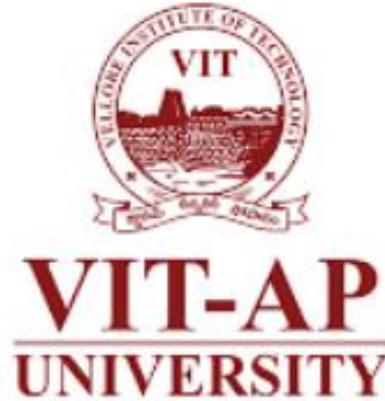**Senior Design Project**



# House Price Prediction Using Ensemble Learning Techniques

Under the guidance of:-
Dr. Manomita Chakraborty

By:-
Raghunath Singh (20BCI7012)
Anuj (20BCE7055)

# AGENDA

PROBLEM DEFINITION

ABOUT DATASET

METHODOLOGY

SUMMARY

REFERENCE

THANK YOU NOTE

# PROBLEM DEFINITION

To design an Machine Learning Algorithm for the accurate prediction of House Price using Random Forest, Support Vector Regression, Boosting algorithm and Ensemble models.

# ABOUT DATASET

The dataset has been taken from Kaggle.

**It has following features:-**
**ID**: Unique identifier for each house listing.
**Date**: Date of the house listing.
**Number of bedrooms**: The count of bedrooms in the house.
**Number of bathrooms:** The count of bathrooms in the house.
**Living area**: Total living area of the house in square feet.
**Lot area**: Total area of the lot in square feet.
**Number of floors**: Number of floors in the house.
**Waterfront present**: Indicates whether the house has a waterfront view (binary: 0 for no, 1 for yes).
**Number of schools nearby**: Number of schools located near the house.
**Distance from the airport:** Distance of the house from the nearest airport in miles.
**Price**: Price of the house.

# ABOUT DATASET

**Number of views**: Number of views the house has received.

**Condition of the house**: Condition rating of the house.

**Grade of the house**: Grade rating of the house.

**Area of the house (excluding basement):** Total area of the house excluding the basement in square feet.

**Area of the basement**: Area of the basement in square feet.

**Built Year**: Year the house was originally built.

**Renovation Year**: Year of the last renovation, if any.

**Postal Code**: Postal code of the house location.

**Latitude**: Latitude coordinate of the house location.

**Longitude**: Longitude coordinate of the house location.

**Living area after renovation**: Total living area of the house after renovation in square feet.

**Lot area after renovation**: Total area of the lot after renovation in square feet.

**Kaggle Link**:- https://www.kaggle.com/datasets/mohamedafsal007/house-price-dataset-of-india

# Methodology

Step 1: Data Exploration:

- Loaded the dataset and checked for missing values.
- Explored data types and basic statistics.
- We are not having any null value in our dataset.
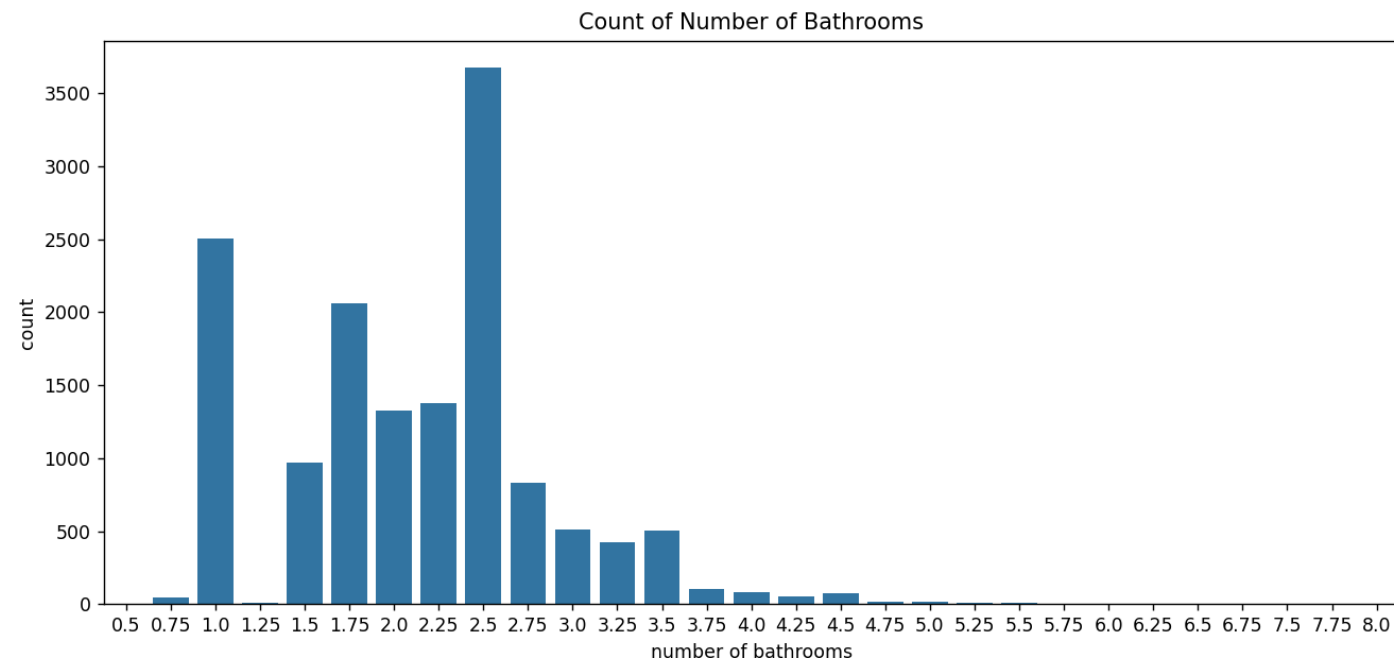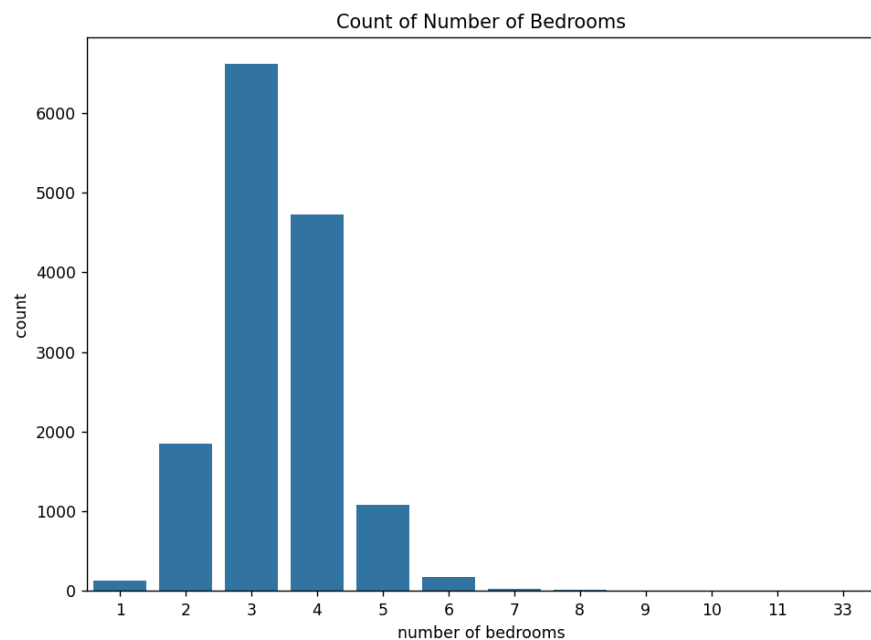
.

```
In [4]: pd.isnull(df).head()
```

Out[4]:

| | number of bathrooms | living area | lot area | number of floors | waterfront present | number of views | condition of the house | ... | Built Year | Renovation Year | Postal Code | Lattitude | Longitude | living_area_renov | lot_area_renov | Number of schools nearby |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | False | False | False | False | False | False | False | ... | False | False | False | False | False | False | False | False |
| | False | False | False | False | False | False | False | ... | False | False | False | False | False | False | False | False |
| | False | False | False | False | False | False | False | ... | False | False | False | False | False | False | False | False |
| | False | False | False | False | False | False | False | ... | False | False | False | False | False | False | False | False |
| | False | False | False | False | False | False | False | ... | False | False | False | False | False | False | False | False |

```
In [5]: df.isna().sum()
```

```
Out[5]: id                                      0
        Date                                    0
        number of bedrooms                      0
        number of bathrooms                     0
        living area                             0
        lot area                                0
        number of floors                        0
        waterfront present                      0
        number of views                         0
        condition of the house                  0
        grade of the house                      0
        Area of the house(excluding basement)   0
        Area of the basement                    0
        Built Year                              0
        Renovation Year                         0
        Postal Code                             0
```
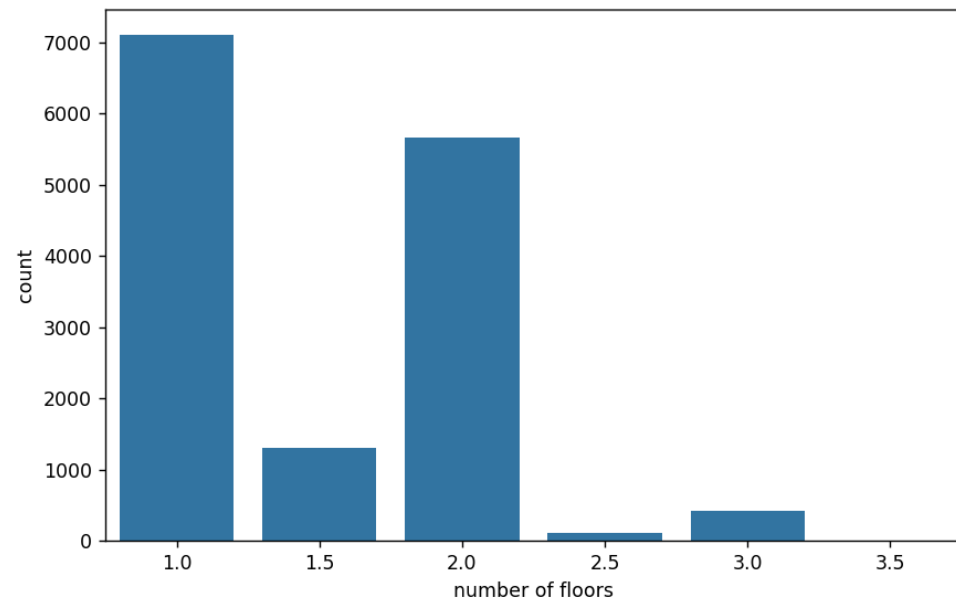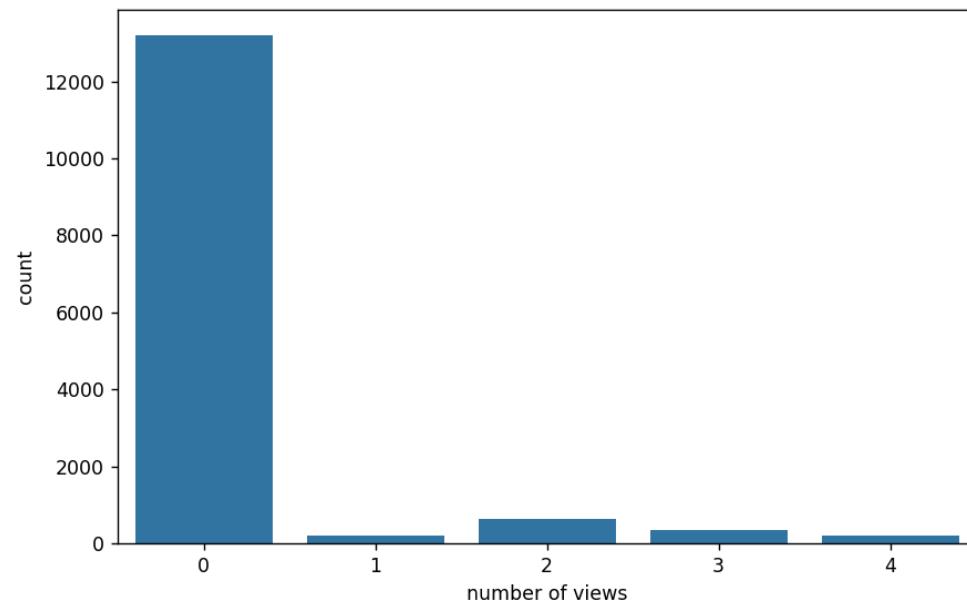
## Step 2: Visualization:

- Plotted count distributions for various features.
- Visualized geographical data and price distributions.
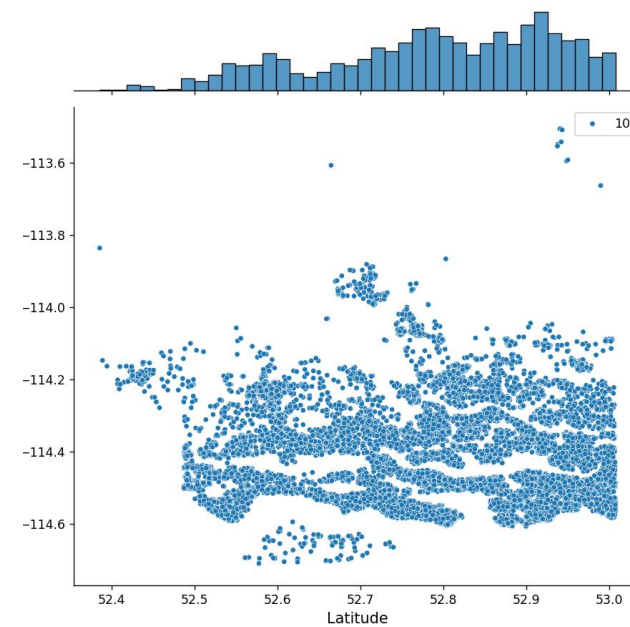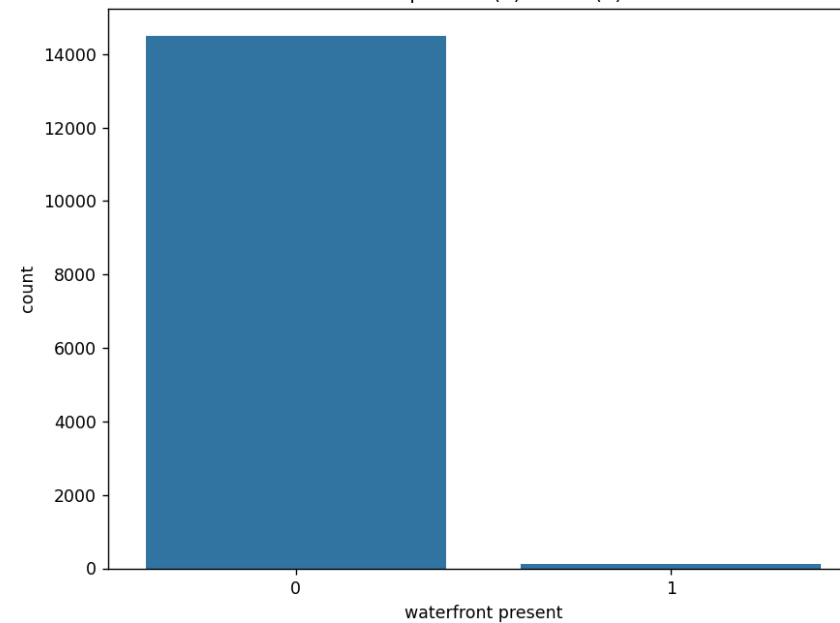- Examined correlations and relationships between variables.



Count of Number of Bedrooms
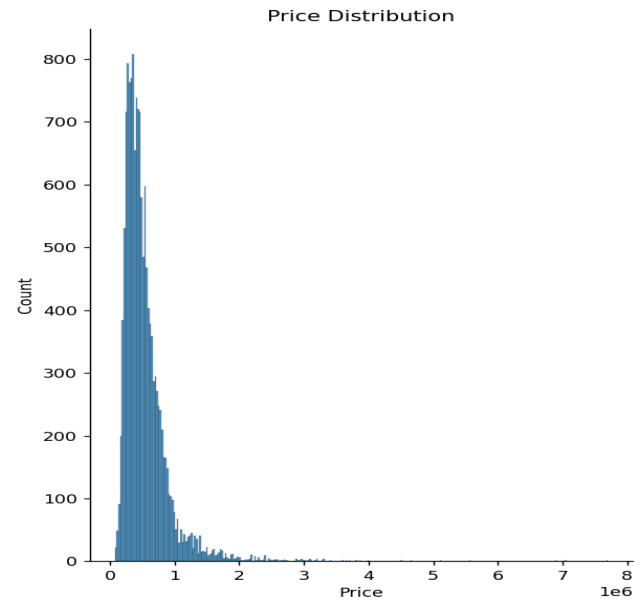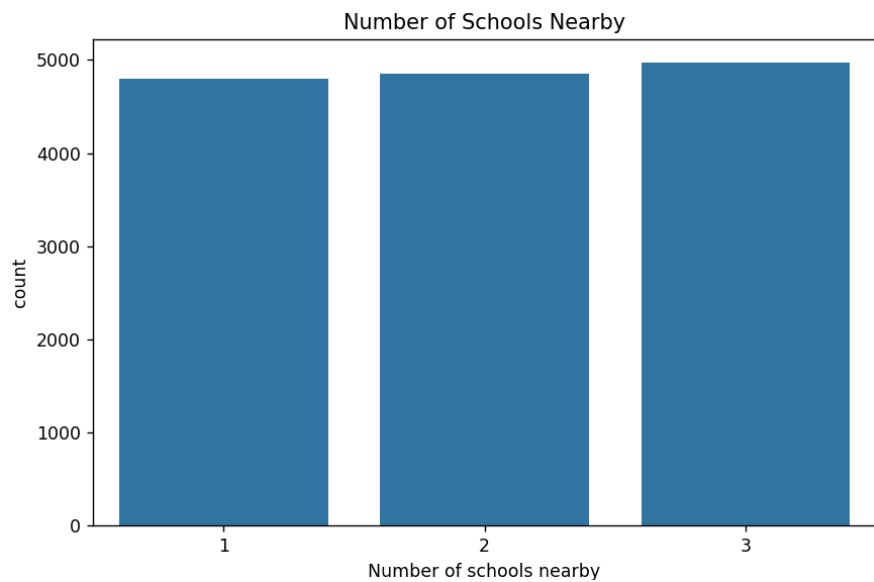
Count of Number of Bathrooms

Pairplot of Numerical Variables

Boxplot of Price by Number of Bedrooms

## Step 3: Feature Engineering:

Created new features like property age and total square footage

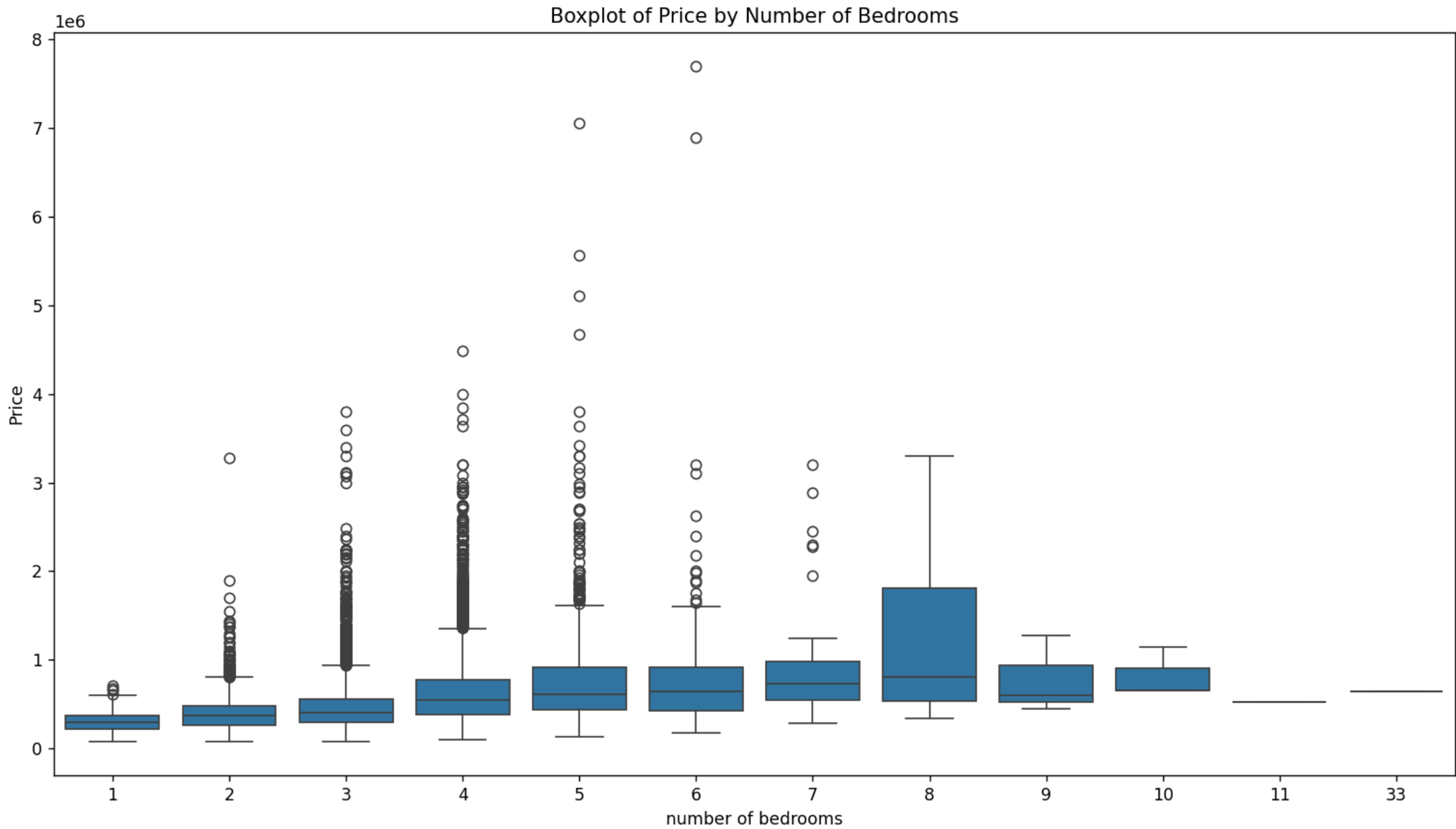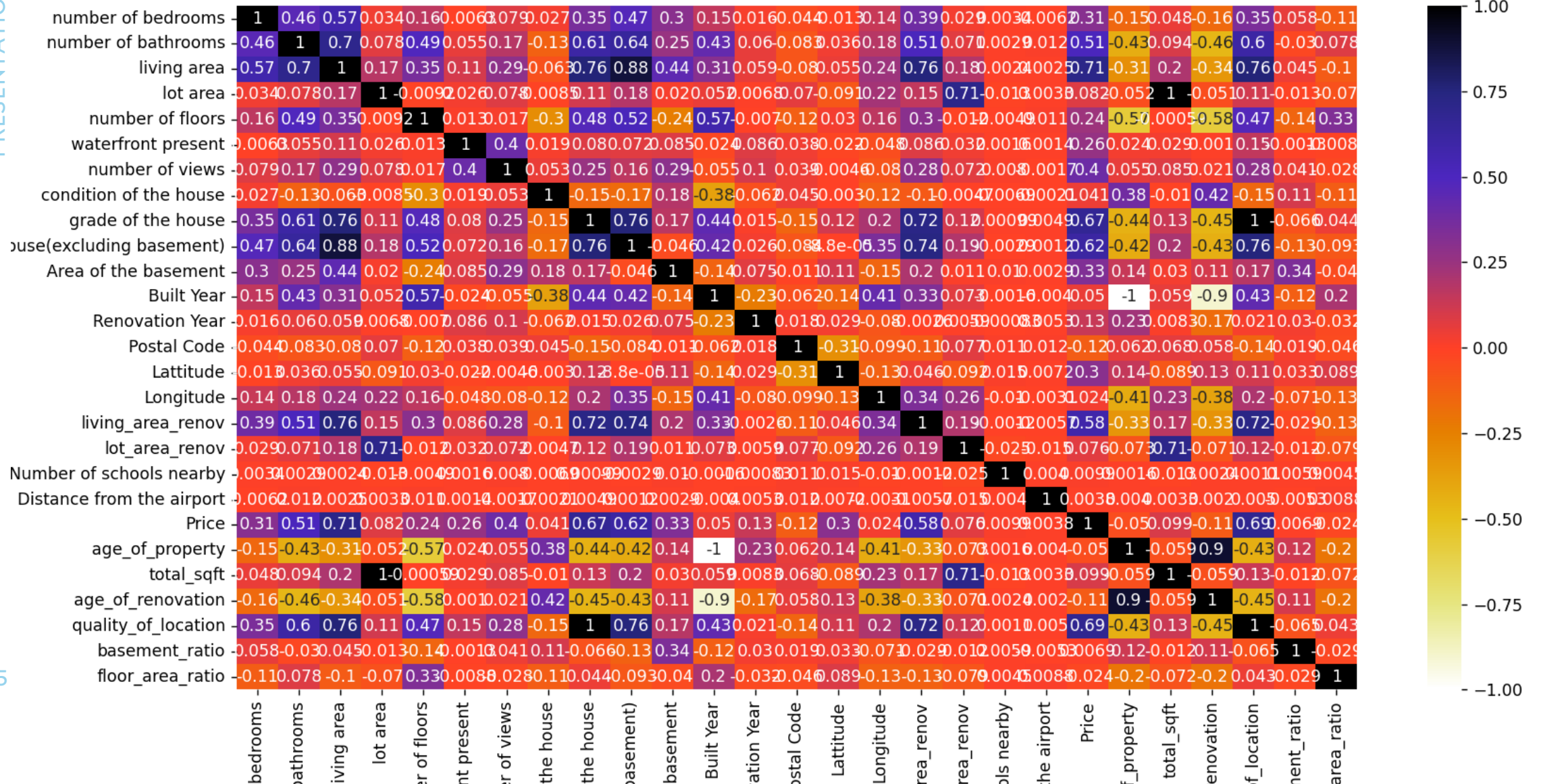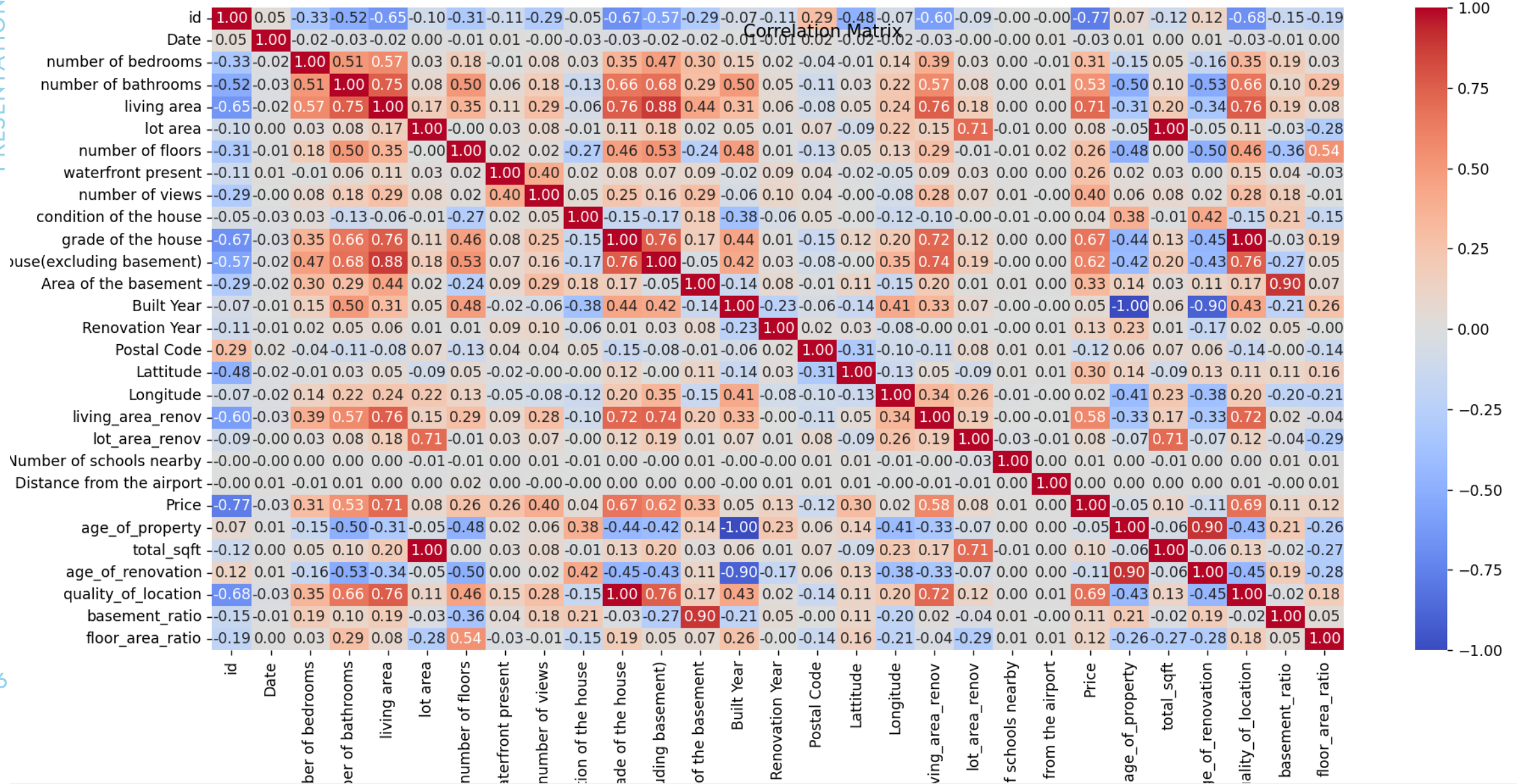## Step 4: Feature Selection:

Identified and removed highly correlated features



Correlation Matrix

## Step 5: Modeling:

- Split data for training and testing.

```python
# Split the data into train and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

- Trained RandomForest, XGBoost, and CatBoost models.

```python
random_forest_model = RandomForestRegressor()
xgboost_model = XGBRegressor()

# Train individual base models
random_forest_model.fit(X_train_scaled, y_train)
xgboost_model.fit(X_train_scaled, y_train)

# Make predictions using individual base models
y_pred_rf = random_forest_model.predict(X_test_scaled)
y_pred_xgb = xgboost_model.predict(X_test_scaled)

y_pred_rf
y_pred_xgb
```

- Evaluated models using R-squared scores.

```
Random Forest Model:
R-squared: 0.6591513191729483


XGBoost Model:
R-squared: 0.6487248781534409
R-squared: 0.6720537148567195
catboost:
R-Squared:  0.6681556299940523
```

- Combined predictions to create an ensemble model.

```
Ensemble Model:
R-squared: 0.6772239409590448
Precision: 0.76878612716763301
Recall: 0.72876712328767121
```

## Step 6: Model Selection:

Selected the best-performing model based on R-squared scores.

```
Best Model: CatBoost
R-squared: 0.6681556299940523
```

## Step 7: Frontend Integeration:

- Saved the selected model for future use.
- Integrated the Machine learning model with the website using flask and saving the executed model file as pickle file.
- Pickel file is used for reusing the compiled model again and again without compiling everytime.
- It increases the speed of website.

# Find the Price of your Dream House

**Number of bedrooms:**

2

**Number of bathrooms:**

1

**Living area:**

2040

**Lot area:**

6090

**Number of floors:**

1

**Waterfront present:**

0

**Number of views:**

0

**Age_of_property:**

30

Submit

# Predicted House price is: ₹587796

# SUMMARY

The project involved analyzing housing data to predict prices in India. After data preprocessing and exploratory data analysis, redundant features were removed, leaving 8 informative features. Three machine learning models (Random Forest, XGBoost, CatBoost) were trained individually and combined using averaging to create an ensemble model. The best model, a combination of the three, achieved high accuracy. A Flask backend was implemented to serve predictions based on user inputs from a frontend interface.
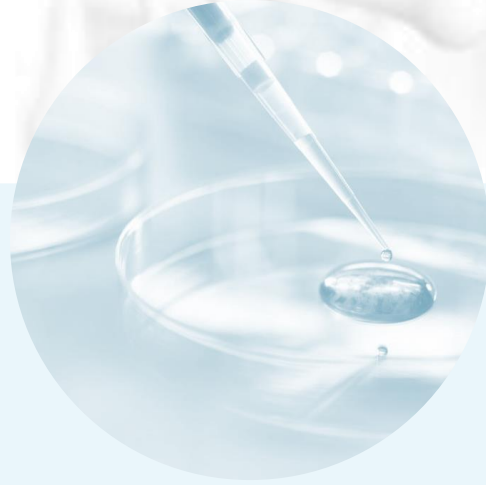
# Reference

https://medium.com/shecodeafrica/predicting-house-prices-gradientboostingregressor-algorithm-ec9d381b0ebc

https://www.mdpi.com/2813-2203/3/1/3

https://gustiyaniz.medium.com/house-price-prediction-with-xgboost-using-r-markdown-d4f891d1f327

https://www.ijraset.com/research-paper/prediction-of-house-prices-using-machine-learning

https://medium.com/hackerdawn/house-prices-prediction-using-random-forest-aa8722347276

# THANK YOU