

✓ Load data

Subtask:

Retry loading the dataset from a different path, including error handling.

Reasoning: Attempt to load the data from the specified path with error handling and display the head and info of the dataframe if successful.

```
import pandas as pd

try:
    df = pd.read_csv('/data/data.csv')
    display(df.head())
    display(df.info())
except FileNotFoundError:
    print("Error: The file was not found at /data/data.csv")
```

Error: The file was not found at /data/data.csv

Toggle Gemini

✓ Load data

Subtask:

Retry loading the dataset from a different path, including error handling.

Reasoning: Import pandas and attempt to load the dataset using a try-except block as instructed.

```
import pandas as pd

try:
    df = pd.read_csv('/content/sample_data/california_housing_train.csv')
    display(df.head())
    display(df.info())
except FileNotFoundError:
    print("Error: The file was not found at /content/sample_data/california_housing_train.csv")
```

| | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income | median_house_value |
|---|-----------|----------|--------------------|-------------|----------------|------------|------------|---------------|--------------------|
| 0 | -114.31 | 34.19 | 15.0 | 5612.0 | 1283.0 | 1015.0 | 472.0 | 1.4936 | 66900.0 |
| 1 | -114.47 | 34.40 | 19.0 | 7650.0 | 1901.0 | 1129.0 | 463.0 | 1.8200 | 80100.0 |
| 2 | -114.56 | 33.69 | 17.0 | 720.0 | 174.0 | 333.0 | 117.0 | 1.6509 | 85700.0 |
| 3 | -114.57 | 33.64 | 14.0 | 1501.0 | 337.0 | 515.0 | 226.0 | 3.1917 | 73400.0 |
| 4 | -114.57 | 33.57 | 20.0 | 1454.0 | 326.0 | 624.0 | 262.0 | 1.9250 | 65500.0 |

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 17000 entries, 0 to 16999
Data columns (total 9 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   longitude              17000 non-null float64
1   latitude               17000 non-null float64
2   housing_median_age     17000 non-null float64
3   total_rooms            17000 non-null float64
4   total_bedrooms        17000 non-null float64
5   population             17000 non-null float64
6   households             17000 non-null float64
7   median_income          17000 non-null float64
8   median_house_value     17000 non-null float64
dtypes: float64(9)
memory usage: 1.2 MB
None
```

✓ Preprocess data

Subtask:

Clean and prepare your data for model training. This might involve handling missing values, encoding categorical features, scaling numerical features, and splitting the data into training, validation, and test sets.

Reasoning: Check for missing values in the DataFrame.

```
display(df.isnull().sum())
```

| | 0 |
|--------------------|---|
| longitude | 0 |
| latitude | 0 |
| housing_median_age | 0 |
| total_rooms | 0 |
| total_bedrooms | 0 |
| population | 0 |
| households | 0 |
| median_income | 0 |
| median_house_value | 0 |

dtype: int64

Reasoning: There are no missing values, so the next steps are to separate the features from the target and then split the data into training, validation, and test sets.

```
from sklearn.model_selection import train_test_split

X = df.drop('median_house_value', axis=1)
y = df['median_house_value']

X_train, X_temp, y_train, y_temp = train_test_split(X, y, test_size=0.3, random_state=42)
X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=42)

display(f'Shape of X_train: {X_train.shape}')
display(f'Shape of X_val: {X_val.shape}')
display(f'Shape of X_test: {X_test.shape}')
display(f'Shape of y_train: {y_train.shape}')
display(f'Shape of y_val: {y_val.shape}')
display(f'Shape of y_test: {y_test.shape}')
```

```
'Shape of X_train: (11900, 8)'
'Shape of X_val: (2550, 8)'
'Shape of X_test: (2550, 8)'
'Shape of y_train: (11900,)'
'Shape of y_val: (2550,)'
'Shape of y_test: (2550,)'
```

✓ Choose a model

Subtask:

Select an appropriate model architecture based on your task (e.g., classification, regression).

Reasoning: Based on the target variable 'median_house_value', which is a continuous numerical value, this is a regression task. An appropriate model for regression from scikit-learn is the Linear Regression model, which is a simple yet effective model for predicting continuous outcomes.

```
from sklearn.linear_model import LinearRegression

# Initialize the Linear Regression model
model = LinearRegression()

# Briefly explain the model choice
print("This is a regression task because the target variable 'median_house_value' is a continuous numerical value.")
print("Choosing Linear Regression: Linear Regression is a suitable choice for this task as it models the linear relationship between the features and the target variable.")

This is a regression task because the target variable 'median_house_value' is a continuous numerical value.
Choosing Linear Regression: Linear Regression is a suitable choice for this task as it models the linear relationship between the features and the target variable.
```

Define hyperparameters

Subtask:

Identify the hyperparameters of your chosen model that you want to tune.

Choose a tuning strategy

Subtask:

Select a method for tuning hyperparameters (e.g., Grid Search, Random Search, Bayesian Optimization).

Implement tuning

Subtask:

Implement the chosen tuning strategy, using a library like scikit-learn, Keras Tuner, or Optuna.

✓ Evaluate model

Subtask:

Evaluate the performance of the tuned model using appropriate metrics on the validation set.

Reasoning: Evaluate the performance of the model using MSE and R-squared on the validation set.

```
from sklearn.metrics import mean_squared_error, r2_score

# Fit the model on the training data
model.fit(X_train, y_train)

# Predict on the validation set
y_pred = model.predict(X_val)

# Calculate MSE and R-squared
mse = mean_squared_error(y_val, y_pred)
r2 = r2_score(y_val, y_pred)

# Print the metrics
print(f'Mean Squared Error (MSE) on validation set: {mse}')
print(f'R-squared score on validation set: {r2}')
```

Mean Squared Error (MSE) on validation set: 4907162398.177931
R-squared score on validation set: 0.6415252340365891

Reasoning: The model was not fitted before attempting to make predictions. Fit the Linear Regression model on the training data before evaluating its performance on the validation set.

```
from sklearn.metrics import mean_squared_error, r2_score

# Fit the model on the training data
model.fit(X_train, y_train)

# Predict on the validation set
y_pred = model.predict(X_val)

# Calculate MSE and R-squared
mse = mean_squared_error(y_val, y_pred)
r2 = r2_score(y_val, y_pred)

# Print the metrics
print(f'Mean Squared Error (MSE) on validation set: {mse}')
print(f'R-squared score on validation set: {r2}')
```

Mean Squared Error (MSE) on validation set: 4907162398.177931
R-squared score on validation set: 0.6415252340365891

✓ Retrain on full data

Subtask:

Train the best performing model with the optimal hyperparameters on the combined training and validation data.

Reasoning: Concatenate the training and validation data, fit the model on the combined data, and then finish the task.

```
import pandas as pd

# Concatenate training and validation data
X_train_val = pd.concat([X_train, X_val])
y_train_val = pd.concat([y_train, y_val])

# Fit the model on the combined data
model.fit(X_train_val, y_train_val)
```

▼ LinearRegression ⓘ ?
LinearRegression()

▼ Final evaluation

Subtask:

Evaluate the final model on the unseen test set to get an unbiased estimate of its performance.

Reasoning: Make predictions on the test set and calculate the MSE and R-squared scores to evaluate the final model performance.

Toggle Gemini

```
from sklearn.metrics import mean_squared_error, r2_score

# Make predictions on the test set
y_pred_test = model.predict(X_test)

# Calculate MSE and R-squared on the test set
mse_test = mean_squared_error(y_test, y_pred_test)
r2_test = r2_score(y_test, y_pred_test)

# Print the metrics
print(f'Mean Squared Error (MSE) on test set: {mse_test}')
print(f'R-squared score on test set: {r2_test}')
```

```
Mean Squared Error (MSE) on test set: 4728547900.164165
R-squared score on test set: 0.6605831628549996
```

Summary:

Data Analysis Key Findings

- The initial attempts to load the dataset from common paths failed due to `FileNotFoundError`. The dataset was successfully loaded from /content/sample_data/california_housing_train.csv.
- The dataset contains no missing values.
- The data was split into training (70% or 11900 samples), validation (15% or 2550 samples), and test (15% or 2550 samples) sets.
- The task is a regression problem as the target variable 'median_house_value' is continuous. A Linear Regression model was chosen as a baseline.
- The standard Linear Regression model in scikit-learn does not have significant hyperparameters for tuning.
- The Linear Regression model, evaluated on the validation set after training on the training set, achieved a Mean Squared Error (MSE) of approximately 4.91×10^9 and an R-squared score of approximately 0.64.
- The model was retrained on the combined training and validation datasets.
- The final evaluation on the unseen test set resulted in a Mean Squared Error (MSE) of approximately 4,728,547,900 and an R-squared score of approximately 0.661.

Insights or Next Steps

- Given the high MSE, explore feature engineering or transforming the target variable to improve model performance.
- Consider using more complex regression models (e.g., Ridge, Lasso, or tree-based models) that have tunable hyperparameters to potentially capture non-linear relationships and improve the R-squared score.

Toggle Gemini