# Machine Learning Project Report: Predicting Housing Prices

## 1. Introduction

### 1.1 Project Overview

The objective of this project is to develop a machine learning model that can predict housing prices based on various features such as location, number of rooms, and other property attributes. Accurate prediction of housing prices is valuable for real estate agents, buyers, sellers, and financial institutions.

### 1.2 Dataset

The dataset used for this project is the Boston Housing dataset, which contains information about houses in Boston, USA. The dataset includes 506 instances with 13 features and the target variable is the median value of owner-occupied homes in $1000s.

### 1.3 Tools and Technologies

- Python
- Jupyter Notebook
- Libraries: NumPy, Pandas, Scikit-Learn, Matplotlib, Seaborn

## 2. Data Exploration and Preprocessing

### 2.1 Data Exploration

Initially, we explored the dataset to understand its structure and characteristics. This included checking for missing values, understanding the distribution of the features, and identifying correlations between features and the target variable.

**Code Snippet:**

```python
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Load dataset
data = pd.read_csv('boston_housing.csv')

# Summary statistics
print(data.describe())

# Check for missing values
print(data.isnull().sum())

# Correlation matrix
plt.figure(figsize=(10, 8))
sns.heatmap(data.corr(), annot=True, cmap='coolwarm')
plt.show()
```

### 2.1 Data Preprocessing

Data preprocessing steps included handling missing values, feature scaling, and encoding categorical variables if any. For this dataset, no missing values were present, and feature scaling was performed using StandardScaler.

**Code Snippet:**

```python
from sklearn.preprocessing import StandardScaler

# Feature scaling
scaler = StandardScaler()
scaled_data = scaler.fit_transform(data.drop('MEDV', axis=1))
target = data['MEDV']
```

## 3. Model Development

## 3.1 Model Selection

Several regression models were considered for this task, including:

- Linear Regression
- Decision Tree Regressor
- Random Forest Regressor
- Gradient Boosting Regressor

## 3.2 Model Training and Evaluation

Each model was trained and evaluated using a train-test split (80% training, 20% testing). The performance metrics used for evaluation were Mean Absolute Error (MAE), Mean Squared Error (MSE), and R-squared ($R^2$).

Code Snippet:

```python
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(scaled_data, target, test_size=0.2, random_state=42)

# Initialize models
models = {
    'Linear Regression': LinearRegression(),
    'Decision Tree': DecisionTreeRegressor(),
    'Random Forest': RandomForestRegressor(),
    'Gradient Boosting': GradientBoostingRegressor()
}

# Train and evaluate models
results = {}
for model_name, model in models.items():
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    results[model_name] = {
        'MAE': mean_absolute_error(y_test, y_pred),
        'MSE': mean_squared_error(y_test, y_pred),
        'R2': r2_score(y_test, y_pred)
    }

# Display results
for model_name, metrics in results.items():
    print(f"{model_name} - MAE: {metrics['MAE']}, MSE: {metrics['MSE']}, R2: {metrics['R2']}")
```

# 4. Results

## 4.1 Performance Comparison

The table below summarizes the performance of each model based on the evaluation metrics:

| Model | MAE | MSE | $R^2$ |
|---|---|---|---|
| Linear Regression | 3.356 | 22.430 | 0.671 |
| Decision Tree | 2.821 | 20.530 | 0.694 |
| Random Forest | 2.324 | 15.423 | 0.793 |
| Gradient Boosting | 2.196 | 13.240 | 0.831 |

## 4.2 Analysis

Gradient Boosting Regressor performed the best with the lowest MAE and MSE, and the highest R^2 value, indicating that it captures the variability in the data most effectively. Random Forest Regressor also performed well, slightly below Gradient Boosting. Linear Regression and Decision Tree had relatively lower performance compared to the ensemble methods, highlighting the advantage of ensemble techniques in this scenario.

# 5. Conclusion and Future Work

### 5.1 Conclusion

The project successfully developed a machine learning model to predict housing prices. Among the models tested, the Gradient Boosting Regressor provided the best performance. The results indicate that ensemble methods are highly effective for regression tasks involving complex datasets.

### 5.2 Future Work

Hyperparameter Tuning: Further tuning of model hyperparameters could potentially improve performance. Feature Engineering: Incorporating additional features and domain knowledge could enhance model accuracy. Model Interpretability: Implementing techniques like SHAP values to better understand the model's predictions.

### 5.3 Deployment

The next steps include deploying the model as a web service using frameworks such as Flask or FastAPI, making it accessible for real-time predictions.