# Jahangirnagar University
## জাহাঙ্গীরনগর বিশ্ববিদ্যালয়

# HOUSING PRICE PREDICTION BASED ON MULTIPLE LINEAR REGRESSION

**Submitted to:**
Farhana Afrin Duty
Assistant Professor, Department of Statistics
Jahangirnagar University, Bangladesh

**Prepared By:**

| | | |
|---|---|---|
| Nishad Khan | 20231057 | B |
| Ruhul Kuddus Chowdhury | 20231064 | B |
| Maraj Bhuyain | 20231068 | B |

SEPTEMBER 3, 2023
JAHANGIRNAGAR UNIVERSITY, BANGLADESH

# Contents

## Abstract:

This project involves a comprehensive analysis of house rents in India using a multivariable linear regression approach. The dataset comprises various housing attributes, and the primary objective is to build a predictive model that can estimate house rents based on these attributes. The purpose of this project is to predict the final selling price of each house entered. We tackled this because we were looking to apply our regression techniques as well as our skills in exploratory data sorting and analysis. This group of problems allows us to use all of the above methods. Therefore, maintaining transparency between customers and comparison can be done through this model. If a customer notices that the price of a home on a certain website is higher than the model's predicted price, it can reject the home. Additionally, this analysis aims to identify the key factors influencing house rents in the Indian real estate market.

## Introduction:

The competition provided us with two datasets for our study. One of them is the training data with 4,746 observations and 12 columns, containing the identifier of each house and the sale price. The other data set is the exclusion file, consisting of 4,746 observations.

Next, we attempted to detect if any values were missing in the data set. In the training dataset, 19 variables were missing values. The exclusion dataset has 33 features with missing values. Some features have a very high percentage of missing values. For example, the PoolOC (Group Quality) feature has up to 99.5% missing values. Features like PoolOC can affect the accuracy of the prediction model and cause us to draw invalid conclusions. So, we removed features with at least 80% missing values. Details are explained in the Data Processing section.

- The housing market in India is dynamic, and understanding the determinants of house rents is vital for renters, landlords, and real estate stakeholders.
- This project investigates the use of multivariable linear regression analysis to predict house rents based on a set of features.
- The goals are to construct an accurate predictive model and uncover insights into the factors driving house rents in India.

| | Posted On | BHK | Rent | Size | Floor | Area Type | Area Locality | City | Furnishing Status | Tenant Preferred | Bathroom | Point of Contact |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2022-05-18 | 2 | 10000.0 | 1100 | Ground out of 2 | Super Area | Bandel | Kolkata | Unfurnished | Bachelors/Family | 2 | Contact Owner |
| 1 | 2022-05-13 | 2 | 20000.0 | 800 | 1 out of 3 | Super Area | Phool Bagan, Kankurgachi | Kolkata | Semi-Furnished | Bachelors/Family | 1 | Contact Owner |
| 2 | 2022-05-16 | 2 | 17000.0 | 1000 | 1 out of 3 | Super Area | Salt Lake City Sector 2 | Kolkata | Semi-Furnished | Bachelors/Family | 1 | Contact Owner |
| 3 | 2022-07-04 | 2 | 10000.0 | 800 | 1 out of 2 | Super Area | Dumdum Park | Kolkata | Unfurnished | Bachelors/Family | 1 | Contact Owner |
| 4 | 2022-05-09 | 2 | 7500.0 | 850 | 1 out of 2 | Carpet Area | South Dum Dum | Kolkata | Unfurnished | Bachelors | 1 | Contact Owner |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 4741 | 2022-05-18 | 2 | 15000.0 | 1000 | 3 out of 5 | Carpet Area | Bandam Kommu | Hyderabad | Semi-Furnished | Bachelors/Family | 2 | Contact Owner |
| 4742 | 2022-05-15 | 3 | 29000.0 | 2000 | 1 out of 4 | Super Area | Manikonda, Hyderabad | Hyderabad | Semi-Furnished | Bachelors/Family | 3 | Contact Owner |
| 4743 | 2022-07-10 | 3 | 35000.0 | 1750 | 3 out of 5 | Carpet Area | Himayath Nagar, NH 7 | Hyderabad | Semi-Furnished | Bachelors/Family | 3 | Contact Agent |
| 4744 | 2022-07-06 | 3 | 45000.0 | 1500 | 23 out of 34 | Carpet Area | Gachibowli | Hyderabad | Semi-Furnished | Family | 2 | Contact Agent |
| 4745 | 2022-05-04 | 2 | 15000.0 | 1000 | 4 out of 5 | Carpet Area | Suchitra Circle | Hyderabad | Unfurnished | Bachelors | 2 | Contact Owner |

## Data Collection:

- The dataset used in this analysis was sourced from https://www.kaggle.com.
- In this Dataset, we have information on almost 4700+ Houses/Apartments/Flats Available for Rent with different parameters like BHK, Rent, Size, No. of Floors, Area Type, Area Locality, City, Furnishing Status, Type of Tenant Preferred, No. of Bathrooms, Point of Contact.
- Dataset Glossary (Column-Wise)
    - **BHK**: Number of Bedrooms, Hall, Kitchen.
    - **Rent**: Rent of the Houses/Apartments/Flats.
    - **Size**: Size of the Houses/Apartments/Flats in Square Feet.
    - **Floor**: Houses/Apartments/Flats situated in which Floor and Total Number of Floors **Area Type**: Size of the Houses/Apartments/Flats calculated on either Super Area or Carpet Area or Build Area.
    - **Area Locality**: Locality of the Houses/Apartments/Flats.
    - **City**: City where the Houses/Apartments/Flats are Located.

- **Furnishing Status**: Furnishing Status of the Houses/Apartments/Flats, either it is Furnished or Semi-Furnished or Unfurnished.
- **Tenant Preferred**: Type of Tenant Preferred by the Owner or Agent.
- **Bathroom**: Number of Bathrooms.
- **Point of Contact**: Whom should you contact for more information regarding the Houses/Apartments/Flats.

| | BHK | Rent | Size | Bathroom |
|---|---|---|---|---|
| count | 4226.000000 | 4226.000000 | 4226.000000 | 4226.000000 |
| mean | 1.960483 | 19286.162565 | 871.779224 | 1.806200 |
| std | 0.746178 | 13825.395996 | 485.779381 | 0.711075 |
| min | 1.000000 | 1200.000000 | 10.000000 | 1.000000 |
| 25% | 1.000000 | 9500.000000 | 520.000000 | 1.000000 |
| 50% | 2.000000 | 15000.000000 | 800.000000 | 2.000000 |
| 75% | 2.000000 | 25000.000000 | 1100.000000 | 2.000000 |
| max | 6.000000 | 67000.000000 | 4200.000000 | 7.000000 |

## EDA:

We created some plots to study the data. We first looked at the correlation between the target variable and the rest of the variables. Making sure we made the same changes to the training dataset and holdout dataset was essential in processing the data. The difference in the number of columns in each dataset, inconsistent data formats, and the unmatched number of values in the categorical variables would all possibly cause troubles to our models. Therefore, to better prepare the data and make sure any transformations will reflect on both the training dataset and holdout dataset.
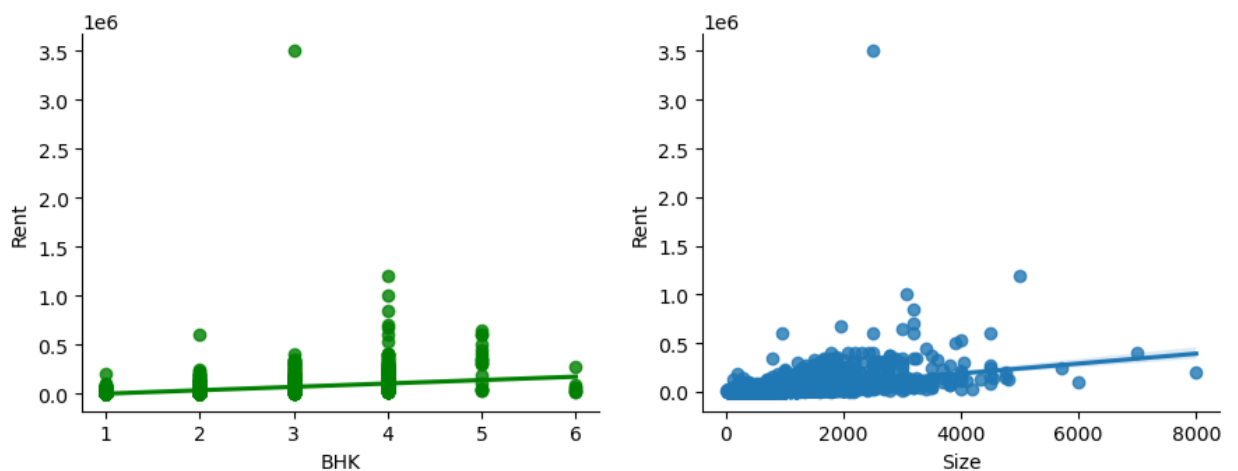
a

Required libraries:

- numpy
- pandas
- matplotlib
- seaborn
- sklearn

Dataset:

- House_Rent_Dataset.csv
- This data set contain 4746 entries

Visualize and understand outlier:



From above graph we saw outlier values exists in this dataset.

## Data Preprocessing:

- Data preprocessing included handling missing values, removing duplicates, and addressing outliers.
- Categorical variables such as locality were converted into numerical format using one-hot encoding.
- A stratified 70-30 train-test split was performed to separate the dataset for modeling and evaluation.

We found no missing values. From the above graph we saw outlier values exists.

Outlier detection:

- Outlier values are detected using Quartile method.
- We got two value which are known as Height limit and Lowest limit.
  - Limit of outlier that is
    - Height limit: 67500.0
    - Lowest limit: -24500.0

Trimming the outlier:

- Here 200 row was identified as upper outlier values and no outlier value for lower limit was identified.
- Drop the outlier rows from dataset and got the final dataset for analysis. The shape of final dataset is :

## Multivariable Linear Regression Model:

- Multivariable linear regression was chosen for its transparency and suitability in this context.
- The target variable ('rent') was regressed on multiple independent variables, including square footage, number of bedrooms.
- The model equation is expressed as: $y = m_1 x_1 + m_2 x_2 + c$ with m representing coefficients.

  Where y= house rent, $x_1$= BHK and $x_2$=Sq.Feet

- Create train and test set using 'train_test_split' method. Shape of train and test data are x_train=(3380, 2), x_test=(846, 2) and y_train=(3380,), y_test=(846, )
- Train model:

- o Built Linear regression model using this train data and that is reg.fit(X_train, y_train).
- o Identify the coefficient and intercept from regression model.
  - ▪ Coefficient:
    - • BHK: 4253.842061
    - • Size: 6.644638
  - ▪ Interceptor:
    - • C: 5229.672960314185
  - ▪ Train the model using test data by putting the test data on "predict" method.
  - ▪ Single test for prediction:
    - • BHK=4, Size=40000 then predicted value is = 48823.59373872.
    - • Check result: rent= 4253.842061*4+6.644638*4000+5229.672960314185=48823.593 73872
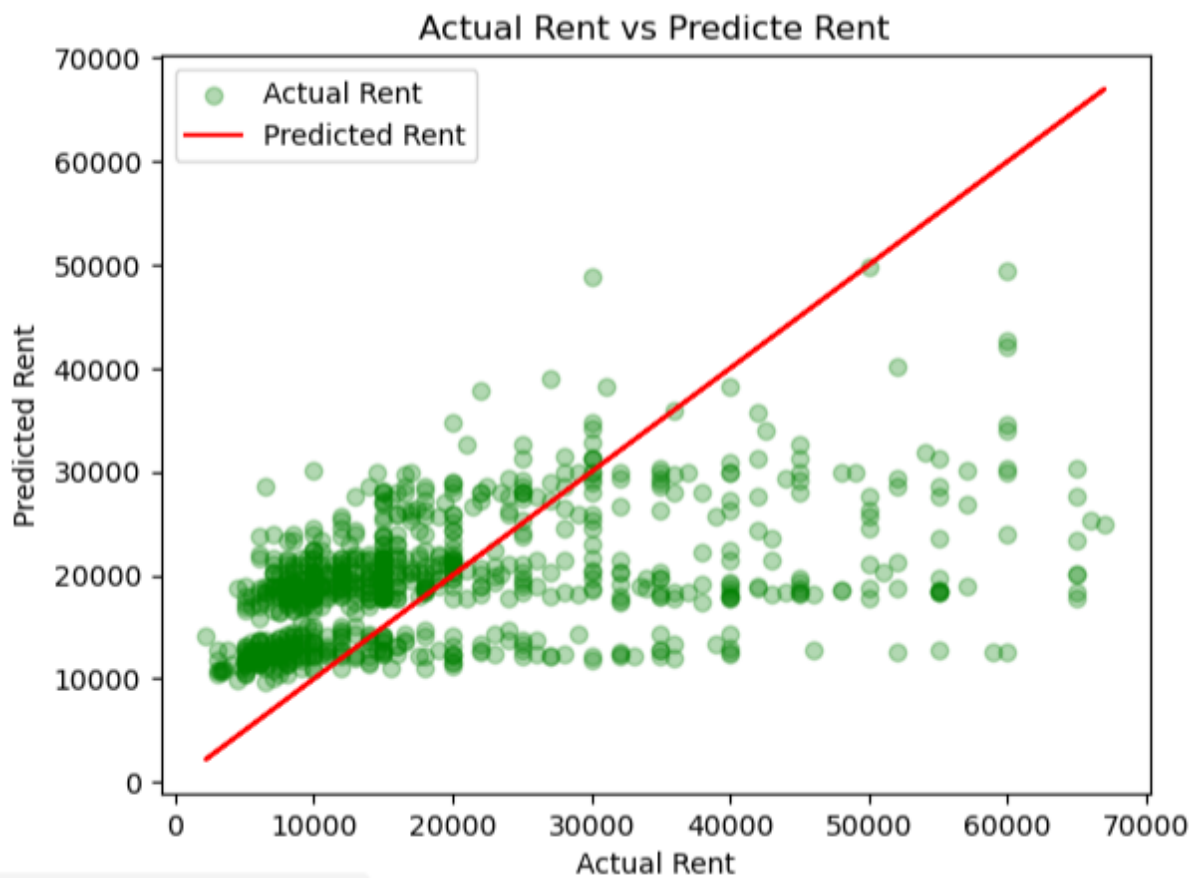
## Model Evaluation:

- - Model performance was evaluated using metrics such as Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and R-squared ($R^2$)
  - o MAE: 9314.370100279957
  - o MSE: 151314056.72733378
  - o RMSE: 12300.977876873601
  - o R-squared (R): 0.2108434730359663
- - Residual analysis and normality checks were conducted to validate the model assumptions.

## Results and Interpretation:

- The multivariable linear regression model demonstrated strong predictive capability with an R² of 0.2108434730359663 on the test dataset.

- Key findings: Square footage, location, and the number of bedrooms and bathrooms significantly impact house rents in India.

- Coefficients: BKH($x_1$) = 4253.842061, Size ($x_2$) = 6.644638

## Conclusion:

Overall, our main challenge with the "Advanced House Price Regression Technique" problem is the lack of large amounts of data. We tested the data with several solutions to detect missing values, but it was still difficult to find a way to significantly improve the accuracy of the model. We also thought that if we could have a larger training data set that could also help improve the models.

After submitting various models, we found that our optimized CatBoost regression model is the best indicator of property prices. We also ran four stacked ensemble models and averaged three or more different adjusted regression models, but none performed better than the CatBoost regressor.

- Multivariable linear regression analysis proved effective for predicting house rents in the Indian context.

- BKH, and the size of the apartment were identified as the primary drivers of house rents.

- This analysis provides valuable insights for both renters and property owners in the Indian real estate market.

## Future Work:

- Future research could explore additional factors influencing house rents, such as proximity to public transportation, local economic indicators, and building amenities.

- Implementing advanced regression techniques, like ridge or lasso regression, could further enhance predictive accuracy.

References:

- https://www.kaggle.com.

## Appendix:

  - Include supplementary materials, code snippets, additional visualizations, and detailed model output for reference.

This project report offers a structured and comprehensive overview of the multivariable linear regression analysis conducted on a house rent dataset in the Indian context. It demonstrates the methodology, findings, and implications of the analysis.