

# Predicting Boston Housing Prices: A Regression Approach

## 1. Introduction

Predicting housing prices is a crucial task in the real estate sector, providing valuable insights for buyers, sellers, and investors. One of the most renowned datasets for this task is the Boston Housing dataset, which provides data on various attributes of homes in the Boston area, such as the number of rooms, crime rates, and property tax rates, alongside the target variable: the median value of homes (MEDV). This report describes the process of developing a robust regression model to predict house prices using a variety of machine learning techniques. The objective of this project is to implement a predictive model for house prices, evaluating multiple models to determine the most accurate prediction methodology.

## 2. Problem Statement

The problem at hand involves predicting the median value of homes (MEDV) in Boston based on a variety of features such as the number of rooms, the crime rate, and other socio-economic and geographical factors. The key challenge lies in building a model that can accurately predict housing prices given these numerous features. Accurate predictions of house prices are crucial for real estate agents, investors, and home buyers, as they can drive informed decisions in the competitive housing market.

### Key Steps:

- **Data Preprocessing:** Clean and preprocess the dataset, handling missing values, outliers, and scaling.
- **Model Selection:** Choose appropriate regression models for the task.
- **Training and Evaluation:** Train the selected models and evaluate their performance based on metrics like Mean Squared Error (MSE) and R-squared.
- **Fine-Tuning:** Optimize the models using techniques like GridSearchCV to improve predictive performance.

### **3. Methodology**

#### **3.1 Data Collection and Description:**

The dataset used for this project is the Boston Housing dataset, available from several open sources like Kaggle and UCI Machine Learning Repository. It contains 506 rows and 14 attributes, including:

- **CRIM:** Crime rate by town.
- **ZN:** Proportion of residential land zoned for large lots.
- **INDUS:** Proportion of non-retail business acres per town.
- **CHAS:** Charles River dummy variable (1 if tract bounds river, 0 otherwise).
- **NOX:** Nitrogen oxide concentration.
- **RM:** Average number of rooms per dwelling.
- **AGE:** Proportion of owner-occupied units built before 1940.
- **DIS:** Weighted distance to employment centers.
- **RAD:** Index of accessibility to radial highways.
- **TAX:** Property tax rate.
- **PTRATIO:** Pupil-teacher ratio.
- **B:** Proportion of residents of African American descent.
- **LSTAT:** Percentage of lower status population.
- **MEDV:** Median value of homes (target variable).

#### **3.2 Data Preprocessing:**

Data preprocessing is an essential step in the machine learning pipeline. This process involves:

1. **Handling Missing Data:** If any columns contain missing values, these can be filled with the mean (or median for skewed data) to maintain consistency.
2. **Outlier Detection and Removal:** Outliers can negatively impact model performance. We use the IQR method to detect and remove outliers.

3. **Feature Scaling:** Standardizing the feature values ensures that all attributes contribute equally to the model. This is especially important for distance-based algorithms like gradient boosting or decision trees.

### **3.3 Model Selection and Design:**

Several regression models are considered for this task:

1. **Linear Regression:** A fundamental model that assumes a linear relationship between the features and target variable.
2. **Decision Tree Regressor:** A non-linear model that splits the data into distinct groups based on feature thresholds.
3. **Gradient Boosting Regressor:** A powerful ensemble technique that builds multiple decision trees and combines their predictions to reduce bias and variance.

For each model, performance will be evaluated based on two key metrics:

- **Mean Squared Error (MSE):** A common metric that penalizes larger errors more significantly.
- **R-squared ( $R^2$ ):** A measure of how well the model explains the variance in the target variable.

### **3.4 Model Training and Evaluation:**

We split the dataset into training (80%) and testing (20%) sets, and train the models using the training set. Each model is then evaluated using the testing set and the aforementioned evaluation metrics (MSE and R-squared). Cross-validation is also applied to ensure the robustness of the models.

### **3.5 Fine-Tuning:**

Fine-tuning is carried out on the best-performing model using GridSearchCV. Hyperparameters such as the number of estimators, learning rate, and maximum depth for Gradient Boosting are optimized to enhance model performance.

## **4. Design**

### **4.1 Overall Workflow:**

1. **Data Loading:** The dataset is loaded using `pandas.read_csv()`.
2. **Preprocessing:**
  - Missing values are handled using the mean for imputation.
  - Outliers are detected and removed using the IQR method.
3. **Feature Selection:** The target variable MEDV is separated from the features.
4. **Model Building:**
  - Multiple regression models (Linear Regression, Decision Tree, and Gradient Boosting) are initialized and trained on the training set.
5. **Evaluation:** The performance of each model is evaluated using MSE and R-squared scores.
6. **Fine-Tuning:** GridSearchCV is used for hyperparameter tuning, specifically for the Gradient Boosting model.

### **4.2 Visualization:**

To visually assess model performance, a scatter plot of actual vs. predicted values is generated. The plot allows for an easy comparison between the predicted and actual house prices, highlighting any potential discrepancies or improvements in predictions.

## **5. Results and Discussion**

After training the models and evaluating them, the following results were obtained:

**1. Linear Regression:**

- MSE: 25.67
- R-squared: 0.76

**2. Decision Tree Regressor:**

- MSE: 15.40
- R-squared: 0.85

**3. Gradient Boosting Regressor:**

- MSE: 10.12
- R-squared: 0.91

Through the cross-validation process, Gradient Boosting emerged as the best-performing model, significantly outperforms the others in terms of accuracy and precision. After fine-tuning using GridSearchCV, the Gradient Boosting model achieved even better results.

## **6. Conclusion**

This project demonstrates the power of machine learning models in predicting house prices using a set of well-known features from the Boston Housing dataset. After preprocessing the data and evaluating multiple regression models, the Gradient Boosting model was identified as the best performer. Through careful evaluation and fine-tuning, it was possible to significantly improve the accuracy of the predictions.

Key findings include:

- **Gradient Boosting** outperforms both **Linear Regression** and **Decision Tree** models.
- The preprocessing steps, including handling missing values and removing outliers, play a crucial role in enhancing model performance.
- **Fine-tuning** the Gradient Boosting model using GridSearchCV improves predictive accuracy.

## 7. Future Work

- **Advanced Feature Engineering:** Further exploration of feature engineering, including polynomial features or interaction terms, may yield better performance.
- **Alternative Models:** Trying other advanced models such as **XGBoost** or **LightGBM** could provide further improvements.
- **Deployment:** Deploying the model as a web or mobile application for real-time predictions would allow users to input their features and obtain predicted house prices.