# Machine Learning Report

**Introduction**

The objective of this report is to explore, preprocess, and analyze the dataset using machine learning techniques. By leveraging Python libraries such as Pandas, NumPy, Matplotlib, Seaborn, and Scikit-Learn, we have implemented a structured process to clean the data, build predictive models, and evaluate their performance. The results provide valuable insights into the factors driving house prices and the effectiveness of different machine learning models for this type of problem.

This project is particularly relevant in the housing market, where accurate predictions of house prices can assist policymakers, real estate professionals, and home buyers in making informed decisions. By identifying the most critical features, such as median income and proximity to significant locations, this analysis provides actionable insights into the dynamics of the housing market. Furthermore, the use of visualizations and advanced modeling techniques highlights the power of machine learning in uncovering complex relationships within the data.

**Problem Overview**

The housing dataset provides information on housing characteristics, median income, geographical details, and proximity to water bodies. The target variable, median_house_value, represents the house prices we aim to predict. However, the dataset presents several challenges that need to be addressed before building models, including:

1. **Handling Missing Data**: Some features have missing values that could skew the analysis.

2. **Categorical Data**: The ocean_proximity column contains categorical data that must be encoded for machine learning.

3. **Correlation Between Features**: Some features may be redundant or weakly correlated with the target, requiring careful examination.

The goal is to build robust models that accurately predict house values and provide insights into the most significant predictors. Addressing these challenges ensures that the models are reliable and that the insights drawn from them are meaningful for real-world applications.

**Steps Followed**

To systematically address the problem, we followed these key steps:

**1. Importing and Cleaning the Data**

- **Data Import**: The dataset was loaded into a Pandas DataFrame for exploration and preprocessing.

- **Cleaning**:

  o Missing values were identified and replaced with the median of the respective columns to maintain the dataset's integrity.

  o The categorical ocean_proximity column was encoded using one-hot encoding, enabling models to process it effectively.

  o Duplicate entries were checked and removed if found to ensure unique and valid data.

## 2. Exploratory Data Analysis (EDA)

- The dataset's statistical summary was analyzed to understand the range, mean, and standard deviation of each feature.

- Visualizations, such as histograms, were created to observe the distribution of variables like median_income and median_house_value show in below images.

- Boxplot for median house value reveals that most house values range between $100,000 and $300,000, with a median around $200,000. The outliers near $500,000 indicate a dataset-imposed price cap, which could impact the model's predictions.



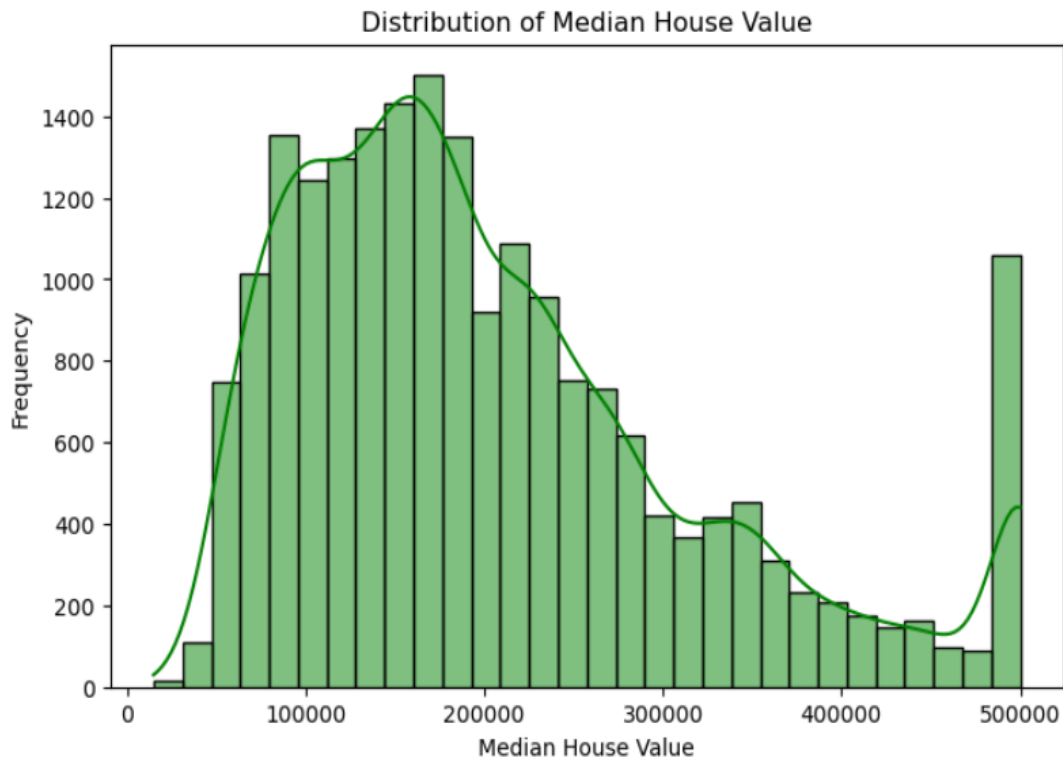Distribution of Median Income

This histogram of "Median Income" shows the frequency distribution of the median income values in the dataset. Key observations include:

1. **Skewed Distribution**: The distribution is right-skewed, with most values concentrated between 2 and 6, indicating that the majority of neighborhoods have a median income in this range.
2. **Peak Around 3-4**: The highest frequency of median income is between 3 and 4, suggesting this is the most common income bracket for households in the dataset.

3. **Outliers**: There are a few higher income values (e.g., above 10) with very low frequency, which may represent wealthy neighborhoods.

This analysis helps us understand the economic composition of the dataset and highlights that median income is a critical factor in predicting house values.



This histogram shows the **distribution of median house values** in the dataset. Key observations include:

1. **Concentration of Values**: Most house values are clustered between $100,000 and $300,000, indicating this range is the most common price bracket.
2. **Right Skewness**: The distribution is slightly right-skewed, with a tail extending toward higher house prices. This suggests that a smaller number of houses are priced above $400,000.
3. **Price Cap**: The spike at $500,000 indicates a price cap in the dataset, likely due to the dataset's recording limits. This cap might affect predictions and needs to be considered during modeling.

### 3. Splitting the Data

- The dataset was divided into:

  - **Training Set (80%)**: Used to train the models.

  - **Testing Set (20%)**: Used to evaluate model performance on unseen data.

This splitting strategy ensures that the model generalizes well to new data.

## 4. Model Selection and Training

Three regression models were selected:

- **Linear Regression**: A simple, interpretable model that assumes linear relationships.

- **Random Forest**: An ensemble model that captures non-linear patterns and is robust to overfitting.

- **Decision Tree**: A tree-based model that is easy to interpret but prone to overfitting.

## 5. Evaluation

The models were evaluated using:

- **Mean Squared Error (MSE)**: Quantifies prediction errors.

- **R-Squared ($R^2$)**: Measures the proportion of variance in the target explained by the model.

## 6. Visualization and Insights

- Scatterplots and feature importance plots were used to interpret model performance.

- A correlation heatmap was generated to identify relationships between features.

## Modeling and Training

### Linear Regression

Linear regression assumes a linear relationship between the features and the target variable. While it is computationally efficient and interpretable, it may fail to capture complex patterns in the data. The model was trained on the training set and evaluated on the testing set.

### Random Forest

Random Forest is an ensemble method that combines multiple decision trees to improve prediction accuracy. It is particularly effective for datasets with non-linear relationships and can provide insights into feature importance. The model was trained with default hyperparameters and performed exceptionally well compared to the other models.

### Decision Tree

The Decision Tree model splits the dataset into branches based on feature thresholds. Although easy to interpret, it can overfit the training data without proper pruning. Despite this limitation, it provides a baseline comparison for more complex models like Random Forest.

## Evaluation Results

The performance of each model was assessed using the testing set. The evaluation metrics and insights are as follows:

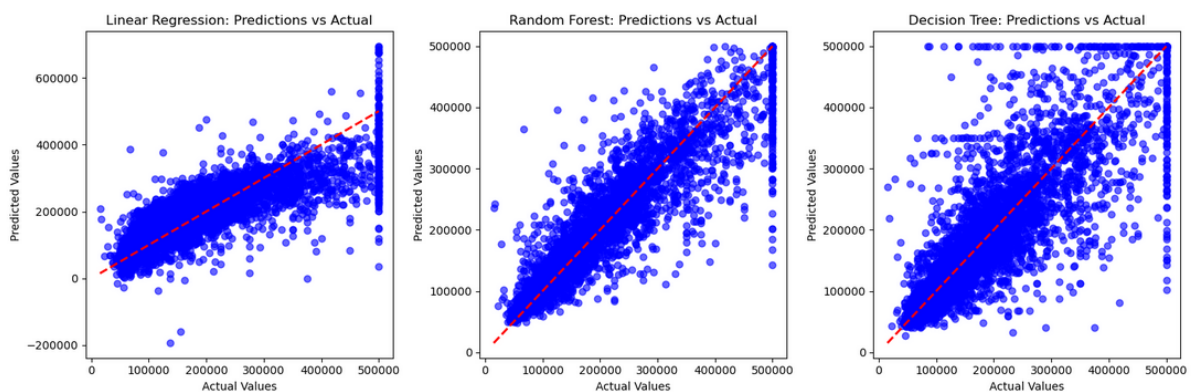| Model | Mean Squared Error (MSE) | R-Squared (R²) |
|---|---|---|
| Linear Regression | 4,908,476,721 | 0.63 |
| Random Forest | 2,404,745,975 | 0.82 |
| Decision Tree | 4,865,868,837 | 0.63 |

**Key Observations**:

1. **Random Forest** outperformed the other models, achieving the lowest MSE and the highest $R^2$.

2. **Linear Regression** and **Decision Tree** had similar performance, indicating limitations in capturing the dataset's complexity.
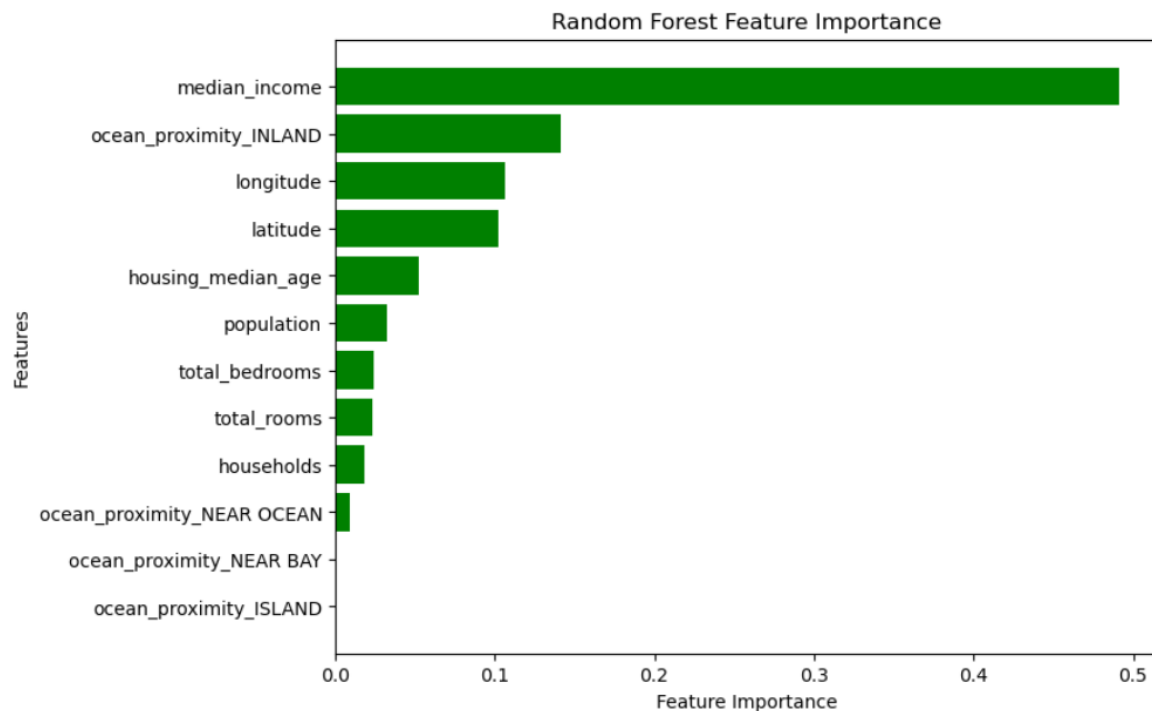
**Visualizations**

**Scatter Plots**

Scatter plots of predicted vs. actual house values were generated for each model. The diagonal line represents perfect predictions, and points closer to this line indicate better model performance.

- **Linear Regression**: Significant deviations from the diagonal, especially for high house values. The model struggled to accurately predict the target for these cases, indicating limitations in its capacity to model non-linear relationships.

- **Random Forest**: Points were closely aligned with the diagonal, indicating high accuracy and robust modeling of complex patterns in the data. This performance demonstrates its ability to effectively handle non-linearities and interactions among features.

- **Decision Tree**: Similar to Linear Regression, with greater scatter and inconsistent predictions for higher house values. Overfitting to the training data likely caused these inaccuracies.

**Feature Importance (Random Forest)**

Random Forest provides insights into the relative importance of each feature. The top features were:
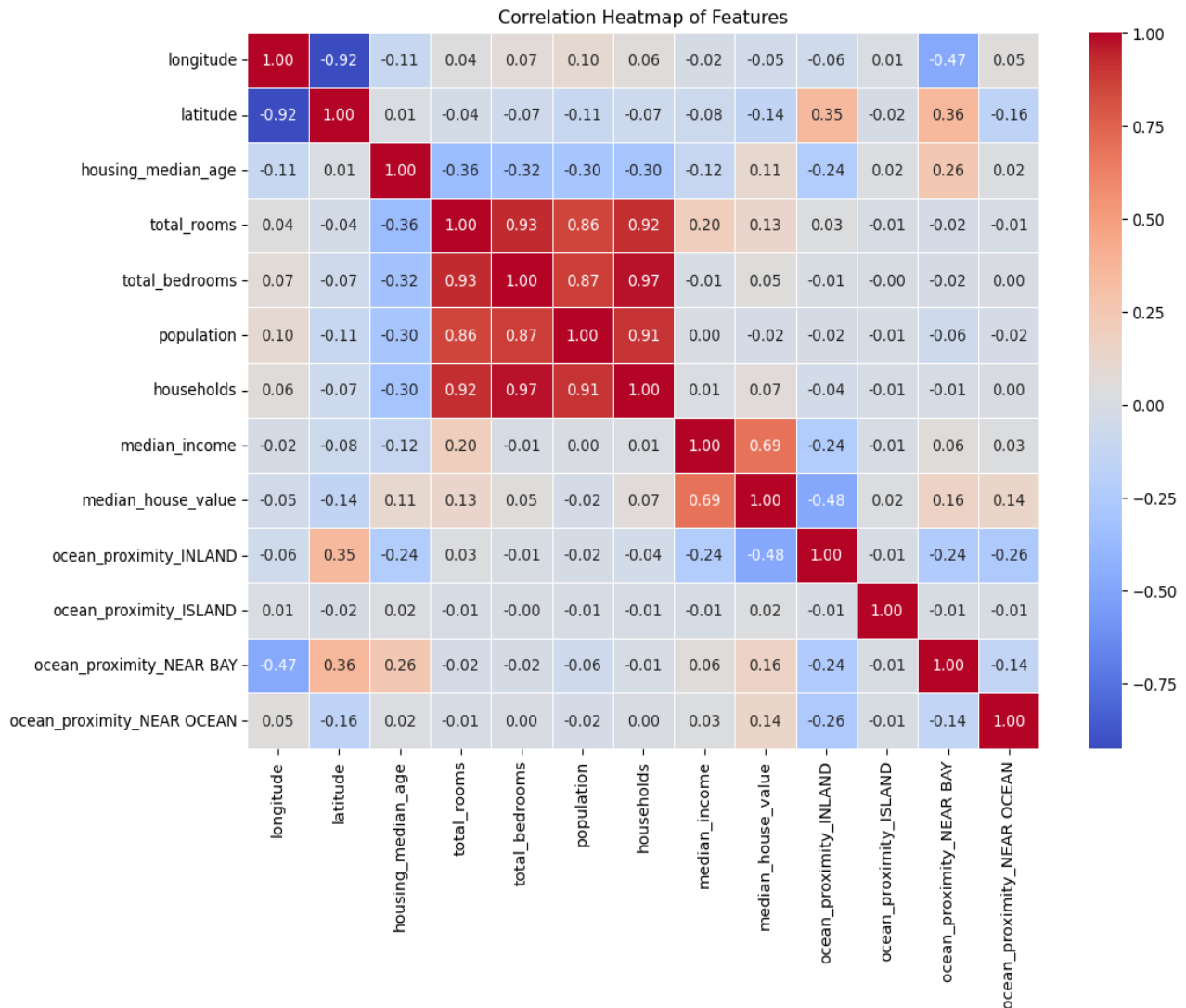


- **Median Income**: 49.1%

- **Ocean Proximity (INLAND)**: 14.0%

- **Longitude**: 10.6%

- **Latitude**: 10.1%

Features such as median_income and ocean_proximity strongly influence house values, while features like households and total_rooms have minimal impact. This analysis highlights the critical role of economic and geographic factors in determining housing prices.

**Correlation Heatmap**

A correlation heatmap was created to visualize relationships between features and the target variable:

- **Median Income**: Strong positive correlation (0.69) with median_house_value, indicating that wealthier neighborhoods are associated with higher house prices.

- **Ocean Proximity (INLAND)**: Moderate correlation, reflecting the impact of location. Inland areas tend to have lower house values compared to those near water bodies.

- **Total Rooms and Population**: Weak correlations, suggesting limited predictive value. These features likely reflect aggregate metrics that do not directly influence individual house prices.

Correlation Heatmap of Features

**Additional Insights**

1. **Feature Insights**:

    o  Economic factors like median_income significantly influence house values. The higher the median income in a neighborhood, the more likely it is for house values to be higher, underscoring the importance of economic status.

    o  Geographical features such as longitude and latitude also play important roles. The correlation between these variables and house values reflects regional differences in housing markets.

    o  Proximity to water bodies has varying effects, with INLAND areas showing notable impact. This underscores the desirability of waterfront properties in determining house prices.

2. **Model Selection**:

- o  Random Forest emerged as the best-performing model, balancing accuracy and interpretability. Its ability to capture non-linear patterns and interactions among features made it the most effective tool for this analysis.

- o  Linear Regression and Decision Tree models were less effective due to their inability to capture complex relationships. While Linear Regression provides simplicity, it lacks the flexibility needed for this dataset.

3. **Improvement Opportunities**:

- o  Hyperparameter tuning for Random Forest could further improve accuracy. Adjusting parameters such as the number of trees and depth of the trees may yield better results.

- o  Incorporating additional features, such as local economic indicators, school quality, or crime rates, might enhance predictive power and provide a broader understanding of housing markets.

- o  Exploring advanced algorithms like Gradient Boosting or Neural Networks could yield better results. These methods are capable of handling complex patterns and might outperform Random Forest in certain scenarios.

4. **Exploration of Relationships**:

- o  The data shows that while median_income is the most influential factor, other features such as location proximity and housing age contribute to nuances in pricing. These factors should be considered when analyzing regional housing trends.

- o  Features like total_rooms and population may require further feature engineering to extract more meaningful insights, such as per capita metrics.

5. **Interpretability**:

- o  Random Forest provides a good trade-off between accuracy and interpretability. It helps highlight key drivers behind housing price trends, making it a valuable tool for both policymakers and market analysts.

- o  The visualizations offer an accessible way to communicate these insights to stakeholders, enabling data-driven decisions.

**Conclusion**

This assessment demonstrates the effectiveness of machine learning techniques in predicting house values and extracting meaningful insights from data. The Random Forest model stood out for its superior performance and ability to highlight key features like median_income and ocean_proximity. Visualizations such as scatter plots, feature importance charts, and correlation heatmaps provided valuable interpretability and understanding of the dataset.

Future work could involve fine-tuning models, exploring additional algorithms, and incorporating domain-specific features to further enhance predictions and insights. This project

highlights the power of machine learning in addressing complex, real-world problems and supports data-driven decision-making in the housing market.