```
1 import pandas as pd
 2 import numpy as np #9
 3 import matplotlib.pyplot as plt #9
 4 from sklearn.model_selection import train_test_split #4
 5 from sklearn.preprocessing import StandardScaler #5
 6 from sklearn.linear_model import LinearRegression #6 & #9
 7 from sklearn.metrics import mean_squared_error #6 & #7
 8 from sklearn.tree import DecisionTreeRegressor #7
 9 from sklearn.ensemble import RandomForestRegressor #8
10
11 # Ignore printing warnings for general readability
12 import warnings
13 warnings.filterwarnings('ignore')
 1 #Making a list of missing value types
 2 missing_values = ["n/a", "na", "--", "...", "NaN"]
 4 # 1. # Read the Excel file using the specified engine and convert to CSV
 5 data = pd.read_excel('1553768847_housing.xlsx', na_values=missing_values)
 6 data.to_csv('1553768847_housing.csv', index=False)
 7 data.head(20)
```

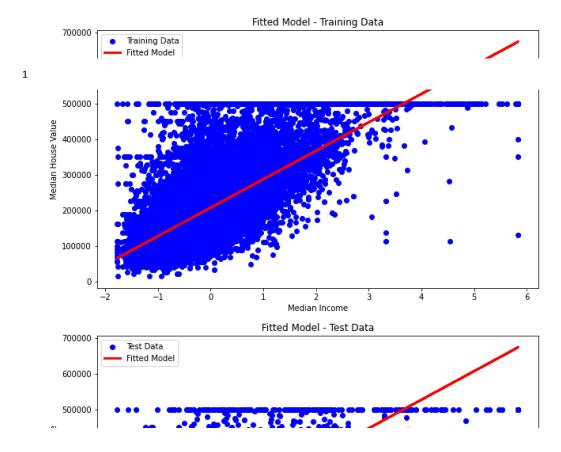


```
longitude latitude housing_median_age total_rooms total_bedrooms population households median_income
            -122.23
                        37.88
      0
                                               41
                                                           088
                                                                         129.0
                                                                                       322
                                                                                                   126
                                                                                                                8.3252
      1
            -122.22
                        37.86
                                               21
                                                          7099
                                                                         1106.0
                                                                                      2401
                                                                                                  1138
                                                                                                                8.3014
      2
            -122.24
                        37.85
                                               52
                                                          1467
                                                                         190.0
                                                                                       496
                                                                                                   177
                                                                                                                7.2574
      3
            -122.25
                        37.85
                                               52
                                                          1274
                                                                         235.0
                                                                                       558
                                                                                                   219
                                                                                                                5.6431
      4
                                                                         280.0
            -122.25
                        37.85
                                               52
                                                          1627
                                                                                       565
                                                                                                   259
                                                                                                                3.8462
 1 # 2. Handle missing values by filling with column mean
 2 data filled = data.fillna(data.mean())
 1 #Identify the data variables which are categorical
 2 #Categorical(classified as strings or integers) columns can be further classified into:
 3 #-nominal(no inherent order/ranking since no specific order)
 4 #-ordinal variables (specific order/ranking) or satisfaction ratings (e.g., low < medium < high)
 6 #Continuous any numerical value(float/int) within a specific range(eg:fractional or decimal values/ ie: age, height, weight, or temperature)
 7 Categorical Columns = []
 8 Numerical Columns = []
 9 for column in data filled.columns:
       if len(data filled[column].unique()) <= 8:</pre>
11
           Categorical Columns.append(column)
12
       else:
13
           Numerical Columns.append(column)
14
15 print("\n Categorical Variables are : " , Categorical Columns)
16 print(" \n Numerical Variables are : " , Numerical Columns)
      Categorical Variables are : ['ocean proximity']
      Numerical Variables are: ['longitude', 'latitude', 'housing_median_age', 'total_rooms', 'total_bedrooms', 'population', 'households', 'median_income',
 1 # 3. Encode categorical data: Convert categorical columns to numerical data
 2 data encoded = pd.get_dummies(data_filled, columns=['ocean_proximity'])
 4 # Display the first few rows of the encoded dataset
 5 data encoded.head()
```

```
longitude latitude housing_median_age total_rooms total_bedrooms population households median_income
           -122.23
                       37.88
                                             41
                                                          880
                                                                        129.0
                                                                                      322
                                                                                                  126
                                                                                                              8.3252
           -122.22
                                                                                                              8.3014
     1
                       37.86
                                              21
                                                         7099
                                                                       1106.0
                                                                                     2401
                                                                                                 1138
           -122.24
                                              52
                                                                                                              7.2574
                       37.85
                                                        1467
                                                                        190.0
                                                                                      496
                                                                                                  177
 1 # Split the dataset into features (X) and target variable (y)
 2 X = data_encoded.drop('median_house_value', axis=1)
 3 y = data_encoded['median_house_value']
 5 #4. Split data into 80% training and 20% testing sets
 6 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
 1 #5. Standardize the data
 2 scaler = StandardScaler()
 3 X train scaled = scaler.fit transform(X train)
 4 X test scaled = scaler.transform(X test)
 1 #6. Perform Linear Regression
 2 model = LinearRegression()
 4 # Train the model on the training data
 5 model.fit(X train scaled, y train)
 7 # Predict the target variable on the test data
 8 y pred = model.predict(X test scaled)
10 # Calculate Root Mean Squared Error (RMSE)
11 rmse = mean_squared_error(y_test, y_pred, squared=False)
12 print(f"Root Mean Squared Error (RMSE): {rmse}")
     Root Mean Squared Error (RMSE): 70031.41991955665
 1 #7. Perform Decision Tree Regression model
 2 tree reg model = DecisionTreeRegressor()
 4 # Fit the model on the training data
 5 tree reg model.fit(X train scaled, y train)
 7 # Predict output for the test dataset
 8 y pred tree = tree reg model.predict(X test scaled)
10 # Calculate Root Mean Squared Error (RMSE) for Decision Tree Regression
11 rmse_tree = mean_squared_error(y_test, y_pred_tree, squared=False)
```

```
12
13 print(f"RMSE (root mean squared error) from Random Forest Regression: {rmse tree}")
     RMSE (root mean squared error) from Random Forest Regression: 69228.49061115789
 1 #8. Perform Random Forest Regression model
 2 forest_reg_model = RandomForestRegressor()
 4 # Fit the model on the training data
 5 forest reg model.fit(X train scaled, y train)
 7 # Predict output for the test dataset
 8 y_pred_forest = forest_reg_model.predict(X_test_scaled)
10 # Calculate Root Mean Squared Error (RMSE) for Random Forest Regression
11 rmse forest = mean squared error(y test, y pred forest, squared=False)
12
13 print(f"Random Forest Regression RMSE: {rmse forest}")
     Random Forest Regression RMSE: 48884.92495518843
 1 #9. Plot the fitted model for test data
 2 # Extract just the median income column
 3 X_train_median_income = X_train_scaled[:, X_train.columns.get_loc('median_income')]
 4 X_test_median_income = X_test_scaled[:, X_test.columns.get_loc('median_income')]
 6 # Reshape the data to fit the model
 7 X_train_median_income = X_train_median_income.reshape(-1, 1)
 8 X_test_median_income = X_test_median_income.reshape(-1, 1)
10 # Initialize Linear Regression model
11 linear_reg_single_var = LinearRegression()
12
13 # Fit the model on the training data
14 linear_reg_single_var.fit(X_train_median_income, y_train)
15
16 # Predict output for the test dataset
17 y pred single var = linear reg single var.predict(X test median income)
18
19 # Plot the fitted model for training data
20 plt.figure(figsize=(10, 6))
21 plt.scatter(X train median income, y train, color='blue', label='Training Data')
22 plt.plot(X_train_median_income, linear_reg_single_var.predict(X_train_median_income), color='red', linewidth=3, label='Fitted Model')
23 plt.xlabel('Median Income')
24 plt.ylabel('Median House Value')
25 plt.title('Fitted Model - Training Data')
26 plt.legend()
27 plt.show()
```

```
28
29 # Plot the fitted model for test data
30 plt.figure(figsize=(10, 6))
31 plt.scatter(X_test_median_income, y_test, color='blue', label='Test Data')
32 plt.plot(X_test_median_income, y_pred_single_var, color='red', linewidth=3, label='Fitted Model')
33 plt.xlabel('Median Income')
34 plt.ylabel('Median House Value')
35 plt.title('Fitted Model - Test Data')
36 plt.legend()
37 plt.show()
```



×