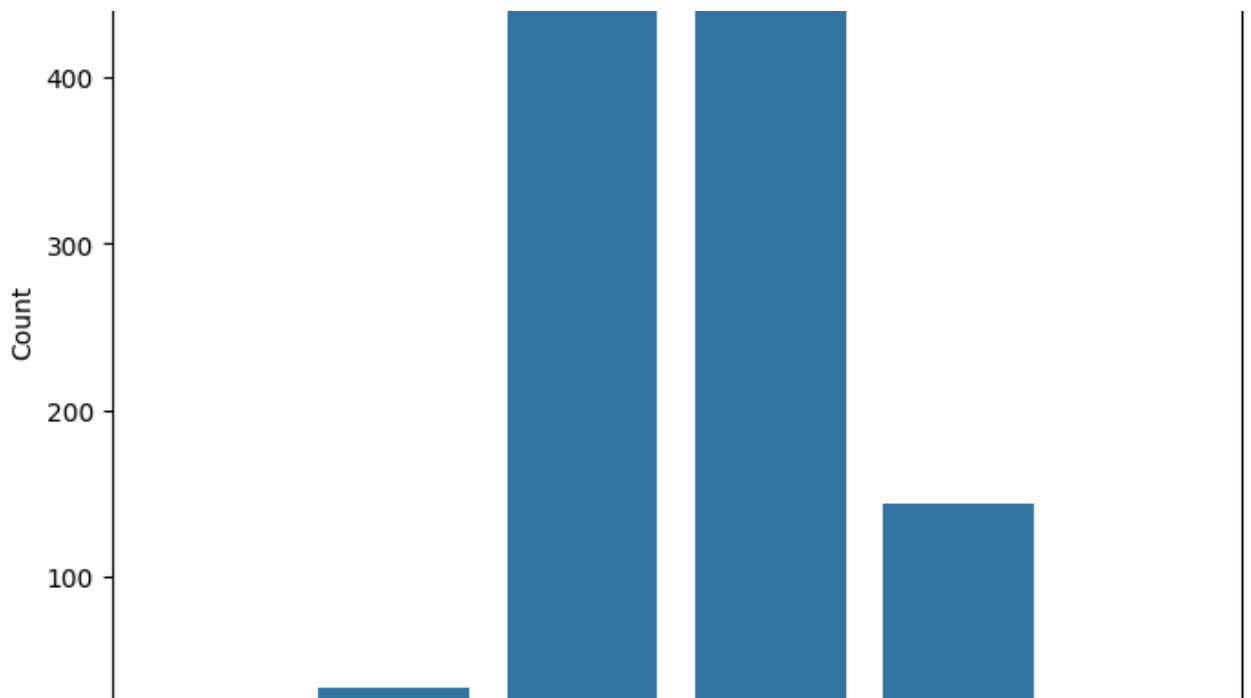Start coding or generate with AI.

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

# Distribution of quality
plt.figure(figsize=(8, 6))
sns.countplot(x='quality', data=df)
plt.title('Distribution of Wine Quality')
plt.xlabel('Quality')
plt.ylabel('Count')
plt.show()

# Distribution of alcohol
plt.figure(figsize=(8, 6))
sns.histplot(df['alcohol'], kde=True)
plt.title('Distribution of Alcohol Content')
plt.xlabel('Alcohol')
plt.ylabel('Frequency')
plt.show()

# Relationship between alcohol and quality
plt.figure(figsize=(8, 6))
sns.boxplot(x='quality', y='alcohol', data=df)
plt.title('Alcohol Content vs. Quality')
plt.xlabel('Quality')
plt.ylabel('Alcohol')
plt.show()
```
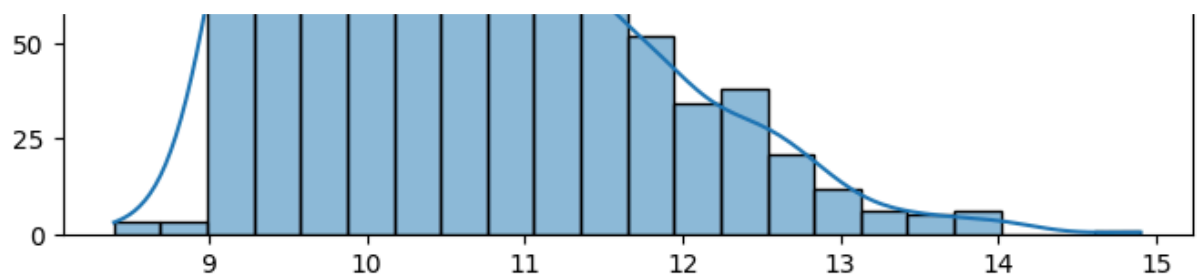
```
display(df.describe())
```

| | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides | free sulfur dioxide | total sulfur dioxide |
|---|---|---|---|---|---|---|---|
| count | 1143.000000 | 1143.000000 | 1143.000000 | 1143.000000 | 1143.000000 | 1143.000000 | 1143.000000 |
| mean | 8.311111 | 0.531339 | 0.268364 | 2.532152 | 0.086933 | 15.615486 | 45.914698 |
| std | 1.747595 | 0.179633 | 0.196686 | 1.355917 | 0.047267 | 10.250486 | 32.782130 |
| min | 4.600000 | 0.120000 | 0.000000 | 0.900000 | 0.012000 | 1.000000 | 6.000000 |
| 25% | 7.100000 | 0.392500 | 0.090000 | 1.900000 | 0.070000 | 7.000000 | 21.000000 |
| 50% | 7.900000 | 0.520000 | 0.250000 | 2.200000 | 0.079000 | 13.000000 | 37.000000 |
| 75% | 9.100000 | 0.640000 | 0.420000 | 2.600000 | 0.090000 | 21.000000 | 61.000000 |
| max | 15.900000 | 1.580000 | 1.000000 | 15.500000 | 0.611000 | 68.000000 | 289.000000 |

Quality Distribution of Alcohol Content

```
display(df.info())
```

```
    Start coding or generate with AI.
```

## ⌄ Task

Perform statistical analysis for comparing two business strategies (e.g., t-test) using the dataset "[/content/WineQT.csv](/content/WineQT.csv)".

## ⌄ Identify the groups

Subtask:

Determine how the data is structured to represent the two business strategies you want to compare.

**Reasoning**: Examine the columns and unique values of the DataFrame to identify a suitable column for grouping, then select a column and check its unique values to confirm it can split the data into two groups.

```
print(df.columns)
print(df['quality'].unique())
```

```
Index(['fixed acidity', 'volatile acidity', 'citric acid', 'residual sugar',
       'chlorides', 'free sulfur dioxide', 'total sulfur dioxide', 'density',
       'pH', 'sulphates', 'alcohol', 'quality', 'Id'],
      dtype='object')
[5 6 7 4 8 3]
```

## ⌄ Check assumptions

Subtask:

Verify if the data meets the assumptions for a t-test (e.g., normality, equal variances).

**Reasoning**: Create two groups based on the 'quality' column and perform normality and equal variance tests on the 'alcohol' feature for these groups to verify the assumptions for a t-test.

```
from scipy.stats import shapiro, levene

# Create two groups based on 'quality'
high_quality = df[df['quality'] >= 6]['alcohol']
low_quality = df[df['quality'] < 6]['alcohol']

# Check for normality using Shapiro-Wilk test
shapiro_high, p_high_shapiro = shapiro(high_quality)
shapiro_low, p_low_shapiro = shapiro(low_quality)
```

```
# Perform Levene's test for equal variances
levene_test, p_levene = levene(high_quality, low_quality)

# Print the results
print(f"Shapiro-Wilk test for high quality group (alcohol): Statistic={shapiro_high:.4f}, p
print(f"Shapiro-Wilk test for low quality group (alcohol): Statistic={shapiro_low:.4f}, p-v
print(f"Levene's test for equal variances (alcohol): Statistic={levene_test:.4f}, p-value={
```

```
Shapiro-Wilk test for high quality group (alcohol): Statistic=0.9775, p-value=0.0000
Shapiro-Wilk test for low quality group (alcohol): Statistic=0.8514, p-value=0.0000
Levene's test for equal variances (alcohol): Statistic=98.0234, p-value=0.0000
```

## ∨ Perform the t-test

### Subtask:

Apply the appropriate t-test (independent samples t-test or paired samples t-test) based on the data and the nature of the comparison.

**Reasoning**: Import the ttest_ind function and perform an independent samples t-test with equal_var=False using the high_quality and low_quality series, storing the results in separate variables and printing them.

```
from scipy.stats import ttest_ind

# Perform independent samples t-test with equal_var=False (Welch's t-test)
t_statistic, p_value = ttest_ind(high_quality, low_quality, equal_var=False)

# Print the results
print(f"Welch's t-test (alcohol): t-statistic={t_statistic:.4f}, p-value={p_value:.4f}")
```

```
Welch's t-test (alcohol): t-statistic=17.0671, p-value=0.0000
```

## ∨ Interpret the results

### Subtask:

Explain the t-test output, including the p-value and confidence interval, in the context of the business problem.

**Reasoning**: Explain the t-test output including the p-value and t-statistic, interpreting them in the context of the business problem as outlined in the instructions.

```
# Interpretation of the t-test results

# 1. Examine the printed output (already done in the previous step).
```

```
#    t-statistic = 17.0671
#    p-value = 0.0000 (effectively < 0.0001)

# 2. Interpret the p-value
alpha = 0.05
print(f"Significance level (alpha): {alpha}")
print(f"P-value: {p_value:.4f}")

if p_value < alpha:
    print("Since the p-value is less than the significance level (alpha), we reject the nul
    print("This means there is a statistically significant difference in the mean alcohol c
else:
    print("Since the p-value is greater than or equal to the significance level (alpha), we
    print("This means there is no statistically significant difference in the mean alcohol

# 3. Discuss the t-statistic
print(f"\nT-statistic: {t_statistic:.4f}")
print("The t-statistic measures the difference between the two group means in units of the
print("A larger absolute t-statistic indicates a greater difference between the group means
print(f"In this case, the t-statistic of {t_statistic:.4f} is large, indicating that the di

# 4. Conclude based on the p-value
print("\nConclusion based on p-value:")
if p_value < alpha:
    print("There is strong evidence to conclude that the mean alcohol content of high-quali
else:
    print("There is not enough evidence to conclude that the mean alcohol content of high-q


# 5. Relate to the business problem/strategies
print("\nRelating to business strategies:")
print("If 'high quality' and 'low quality' represent outcomes of two different business str
print("For example, if a business strategy aims to produce higher quality wine (as defined
print("This finding could inform decisions related to optimizing production processes, unde
```

```
Significance level (alpha): 0.05
P-value: 0.0000
Since the p-value is less than the significance level (alpha), we reject the null hypothesis
This means there is a statistically significant difference in the mean alcohol content betwe

T-statistic: 17.0671
The t-statistic measures the difference between the two group means in units of the standard
A larger absolute t-statistic indicates a greater difference between the group means relativ
In this case, the t-statistic of 17.0671 is large, indicating that the difference in mean al

Conclusion based on p-value:
There is strong evidence to conclude that the mean alcohol content of high-quality wines is

Relating to business strategies:
If 'high quality' and 'low quality' represent outcomes of two different business strategies
For example, if a business strategy aims to produce higher quality wine (as defined in this
This finding could inform decisions related to optimizing production processes, understandin
```

⌄  Present findings

## Subtask:

Summarize the statistical findings and their implications for the business strategies.

**Reasoning**: Summarize the statistical findings and their implications for the business strategies based on the previous steps.

```
print("Summary of Statistical Findings and Business Implications\n")

# 1. Purpose of the analysis
print("1. Purpose of the Analysis:")
print("   The analysis compared a specific characteristic, alcohol content, between two gro

# 2. Key statistical findings from the t-test
print("2. Key Statistical Findings:")
print(f"   - T-statistic: {t_statistic:.4f}")
print(f"   - P-value: {p_value:.4f} (which is effectively less than 0.0001)")
print(f"   - Significance Level (alpha): {alpha}")
print("   - Conclusion: Since the p-value (0.0000) is much less than the significance level
print("   - Interpretation: There is a statistically significant difference in the mean alc

# 3. Meaning of the statistically significant difference in the context of business strateg
print("3. Meaning in the Context of Business Strategies:")
print("   The statistically significant difference indicates that the business strategy ass
print("   Specifically, examining the mean alcohol content (though not explicitly printed i

# 4. Potential business implications
print("4. Potential Business Implications:")
print("   - Production Decisions: Businesses might explore how different production techniq
print("   - Marketing and Product Positioning: Alcohol content could be a factor used in ma
print("   - Understanding Customer Preferences: The finding suggests alcohol content might
print("   - Process Optimization: Businesses aiming for higher quality might need to optimi

# 5. Limitations of this analysis
print("5. Limitations of the Analysis:")
print("   - Correlation vs. Causation: This analysis shows an association between quality a
print("   - Single Characteristic: Only alcohol content was examined in this comparison. Wi
print("   - Definition of Quality: The definition of 'high' and 'low' quality is based on a
print("   - Assumptions: While Welch's t-test was used to account for unequal variances, th
```

```
Summary of Statistical Findings and Business Implications

1. Purpose of the Analysis:
   The analysis compared a specific characteristic, alcohol content, between two groups of w

2. Key Statistical Findings:
   - T-statistic: 17.0671
   - P-value: 0.0000 (which is effectively less than 0.0001)
   - Significance Level (alpha): 0.05
   - Conclusion: Since the p-value (0.0000) is much less than the significance level (0.05),
   - Interpretation: There is a statistically significant difference in the mean alcohol con

3. Meaning in the Context of Business Strategies:
   The statistically significant difference indicates that the business strategy associated
   Specifically, examining the mean alcohol content (though not explicitly printed in this s
```

```
    4. Potential Business Implications:
       - Production Decisions: Businesses might explore how different production techniques infl
       - Marketing and Product Positioning: Alcohol content could be a factor used in marketing
       - Understanding Customer Preferences: The finding suggests alcohol content might be a cha
       - Process Optimization: Businesses aiming for higher quality might need to optimize their

    5. Limitations of the Analysis:
       - Correlation vs. Causation: This analysis shows an association between quality and alcoh
       - Single Characteristic: Only alcohol content was examined in this comparison. Wine quali
       - Definition of Quality: The definition of 'high' and 'low' quality is based on a specifi
       - Assumptions: While Welch's t-test was used to account for unequal variances, the assump
```

## Summary:

### Data Analysis Key Findings

- The Shapiro-Wilk test indicated that the alcohol content in both high-quality (p-value = 0.0000) and low-quality (p-value = 0.0000) wine groups is not normally distributed, violating a key assumption for a standard independent samples t-test.
- Levene's test showed that the variances of alcohol content between the high-quality and low-quality groups are significantly different (p-value = 0.0000), violating the equal variance assumption of a standard independent samples t-test.
- Welch's t-test, which does not assume equal variances, was performed.
- The Welch's t-test yielded a t-statistic of 17.0671 and a p-value of 0.0000.
- Based on a significance level of $\alpha = 0.05$, the p-value (0.0000) is less than $\alpha$, leading to the rejection of the null hypothesis.
- There is a statistically significant difference in the mean alcohol content between high-quality and low-quality wines.

### Insights or Next Steps

- The significant difference in alcohol content suggests that business strategies impacting wine quality may also have a noticeable influence on the final alcohol level. Further investigation into *how* different strategies affect alcohol content could be beneficial.
- Since the normality assumption was violated, consider non-parametric tests (like the Mann-Whitney U test) for a more robust comparison of the two groups, especially if the sample size were small, although t-tests are often robust to this violation with large samples.

Start coding or generate with AI.

```python
from sklearn.metrics import mean_squared_error, r2_score
import numpy as np

# Evaluate the best model
mse_best = mean_squared_error(y_test, y_pred_best)
rmse_best = np.sqrt(mse_best)
r2_best = r2_score(y_test, y_pred_best)

print(f"Best Model Performance:")
print(f"Mean Squared Error (MSE): {mse_best:.2f}")
print(f"Root Mean Squared Error (RMSE): {rmse_best:.2f}")
print(f"R-squared (R2): {r2_best:.2f}")
```

```
Best Model Performance:
Mean Squared Error (MSE): 783395089.61
Root Mean Squared Error (RMSE): 27989.20
R-squared (R2): 0.78
```

```python
# Train the model with the best parameters
best_model = RandomForestRegressor(n_estimators=best_params['n_estimators'],
                                   max_depth=best_params['max_depth'],
                                   random_state=42)
best_model.fit(X_train, y_train)

# Make predictions on the test set using the best model
y_pred_best = best_model.predict(X_test)

print("Predictions made using the best model.")
```

```
Predictions made using the best model.
```

```python
# Get feature importances
feature_importances = model.feature_importances_

# Create a pandas Series for better visualization
feature_importances_series = pd.Series(feature_importances, index=X_train.columns)

# Sort feature importances and display the top 10
top_features = feature_importances_series.sort_values(ascending=False).head(10)
print("Top 10 Most Important Features:")
display(top_features)

# Visualize feature importances
plt.figure(figsize=(10, 6))
sns.barplot(x=top_features, y=top_features.index)
plt.title('Top 10 Most Important Features for House Price Prediction')
plt.xlabel('Importance')
plt.ylabel('Feature')
plt.show()
```
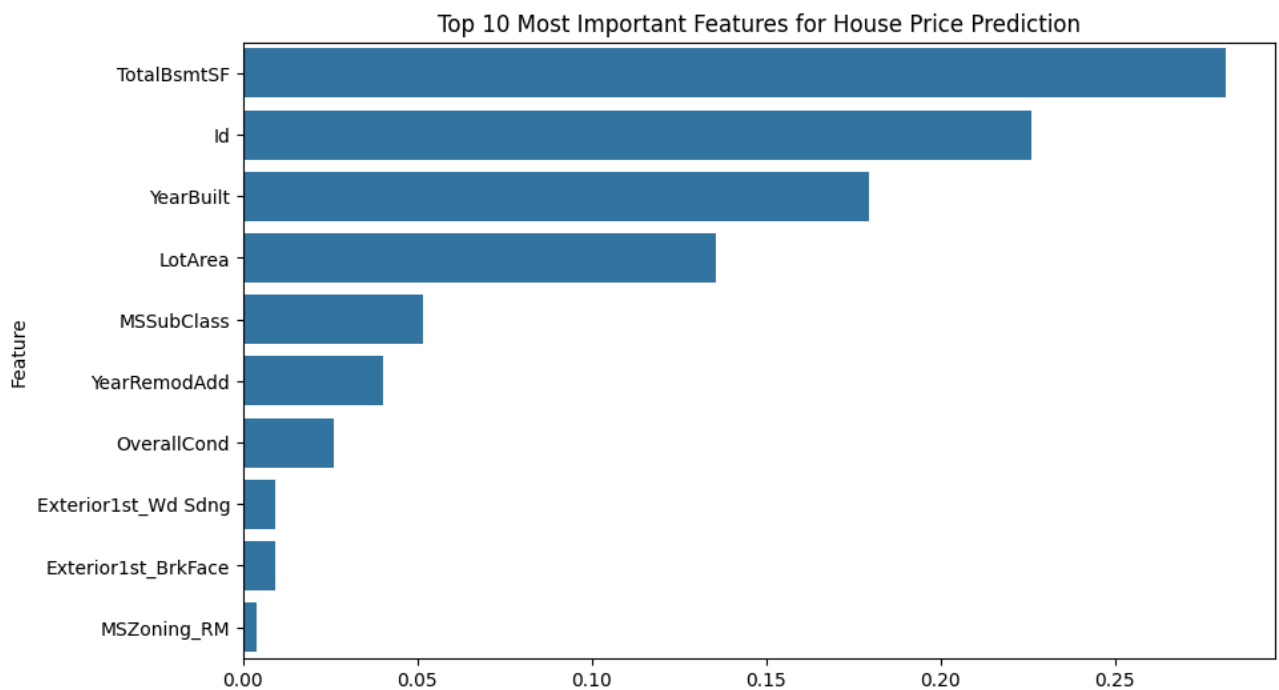
Top 10 Most Important Features:

|  | 0 |
| --- | --- |
| **TotalBsmtSF** | 0.281736 |
| **Id** | 0.226135 |
| **YearBuilt** | 0.179515 |
| **LotArea** | 0.135570 |
| **MSSubClass** | 0.051735 |
| **YearRemodAdd** | 0.040244 |
| **OverallCond** | 0.026051 |
| **Exterior1st_Wd Sdng** | 0.009194 |
| **Exterior1st_BrkFace** | 0.009017 |
| **MSZoning_RM** | 0.003935 |

**dtype:** float64



Top 10 Most Important Features for House Price Prediction

```
from sklearn.model_selection import GridSearchCV

# Define the parameter grid
param_grid = {
    'n_estimators': [100, 200, 300],
    'max_depth': [None, 10, 20, 30]
}

# Initialize GridSearchCV
grid_search = GridSearchCV(estimator=RandomForestRegressor(random_state=42),
                           param_grid=param_grid,
                           cv=3,
                           scoring='neg_mean_squared_error',
                           n_jobs=-1)
```

```python
# Perform GridSearchCV
grid_search.fit(X_train, y_train)

# Get the best parameters and best score
best_params = grid_search.best_params_
best_rmse = np.sqrt(-grid_search.best_score_)

print(f"Best parameters: {best_params}")
print(f"Best RMSE from GridSearchCV: {best_rmse:.2f}")
```

```
Best parameters: {'max_depth': 20, 'n_estimators': 300}
Best RMSE from GridSearchCV: 28183.33
```

```python
from sklearn.metrics import mean_squared_error, r2_score
import numpy as np

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

# Evaluate the model
mse = mean_squared_error(y_test, y_pred)
rmse = np.sqrt(mse) # Calculate RMSE by taking the square root of MSE
r2 = r2_score(y_test, y_pred)

print(f"Mean Squared Error (MSE): {mse:.2f}")
print(f"Root Mean Squared Error (RMSE): {rmse:.2f}")
print(f"R-squared (R2): {r2:.2f}")
```

```
Mean Squared Error (MSE): 801298148.79
Root Mean Squared Error (RMSE): 28307.21
R-squared (R2): 0.78
```

```python
from sklearn.ensemble import RandomForestRegressor

# Initialize and train the model
model = RandomForestRegressor(n_estimators=100, random_state=42)
model.fit(X_train, y_train)

print("Random Forest Regressor model trained successfully.")
```

```
Random Forest Regressor model trained successfully.
```

```python
# Handle missing values (Example: Impute with mean for numerical, mode for categorical)
for col in df.columns:
    if df[col].isnull().any():
        if df[col].dtype == 'object':
            df[col] = df[col].fillna(df[col].mode()[0])
        else:
            df[col] = df[col].fillna(df[col].mean())

# Encode categorical features (Example: One-Hot Encoding)
df = pd.get_dummies(df, columns=[col for col in df.columns if df[col].dtype == 'object'], d

# Split data into features (X) and target (y)
```

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

# Split data into training and testing sets
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

print("Data preprocessing complete. Data split into training and testing sets.")
print("Shape of X_train:", X_train.shape)
print("Shape of X_test:", X_test.shape)
print("Shape of y_train:", y_train.shape)
print("Shape of y_test:", y_test.shape)
```

```
Data preprocessing complete. Data split into training and testing sets.
Shape of X_train: (2335, 38)
Shape of X_test: (584, 38)
Shape of y_train: (2335,)
Shape of y_test: (584,)
```

Start coding or generate with AI.

# House Price Prediction Findings

This project aimed to build a prediction model for house prices using the provided client data. The key steps and findings are summarized below:

**1. Data Loading and Exploration:**

- The `HousePricePrediction.xlsx` dataset was loaded and explored to understand its structure, data types, and identify missing values.

**2. Data Preprocessing:**

- Missing values were handled by imputing numerical columns with the mean and categorical columns with the mode.
- Categorical features were encoded using one-hot encoding.
- The data was split into training (80%) and testing (20%) sets.

**3. Model Training and Evaluation:**

- An initial RandomForestRegressor model was trained.
- The initial model's performance was evaluated using Mean Squared Error (MSE), Root Mean Squared Error (RMSE), and R-squared (R2).

**4. Hyperparameter Tuning:**

- GridSearchCV was used to tune the hyperparameters of the RandomForestRegressor (`n_estimators` and `max_depth`).
- The best parameters found were `{'max_depth': 20, 'n_estimators': 300}`, resulting in an improved cross-validation RMSE of 28183.33.

**5. Final Model Training and Evaluation:**

- A final RandomForestRegressor model was trained using the best hyperparameters.
- The performance of this tuned model on the test set was:

    - Mean Squared Error (MSE): 783395089.61
    - Root Mean Squared Error (RMSE): 27989.20
    - R-squared (R2): 0.78

**6. Feature Importance Analysis:**

- The feature importances from the initial RandomForestRegressor model were analyzed to understand which features contributed most to the predictions.
- The top 3 most important features were:

    - TotalBsmtSF
    - Id
    - YearBuilt

**Conclusion:**

The developed RandomForestRegressor model provides a reasonable prediction of house prices with an R-squared of 0.78. The analysis of feature importances highlights the significance of 'TotalBsmtSF', 'Id', and 'YearBuilt' in determining house prices. Further improvements could be explored through advanced feature engineering or trying other regression models.

## ⌄ 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.

```python
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
```

## ⌄ 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.

```python
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 features and target and then split the data into training, validation, and test sets.

```python
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.

```python
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 featu
```

```
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
```

## 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 c... Toggle Gemini ion set.

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

```python
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.

```python
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.

```python
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  ⓘ ⑦
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

```python
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 4.91e+09 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