# Predicting the housing prices in household surveys in Uganda

By Kartikay Dhall, Malhaar Khanna, Rishiraj Datta

## Summary

In order to estimate the rental value of homes in low-income nations like Tanzania, Uganda, and Malawi, the research article compares the performance of conventional hedonic pricing models versus machine learning methods.

The study developed hedonic pricing models and machine learning models using information from household surveys carried out in these nations. The machine learning models included a variety of input variables, including demographic and socioeconomic traits of tenants, whereas the hedonic pricing models used standard input variables like location, size, and quality of housing.

The study's findings demonstrated that machine learning models performed better at forecasting rental values in all three nations than hedonic pricing models. The predictions made using the machine learning models were more accurate because they could capture the subtler correlations between the input variables and rental values.

The ability of machine learning models to manage large levels of housing market heterogeneity is one of the key benefits of employing them to forecast rental values in low-income nations. This is crucial in many nations since housing conditions and quality can differ greatly even within relatively small geographic areas.

Overall, the research offers insightful information about the potential of machine learning models to enhance our comprehension of the housing market in low-income countries and to guide housing policy decisions.

## Introduction

The purpose of this study is to determine whether hedonic pricing and machine learning models can accurately forecast the rental value of homes in low-income nations, particularly Tanzania, Uganda, and Malawi.

The study makes use of a number of variables that are pertinent to forecasting the rental value of homes in order to answer this research issue. These elements are:

1.Location: The geographic location of the home is a factor that can significantly affect how much it rents for. Latitude and longitude coordinates are used in the study to determine location.

2. Size: The size of the house, measured in square metres, is a key component in determining its rental value.

3. Housing quality: The physical state of the home, including its age, construction methods, and general quality, is referred to by this variable.

4. Age, income, level of education, and employment are just a few examples of the tenants' demographic and socioeconomic variables that are included in this variable.

The study develops hedonic pricing models and machine learning models using these characteristics to forecast rental values in the three nations. The machine learning models use a more varied set of input factors, including demographic and socioeconomic features of tenants, whereas the hedonic pricing models use more conventional input variables, such as location, size, and quality of housing.

## Methodology

The research predicts the rental value of homes using two machine learning algorithms, Random Forest and Gradient Boosted Tree. These methods are contrasted with conventional hedonic pricing models, such as the Tobit and Ordinary Least Squares (OLS) models.

Both the Random Forest and Gradient Boosted Tree ensemble learning techniques mix different decision trees to get predictions that are more accurate. Due to their capacity for handling complex data and capturing non-linear correlations between variables, these algorithms are frequently utilised in machine learning.

One would need to collect and clean pertinent data from household surveys in Tanzania, Uganda, and Malawi, similar to what was done in the paper, to replicate the ML/econometric exercise outlined in the paper. Then, based on a set of input factors, one would need to construct and train Random Forest and Gradient Boosted Tree models to forecast the rental value of homes.

These models' accuracy can be compared to more established hedonic pricing models like the OLS and Tobit models. It is crucial to remember that the quality and amount of the input variables utilised, as well as the specific parameters and methods employed in the modeling process, will all have an impact on how well the machine learning models perform.

Overall, a thorough understanding of data cleaning, data analysis, and machine learning algorithms would be necessary to replicate the ML/econometric exercise reported in the research utilizing pertinent machine learning approaches.

## Results from the analysis

The findings from the econometric and machine learning analyses presented in this study can be used to make predictions about the rental values of homes in Tanzania, Uganda, and Malawi.

First, the study discovered that in terms of prediction accuracy, both machine learning models—Random Forest and Gradient Boosted Tree—performed better than more established hedonic pricing models, such as OLS and Tobit models. This means that using machine learning techniques to estimate rental values in these nations may be more successful.

Second, the study discovered that factors like location, size, quality, age, building type, income, and education level were some of the most important predictors of rental value in these nations. These results support earlier studies and offer more proof that these elements have a significant role in determining rental value in the area.
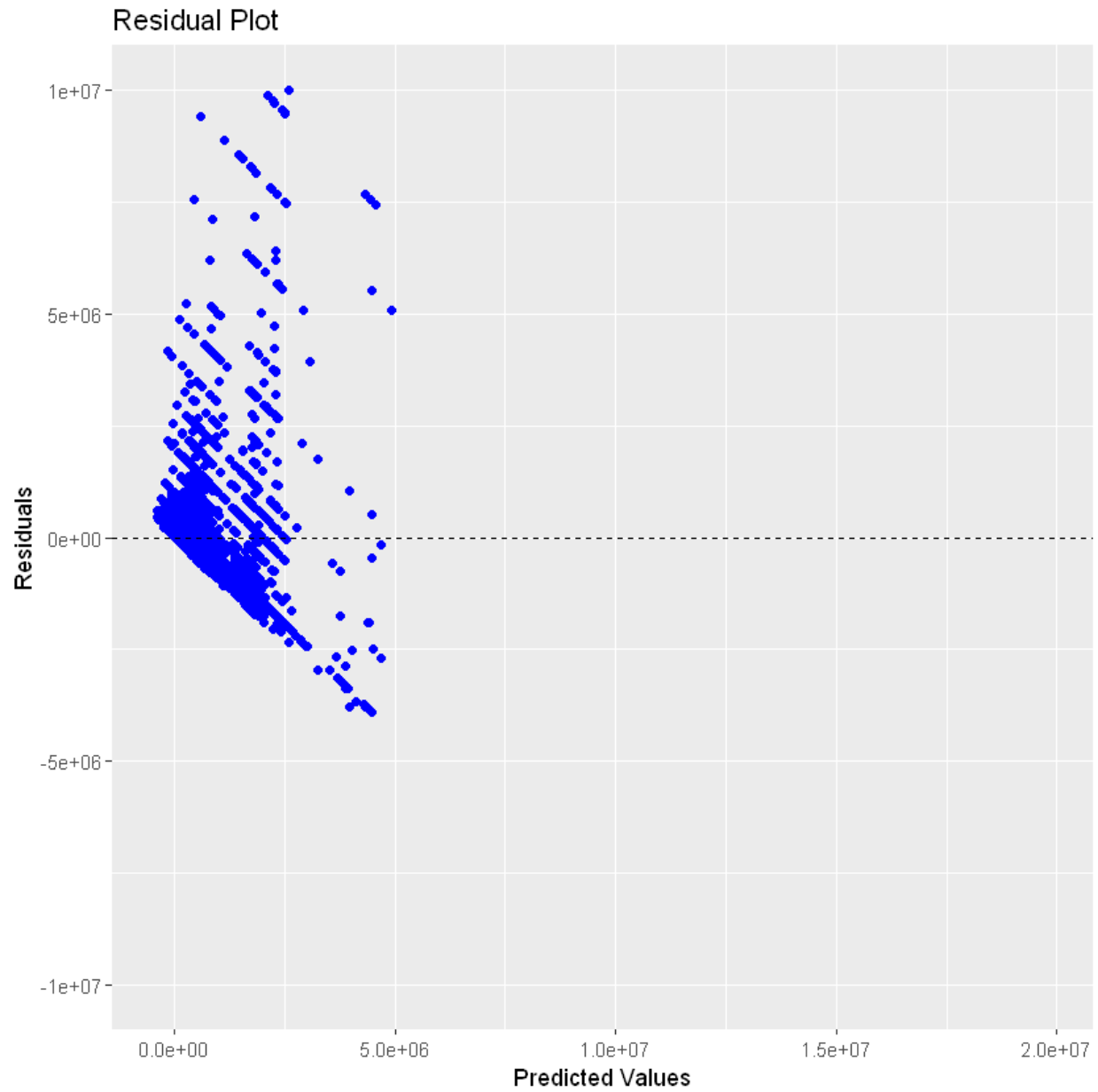
Last but not least, the study emphasises the potential utility of leveraging data from household surveys to forecast rental values, as this strategy may be more practical and affordable than other approaches like property valuation or real estate market analysis. It is crucial to remember that the quality and representativeness of the survey data, as well as other elements like the local political and economic situation, may have an impact on how accurate these predictions are.

Overall, the findings of this study can help inform policymakers and real estate stakeholders in Tanzania, Uganda, and Malawi on how to better predict and understand rental values, which can ultimately lead to more effective policy and investment decisions.
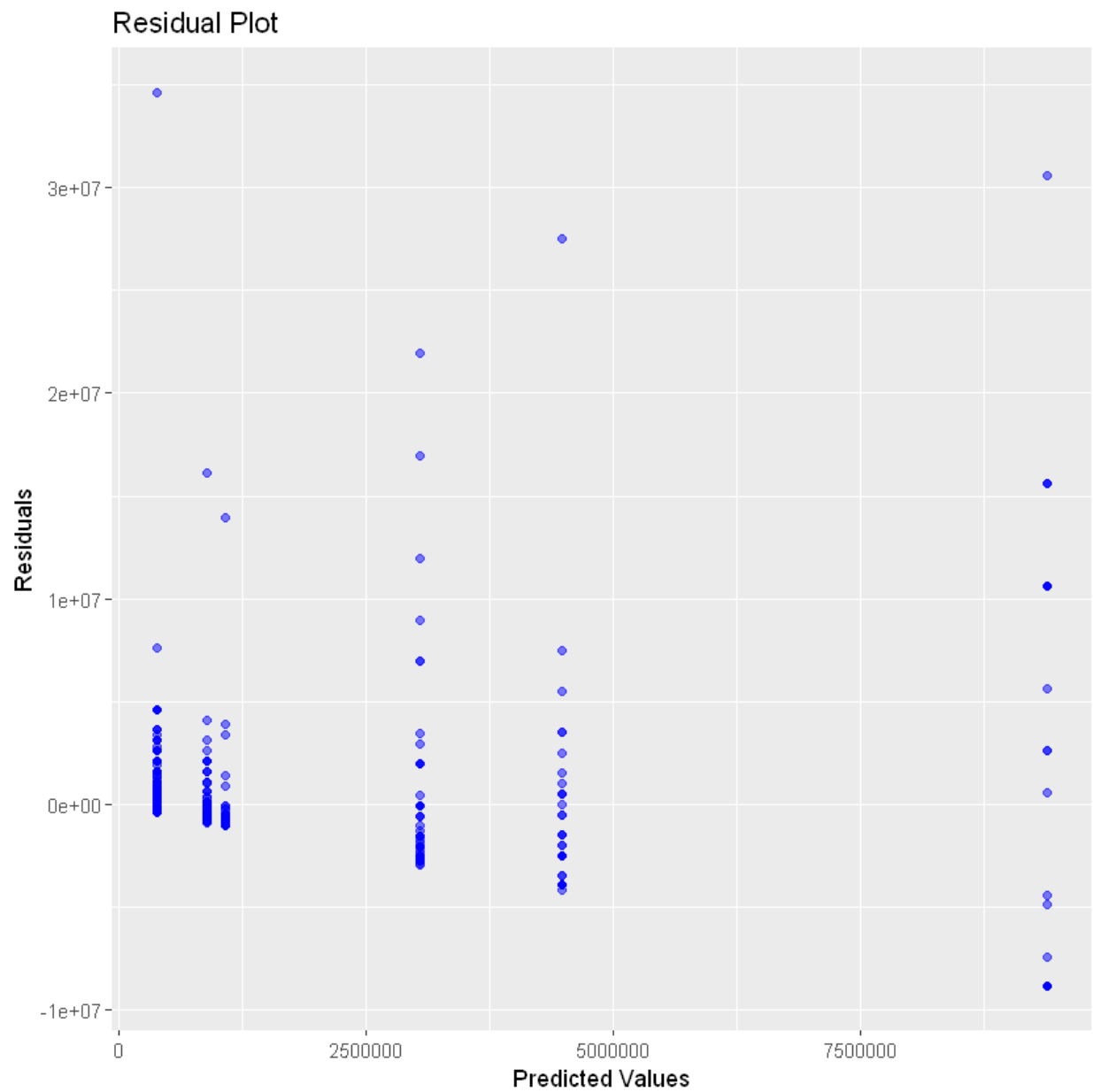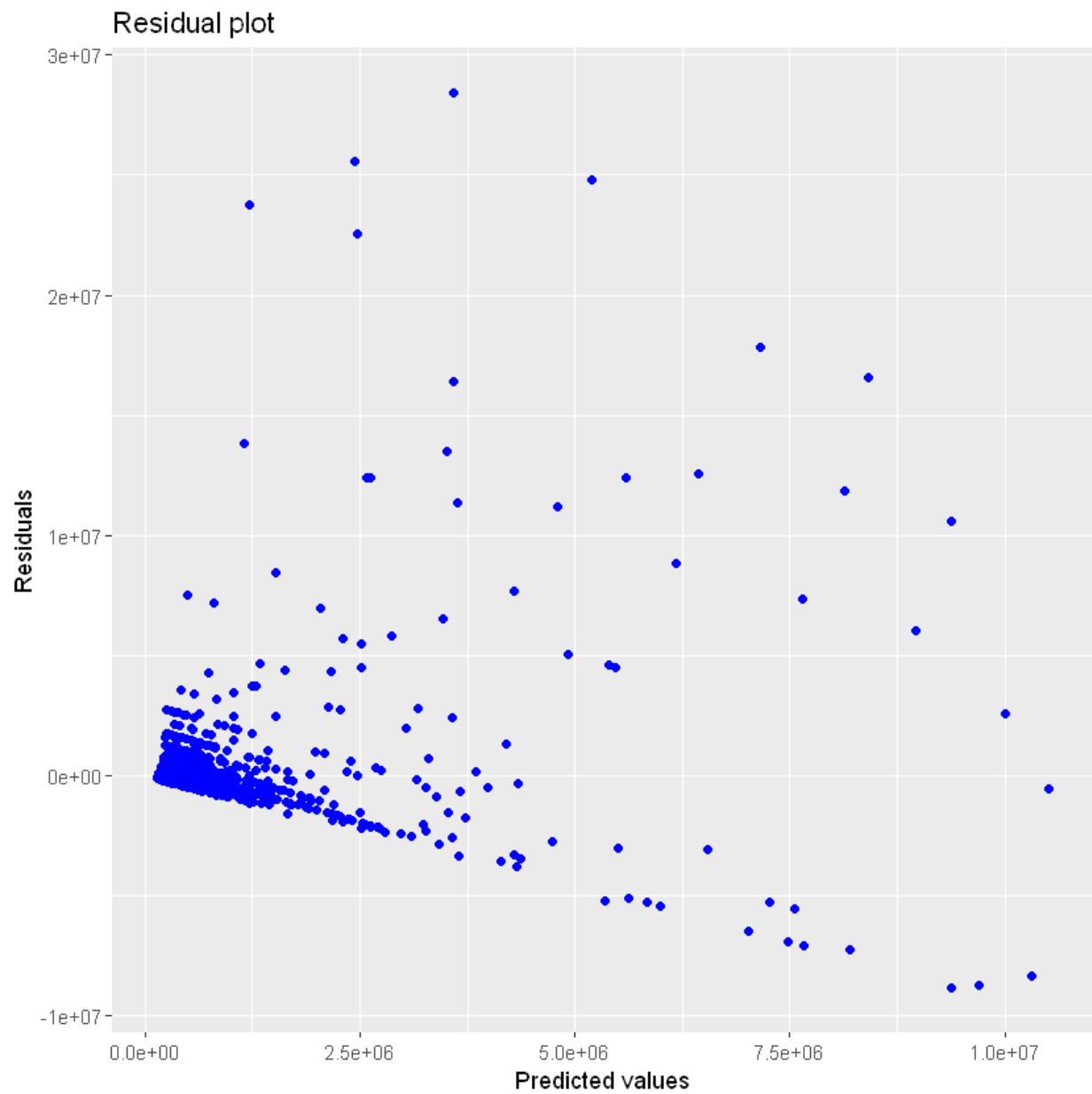
**Graphs**

- **Residual Plot for OLS Regression Model**
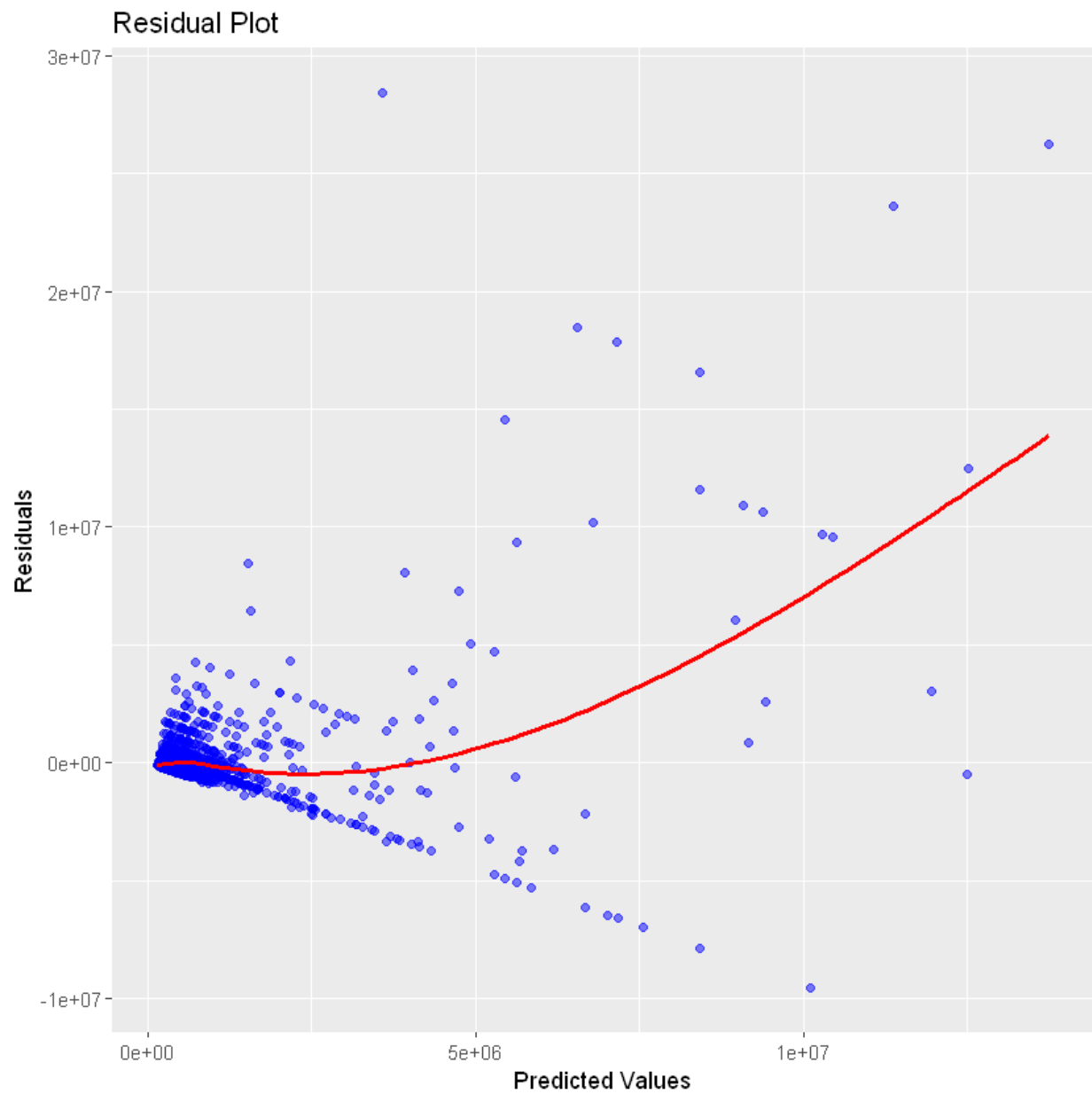


Residual Plot

**Tree based Regression Model**



Residual Plot

**Random Forest Regression Model**



Residual plot

**Gradient Boosting Method**



Residual Plot

**Table for Root Mean Squared Error Values for Various Econometric Methods Used -**

| Method | RMSE Value | Normalized Error |
|---|---|---|
| OLS | 1558824 | 1.00 |
| Ridge | 1838892 | 1.18 |
| Tree | 1504464 | 0.96 |
| Random Forest | 1439510 | 0.92 |
| Gradient Boost | 1253958 | 0.80 |

It is clearly evident that Gradient Boost Algorithm gives the least error out of all the ML methods we have implemented to build models on our dataset.

| | Uganda | |
|---|---|---|
| | **2010** | **2012** |
| OLS[+] | 1.00 | 1.00 |
| Ridge | 0.94 | 0.98 |
| LASSO[+ +] | 0.88 | 0.96 |
| Tree | 0.83 | 0.86 |
| Bagging | 0.78 | 0.83 |
| Forest | 0.82 | 0.91 |
| Boosting | 0.87 | 0.88 |

Comparative Error values (as given in the Research Paper)

Error Comparison in various Algorithms

# Interpretation of the coefficients and their relevance in the context of the research paper.

| | housing_price | Rent | area_property | cooking_fuel | electricity | electricity_ESCOM | MTLtelephone | source_drinkingWater | toilet_facility | disposal_facility | use_bednet |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **housing_price** | 1.00000000 | 0.145637542 | 0.01491336 0 | 0.200266185 | -0.295887837 | 0.018338 5335 | -0.1221852 71 | -0.244048322 | -0.229374 800 | -0.1091663 26 | -0.03182363 01 |
| **Rent** | 0.14563754 | 1.000000000 | 0.001198433 | 0.067299877 | -0.117298602 | 0.008759 2067 | -0.0279825 57 | -0.096144964 | -0.101204 239 | -0.0445885 04 | -0.01245435 97 |
| **area_property** | 0.01491336 | 0.001198433 | 1.00000000 0 | -0.001844503 | -0.007906939 | 0.036318 7542 | 0.0027569 66 | 0.005288152 | 0.006577 243 | -0.0023012 72 | -0.00181632 60 |
| **cooking_fuel** | 0.20026619 | 0.067299877 | -0.001844503 | 1.000000000 | -0.512692772 | 0.027241 7699 | -0.0713864 80 | -0.494255776 | -0.294385 356 | -0.176811547 | -0.04342296 57 |
| **electricity** | -0.295887884 | -0.117298602 | -0.007906939 | -0.512692772 | 1.000000000 | -0.2403 458222 | 0.1437319 76 | 0.528477996 | 0.421619 480 | 0.2052873 00 | 0.06296557 60 |
| **electricity_ESCOM** | 0.01833853 | 0.008759207 | 0.036318754 | 0.027241770 | -0.240345822 | 1.000000000 | -0.0108961 96 | -0.023097193 | -0.049732 025 | -0.0403304 04 | -0.00077387 47 |
| **MTLtelephone** | -0.122185277 | -0.027982557 | 0.002756966 | -0.071386480 | 0.143731976 | -0.0108 961960 | 1.0000000 00 | 0.106836397 | 0.098066 153 | 0.0334716 43 | 0.01929333 04 |
| **source_drinkingWater** | -0.244048322 | -0.096144964 | 0.005288152 | -0.494255776 | 0.528477996 | -0.0230 971930 | 0.1068363 97 | 1.000000000 | 0.339072 709 | 0.1682731 12 | 0.03320159 80 |

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **toilet_facility** | -0.2293748 0 | -0.10 1204 239 | 0.006 5772 43 | -0.29 4385 356 | 0.42 1619 480 | -0.0497 320248 | 0.098 06615 3 | 0.3390727 09 | 1.00 0000 000 | 0.2434 28172 | 0.085 10688 74 |
| **disposal_ facility** | -0.10 91663 3 | -0.04 4588 504 | -0.00 2301 272 | -0.17 6811 547 | 0.20 5287 300 | -0.0404 033037 | 0.033 47164 3 | 0.1682731 12 | 0.24 3428 172 | 1.0000 00000 | 0.087 71428 54 |
| **use_bedn et** | -0.03 18236 3 | -0.01 2454 360 | -0.00 1816 326 | -0.04 3422 966 | 0.06 2965 576 | -0.0007 738747 | 0.019 29333 0 | 0.0332015 98 | 0.08 5106 887 | 0.0877 14285 | 1.000 00000 00 |

The above matrix represents the correlation between different variables in our dataset, where we can say that -

1. Correlation between area_property and hous_price is very low , indicating that social factors like electricity, source of water, toilet facilities, etc. have a bigger role to play.
2. Some of the categoricals columns like Electricity have negative correlation with housing price that is because we consider our binary category as (1 : Yes 2 : No)
3. The above explaination can beextended for other categorical columns such as toilet facility, disposal facility, etc.

## Extra points

- **Dealing with Null Values** - In our model, we have replaced the Null Values in the dataset with the mean/mode of the other entries in that column respectively. We didn't use Zero(0) to fill these Null values as the errors might be increased when such an assumption is taken.

- **Bootstrapping** - We could use bootstrapping to repeatedly draw sample data in order to avoid the problems with any outliers that might be present. Bootstrapping from the dataset can also help us to build and train models on relatively smaller chunks of data as compared to the original methods.

- **K Fold Cross Validation** - If we were to implement Bootstrapping technique to build and train models on the dataset, we could also use K-Fold Cross validation method in order to get an average value for the Root Mean Square Errors across various algorithms implemented. This reduces dependency on the subset of dataset sampled in Bootstrapping.

- **KNN** - By using KNN for imputation, we can create a more complete dataset that can be used for further analysis. KNN imputation is especially useful when there are only a small number of null values in the dataset, as it can be computationally expensive for large datasets.