



# Housing Price Prediction :

This project delves into predicting housing prices using advanced regression techniques, focusing on machine learning models to achieve accurate price estimations. We explore various housing features that significantly influence pricing.

Presented by :P.Dhanush

Date :03/08/2025

Course :AI&DS

# Abstract :

Our primary objective is to accurately predict the sales price of houses based on a comprehensive set of features. We utilize a range of regression models, including sophisticated techniques like XGBoost, alongside traditional linear regression. This analysis provides invaluable insights for potential buyers, sellers, and real estate investors, guided by meticulous data preprocessing, insightful visualization, and rigorous model evaluation.

# Introduction :

Accurate housing price prediction is paramount for robust real estate market analysis. Machine learning offers powerful tools to model the complex relationships within housing data. Our dataset incorporates critical features such as house size, number of bedrooms, bathrooms, specific location attributes, and available amenities. We address common challenges including multicollinearity and the effective handling of missing data to ensure model reliability.



# Project Objectives



## Model Accuracy

Develop highly accurate regression models capable of robust house price prediction.



## Data Readiness

Explore and preprocess housing datasets to ensure optimal model readiness and performance.



## Algorithm Comparison

Rigorously compare various regression algorithms to identify the best performers for this task.

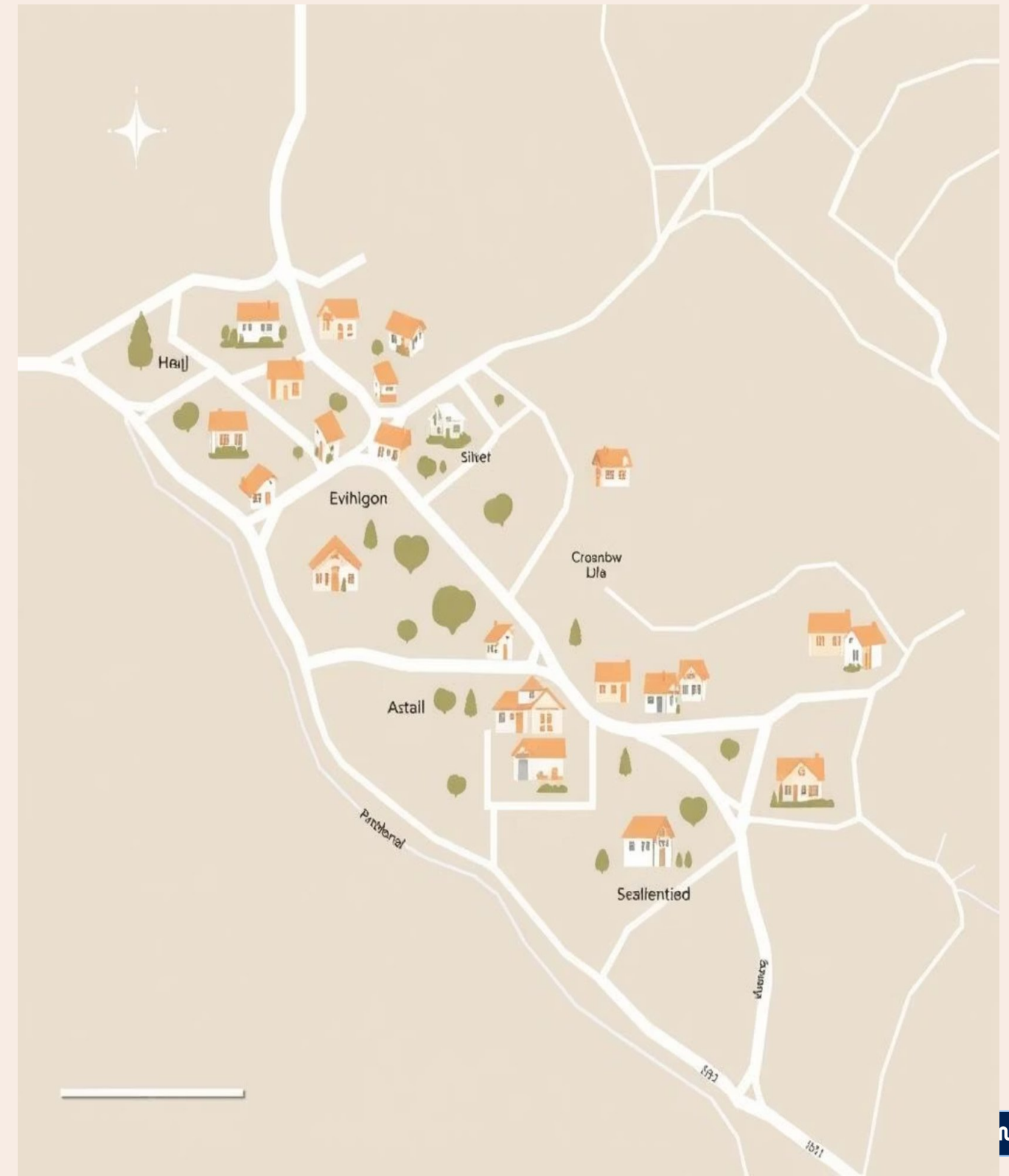


## Actionable Insights

Provide clear, actionable insights derived from comprehensive model evaluation metrics.

# Dataset Overview :

- **Example Dataset:** Utilizes standard datasets like California Housing or the Kaggle Housing Prices dataset.
- **Key Features:** Includes attributes such as lot area, house age, number of bedrooms and bathrooms, stories, and various amenities.
- **Dataset Size:** Varies from approximately 500 samples to over 30,000+, depending on the specific source.
- **Target Variable:** The primary prediction target is the "SalePrice" or "Median House Value."





# Data Preprocessing & Visualization

## Missing Value Handling

Address missing values through strategic deletion or imputation methods like mean/mode.

## Exploratory Data Analysis (EDA)

Conduct comprehensive EDA with correlation heatmaps and scatter plots to reveal data relationships.

## Categorical Encoding

Transform categorical variables into numerical formats using techniques such as OneHotEncoder.

## Train-Test Split

Implement train-test splits (e.g., 70/30 or 60/40) to rigorously evaluate model generalization.

# Regression Models Employed

We explored a diverse range of regression models to capture various complexities in the data:



## Linear Regression

Serves as a foundational baseline model for identifying linear relationships within the data.



## XGBoost Regression

An advanced gradient boosting algorithm, highly effective for modeling complex, non-linear patterns.



## Random Forest Regression

An ensemble method leveraging multiple decision trees for enhanced robustness and accuracy.



## Decision Tree & Gradient Boosting

Additional models explored to provide a comprehensive comparative analysis of performance.

# Model Evaluation Metrics

To assess the efficacy of our predictive models, we utilized several key evaluation metrics:

## R-squared ( $R^2$ )

Quantifies the proportion of variance in the dependent variable explained by the model, indicating goodness of fit.

## Mean Absolute Error (MAE)

Represents the average absolute difference between predicted and actual values, offering a clear measure of error magnitude.

## Root Mean Squared Error (RMSE)

Penalizes larger errors more significantly, providing a robust measure of overall model accuracy.

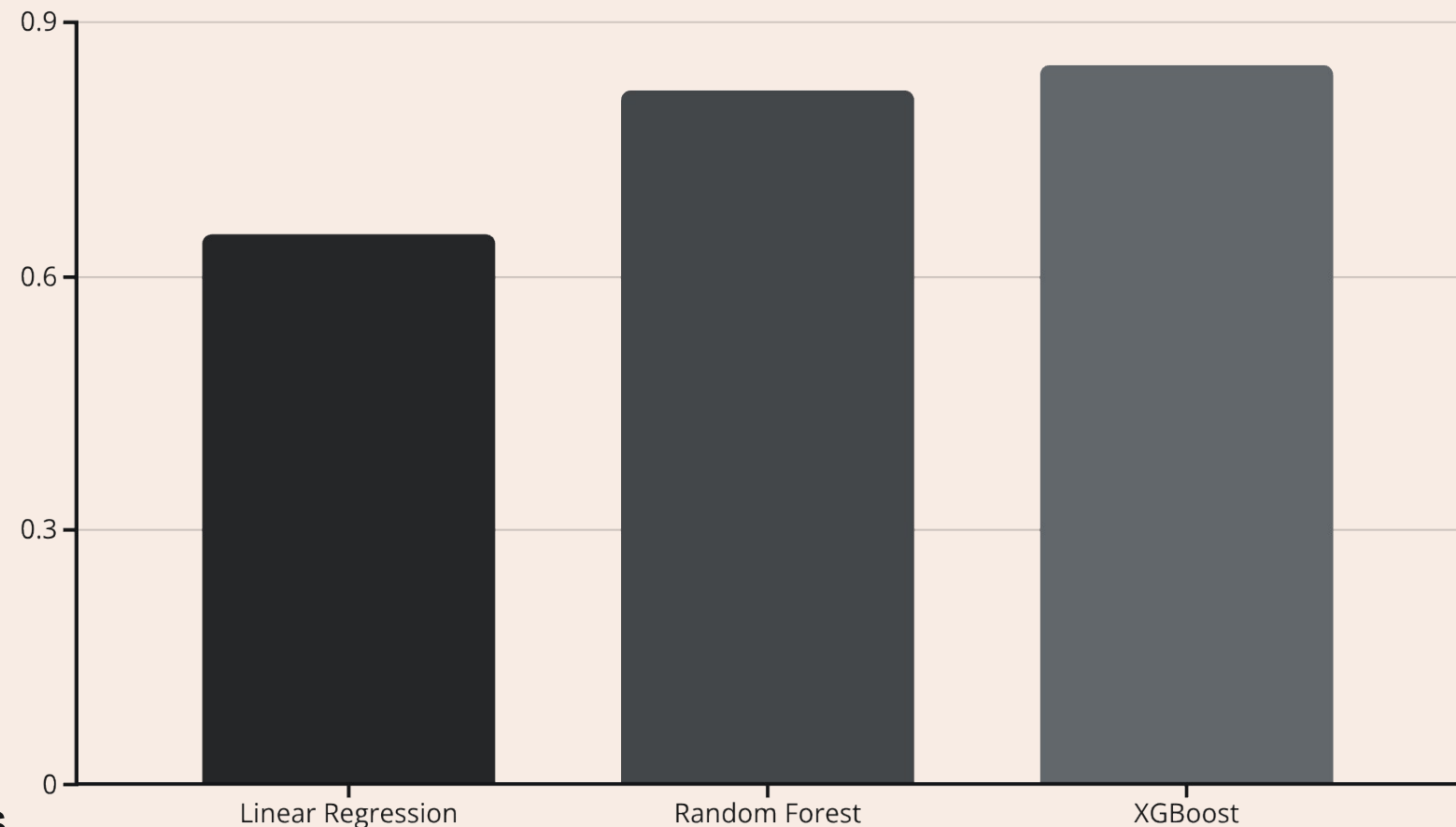
## Visual Comparison

Scatter plots comparing predicted versus actual prices provide intuitive insights into model performance and distribution of errors.



# Results Summary :

- **Performance Leaders:** XGBoost and Random Forest models consistently outperformed linear regression, demonstrating superior predictive power.
- **Accuracy Achieved:** Our best models achieved impressive  $R^2$  scores ranging from 0.7 to 0.85 on unseen test data, indicating high explanatory power.
- **Key Drivers:** Feature importance analysis revealed that location and house size are among the most significant drivers of housing prices.
- **Refinement:** Continuous model tuning and rigorous cross-validation techniques were instrumental in further enhancing prediction accuracy.





## Conclusion :

Our project successfully demonstrates that advanced regression techniques are highly effective for predicting housing prices with considerable accuracy. The journey from raw data to actionable insights hinges on meticulous data preprocessing and innovative feature engineering. Ultimately, robust model evaluation guides the selection of the most powerful predictive approach, empowering informed decision-making across real estate markets.