

This process will include data loading, cleaning, visualization, and model creation.

1. Define the Problem and Objective

- **Objective:** Predict the quality of wine based on various chemical properties such as acidity, alcohol content, sugar, and others.
 - **Goal:** Perform EDA to understand the relationships between features and the target variable (`quality`), which represents the quality of the wine on a scale from 0 to 10.
-

2. Load the Dataset

```
import pandas as pd

# Load Wine Quality dataset
df = pd.read_csv('wine_quality.csv')

# Display the first few rows of the dataset
df.head()
```

3. Data Exploration (EDA)

Check the Shape and Data Types

```
# Check dataset shape
print(f"Dataset Shape: {df.shape}")

# Data types and missing values
print(df.info())

# Summary statistics for numerical columns
print(df.describe())
```

Visualize Distributions

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

# Distribution of wine quality
sns.countplot(x='quality', data=df)
plt.title("Wine Quality Distribution")
plt.show()

# Distribution of Alcohol content
sns.histplot(df['alcohol'], kde=True)
plt.title("Alcohol Content Distribution")
plt.show()
```

4. Data Cleaning

Handle Missing Values

```
# Check for missing values
print(df.isnull().sum())

# Handle missing values (if any)
# In this case, assume there are no missing values in the dataset.
```

Handle Duplicates

```
# Remove duplicate rows if any
df.drop_duplicates(inplace=True)
```

5. Feature Engineering (Optional)

Create New Features (if needed)

```
# Example: Add a new feature "high_quality" indicating if quality is above a
certain threshold (e.g., 7)
df['high_quality'] = df['quality'].apply(lambda x: 1 if x >= 7 else 0)

# Display unique high_quality values
print(df['high_quality'].unique())
```

6. Data Preprocessing

Scale Numerical Features

```
from sklearn.preprocessing import StandardScaler

# Initialize the scaler
scaler = StandardScaler()

# Scale the numerical features (e.g., alcohol, volatile acidity, etc.)
scaled_features = df.drop(columns=['quality', 'high_quality']) # Exclude
target and new feature
df[scaled_features.columns] = scaler.fit_transform(scaled_features)
```

Split the Data into Features (X) and Target (y)

```
# Define features (X) and target (y)
X = df.drop(columns=['quality', 'high_quality'])
y = df['quality'] # or you could use 'high_quality' as the target
```

7. Split the Data into Training and Testing Sets

```
from sklearn.model_selection import train_test_split

# Split data into training and testing sets (80% train, 20% test)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)
```

8. Model Selection

Choose the Model

```
from sklearn.ensemble import RandomForestClassifier

# Initialize the Random Forest model
model = RandomForestClassifier(n_estimators=100, random_state=42)
```

9. Model Training

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

10. Model Evaluation

Predictions and Evaluation Metrics

```
from sklearn.metrics import accuracy_score, classification_report,
confusion_matrix

# Make predictions
y_pred = model.predict(X_test)

# Evaluate the model
print(f'Accuracy: {accuracy_score(y_test, y_pred)}')
print("\nClassification Report:")
print(classification_report(y_test, y_pred))
print("\nConfusion Matrix:")
print(confusion_matrix(y_test, y_pred))
```

Confusion Matrix Visualization

```
# Visualize confusion matrix
sns.heatmap(confusion_matrix(y_test, y_pred), annot=True, fmt='d',
cmap='Blues', cbar=False)
plt.title("Confusion Matrix")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.show()
```

11. Model Interpretation and Insights

Feature Importance

```
# Get feature importances from the trained model
importances = model.feature_importances_

# Create a DataFrame to display features and their importance
features = X.columns
feature_importance_df = pd.DataFrame({'Feature': features, 'Importance':
importances})

# Sort by importance
feature_importance_df = feature_importance_df.sort_values(by='Importance',
ascending=False)

# Plot feature importances
sns.barplot(x='Importance', y='Feature', data=feature_importance_df)
plt.title("Feature Importance")
plt.show()

print(feature_importance_df)
```

12. Fine-Tuning and Optimization (Optional)

Hyperparameter Tuning

```
from sklearn.model_selection import GridSearchCV

# Set up grid search parameters for Random Forest
param_grid = {
    'n_estimators': [100, 200, 300],
    'max_depth': [None, 10, 20],
    'min_samples_split': [2, 5, 10]
}

# Initialize GridSearchCV
grid_search = GridSearchCV(estimator=model, param_grid=param_grid, cv=3,
n_jobs=-1, verbose=2)

# Fit grid search
grid_search.fit(X_train, y_train)

# Get best parameters
print("Best Parameters:", grid_search.best_params_)
```

13. Final Reporting

Summarize Insights

- **EDA Insights:** Discuss the distribution of features like alcohol content, acidity, and others. Highlight any trends in wine quality and any correlations between features and target variable.
 - **Model Evaluation:** Summarize the accuracy, classification report, and confusion matrix.
 - **Important Features:** Identify the most important features influencing wine quality predictions.
 - **Next Steps:** Explore other models, optimize hyperparameters, or improve feature engineering.
-

Additional Notes:

- Start with a **simple model** like Random Forest and evaluate its performance.
 - Perform **cross-validation** to assess model robustness.
 - Explore feature interactions and correlations to refine features.
 - **Fine-tune hyperparameters** using GridSearchCV to get better results.
-