

Assignment- 24 March 2024

Q1. What are the key features of the wine quality data set? Discuss the importance of each feature in predicting the quality of wine.

Aa Key Features of the Wine Quality Dataset

The wine quality dataset is a widely used dataset for regression and classification tasks, containing physicochemical properties (numerical variables) of red or white wine samples and their quality ratings (target variable). Each feature contributes uniquely to predicting wine quality, which is rated on a scale, typically from 0 to 10.

Features and Their Importance

1. Fixed Acidity

- **Description:** Tartaric acid present in wine, contributing to its acidity.
- **Importance:**
 - Affects the perception of freshness and taste.
 - Optimal acidity levels contribute to a balanced taste.

2. Volatile Acidity

- **Description:** Acetic acid present in wine, leading to a vinegar-like taste if in high concentration.
- **Importance:**
 - High levels negatively affect quality ratings.
 - Strongly correlated with customer preferences.

3. Citric Acid

- **Description:** Adds freshness to the wine.
- **Importance:**
 - Low levels make wine flat.
 - Contributes to acidity balance and stability.

4. Residual Sugar

- **Description:** Sugar remaining after fermentation.

- **Importance:**
 - Influences sweetness.
 - Higher residual sugar is preferred in dessert wines.

5. Chlorides

- **Description:** Represents salt content in wine.
- **Importance:**
 - High chlorides can impart a salty taste, reducing quality.
 - Reflects possible contamination.

6. Free Sulfur Dioxide

- **Description:** Sulfur dioxide in wine that is not chemically bound and is active as an antimicrobial agent.
- **Importance:**
 - Helps preserve wine.
 - Excessive levels may affect taste.

7. Total Sulfur Dioxide

- **Description:** Total sulfur dioxide in wine, both free and bound.
- **Importance:**
 - High levels can mask flavors and reduce quality.
 - Regulated in the wine industry.

8. Density

- **Description:** Mass per unit volume, related to sugar and alcohol content.
- **Importance:**
 - Higher density often correlates with sweetness.
 - Used to infer residual sugar levels.

9. pH

- **Description:** Measures acidity/basicity.

- **Importance:**
 - Impacts microbial stability and sensory perception.
 - Wines with balanced pH levels are often rated higher.

10. Sulphates

- **Description:** Contribute to wine's preservative properties and flavor.
- **Importance:**
 - Enhances antioxidant properties.
 - Impacts the perception of quality positively when in balance.

11. Alcohol

- **Description:** Percentage of ethanol content.
- **Importance:**
 - Strong positive correlation with wine quality.
 - Impacts aroma, flavor, and overall sensory experience.

12. Quality (Target Variable)

- **Description:** Sensory rating of the wine.
- **Importance:**
 - Determined by a combination of all other features.
 - Typically rated by human tasters.

In Short- Key Features:

The wine quality dataset typically includes features such as:

1. **Fixed Acidity:** Affects the freshness and tartness of wine.
2. **Volatile Acidity:** Impacts the vinegar-like taste; high levels can reduce quality.
3. **Citric Acid:** Adds freshness and flavor.
4. **Residual Sugar:** Determines sweetness.
5. **Chlorides:** Affects salinity.

6. **Free Sulfur Dioxide:** Prevents microbial growth; excessive levels may affect taste.
7. **Total Sulfur Dioxide:** Affects preservation; must be balanced.
8. **Density:** Correlates with alcohol content.
9. **pH:** Influences acidity perception.
10. **Sulphates:** Acts as a preservative and affects flavor.
11. **Alcohol:** Significant for overall perception and quality.

Steps to Analyze the Importance of Each Feature

1. Load the Dataset

Import the dataset to explore and understand its structure.

2. Exploratory Data Analysis (EDA)

Analyze the relationship between features and the target variable (quality). Use statistical and visual tools to identify trends and patterns.

3. Correlation Analysis

Measure how strongly each feature correlates with quality using a correlation matrix.

4. Feature Importance with Machine Learning

Use machine learning models (e.g., Random Forest) to rank the importance of features.

Code Implementation

Here's how to analyze the key features of the wine quality dataset:

```
python
```

```
CopyEdit
```

```
import pandas as pd
```

```
import seaborn as sns
```

```
import matplotlib.pyplot as plt
```

```
from sklearn.ensemble import RandomForestRegressor
```

```
# Step 1: Load the Dataset
```

```
wine_data_path = '/mnt/data/winequality-red.csv' # Replace with your dataset path
wine_data = pd.read_csv(wine_data_path)
```

```
# Step 2: Correlation Analysis
```

```
correlation_matrix = wine_data.corr()
plt.figure(figsize=(10, 8))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm')
plt.title('Correlation Matrix')
plt.show()
```

```
# Step 3: Feature Importance with Random Forest
```

```
X = wine_data.drop(columns=['quality'])
y = wine_data['quality']
```

```
# Train Random Forest Regressor
```

```
rf_model = RandomForestRegressor(random_state=42)
rf_model.fit(X, y)
```

```
# Feature Importance
```

```
feature_importances = rf_model.feature_importances_
importance_df = pd.DataFrame({'Feature': X.columns, 'Importance':
feature_importances})
importance_df = importance_df.sort_values(by='Importance', ascending=False)
```

```
# Plot Feature Importance
```

```
plt.figure(figsize=(10, 6))
```

```
sns.barplot(x='Importance', y='Feature', data=importance_df)

plt.title('Feature Importance')

plt.show()
```

```
print("\nFeature Importance Ranking:\n", importance_df)
```

Results and Insights

1. Top Features:

- Alcohol: Strongly correlated with quality; higher levels generally indicate better quality.
- Volatile Acidity: High levels negatively impact quality.
- Sulphates: Contributes positively to quality when in balance.

2. Feature Relationships:

- Residual Sugar and Density: Correlated, with density often serving as a proxy for sugar levels.
- Free Sulfur Dioxide and Total Sulfur Dioxide: Reflect preservation properties.

3. Actionable Insights:

- Winemakers can focus on controlling volatile acidity and balancing alcohol content to improve quality.
- Excessive use of sulphates and chlorides should be avoided.

Advantages of Feature Importance Analysis

- Helps identify which features significantly influence wine quality.
- Enables winemakers to adjust production processes.
- Reduces model complexity by focusing on important features.

Key Takeaways

- Features like Alcohol, Volatile Acidity, and Sulphates are critical for predicting wine quality.

- Feature importance analysis helps streamline data and improve model performance.
- Understanding the role of each feature can guide practical improvements in winemaking processes.

In Short- Importance:

Each feature contributes differently to predicting wine quality. For example:

- **Alcohol** and **Sulphates** are often highly correlated with higher quality ratings.
- **Volatile Acidity** usually negatively affects quality.

Q2. How did you handle missing data in the wine quality data set during the feature engineering process? Discuss the advantages and disadvantages of different imputation techniques.

Aa

Handling Missing Data in the Wine Quality Dataset

Dealing with missing data is a crucial step in data preprocessing. The presence of missing values can significantly impact the quality of insights and the performance of machine learning models. The wine quality dataset typically contains features such as acidity, alcohol, and residual sugar, which might have missing entries that need to be addressed before analysis or modeling.

Steps for Handling Missing Data

1. Understand the Extent of Missing Data

- **Initial Analysis:** Identify the columns with missing values and the proportion of missing data.
- **Visualization:** Use heatmaps or bar charts to visually explore missing data patterns.

2. Analyze Missing Data Mechanism

- **Missing Completely at Random (MCAR):** No relationship between missing data and other values.
- **Missing at Random (MAR):** Missingness depends on observed data.

- **Missing Not at Random (MNAR):** Missingness depends on unobserved values.

3. Choose an Imputation Strategy

Based on the nature and extent of missing data, apply one or more of the following techniques.

Common Imputation Techniques

1. Mean Imputation

- Replace missing values with the mean of the respective column.

Advantages:

- Simple and quick.
- Retains the mean of the original data.

Disadvantages:

- Reduces variability, which can lead to biased estimates.

Code Example:

```
from sklearn.impute import SimpleImputer
```

```
# Mean imputation
```

```
mean_imputer = SimpleImputer(strategy='mean')
```

```
wine_data['fixed_acidity'] =  
mean_imputer.fit_transform(wine_data[['fixed_acidity']])
```

2. Median Imputation

- Replace missing values with the median of the column, often used for skewed distributions.

Advantages:

- Robust to outliers.

Disadvantages:

- Can distort relationships for normally distributed data.

Code Example:

```
median_imputer = SimpleImputer(strategy='median')  
wine_data['chlorides'] = median_imputer.fit_transform(wine_data[['chlorides']])
```

3. Mode Imputation

- Replace missing values with the mode (most frequent value) of the column, suitable for categorical data.

Advantages:

- Simple and effective for categorical variables.

Disadvantages:

- Not suitable for continuous data.

Code Example:

```
mode_imputer = SimpleImputer(strategy='most_frequent')  
wine_data['quality'] = mode_imputer.fit_transform(wine_data[['quality']])
```

4. K-Nearest Neighbors (KNN) Imputation

- Use the mean/median of the k-nearest neighbors to impute missing values.

Advantages:

- Retains relationships between features.
- Suitable for small datasets with limited missing values.

Disadvantages:

- Computationally expensive for large datasets.
- Sensitive to the choice of k and scaling.

In Short Common Imputation Techniques:

1. Mean/Median Imputation:

- Replace missing values with the mean or median of the column.
- **Advantage:** Simple and quick.
- **Disadvantage:** Reduces variability in data.

2. K-Nearest Neighbors (KNN) Imputation:

- Fills missing values using the mean or median of the k-nearest rows.
- **Advantage:** Retains relationships between variables.
- **Disadvantage:** Computationally expensive for large datasets.

3. Regression Imputation:

- Predict missing values based on other features.
- **Advantage:** Captures relationships among features.
- **Disadvantage:** Can introduce bias.

Code Example:

```
from sklearn.impute import KNNImputer
```

```
knn_imputer = KNNImputer(n_neighbors=5)
```

```
wine_data = pd.DataFrame(knn_imputer.fit_transform(wine_data),  
columns=wine_data.columns)
```

5. Regression Imputation

- Predict missing values using a regression model trained on non-missing values.

Advantages:

- Leverages relationships between features.
- Can be very accurate.

Disadvantages:

- Computationally intensive.
- Can introduce bias if assumptions about the relationship are incorrect.

Code Example:

```
from sklearn.linear_model import LinearRegression
```

```
# Example: Imputing 'alcohol' based on other features

train_data = wine_data.dropna(subset=['alcohol'])

missing_data = wine_data[wine_data['alcohol'].isnull()]

model = LinearRegression()

model.fit(train_data.drop('alcohol', axis=1), train_data['alcohol'])

wine_data.loc[wine_data['alcohol'].isnull(), 'alcohol'] =
model.predict(missing_data.drop('alcohol', axis=1))
```

6. Forward/Backward Fill

- Fill missing values with the previous or next value in the column, useful for time-series data.

Advantages:

- Simple and intuitive for time-series datasets.

Disadvantages:

- Not suitable for non-sequential data.

Code Example:

```
wine_data['density'] = wine_data['density'].fillna(method='ffill')
```

Detailed Implementation

```
import pandas as pd
```

```
from sklearn.impute import SimpleImputer, KNNImputer
```

```
# Step 1: Load the Dataset
```

```
wine_data_path = '/mnt/data/winequality-red.csv'
```

```
wine_data = pd.read_csv(wine_data_path)
```

```
# Step 2: Identify Missing Data
```

```
missing_summary = wine_data.isnull().sum()

print("Missing Data Summary:\n", missing_summary)
```

Step 3: Apply Imputation Techniques

1. Mean Imputation for fixed acidity

```
mean_imputer = SimpleImputer(strategy='mean')

wine_data['fixed_acidity'] =
mean_imputer.fit_transform(wine_data[['fixed_acidity']])
```

2. Median Imputation for chlorides

```
median_imputer = SimpleImputer(strategy='median')

wine_data['chlorides'] = median_imputer.fit_transform(wine_data[['chlorides']])
```

3. KNN Imputation for remaining features

```
knn_imputer = KNNImputer(n_neighbors=5)

wine_data = pd.DataFrame(knn_imputer.fit_transform(wine_data),
columns=wine_data.columns)
```

Step 4: Verify Imputation

```
print("\nAfter Imputation Missing Data Summary:\n", wine_data.isnull().sum())
```

Advantages and Disadvantages of Imputation Techniques

Technique	Advantages	Disadvantages
Mean Imputation	Quick, retains mean of data	Reduces variability, biased estimates
Median Imputation	Handles outliers well	Can distort relationships for normal distributions

Mode Imputation	Effective for categorical data	Not suitable for continuous data
KNN Imputation	Preserves feature relationships	Computationally expensive, scaling sensitive
Regression Imputation	Leverages feature relationships, accurate	Computationally intensive, risk of bias
Forward/Backward Fill	Intuitive for sequential data	Assumes temporal relationships, not generalizable

Key Takeaways

1. **Selection of Method:** The choice of imputation method depends on the nature of the data and the extent of missingness.
2. **Trade-offs:** Simpler methods like mean/median imputation are computationally efficient but can reduce variability, while advanced methods like KNN or regression are computationally expensive but retain relationships.
3. **Validation:** Always validate the effectiveness of the imputation by analyzing the impact on feature distributions and model performance.

Example-2 Code:

```
import pandas as pd

from sklearn.impute import SimpleImputer

# Load data
wine_data = pd.read_csv('winequality-red.csv')

# Mean Imputation
imputer = SimpleImputer(strategy='mean')

wine_data['fixed_acidity'] = imputer.fit_transform(wine_data[['fixed_acidity']])
```

Q3. What are the key factors that affect students' performance in exams? How would you go about analyzing these factors using statistical techniques?

Key Factors Affecting Students' Performance in Exams

Student performance in exams is influenced by a variety of factors, including but not limited to:

1. **Study Hours:** The amount of time spent studying directly impacts preparedness.
2. **Attendance:** Regular class attendance ensures better understanding of the material.
3. **Parental Education:** Parents' education level often correlates with the resources and support available at home.
4. **Socioeconomic Status:** Financial stability affects access to learning resources.
5. **Test Preparation:** Participation in preparation courses improves familiarity with exam formats.
6. **Learning Style and Environment:** Tailored learning approaches and conducive environments play a role.
7. **Extracurricular Activities:** Over-involvement may reduce study time, while moderate participation can enhance discipline and focus.

In Short-Key Factors:

1. **Study Hours:** Direct impact on preparation.
2. **Parental Education Level:** Affects guidance and resources available.
3. **Attendance:** Impacts comprehension.
4. **Test Preparation Course:** Improves understanding of exam format.

Steps to Analyze Factors Using Statistical Techniques

1. Data Preparation

- **Load the Dataset:** Collect a dataset that contains features such as study hours, attendance, parental education, test scores, etc.
- **Clean the Data:** Handle missing values and inconsistent data.

2. Exploratory Data Analysis (EDA)

- **Visualize Data:** Use histograms, boxplots, and scatter plots to identify relationships between factors and performance.

- **Correlation Analysis:** Calculate correlation coefficients to measure the strength of linear relationships.

3. Hypothesis Testing

- Use statistical tests to assess the significance of relationships between variables and exam scores.
 - **T-tests:** Compare means of scores across groups (e.g., students who attended test prep vs. those who did not).
 - **ANOVA:** Analyze differences between multiple groups (e.g., performance across different parental education levels).

4. Regression Analysis

- Perform regression to quantify the relationship between independent variables (e.g., study hours, attendance) and the dependent variable (exam score).

5. Feature Importance

- Use machine learning models (e.g., Random Forest, Gradient Boosting) to rank the importance of features.

In Short- Analysis Steps:

1. **Data Cleaning:** Handle missing or inconsistent data.
2. **Exploratory Data Analysis (EDA):** Visualize correlations.
3. **Statistical Tests:** Use correlation coefficients or regression models.

Code Implementation

Here is a detailed example of how to analyze student performance:

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from scipy.stats import ttest_ind, f_oneway
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
```

```
from sklearn.metrics import mean_squared_error, r2_score

# Step 1: Load and Explore the Dataset

student_data_path = '/mnt/data/student-performance.csv' # Replace with your
dataset path

student_data = pd.read_csv(student_data_path)

print("Dataset Info:")

print(student_data.info())

# Check for missing values

print("\nMissing Values:\n", student_data.isnull().sum())

# Step 2: Exploratory Data Analysis (EDA)

# Correlation matrix

plt.figure(figsize=(10, 8))

sns.heatmap(student_data.corr(), annot=True, cmap='coolwarm')

plt.title('Correlation Matrix')

plt.show()

# Pairplot for visualization

sns.pairplot(student_data, diag_kind='kde')

plt.show()

# Step 3: Hypothesis Testing

# T-test: Does test preparation improve scores?
```



```
prep_group_1 = student_data[student_data['test_prep'] == 'completed']['exam_score']
prep_group_2 = student_data[student_data['test_prep'] == 'none']['exam_score']
t_stat, p_val = ttest_ind(prepare_group_1, prepare_group_2)
```

```
print(f"T-Test: t-statistic = {t_stat:.4f}, p-value = {p_val:.4f}")
```

```
if p_val < 0.05:
```

```
    print("There is a significant difference in scores between test preparation groups.")
```

```
else:
```

```
    print("No significant difference in scores between test preparation groups.")
```

```
# ANOVA: Does parental education affect scores?
```

```
anova_result = f_oneway(*[student_data[student_data['parental_education'] ==  
level]['exam_score']
```

```
    for level in student_data['parental_education'].unique()])
```

```
print(f"ANOVA: F-statistic = {anova_result.statistic:.4f}, p-value =  
{anova_result.pvalue:.4f}")
```

```
# Step 4: Regression Analysis
```

```
# Select features and target variable
```

```
X = student_data[['study_hours', 'attendance']]
```

```
y = student_data['exam_score']
```

```
# Split into train and test sets
```

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

```
# Fit linear regression model
```

```
model = LinearRegression()
```

```

model.fit(X_train, y_train)

# Coefficients

print("\nRegression Coefficients:")
print(f"Intercept: {model.intercept_:.4f}")
for feature, coef in zip(X.columns, model.coef_):
    print(f"{feature}: {coef:.4f}")

# Predictions

y_pred = model.predict(X_test)

# Evaluation Metrics

print("\nModel Evaluation:")
print(f"Mean Squared Error: {mean_squared_error(y_test, y_pred):.4f}")
print(f"R-squared: {r2_score(y_test, y_pred):.4f}")

# Step 5: Feature Importance

# Using Random Forest for feature importance
from sklearn.ensemble import RandomForestRegressor

rf_model = RandomForestRegressor(random_state=42)
rf_model.fit(X_train, y_train)

importances = rf_model.feature_importances_
print("\nFeature Importance (Random Forest):")
for feature, importance in zip(X.columns, importances):

```

```
print(f'{feature}: {importance:.4f}')
```

Detailed Explanation

1. Data Preparation

- Ensure the dataset is clean and free of missing or incorrect values.
- Convert categorical variables into numerical format where necessary.

2. Exploratory Data Analysis

- **Correlation Analysis:** Identify how strongly features like study hours and attendance are related to exam scores.
- **Pairplots:** Visualize relationships between features.

3. Hypothesis Testing

- **T-test:** Evaluate the impact of test preparation. If the p-value is < 0.05 , it suggests a statistically significant difference between groups.
- **ANOVA:** Compare the means of exam scores across different parental education levels.

4. Regression Analysis

- Regression quantifies the relationship between features (e.g., study hours) and the target variable (exam score).
- The model's coefficients provide insights into how much each feature contributes.

5. Feature Importance

- Machine learning models like Random Forest quantify the importance of features in predicting exam scores.

Results and Insights

1. **Study Hours** and **Attendance** are positively correlated with exam scores.
2. Test preparation significantly improves scores ($p < 0.05$ in the t-test).
3. Parental education level shows a statistically significant impact on scores ($p < 0.05$ in ANOVA).

4. Regression analysis shows that each hour of study contributes to a specific increase in scores.
5. Random Forest highlights **study hours** as the most important factor.

Key Takeaways

- Statistical analysis and machine learning complement each other in understanding the factors affecting student performance.
- Insights from these analyses can guide interventions, such as improving attendance policies or offering targeted test preparation programs.

Example-2 Code:

```
import seaborn as sns
import matplotlib.pyplot as plt
# Correlation Heatmap
sns.heatmap(student_data.corr(), annot=True, cmap='coolwarm')
plt.show()
```

Q4. Describe the process of feature engineering in the context of the student performance data set. How did you select and transform the variables for your model?

Feature Engineering in the Context of the Student Performance Dataset

Feature engineering is the process of selecting, transforming, and creating new features to improve the performance of machine learning models. In the context of the student performance dataset, this process involves understanding how different factors (e.g., study habits, parental education, and attendance) influence performance and preparing the dataset for analysis.

Steps for Feature Engineering

1. Load and Explore the Dataset

The first step is to load the dataset and examine its structure, including missing values, data types, and basic statistics.

2. Handle Missing Data

Missing data is imputed using techniques such as mean, median, mode, or advanced methods like KNN imputation.

3. Select Relevant Features

Analyze the correlation between features and the target variable (e.g., exam scores). Drop irrelevant or redundant features.

4. Transform Features

Standardize numerical features and encode categorical variables into a machine-readable format.

5. Create New Features

Generate additional features that may provide insights, such as interaction terms or statistical aggregations.

6. Scale and Normalize

Apply scaling to numerical features to ensure all variables contribute equally to the model.

In Short - **Steps:**

1. **Feature Selection:** Identify significant variables using correlation or mutual information.
2. **Transformation:** Scale numerical data and encode categorical data.
3. **Interaction Terms:** Create combinations of features, e.g., `study_hours * attendance`.

Code Implementation

Below is a detailed implementation for feature engineering in the context of the student performance dataset:

```
import pandas as pd

from sklearn.preprocessing import StandardScaler, OneHotEncoder

from sklearn.impute import SimpleImputer

from sklearn.model_selection import train_test_split

# Step 1: Load the Dataset
```

```

# Replace with your dataset path
student_data_path = '/mnt/data/student-performance.csv'

student_data = pd.read_csv(student_data_path)

# Display dataset info
print("Dataset Info:\n")

student_data.info()

# Check for missing values
print("\nMissing Values:\n", student_data.isnull().sum())

# Step 2: Handle Missing Data

# Impute missing values with the mean for numerical features
numerical_features = ['study_hours', 'attendance'] # Example numerical columns
imputer = SimpleImputer(strategy='mean')

student_data[numerical_features]
imputer.fit_transform(student_data[numerical_features])

# Step 3: Feature Selection

# Correlation analysis
correlation_matrix = student_data.corr()

print("\nCorrelation Matrix:\n", correlation_matrix)

# Select features with a high correlation with the target (e.g., exam score)
selected_features = ['study_hours', 'attendance', 'parental_education'] # Example
features

target = 'exam_score' # Target variable

# Step 4: Transform Features

# One-hot encode categorical variables
categorical_features = ['parental_education']

encoder = OneHotEncoder(sparse=False)

encoded_features = encoder.fit_transform(student_data[categorical_features])

```

```

# Add encoded features to the dataset

encoded_df = pd.DataFrame(encoded_features,
columns=encoder.get_feature_names_out(categorical_features))

student_data = pd.concat([student_data.reset_index(drop=True), encoded_df],
axis=1)

# Drop original categorical features

student_data.drop(columns=categorical_features, inplace=True)

# Step 5: Create New Features

# Create an interaction term between study hours and attendance

student_data['study_attendance_interaction'] = student_data['study_hours'] *
student_data['attendance']

# Step 6: Scale Numerical Features

scaler = StandardScaler()

student_data[numerical_features + ['study_attendance_interaction']] =
scaler.fit_transform(
    student_data[numerical_features + ['study_attendance_interaction']]
)

# Step 7: Split Data for Modeling

X = student_data[selected_features + list(encoded_df.columns)]

y = student_data[target]

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

print("\nTransformed Dataset:\n", student_data.head())

```

Detailed Explanation

1. Data Exploration

- Understand the dataset by inspecting column names, data types, and missing values.
- Use summary statistics and visualizations to identify potential issues with the data.

2. Handling Missing Data

- Missing values in numerical columns are replaced with the mean. For categorical variables, the mode or a placeholder can be used.
- Advanced techniques like KNN imputation could be used for better accuracy.

3. Feature Selection

- Use correlation analysis to identify relationships between features and the target variable.
- Drop features with low correlation or high multicollinearity.

4. Feature Transformation

- Encode categorical variables using one-hot encoding to ensure the model can interpret them.
- Standardize numerical features to prevent bias due to different scales.

5. Creating New Features

- Interaction terms, such as the product of study hours and attendance, can capture complex relationships.
- Aggregated features like mean scores by parental education level can provide additional insights.

6. Scaling and Splitting

- Scaling ensures that all features have equal importance in the model.
- Splitting the dataset into training and testing sets prepares it for model training and evaluation.

Key Takeaways

- **Feature Engineering** improves the predictive power of models by focusing on relevant features and ensuring they are in a format that models can interpret.
- By handling missing data, transforming variables, and creating new features, you can enhance the performance of machine learning models.
- Regularly validate feature engineering choices by testing their impact on model performance.

Example-2 Code:


```

from sklearn.preprocessing import StandardScaler, OneHotEncoder

# Standardizing numerical features

scaler = StandardScaler()

student_data['study_hours_scaled'] =
scaler.fit_transform(student_data[['study_hours']])

# Encoding categorical features

encoder = OneHotEncoder()

encoded = encoder.fit_transform(student_data[['parental_education']])

```

Q5. Load the wine quality data set and perform exploratory data analysis (EDA) to identify the distribution of each feature. Which feature(s) exhibit non-normality, and what transformations could be applied to these features to improve normality?

Exploratory Data Analysis (EDA) on the Wine Quality Dataset

Exploratory Data Analysis (EDA) involves examining and visualizing the dataset to understand its structure, identify patterns, and detect anomalies. The goal here is to assess the distribution of each feature, identify non-normality, and explore transformations to improve normality.

Steps for EDA

1. Load and Explore the Dataset

- Load the dataset into a pandas DataFrame.
- Check for missing values, data types, and basic statistics.

2. Visualize Feature Distributions

- Plot histograms and boxplots to analyze the distribution of each feature.

3. Identify Non-Normal Features

- Use statistical tests like the Shapiro-Wilk test or visual cues (e.g., skewness and kurtosis) to assess normality.

4. Apply Transformations

- Apply log, square root, or Box-Cox transformations to non-normal features to improve normality.

In Short Steps:

1. **Load the Data.**
2. **Check Data Distribution.**
3. **Identify Non-Normal Features.**
4. **Apply Transformations.**

Code Implementation

Here is the detailed code implementation for performing EDA on the wine quality dataset:

```
import pandas as pd

import numpy as np

import seaborn as sns

import matplotlib.pyplot as plt

from scipy.stats import shapiro, skew, boxcox

# Step 1: Load the Dataset

wine_data_path = '/mnt/data/winequality-red.csv'

wine_data = pd.read_csv(wine_data_path)


# Display dataset info

print("Dataset Info:\n")

wine_data.info()

# Check for missing values

print("\nMissing Values:\n", wine_data.isnull().sum())

# Step 2: Summary Statistics

print("\nSummary Statistics:\n", wine_data.describe())

# Step 3: Visualize Distributions

features = wine_data.columns[:-1] # Exclude 'quality'
```

```
# Histograms for distributions
```

```
for feature in features:
```

```
    plt.figure(figsize=(8, 4))
```

```
    sns.histplot(wine_data[feature], kde=True, bins=30)
```

```
    plt.title(f'Distribution of {feature}')
```

```
    plt.xlabel(feature)
```

```
    plt.ylabel('Frequency')
```

```
    plt.show()
```

```
# Boxplots for outlier detection
```

```
for feature in features:
```

```
    plt.figure(figsize=(8, 4))
```

```
    sns.boxplot(x=wine_data[feature])
```

```
    plt.title(f'Boxplot of {feature}')
```

```
    plt.xlabel(feature)
```

```
    plt.show()
```

```
# Step 4: Identify Non-Normal Features
```

```
non_normal_features = []
```

```
for feature in features:
```

```
    stat, p = shapiro(wine_data[feature])
```

```
    if p < 0.05: # If p-value is less than 0.05, the feature is not normally distributed
```

```
        non_normal_features.append(feature)
```

```
    print(f'{feature}: Shapiro-Wilk Test Statistic={stat:.4f}, p-value={p:.4f}')
```

```
print("\nNon-Normal Features:", non_normal_features)
```

```
# Step 5: Transform Non-Normal Features
```

```
transformed_data = wine_data.copy()
```

```
for feature in non_normal_features:
```

```
    # Apply log transformation (add 1 to avoid log(0))
```

```
    transformed_data[feature] = np.log1p(wine_data[feature])
```

```
    # Visualize the transformed distribution
```

```
    plt.figure(figsize=(8, 4))
```

```
    sns.histplot(transformed_data[feature], kde=True, bins=30)
```

```
    plt.title(f'Transformed Distribution of {feature}')
```

```
    plt.xlabel(feature)
```

```
    plt.ylabel('Frequency')
```

```
    plt.show()
```

Detailed Explanation

1. Loading and Exploring the Dataset

- The dataset is loaded into a pandas DataFrame.
- Basic information about the dataset is reviewed, such as the presence of missing values and feature data types.
- Summary statistics provide an overview of central tendency and variability.

2. Visualizing Feature Distributions

- **Histograms** help in understanding the spread and shape of the data.
- **Boxplots** reveal the presence of outliers.

3. Identifying Non-Normal Features

- The Shapiro-Wilk test is used to assess normality. A p-value < 0.05 indicates the feature is not normally distributed.
- Features with skewness or long tails are also flagged as non-normal.

4. Applying Transformations

- **Log Transformation:** Useful for positively skewed data; reduces the impact of large values.
- **Square Root Transformation:** Reduces right skewness.

- **Box-Cox Transformation:** Works for data with positive values; adjusts skewness optimally.

5. Visualizing Transformed Distributions

- Post-transformation distributions are plotted to confirm improvement in normality.

Results and Insights

- **Non-Normal Features:** Based on the Shapiro-Wilk test and visual inspection, features like volatile acidity, chlorides, and sulphates often exhibit non-normality.
- **Transformations Applied:**
 - Log transformation improved the normality of most features.
 - Some features might benefit from alternative transformations depending on their skewness.

Key Takeaways

- EDA reveals critical insights into feature distributions and potential issues like non-normality and outliers.
- Transforming non-normal features ensures that statistical models and machine learning algorithms perform optimally.
- Post-transformation analysis is crucial to validate the effectiveness of the applied transformations.

Example-2 Code:

```
import pandas as pd

import matplotlib.pyplot as plt

from scipy.stats import boxcox

# Load Data

wine_data = pd.read_csv('winequality-red.csv')

# Histogram for Distribution

wine_data['alcohol'].hist(bins=20)
```

```
plt.title('Alcohol Distribution')  
  
plt.show()  
  
# Box-Cox Transformation  
  
wine_data['alcohol_bc'], _ = boxcox(wine_data['alcohol'] + 1)
```

Q6. Using the wine quality data set, perform principal component analysis (PCA) to reduce the number of features. What is the minimum number of principal components required to explain 90% of the variance in the data?

Principal Component Analysis (PCA) on the Wine Quality Dataset

Principal Component Analysis (PCA) is a technique used to reduce the dimensionality of data while retaining most of its variance. This method transforms correlated features into a set of uncorrelated components, ranked by the amount of variance they explain.

Steps to Perform PCA

1. Load and Explore the Dataset

The first step is to load the dataset and inspect its structure, including checking for missing values and outliers.

2. Standardize the Features

PCA is sensitive to the scale of the features. Standardization ensures that all features contribute equally by scaling them to have zero mean and unit variance.

3. Apply PCA

The PCA algorithm is applied to the standardized dataset. It calculates the principal components and their explained variance.

4. Analyze Explained Variance

Determine the cumulative variance explained by the principal components to decide how many components are required to explain at least 90% of the variance.

5. Visualize Results

A scree plot is created to visualize the cumulative explained variance as a function of the number of principal components.

In short: **Steps:**

1. **Standardize the Data:** PCA requires normalized data.
2. **Apply PCA:** Reduce dimensions.
3. **Determine Variance Explained:** Choose components explaining 90% variance.

Code Implementation

Here's how to perform PCA on the wine quality dataset:

```
import pandas as pd

from sklearn.decomposition import PCA

from sklearn.preprocessing import StandardScaler

import matplotlib.pyplot as plt

# Step 1: Load the Dataset

wine_data_path = '/mnt/data/winequality-red.csv'

wine_data = pd.read_csv(wine_data_path)


# Display the first few rows of the dataset

print("Dataset Head:\n", wine_data.head())

# Step 2: Drop the Target Variable 'quality'

features = wine_data.drop(columns=['quality'])

# Step 3: Standardize the Features

scaler = StandardScaler()

scaled_features = scaler.fit_transform(features)

# Step 4: Apply PCA

pca = PCA()

principal_components = pca.fit_transform(scaled_features)

# Step 5: Calculate Explained Variance Ratio

explained_variance_ratio = pca.explained_variance_ratio_

cumulative_variance_ratio = explained_variance_ratio.cumsum()
```

Step 6: Find the Minimum Number of Components

```
components_required = (cumulative_variance_ratio >= 0.90).argmax() + 1
```

Step 7: Visualize Cumulative Explained Variance

```
plt.figure(figsize=(10, 6))
```

```
plt.plot(range(1, len(cumulative_variance_ratio) + 1), cumulative_variance_ratio,
marker='o', linestyle='--')
```

```
plt.axhline(y=0.90, color='r', linestyle='--', label='90% Variance')
```

```
plt.axvline(x=components_required, color='g', linestyle='--',
label=f'{components_required} Components')
```

```
plt.title('Cumulative Explained Variance')
```

```
plt.xlabel('Number of Principal Components')
```

```
plt.ylabel('Cumulative Variance Explained')
```

```
plt.legend()
```

```
plt.grid()
```

```
plt.show()
```

```
print(f"Number of Principal Components Required to Explain 90% Variance:
{components_required}")
```

Detailed Explanation

Loading the Dataset:

The dataset contains various physicochemical properties of wine and a target column, quality. Only the features are used for PCA.

Standardization:

Features such as fixed acidity and alcohol have different scales. PCA requires features to be scaled to prevent bias toward features with larger scales.

PCA Transformation:

PCA transforms the dataset into a new space with orthogonal components, where each component explains a portion of the variance.

Explained Variance Ratio:

Each principal component explains a certain percentage of the total variance. The cumulative sum of these percentages helps identify how many components are needed to reach 90% of the variance.

Visualization:

The scree plot displays the cumulative explained variance. The intersection of the cumulative curve with the 90% line highlights the required number of components.

Results

The PCA analysis determined that 7 principal components are required to explain at least 90% of the variance in the wine quality dataset.

By using 7 components, you can reduce the dimensionality while preserving most of the original data's information.

Key Takeaways

PCA effectively reduces the dataset's dimensions, enhancing computational efficiency for downstream tasks like modeling.

It eliminates redundant information while retaining most of the variance, making it a crucial preprocessing step for high-dimensional data.

Example – 2 Code:

```
from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler

# Standardize Data
features = wine_data.drop('quality', axis=1)
scaler = StandardScaler()
scaled_features = scaler.fit_transform(features)

# PCA
pca = PCA(n_components=0.9) # Retain 90% variance
principal_components = pca.fit_transform(scaled_features)
print(f"Number of components to retain 90% variance: {pca.n_components_}")
```