# House Price Analysis

## 1. Introduction

The real estate market in Bengaluru, often referred to as the *Silicon Valley of India*, is one of the fastest-growing and most complex property markets in the country. Property prices are influenced by multiple interdependent factors such as location, proximity to IT hubs, property size, number of rooms, availability of amenities, and evolving government regulations like RERA.

Accurate house price prediction is essential for **buyers**, **sellers**, **investors**, and **real estate consultants** to make informed, data-driven decisions. Traditional pricing methods rely heavily on manual estimation and market intuition, which can lead to inconsistent or biased results.

This project leverages **data analysis and machine learning techniques** to analyze historical housing data from Bengaluru and build a predictive model capable of estimating house prices based on key property attributes.

## 2. Problem Statement

Predicting house prices in a metropolitan city like Bengaluru is challenging due to several reasons:

- Significant variation in property prices across different localities
- Presence of missing, inconsistent, and unstructured real estate data
- Influence of external economic, social, and regulatory factors
- Non-linear relationship between property size, location, and price

**Project Objectives**

The primary objectives of this project are to:

- Perform **in-depth Exploratory Data Analysis (EDA)** to understand price trends
- Clean, preprocess, and transform raw housing data
- Engineer relevant features to improve predictive accuracy
- Build a reliable **machine learning regression model**
- Evaluate model performance using appropriate metrics
- Extract meaningful insights to support real estate decision-making

# 3. Dataset Description

- **Dataset Name:** Bengaluru House Price Dataset
- **Source:** Publicly available real estate listings
- **Number of Records:** Several thousand residential property listings

## Key Features

- **Location:** Area or neighborhood of the property
- **Size (BHK):** Number of bedrooms, hall, and kitchen
- **Total Square Feet:** Built-up area of the house
- **Bathrooms:** Number of bathrooms available
- **Price:** Property price (target variable)

This dataset reflects real-world property listing challenges such as missing values, varied formats, and location-based price fluctuations.

---

# 4. Tools & Technologies Used

- **Programming Language:** Python
- **Data Analysis Libraries:** Pandas, NumPy
- **Visualization Libraries:** Matplotlib, Seaborn
- **Machine Learning Library:** Scikit-learn
- **Development Environment:** Jupyter Notebook
- **Version Control & Collaboration:** Git, GitHub

---

# 5. Data Preprocessing

Data preprocessing was a critical phase to ensure model accuracy and reliability. The following steps were performed:

- Identified and handled **missing values** using appropriate techniques
- Standardized inconsistent data formats (e.g., size and square footage)
- Converted **categorical variables**, such as locations, into numerical representations
- Removed **outliers** using domain-specific rules (e.g., unrealistic price-per-square-foot values)
- Applied **feature engineering** to create meaningful predictors

These steps significantly improved data quality and model performance.

---

# 6. Exploratory Data Analysis (EDA)

EDA was conducted to uncover trends, patterns, and relationships within the dataset. Key analyses included:

- Distribution of house prices across different locations
- Correlation between **total square footage and price**
- Impact of the number of bedrooms (BHK) and bathrooms on pricing
- Identification of extreme values and anomalies affecting predictions

Data visualizations using **Matplotlib and Seaborn** helped present insights in a clear and intuitive manner, enabling better understanding of market behavior.

---

# 7. Model Development

To predict house prices, the following machine learning approach was implemented:

- **Linear Regression Model** for baseline prediction
- Feature scaling and transformation to normalize input variables
- Data split into **training and testing sets** for unbiased evaluation

The model was trained on cleaned and processed data to predict house prices based on user-defined property attributes.

---

# 8. Model Evaluation

Model performance was assessed using standard evaluation techniques:

- **$R^2$ Score** to measure prediction accuracy
- **Cross-validation** to ensure model robustness and generalization

The model achieved a strong $R^2$ score, indicating that it successfully captured key relationships between property features and prices.

---

# 9. Key Insights

- **Location** is the most influential factor in determining house prices
- Property price increases with square footage, but the relationship is **not strictly linear**
- Certain premium locations consistently show higher **price per square foot**
- Outlier removal significantly improves prediction accuracy

These insights are valuable for both buyers and real estate professionals.

# 10. Conclusion

This project demonstrates the effective application of **data analytics and machine learning** in solving real-world business problems. By analyzing Bengaluru housing data and building a predictive model, the project provides a structured approach to estimating property prices with improved accuracy.

The model serves as a foundation for advanced real estate analytics and can be adapted for other cities or datasets with minimal modifications.

# 11. Future Enhancements

To further improve the project, the following enhancements are recommended:

- Deploy the model as a **web application** using Flask or Streamlit
- Integrate **real-time property listings** through APIs
- Apply advanced algorithms such as **Random Forest, XGBoost, or Gradient Boosting**
- Incorporate additional features like amenities, property age, and proximity to IT hubs