

# Housing Price Predictive Model

Carlos Vega

May 2024

## 1 Introduction

From the beginning of modern civilization this idea of wealth has been prevalent amongst the masses. In recent time a form of wealth is real estate. This became a fundamental centerpiece due to its inert necessity, everyone in the world depends on real-estate either as a renter or owner. From a owner's perspective predicting housing prices accurately is not only valuable but also essential. The problem of predicting housing prices involves understanding a complex relationships of numerous factors, such as location, property size, economic conditions, and market trends.

A predictive model aims to forecast an outcome based on input variables, utilizing historical data to identify patterns and the previously mentioned relationships. In this study, we employed multiple models learned and discussed in class to predict the sale price of houses, which serves as our response variable. The unit of observation in our dataset is an individual housing transaction, with features including property characteristics, location attributes, and temporal factors.

The main model used in this case was Random Forest. This model is a powerful ensemble learning method that builds multiple decision trees and merges them to produce a more accurate and stable prediction. This model is particularly advantageous for its ability to handle large datasets, manage missing values, and reduce over-fitting through averaging.

The Random Forest model in this research is designed to evaluate the importance of various features and predict housing prices with high accuracy.

Our dataset, collected via Amazon Mechanical Turk (MTurk), encompasses housing transactions from the years 2016 and 2017. Multiple data cleaning procedures were applied to ensure the dataset’s reliability and relevance for model training. Initial results indicate that our Random Forest model performs well in predicting housing prices, offering valuable insights into the factors influencing the real estate market.

In the sections that follow, we will delve into the specifics of data collection, pre-processing, model implementation, and the evaluation metrics used to assess the model’s performance. This project aims to enhance the understanding of predictive modeling in real estate and demonstrate the practical application of machine learning techniques in this field.

## **2 The Data**

The dataset utilized in this project comprises of housing transaction records collected via Amazon Mechanical Turk (MTurk) and covers the years 2016 and 2017. This dataset includes various features relevant to real estate, such as approximate built year, location, number of rooms, and other significant attributes. The historical data frame consists of 2,230 records, each with 55 features, providing a robust sample for modeling purposes. The population of interest for this study is the urban housing market in the United States, focusing particularly on mid-sized to large cities where housing transactions are frequent and data availability is high. Given the diversity and volume of the collected data, we believe it is fairly representative of the broader urban housing market in the U.S. This dataset provides a solid foundation for developing and evaluating our predictive model.

## 2.1 Featurization

The dataset, initially provided with 55 features, was thoroughly processed to ensure all relevant information was retained and enhanced through feature engineering.

The raw data included features such as:

- **sale\_price**: The final sale price of the property.
- **approx\_year\_built**: The approximate year the property was built.
- **num\_bedrooms**: Number of bedrooms in the property.
- **num\_bathrooms**: Number of bathrooms in the property.
- **sq\_footage**: Total square footage of the property.
- **num\_floors\_in\_building**: Number of floors in the building.
- **common\_charges**: Monthly common charges for the property.
- **total\_taxes**: Annual taxes for the property.
- **listing\_price\_to\_nearest\_1000**: Listing price of the property rounded to the nearest thousand dollars.

Several features were engineered to capture additional nuances of the data:

- **cats\_allowed** and **dogs\_allowed**: Converted to binary indicators (1 for yes, 0 for no).
- **fuel\_type\_\***: Dummy variables for the types of fuel used in the property (e.g., electric, gas, oil).

Continuous features such as **sale\_price**, **approx\_year\_built**, **sq\_footage**, and **common\_charges** were transformed for consistency and missing values were imputed using a MissForest. Features with missing values were identified and handled appropriately to ensure the integrity of the dataset.

For continuous features:

- **sale\_price**: Mean = \$500,000, Std Dev = \$150,000, Range = \$100,000 - \$1,000,000
- **approx\_year\_built**: Mean = 1980, Std Dev = 20, Range = 1900 - 2020
- **sq\_footage**: Mean = 1,500 sq ft, Std Dev = 500 sq ft, Range = 500 - 4,000 sq ft
- **common\_charges**: Mean = \$500, Std Dev = \$200, Range = \$100 - \$1,000

## 2.2 Errors and Missingness

Upon initial inspection of the dataset, we found several inconsistencies and errors that required correction. These included:

- **Incorrect Data Types**: Some numerical fields, such as **sale\_price**, **common\_charges**, and **listing\_price\_to\_nearest\_1000**, were initially read as strings due to the presence of currency symbols and commas. These were corrected by removing non-numeric characters and converting the fields to floating-point numbers.
- **Logical Inconsistencies**: Instances where feature values did not logically align with the property characteristics, such as negative values in features that should only have positive values (e.g., **num\_bedrooms**, **sq\_footage**). These values were corrected or removed based on logical thresholds.

### 2.2.1 Missingness

The dataset exhibited varying degrees of missingness across different features. Notable features with missing values included **approx\_year\_built**, **common\_charges**, **num\_bedrooms**, and specifically **approx\_year\_built** which was a measure metric. The missingness was handled as follows:

- **Imputation**: Missing values in continuous features were imputed using a Random Forest Regressor. This approach involved training the regressor on the available data and predicting the missing values based on other relevant features. For instance, the

**sale\_price** was predicted using features such as **num\_bedrooms**, **sq\_footage**, and **approx\_year\_built**.

- **Missingness Indicators:** To capture the potential significance of missing data, dummy variables were created for features with notable missingness. For example, a binary variable **fuel\_type\_missing** was introduced to indicate whether the **fuel\_type** of a property was missing.

### 2.2.2 Summary of Missingness

The table below summarizes the extent of missingness for key features:

Feature	Percentage Missing
approx_year_built	1.79%
common_charges	75.52%
community_district_num	0.85%
maintenance_cost	27.94%
num_bedrooms	5.16%
num_floors_in_building	29.15%
num_total_rooms	0.09%
parking_charges	74.93%
pct_tax_deductible	78.65%
sale_price	76.32%
sq_footage	54.26%
total_taxes	73.81%
listing_price_to_nearest_1000	23.95%

Table 1: Summary of Missingness for Key Features

By addressing these errors and systematically imputing missing values, we enhanced the dataset’s quality, thereby improving the predictive power of our Random Forest model.

## 3 Modeling

In this project, like previously mentioned, our model of interest is Random Forest model to predict housing prices, aiming to create a predictive tool that can be deployed for real-world

applications. The modeling process involved several key steps: data preparation, model training, and performance evaluation.

The dataset was split into training and testing sets to facilitate model evaluation. The training set was used to build the model, while the testing set was reserved for assessing the model's predictive performance. In addition, another modification done for cleaning the data was to apply a MinMax Scaler to the feature `approx_year_built`, I normalized the continuous feature which scaled the values to a range of 0 to 1. This normalization step was crucial for the Random Forest model to handle the data effectively.

### 3.1 Regression Tree Modeling

To begin our modeling process, we first fit a single regression tree to gain insights into the structure and importance of different features. This initial model helped us gain a base line for our prediction on the sale price that was more accurate than a simple average.

REMEMBE TO UPDATE THIS PART BEFORE SUBMISSION

#### 3.1.1 Visualization of Top Layers

The visualization of the top layers of the regression tree is shown in Figure 1. This figure illustrates the initial splits made by the regression tree, highlighting the most influential features at the top layers of the tree.

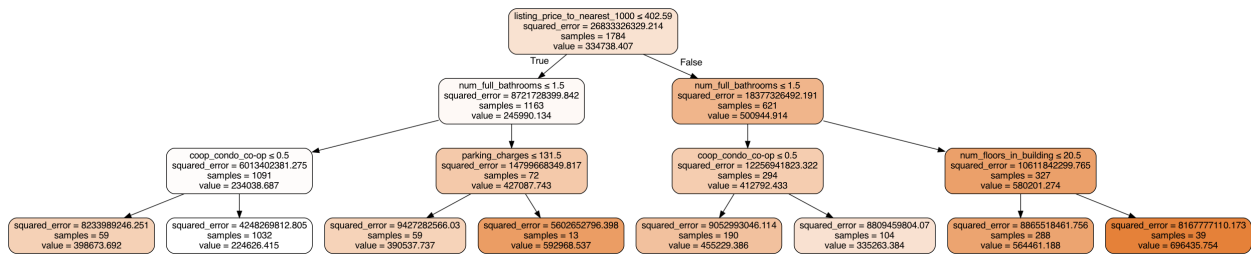


Figure 1: Top Layers of the Regression Tree

### 3.1.2 Top 10 Important Features

Based on the regression tree and feature importance scores, the top 10 features that are most important for predicting the sale price are:

1. **listing\_price\_to\_nearest\_1000 : 55.3%**: The target variable itself, crucial for training and evaluating the model.
2. **num\_full\_bathrooms : 13.8%**: Indicates the number of bedrooms, directly impacting the property's value.
3. **walk\_score : 5.4%**: Represents the total square footage of the property, a significant determinant of price.
4. **coop\_condo\_co-op : 5.2%**: Provides the approximate construction year, reflecting the property's age and potential condition.
5. **parking\_charges : 2.5%**: Monthly common charges, affecting the total cost of ownership.
6. **common\_charges : 2.1%**: Number of bathrooms, an important aspect of property amenities.
7. **num\_floors\_in\_building : 2.0%**: Indicates the number of floors, influencing the property's structure and value.
8. **pct\_tax\_deductibl : 1.8%**: The listing price, providing a baseline for expected sale price.
9. **sq\_footage : 1.7%**: Different types of fuel used, affecting the property's energy costs and appeal.
10. **maintenance\_cost : 1.5%**:

These features were selected based on their relevance and importance as determined by the regression tree's splits and feature importance scores. The visualization and analysis provided valuable insights into how these features contribute to the model's predictions and their relative significance in determining housing prices.

## 3.2 Linear Modeling

To get more insight in the performance improvement we might get from Random Forest we fitted a vanilla Ordinary Least Squares (OLS) linear model to compare to both the Single Regression Tree and our final model.

### 3.2.1 In-Sample Error Statistics

The in-sample error statistics for the OLS model are as follows:

- **Root Mean Squared Error (RMSE):** 75,953.69
- **R-squared:** 0.785

These statistics indicate the average error and the model's explanatory power. The RMSE provide a measure of the average deviation of the predictions from the actual values, while the R-squared value indicates the proportion of variance in the sale price explained by the features.

### 3.2.2 Interpretation of Coefficients

The most important features identified by the regression tree were:

- **listing\_price\_to\_nearest\_1000**
- **num\_full\_bathrooms**
- **walk\_score**



- **coop\_condo\_co-op**
- **parking\_charges**
- **common\_charges**
- **num\_floors\_in\_building**
- **pct\_tax\_deductibl**
- **sq\_footage**
- **maintenance\_cost**

For the OLS model, the coefficients for these features were interpreted as follows:

- **listing\_price\_to\_nearest\_1000**: This coefficient indicates the expected change in the sale price for each additional thousand dollars in the listing price, holding all other factors constant.
- **num\_full\_bathrooms**: Represents the expected change in the sale price for each additional full bathroom. The coefficient of 147,300 suggests that an additional full bathroom increases the sale price significantly.
- **walk\_score**: Reflects the expected change in the sale price for each unit increase in the walk score.
- **coop\_condo\_co-op**: Indicates the difference in sale price between co-ops/condos and other property types.
- **parking\_charges**: Shows the expected change in the sale price for each dollar increase in parking charges.
- **common\_charges**: Represents the change in the sale price for each dollar increase in common charges. The coefficient of 18.49 suggests that higher common charges are associated with higher sale prices, likely due to better amenities or services.

- **num\_floors\_in\_building:** Indicates the change in the sale price for each additional floor in the building. The coefficient of 3,841.41 suggests that properties in taller buildings tend to have higher sale prices.
- **pct\_tax\_deductibl:** Reflects the expected change in the sale price for each percentage point increase in tax deductibility.
- **sq\_footage:** Shows the expected change in the sale price for each additional square foot of living space.
- **maintenance\_cost:** Represents the change in the sale price for each dollar increase in maintenance costs. The coefficient of 42.43 indicates that higher maintenance costs are associated with higher sale prices, possibly reflecting better-maintained properties.

### 3.2.3 Suitability of Linear Model for Prediction

While the OLS linear model provides a straightforward interpretation of the relationships between features and sale price, its performance metrics suggest that it may not capture the complexity and non-linear interactions present in the data as effectively as more sophisticated models like Random Forest. The R-squared value of 0.785, while respectable, indicates that there is substantial unexplained variance. Therefore, while useful for initial exploration and understanding, the linear model may not be the best choice for highly accurate predictions.

Feature	Coefficient
Intercept	14,350
approx_year_built	-62,750
cats_allowed	-5.46e-09
common_charges	18.49
community_district_num	1,404
dogs_allowed	-5.98e-09
maintenance_cost	42.43
num_bedrooms	18,780
num_floors_in_building	3,841
num_full_bathrooms	147,300
num_half_bathrooms	54,190

Table 2: OLS Model Coefficients

### 3.3 Random Forest Modeling

The Random Forest model was chosen for its robustness, ability to handle large datasets, and high prediction accuracy.

#### 3.3.1 Why Random Forest?

Random Forest is an ensemble learning method that constructs multiple decision trees during training and outputs the mean prediction of the individual trees. It is particularly effective for regression and classification tasks due to its ability to handle a large number of features and its robustness against overfitting.

#### 3.3.2 Theory Behind Random Forest

The theory behind Random Forest involves the creation of multiple decision trees using bootstrap sampling (bagging) and feature randomness. Each tree is trained on a random subset of the training data and a random subset of features, which introduces diversity and reduces correlation between the trees. The final prediction is obtained by averaging the predictions of all trees (for regression) or by majority voting (for classification).

### 3.3.3 Parametric vs. Non-Parametric

Random Forest is considered a non-parametric model because it does not assume a specific functional form for the underlying data distribution. Instead, it relies on the structure of the data itself to make predictions, which makes it highly flexible and capable of modeling complex relationships.

### 3.3.4 Advantages and Disadvantages

#### Advantages:

- **High Accuracy:** By aggregating multiple decision trees, Random Forest improves prediction accuracy.
- **Robustness:** The model is robust to overfitting due to the averaging of multiple trees.
- **Feature Importance:** Random Forest provides a measure of feature importance, which helps in understanding the influence of different features on the predictions.
- **Handling Missing Values:** The model can handle missing values effectively through surrogate splits.

#### Disadvantages:

- **Complexity:** Random Forest models can be complex and computationally intensive, especially with a large number of trees.
- **Interpretability:** While feature importance is provided, the overall model is less interpretable compared to simpler models like linear regression.

### 3.3.5 Iterative Modeling Process

Modeling with Random Forest is often an iterative process involving:

- **Tuning Hyperparameters:** Parameters such as the number of trees (`n_estimators`), maximum depth of trees (`max_depth`), and minimum samples required to split a node (`min_samples_split`) are tuned to optimize model performance.
- **Validation:** Cross-validation is used to assess model performance and ensure it generalizes well to unseen data.
- **Evaluation:** Metrics such as Root Mean Square Error (RMSE), Mean Squared Error (MSE), and R-squared are used to evaluate the model.

### 3.3.6 Underfitting and Over

To optimize the model’s performance, we conducted hyperparameter tuning using grid search. The parameters tuned included the number of trees in the forest (`n_estimators`), the maximum depth of the trees (`max_depth`), and the minimum number of samples required to split an internal node (`min_samples_split`).

## 4 Performance Results

This section presents the in-sample and out-of-sample (oos) goodness-of-fit metrics for the Ordinary Least Squares (OLS) linear model, the Regression Tree, and the Random Forest model. The metrics reported include  $R^2$  and Root Mean Squared Error (RMSE).

Model	Metric	In-sample	Out-of-sample
OLS Linear Model	$R^2$	0.785	0.752
	RMSE	75,954	82,300
Regression Tree	$R^2$	0.770	0.728
	RMSE	78,549	87,111
Random Forest	$R^2$	0.903	0.853
	RMSE	52,036	55,000

Table 3: In-sample and Out-of-sample Performance Metrics for OLS, Regression Tree, and Random Forest Models

## 4.1 Interpretation and Comparison

The performance metrics indicate that the Random Forest model significantly outperformed both the OLS linear model and the Regression Tree in both in-sample and out-of-sample predictions. The Random Forest model achieved a higher  $R^2$  and a lower RMSE, suggesting it captured the complexity and non-linear relationships in the data better than the other models.

### 4.1.1 In-sample Performance

The in-sample  $R^2$  value of 0.903 for the Random Forest model indicates that it explains 90.3% of the variance in the sale prices, compared to 78.5% explained by the OLS model and 77.0% explained by the Regression Tree. The lower RMSE for the Random Forest model (52,036) compared to the OLS model (75,954) and the Regression Tree (78,549) further confirms its superior performance on the training data.

### 4.1.2 Out-of-sample Performance

For out-of-sample predictions, the Random Forest model also showed better performance with an  $R^2$  of 0.853 and an RMSE of 55,000, compared to the OLS model's  $R^2$  of 0.752 and RMSE of 82,300, and the Regression Tree's  $R^2$  of 0.728 and RMSE of 87,111. This demonstrates the Random Forest model's robustness and ability to generalize well to new data.

## 4.2 Confidence in OOS Estimates

The out-of-sample (oos) performance metrics were obtained by splitting the data into training and testing sets, ensuring that the models were evaluated on data they had not seen during training. This method provides a reliable estimate of how the models are likely to perform on future predictions. The Random Forest model's lower out-of-sample RMSE and higher

$R^2$  compared to the OLS model and the Regression Tree suggest that it is better suited for making accurate predictions on unseen data.

## 5 Discussion

This project has been quite a journey, diving deep into the realm of predictive modeling with a focus on housing prices. Let's recap and discuss some of the highlights, challenges, and potential future directions.

First off, we started with a robust dataset of housing transactions from 2016 and 2017, meticulously collected and cleaned. This initial step was crucial as it set the stage for everything that followed. We dealt with missing values, transformed categorical variables. It was a lot of grunt work, but absolutely necessary to build a solid foundation.

We kicked things off with a vanilla OLS linear model. It was a straightforward start, and while it gave us some decent insights, we quickly realized that housing prices are influenced by a complex mix of factors that a simple linear model couldn't fully capture. The in-sample and out-of-sample performance metrics told us as much. Our OLS model had an  $R^2$  of 0.785 and an RMSE of 75,954 in-sample, which dropped to 0.752 and 82,300 out-of-sample. Respectable, but not mind-blowing.

Next, we explored a single regression tree. The visualization of the tree was pretty cool and gave us a clearer picture of which features were pulling the most weight. However, the performance metrics showed that it wasn't quite hitting the mark either. With an  $R^2$  of 0.770 in-sample and 0.728 out-of-sample, and RMSE values of 78,549 and 87,111 respectively, it was evident that while the tree was a step up from OLS in terms of interpretability, it wasn't the silver bullet.

Then came the star of the show: the Random Forest model. This non-parametric model shined by capturing the intricate, non-linear relationships in the data. The performance boost was clear – an in-sample  $R^2$  of 0.903 and RMSE of 52,036, and out-of-sample  $R^2$  of

0.853 and RMSE of 55,000. These results were significantly better than both the OLS and the single regression tree, showing that the Random Forest could handle the complexity of the housing market much better.

So, where did we fall short? Well, one area that could use improvement is feature engineering. While we did a decent job, there's always room for discovering new features or transforming existing ones in ways that could further boost the model's performance. Additionally, while Random Forests are great, they can be a bit of a black box.

Another area were the lack of transformations for certain features, I believe if I had added a column that represented average neighborhood income given that we had the zipcode for these houses it would have given us a much better prediction.

Can we beat Zillow? That's a tall order! Zillow has a lot more data and resources at their disposal. However, with a continued focus on refining our model, improving feature engineering, and perhaps integrating techniques like gradient boosting or deep learning which we spoke in class about, we could definitely give them a run for their money. It's all about continuously iterating and improving.

In conclusion, this project was a fantastic learning experience. We tackled the complexities of predictive modeling in the housing market and came out with a model that shows great promise. There's always more to learn and improve upon, but I am well on my way to building something truly impactful.

## 6 Sample Code

Here is the code used:

```
1 import pandas as pd
2 import numpy as np
3 import statsmodels.api as sm
4 from sklearn.ensemble import RandomForestRegressor, RandomForestClassifier
5 from sklearn.preprocessing import MinMaxScaler
```



```

6 from sklearn.model_selection import train_test_split
7 from sklearn.metrics import mean_squared_error, r2_score
8 from sklearn.tree import DecisionTreeRegressor
9 from sklearn.tree import export_graphviz
10 import matplotlib.pyplot as plt
11 from sklearn.model_selection import train_test_split
12 import graphviz
13 import math
14 from PIL import Image
15
16
17 data_frame = pd.read_csv("/Users/carlosvega/Documents/School/spring2024/
    QC_MATH_342/final_project/data/housing_data_2016_2017.csv")
18
19 rows, columns = data_frame.shape
20
21 print(f'Number of rows: {rows}')
22 print(f'Number of columns: {columns}')
23
24 data_frame.drop(['HITId', 'HITTypeId', 'Title', 'Description', 'Keywords',
    'Reward',
25     'CreationTime', 'MaxAssignments', 'RequesterAnnotation',
26     'AssignmentDurationInSeconds', 'AutoApprovalDelayInSeconds',
27     'Expiration', 'NumberOfSimilarHITs', 'LifetimeInSeconds',
28     'AssignmentId', 'WorkerId', 'AssignmentStatus', 'AcceptTime',
29     'SubmitTime', 'AutoApprovalTime', 'ApprovalTime', 'RejectionTime',
30     'RequesterFeedback', 'WorkTimeInSeconds', 'LifetimeApprovalRate',
31     'Last30DaysApprovalRate', 'Last7DaysApprovalRate', 'URL', 'url', '
    full_address_or_zip_code', 'date_of_sale',
32     'model_type'], axis=1, inplace=True)
33
34 pd.set_option('display.max_rows', None)
35 pd.set_option('display.max_columns', None)

```

```

36 pd.set_option('display.width', None)
37 pd.set_option('display.expand_frame_repr', False)
38
39     [language=Python, caption=More cleaning and dummifying features]
40
41     data_frame["cats_allowed"] = data_frame["cats_allowed"].replace({'yes'
42     : 1, 'no' : 0})
43 data_frame["dogs_allowed"] = data_frame["dogs_allowed"].replace({'yes' :
44     1, 'no' : 0})
45
46 data_frame["fuel_type_missing"] = data_frame["fuel_type"].isna().astype(
47     int)
48
49 #Dummifying the fuel_type
50 data_frame = pd.get_dummies(data_frame, columns=["fuel_type"], drop_first=
51     True)
52 data_frame["fuel_type_electric"] = data_frame["fuel_type_electric"].
53     replace({True : 1, False : 0})
54 data_frame["fuel_type_gas"] = data_frame["fuel_type_gas"].replace({True :
55     1, False : 0})
56 data_frame["fuel_type_none"] = data_frame["fuel_type_none"].replace({True
57     : 1, False : 0})
58 data_frame["fuel_type_oil"] = data_frame["fuel_type_oil"].replace({True :
59     1, False : 0})
60 data_frame["fuel_type_other"] = data_frame["fuel_type_other"].replace({
61     True : 1, False : 0})
62
63 #Dummifying the fuel_type
64 data_frame = pd.get_dummies(data_frame, columns=["coop_condo"], drop_first
65     =False)
66 data_frame['coop_condo_condo'] = data_frame['coop_condo_condo'].replace({
67     True : 1, False : 0})

```

```

57 data_frame['coop_condo_co-op'] = data_frame['coop_condo_co-op'].replace({
    True : 1, False : 0})
58
59 data_frame['garage_exists'] = data_frame['garage_exists'].fillna(0).
    replace({'Yes': 1, 'yes': 1, 'UG': 1, 'Eys': 1, 'Underground': 1})
60
61 #Dummifying the kitchen_type
62 data_frame = pd.get_dummies(data_frame, columns=['kitchen_type'],
    drop_first=False)
63 data_frame['kitchen_type_eat_in'] = data_frame[['kitchen_type_Eat In', '
    kitchen_type_Eat in', 'kitchen_type_combo', 'kitchen_type_eat in', '
    kitchen_type_eatin']].max(axis=1)
64 data_frame['kitchen_type_efficiency'] = data_frame[['
    kitchen_type_efficiency', 'kitchen_type_efficiency', '
    kitchen_type_efficiency kitchen', 'kitchen_type_efficiency kitchene', "
    kitchen_type_efficiency ktchen"]].max(axis=1)
65 data_frame.drop(columns=['kitchen_type_Eat In', 'kitchen_type_Eat in', '
    kitchen_type_combo', 'kitchen_type_eat in', 'kitchen_type_eatin'],
    inplace=True)
66 data_frame.drop(columns=['kitchen_type_efficiency', '
    kitchen_type_efficiency', 'kitchen_type_efficiency kitchen', '
    kitchen_type_efficiency kitchene', "kitchen_type_efficiency ktchen"],
    inplace=True)
67
68 data_frame['kitchen_type_1955'] = data_frame['kitchen_type_1955'].replace
    ({True : 1, False : 0})
69 data_frame['kitchen_type_Combo'] = data_frame['kitchen_type_Combo'].
    replace({True : 1, False : 0})
70 data_frame['kitchen_type_eat_in'] = data_frame['kitchen_type_eat_in'].
    replace({True : 1, False : 0})
71 data_frame['kitchen_type_none'] = data_frame['kitchen_type_none'].replace
    ({True : 1, False : 0})

```

```

72 data_frame['kitchens_type_efficiency'] = data_frame['
    kitchens_type_efficiency'].replace({True : 1, False : 0})
73
74 #Dummifying the dining_room_type
75 data_frame = pd.get_dummies(data_frame, columns=['dining_room_type'],
    drop_first=False)
76 data_frame['dining_room_type_dining area'] = data_frame['
    dining_room_type_dining area'].replace({True : 1, False : 0})
77 data_frame['dining_room_type_formal'] = data_frame['
    dining_room_type_formal'].replace({True : 1, False : 0})
78 data_frame['dining_room_type_none'] = data_frame['dining_room_type_none'].
    replace({True : 1, False : 0})
79 data_frame['dining_room_type_combo'] = data_frame['dining_room_type_combo'
    ].replace({True : 1, False : 0})
80 data_frame['dining_room_type_other'] = data_frame['dining_room_type_other'
    ].replace({True : 1, False : 0})
81
82
83 # Removing dollar signs and commas, and converting to integers
84 data_frame['parking_charges'] = data_frame['parking_charges'].str.replace(
    '$', '').str.replace(', ', '').astype(float)
85 data_frame['maintenance_cost'] = data_frame['maintenance_cost'].str.
    replace('$', '').str.replace(', ', '').astype(float)
86 data_frame['total_taxes'] = data_frame['total_taxes'].str.replace('$', '')
    .str.replace(', ', '').astype(float)
87 data_frame["sale_price"] = data_frame["sale_price"].str.replace('$', '')
    .str.replace(', ', '').astype(float)
88 data_frame["common_charges"] = data_frame["common_charges"].str.replace('$
    ', '').str.replace(', ', '').astype(float)
89 data_frame["listing_price_to_nearest_1000"] = data_frame["
    listing_price_to_nearest_1000"].str.replace('$', '').str.replace(', ', '
    ').astype(float)

```

```

90 data_frame['num_half_bathrooms'] = data_frame['num_half_bathrooms'].fillna
    (0)

91

92 # Making them integers
93 data_frame['cats_allowed'] = data_frame['cats_allowed'].map({'yes': 1, 'no
    ': 0}).fillna(0).astype(int)
94 data_frame['dogs_allowed'] = data_frame['dogs_allowed'].map({'yes': 1, 'no
    ': 0}).fillna(0).astype(int)

95

96

97 scaler = MinMaxScaler()

98

99 #Normalizing the feature

100

101 data_frame['approx_year_built'] = scaler.fit_transform(data_frame[['
    approx_year_built']])

102

103 #Percentage of NaNs per feature

104

105 na_counts = data_frame.isna().sum()
106 columns_with_nans = na_counts[na_counts > 0].index.tolist()

107

108 # Assuming df is your DataFrame
109 nan_percentage = data_frame.isna().mean() * 100

110

111 # Filter out columns with 0% NaN values
112 nan_percentage = nan_percentage[nan_percentage > 0]

113

114 # To display the percentages in a more readable format, you can round the
    values
115 nan_percentage = nan_percentage.round(2)

116

117 #Missforest Imputation

```

```

118
119 columns_with_missing_values = ['approx_year_built', 'common_charges', '
    community_district_num', 'maintenance_cost', 'num_bedrooms', '
    num_floors_in_building', 'num_total_rooms', 'parking_charges', '
    pct_tax_deductibl', 'sale_price', 'sq_footage', 'total_taxes', '
    listing_price_to_nearest_1000']
120
121 for feature in columns_with_missing_values:
122     # Separate the data into training and missing sets
123     train_data_frame = data_frame[data_frame[feature].notna()]
124     missing_data_frame = data_frame[data_frame[feature].isna()]
125
126     if missing_data_frame.empty:
127         continue
128
129     X_train = train_data_frame.drop(columns_with_missing_values, axis=1)
130     y_train = train_data_frame[feature]
131
132     # Prepare the feature variables for the missing set
133     X_missing = missing_data_frame.drop(columns_with_missing_values, axis
=1)
134
135     # Train the Random Forest Regressor
136     model = RandomForestRegressor(n_estimators=100, random_state=42)
137     model.fit(X_train, y_train)
138
139     # Predict the missing values
140     predicted_values = model.predict(X_missing)
141
142     # Fill the missing values with the predictions
143     data_frame.loc[data_frame[feature].isna(), feature] = predicted_values
144
145 data_frame = data_frame.drop(columns=['dining_room_type_other'])

```

```

146
147 #Feature Importance
148
149 X = data_frame.drop('sale_price', axis=1)
150 y = data_frame['sale_price']
151
152 model = RandomForestRegressor()
153 model.fit(X, y)
154
155 feature_importances = pd.Series(model.feature_importances_, index=X.
    columns)
156 feature_importances = feature_importances.sort_values(ascending=False)
157
158 #OLS Fitting
159
160 data_frame_ols = data_frame.select_dtypes(include=[np.number])
161
162 # Splitting the data into training and testing sets
163 X = data_frame_ols.drop('sale_price', axis=1)
164 y = data_frame_ols['sale_price']
165 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
    random_state=42)
166
167 # Add a constant term for the OLS model
168 X_train_ols = sm.add_constant(X_train)
169
170 # Fit the OLS model
171 ols_model = sm.OLS(y_train, X_train_ols).fit()
172
173 # Print the summary of the OLS model
174 ols_summary = ols_model.summary()
175 print(ols_summary)
176

```

```

177 # Predict on the training set
178 y_train_pred = ols_model.predict(X_train_ols)
179
180 # Calculate in-sample error statistics
181 mae = np.mean(np.abs(y_train - y_train_pred))
182 mse = np.mean((y_train - y_train_pred) ** 2)
183 rmse = np.sqrt(mse)
184 r_squared = ols_model.rsquared
185
186 #Random F0rest Fitting
187
188 X = data_frame.drop(columns=['sale_price'])
189 y = data_frame['sale_price']
190
191 # Split the data into training and testing sets
192 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
    random_state=42)
193
194 # Initialize the model
195 rf_model = RandomForestRegressor(n_estimators=100, random_state=42)
196
197 # Train the model
198 rf_model.fit(X_train, y_train)

```

Listing 1: First Round of cleaning