trends: which house attributes contribute to a higher priced home?
- Simplify model: group low-impact categorical variables. Identify categorical variables that have a minimal impact on the predictions and explore options to group them together or remove them from the model.