

Rent Price Prediction Model For India Using Linear Regression

Machine Learning CS-ELEC2C First Laboratory Activity

Lorenzo, Lex Zedrick M.

I. INTRODUCTION

One of the most common problems in real life is predicting a variable's value (the dependent variable) based on other factors of interest (the independent variables). This problem is solved by applying Linear Regression analysis. This method aims to minimize the discrepancies between the predicted and actual outcomes by fitting the variables in a straight line.

It is essential to remember that some facts in our case—the challenge of using linear regression to predict rent prices based on other provided data—must be more quantitative. As a result, we need first to transform the mentioned data into a quantitative format. Furthermore, there are disadvantages to linear regression despite its widespread usage and robustness. This is mainly observed when it is applied to resolve practical problems. Data collected in real-world situations occasionally deviates from a straight line, introducing complexity and variables outside the fundamental assumptions of the model.

II. METHODOLOGY

The dataset used in the study consists of 4,747 rows that are kept in a CSV file, each containing 12 characteristics. The posting date, number of bedrooms, halls, and kitchens (BHK), rental price, size in square feet, floor location, area type, area of locality, city, furnished status, preferred tenants, number of bathrooms, and point of contact are all included in these aspects.

Posted On	BHK	Size	Floor	Area Type	Area Locality	City	Furnishing	Tenant Profile	Bathrooms	Point of Contact	
2022-06-06	1	100000	8000	Ground out of 4	Super Area	Banarganga, Nc	Hybridbed	Unfurnished	Bachelors/Family	10	Contact Owner
2022-07-06	4	400000	7000	Lower Basement	Carpet Area	Jadhav Hills	Hybridbed	Semi-Furnished	Bachelors/Family	6	Contact Agent
2022-06-24	5	100000	8000	Floor of 12	Super Area	Nariman	Chennai	Semi-Furnished	Bachelors/Family	4	Contact Agent
2022-09-31	4	200000	5700	2d out of 25	Carpet Area	whitefield	Bangalore	Semi-Furnished	Bachelors/Family	4	Contact Agent
2022-06-01	4	1000000	5000	4 out of 15	Carpet Area	Juhu	Mumbai	Semi-Furnished	Bachelors/Family	4	Contact Agent
2022-07-02	4	130000	4000	11 out of 19	Carpet Area	whitefield	Bangalore	Semi-Furnished	Bachelors/Family	4	Contact Agent
2022-07-10	4	100000	4710	1st out of 19	Super Area	Kalyanpada, Nc	Hybridbed	Semi-Furnished	Family	5	Contact Agent
2022-06-11	4	200000	4710	1 out of 35	Carpet Area	K P H B Phase 4	Hybridbed	Semi-Furnished	Bachelors	6	Contact Owner
2022-06-18	4	150000	4510	3 out of 5	Super Area	Indira Parkstation	Bangalore	Semi-Furnished	Family	5	Contact Agent
2022-02-27	5	400000	4500	7 out of 20	Carpet Area	Baroda West	Mumbai	Furnished	Bachelors	5	Contact Agent
2022-06-10	4	100000	4500	8 out of 6	Carpet Area	Langford Garden	Bangalore	Unfurnished	Bachelors	4	Contact Agent
2022-06-18	6	200000	4500	Ground out of 1	Carpet Area	Raja Anandam	Chennai	Semi-Furnished	Bachelors	5	Contact Agent
2022-06-30	4	200000	4500	Ground out of 2	Carpet Area	Jadhav Hills	Hybridbed	Semi-Furnished	Family	4	Contact Agent
2022-07-06	6	80000	4500	1 out of 2	Super Area	Kalyanpada Nagar	Hybridbed	Semi-Furnished	Bachelors/Family	6	Contact Owner
2022-07-06	4	100000	4410	3 out of 15	Carpet Area	Aravali	Bangalore	Semi-Furnished	Bachelors/Family	4	Contact Agent
2022-07-06	6	30000	4200	Ground out of 2	Super Area	Kuribara	Hybridbed	Semi-Furnished	Bachelors/Family	5	Contact Owner
2022-06-08	4	100000	4105	11 out of 25	Carpet Area	Hobli	Bangalore	Semi-Furnished	Bachelors/Family	4	Contact Agent
2022-07-02	5	200000	4020	Ground out of 5	Super Area	Banarganga High	Hybridbed	Unfurnished	Bachelors/Family	4	Contact Owner
2022-07-04	4	30000	4000	Ground out of 3	Super Area	Sahayra Sadar, Nc	Kolkata	Furnished	Bachelors/Family	3	Contact Agent
2022-06-08	4	100000	4000	Ground out of 2	Super Area	Koramangala	Bangalore	Furnished	Bachelors/Family	5	Contact Owner
2022-06-11	4	100000	4000	3 out of 4	Carpet Area	Aravali Station	Delhi	Semi-Furnished	Bachelors/Family	5	Contact Agent
2022-06-18	4	100000	4000	3 out of 3	Carpet Area	West End	Delhi	Semi-Furnished	Bachelors	7	Contact Agent
2022-07-08	4	100000	4000	8 out of 9	Super Area	Heysgarth	Chennai	Semi-Furnished	Bachelors	4	Contact Agent
2022-07-10	4	200000	4000	1 out of 2	Super Area	Poon Garden, Nc	Chennai	Semi-Furnished	Bachelors/Family	4	Contact Agent
2022-07-09	5	100000	3900	8 out of 6	Carpet Area	Atturam Road	Mumbai	Furnished	Bachelors/Family	5	Contact Agent

Figure 2.1 Data Set used for the model

As a part of feature engineering and preprocessing of the data, It has been decided that 10 of the features will be used for the model. The "Posted on," "Area of Locality," "Bathroom," and "Area Type" have been excluded because, during the experiment, they either lessened the model score or were deemed insignificant. During data exploration, it can be seen that some features had inconsistent rows, while some features, like the City and Furnishing status, can not be directly used in the model. Therefore, some data manipulation is needed to be used in the model.

Data cleaning was primarily done in the "Floor" column of the dataset. Data inconsistency was addressed by correcting the naming convention of the ground and basement floors, switching the house floor and the total number of floors, and adapting the format to suit the model.

One hot encoding addresses the features that need to be quantified. This method creates separate columns for the unique instances of the feature that need to be manipulated for it to be inputted in the data model. The separate columns will then be filled with boolean values to identify each row uniquely.

Outlier detection and removal were also done on the target (the rent). Removal of outliers is necessary because it prevents the outliers from affecting the estimation of the model. The detection and removal were implemented using z-scores to calculate which data points deviate too far from the mean.

III. EXPERIMENTS

3.1 Feature Engineering

Feature engineering is transforming raw data from datasets into features suitable for machine learning models—in other words, adding, deleting, creating, and transforming the most suitable features from the data at hand to make the machine learning models more accurate and efficient.

3.1.1 Selecting features

This machine learning model's features were selected mainly by trial and error. Adding and deleting features that significantly change the R^2 and Mean Squared Error.

The decision to exclude the "Posted On" column is mainly because of the difficulty and insignificance of the date the house is posted on to the prediction of house prices. The "Bathroom," "Area Type," and "Area Locality" are excluded because these features did not have that many significant changes in the overall test score of the model. It just made the model have to calculate complex, unnecessary relationships. The remaining eight features proved to help maintain a balance of high training and test scores.

3.2 Data cleaning for the selected features

3.2.1 Data cleaning for floor column

In the data cleaning process, the "Floor" column is the one that went with the most refinement. During data exploration, inconsistencies in the naming convention for the ground and basement floors were the easiest to spot, and things like there is no total floor in some data points and some points where the house floor is higher than the total floor were the ones that needed proper searching.

The representation of the data points where a number does not represent the house floor needs to be equated with a proper number equivalent to the house floor. The first thought was to equate it with negative values, which resulted in the evaluation being biased against the basement floors. Then, I decided to count it as part of the total floors and gave the upper basement a value of 2 and the ground and lower basement a value of 1. This lessened the bias against the basement floor while removing the negative values in the evaluation.

0.5	, 0.33333333,	1.	, 0.25	, 0.66666667,
0.8	, 0.4	, 0.28571429,	0.875	, 0.75
0.2	, 0.625	, 0.16666667,	-0.25	, 0.6
0.57894737,	0.78571429,	-1.	, 0.7	, 0.15384615,
0.85714286,	0.57142857,	0.55128205,	0.72222222,	0.41666667,
0.42857143,	0.5483871	, 0.52380952,	0.36842105,	0.60869565,
0.45	, -0.11111111,	0.79166667,	0.14285714,	0.04545455,
0.06896552,	-0.0625	, 0.90909091,	0.70833333,	0.35714286,

Figure 3.1 Negative evaluation of "Floor" values

```
array([0.5, 0.33333333, 1., 0.25, 0.66666667,
0.8, 0.4, 0.28571429, 0.875, 0.75,
0.2, 0.625, 0.16666667, 0.6, 0.57894737,
0.78571429, 0.7, 0.15384615, 0.85714286, 0.57142857,
0.55128205, 0.72222222, 0.41666667, 0.42857143, 0.5483871,
0.52380952, 0.36842105, 0.60869565, 0.45, 0.22222222,
0.79166667, 0.14285714, 0.04545455, 0.06896552, 0.125,
0.90909091, 0.70833333, 0.35714286, 0.35, 0.71428571,
0.40909091, 0.59090909, 0.5952381, 0.6097561, 0.67948718,
0.65, 0.69565217, 0.55555556, 0.76190476, 0.3125,
0.7804878, 0.1, 0.61904762, 0.31034483, 0.52808989,
0.17073171, 0.93333333, 0.86666667, 0.26315789, 0.27272727,
0.4047619, 0.83333333, 0.63636364, 0.52941176, 0.26666667,
0.46666667, 0.05, 0.27777778, 0.5862069, 0.63157895,
0.375, 0.77272727, 0.53333333, 0.48387097, 0.39285714,
0.45454545, 0.78947368, 0.9, 0.58333333, 0.84615385,
0.47368421, 0.7254902, 0.54545455, 0.55, 0.43478261,
0.30434783, 0.36363636, 0.39534884, 0.31818182, 0.77777778,
0.24137931, 0.7826087, 0.92857143, 0.06666667, 0.51612903,
0.92307692, 0.21568627, 0.91666667, 0.38461538, 0.68,
0.64285714, 0.2962963, 0.85, 0.81818182, 0.3,
0.42307692, 0.62962963, 0.85526316, 0.28333333, 0.72,
0.88235294, 0.65217391, 0.29411765, 0.10714286, 0.20833333,
0.95454545, 0.53846154, 0.46875, 0.9375, 0.44444444,
0.68181818, 0.81081081, 0.09090909, 0.55172414, 0.64,
...
0.76470588, 0.07142857, 0.77419355, 0.20689655, 0.90322581,
0.61290323, 0.78125, 0.6875, 0.64705882, 0.05263158,
0.23333333, 0.84210526, 0.08, 0.76666667, 0.52,
0.65714286, 0.14705882, 0.12903226, 0.23529412, 0.02857143,
0.31428571, 0.03703704, 0.67647059])
```

Figure 3.2 Solved the issue of negative values for "Floor"

The next problem I addressed was that some data points had no total floors; it was just a number or "Ground." In this case, I assumed that this instance means that the house floor is the top floor. Therefore, I replaced these data points to be equivalent to 1.

```
Rows without total floors:
3
Ground
1
1
```

Figure 3.3 List of data points without total floors

The last problem addressed for the "Floor" column is that some data points on the house floor had higher values than the total floor. I assumed this was caused by human error; therefore, I switched the places of the two numbers.

```
Before manipulation: 8 out of 5
After manipulation: 5 out of 8
Before manipulation: 2 out of 1
After manipulation: 1 out of 2
```

Figure 3.4 Switched the places of house floor and total floor

Now that the data of the "Floor" column had been cleaned, I evaluated each data point by replacing " out of " with "/" and then evaluated the string using a function.

3.2.3 Removing Outliers

Outliers are typically removed for machine learning models because these values significantly affect the model's prediction. In this prediction model, the rent prices are scored using a z-score, which measures the standard deviation of each date from the mean. The threshold set for the measurement is 1.7. This allows for some leniency while also reducing the number of outliers.

```
Original DataFrame Shape: (4746, 15)
DataFrame Shape after Removing Outliers: (4595, 15)
```

Figure 3.5 Results of removing outliers

3.3 Splitting of data

The splitting proportion for this model is 80% for the training set and 20% for the test set. This data division allows the model to learn more about the relationships between the features while maintaining a significant number for the test set.

3.4 Regularization

Regularization is done so that the model can generalize better and prevent it from overfitting. This model uses a data Library called Ridge Cross-Value to find the best alpha/penalty value multiplied by the

model's weight squared. The optimal alpha will then be inputted into a Regularization model Library called Ridge.

IV. RESULTS AND ANALYSIS

4.1 Final MSE, training, and R² scores

Mean Squared Error	153,225,096.37
Training Score	0.75
R ² score	0.80

Figure 4.1 Final results of the model

After all the procedures, the model results are 153,225,096.37 for the Mean Squared Error (MSE), 0.75 for the training set, and 0.80 for the test set. Although the MSE is relatively high, the training and test set scores are considerably high but could be better. The closeness of the training and test scores means there is no overfitting, and the model has done quite well in predicting the prices for those sets that it had not seen.

4.2 Scatter Plot Analysis



Figure 4.2

Figure 4.2 shows that the predictions on the lower left quadrant are closer to the actual rent prices than those of the upper right quadrant. This entails that the model can

more accurately predict the prices of those low-value listings compared to high-value ones. This may be because of the nature of the given dataset.

V. CONCLUSIONS AND RECOMMENDATIONS

5.1 Conclusion

While the final R² score of 0.80 indicates that the model has a high accuracy in predicting rent prices, it is also essential to consider the MSE, which is significantly high. This means that the model still needs improvement, and improvements can still happen. The model's tendency to underestimate the prices of high-value houses can cause problems when predicting the prices of those particular houses.

During the experimentation stage, it was found that not all features proved to be helpful. For instance, the "Area Locality" and "Area Type" proved to have lessened the test scores. Moreover, seemingly insignificant features like "Point of Contact" had one of the most significant impacts on prices. This may be because agents have to increase the prices for their commission, while the owners do not have to do this because they will pocket the rent for themselves.

5.2 Recommendations

To further enhance this model's accuracy, adding factors that may affect rent prices, like the neighborhood status, whether the area around the house is safe, or if amenities are offered, is recommended. The house's proximity can also be considered for this model, whether it is close to a public

transportation network or if it is close to schools, office areas, and such.

VI. REFERENCES

- The Magic of Linear Regression Model.
(n.d.). Wwww.linkedin.com. Retrieved
February 4, 2024, from
<https://www.linkedin.com/pulse/magic-linear-regression-model-bhagyashree-ghosh/>
- What is One-hot Encoding. (n.d.).
Deepchecks. Retrieved February 4,
2024, from
<https://deepchecks.com/glossary/one-hot-encoding/#:~:text=One-hot%20encoding%20in%20machine>
- What is Feature Engineering? (2023, March 20). GeeksforGeeks.
<https://www.geeksforgeeks.org/what-is-feature-engineering/>
- Barkved, K. (2022, March 9). How To
Know if Your Machine Learning Model
Has Good Performance | Obviously AI.
Wwww.obviously.ai.
<https://www.obviously.ai/post/machine-learning-model-performance>