

Here is the markdown version of the analysis:

Step-by-Step Analysis Structure

Step 1: Selecting the Dataset

****Question Before Dataset Selection**:**

– Why did I choose the ****Student Performance**** dataset for my analysis?

****Answer**:**

– I chose the ****Student Performance**** dataset because it contains a rich set of features

Step 2: Data Understanding

****Questions Before Data Understanding**:**

1. What are the key attributes of the dataset, and what do they represent?
2. Are there any missing values or anomalies in the data?

****Answers**:**

1. The dataset contains 33 attributes, including demographic information (e.g., age, gender).
2. There are no missing values in the dataset; however, some attributes might require further exploration.

Step 3: Data Wrangling

****Questions Before Data Wrangling**:**

1. How can we handle categorical variables for machine learning models?
2. Are there any irrelevant or redundant features that should be removed?
3. How should we deal with outliers in the data?

****Answers**:**

1. Categorical variables can be handled using techniques like one-hot encoding, which converts categories into numerical form.
2. Some features like `school` (if all students belong to the same school) might be redundant.
3. Outliers can be detected using box plots or z-score analysis and handled by either removing them or using robust statistical methods.

Step 4: Feature Engineering

****Questions Before Feature Engineering**:**

1. Which features are the most important in predicting student performance?
2. Can we create new features that might better represent the underlying patterns?

****Answers**:**

1. Features like `study time`, `parental education`, `absences`, and `previous grades (G1`
2. Yes, new features can be created by combining existing ones. For instance, combining `

Step 5: Model Building****Questions Before Model Building**:**

1. Which machine learning models are suitable for this type of prediction problem?
2. How should we split the dataset into training and testing sets to ensure robust evaluation?

****Answers**:**

1. Models like Linear Regression, Decision Trees, Random Forest, and Gradient Boosting are
2. The dataset should be split into training and testing sets (e.g., 80% training, 20% test).

Step 6: Model Evaluation****Questions Before Model Evaluation**:**

1. What evaluation metrics should be used to assess the model's performance?
2. How can we improve model performance if the initial results are not satisfactory?

****Answers**:**

1. For regression tasks, metrics like Mean Absolute Error (MAE), Mean Squared Error (MSE)
2. Model performance can be improved by tuning hyperparameters, using cross-validation, for

Conclusion****Question After Full Process**:**

- What insights have you gained from the analysis of the ****Student Performance**** dataset?

****Answer**:**

- The analysis revealed that certain socio-economic and behavioral factors, such as `study

You can use this Markdown code to format the content as needed.

In-Depth Data Wrangling for Student Performance Dataset

The data wrangling process involves cleaning and preparing the dataset to ensure it is ready for

Step 1: Data Loading and Initial Inspection

First, let's load the dataset and inspect its structure, including the first few rows, da

```
```python
import pandas as pd

Load the dataset
student_mat = pd.read_csv('student-mat.csv')
student_por = pd.read_csv('student-por.csv')

Initial inspection of the data
print(student_mat.head())
print(student_mat.info())
print(student_mat.describe(include='all'))
```

### Key Observations:

- We have 33 columns in each dataset.
- The data types include both numeric and categorical variables.
- No missing values are reported, but further checks are needed.

## ✓ Step 2: Handling Missing Values

Even though the initial inspection shows no missing values, it is good practice to confirm by checking for any potential anomalies that may not be immediately visible.

```
Check for missing values
print(student_mat.isnull().sum())
print(student_por.isnull().sum())
```

### Handling Missing Values:

- If there were any missing values, we'd decide on strategies like imputation (using mean, median, or mode) or removing rows/columns based on the context. However, in this dataset, there are no missing values.

## Step 3: Encoding Categorical Variables

The dataset contains several categorical variables (e.g., school, sex, address, famsize). Machine learning models typically require numeric input, so we need to encode these variables.

```
List of categorical columns
categorical_columns = ['school', 'sex', 'address', 'famsize', 'Pstatus', 'Mjob', 'Fjob',
 'schoolsup', 'famsup', 'paid', 'activities', 'nursery', 'higher',
```

```
Use one-hot encoding for categorical variables
student_mat_encoded = pd.get_dummies(student_mat, columns=categorical_columns, drop_first=True)
student_por_encoded = pd.get_dummies(student_por, columns=categorical_columns, drop_first=True)

print(student_mat_encoded.head())
print(student_por_encoded.head())
```

### Encoding Strategy:

- We use **one-hot encoding** for categorical variables to convert them into binary variables (0 or 1). The `drop_first=True` parameter avoids multicollinearity by dropping the first category.

## Step 4: Handling Outliers

Outliers can skew the model and reduce performance. We will use visualization and statistical methods to detect outliers in continuous variables like `age`, `absences`, and grades (`G1`, `G2`, `G3`).

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

Visualize outliers using box plots
sns.boxplot(data=student_mat[['age', 'absences', 'G1', 'G2', 'G3']])
plt.show()

Detect outliers using the Z-score method
from scipy import stats

z_scores = stats.zscore(student_mat[['age', 'absences', 'G1', 'G2', 'G3']])
abs_z_scores = abs(z_scores)
filtered_entries = (abs_z_scores < 3).all(axis=1)
student_mat_cleaned = student_mat[filtered_entries]

print(f"Number of rows before cleaning: {student_mat.shape[0]}")
print(f"Number of rows after removing outliers: {student_mat_cleaned.shape[0]}")
```

### Outlier Detection and Treatment:

- **Box Plots:** Visualize the spread of data and identify potential outliers.
- **Z-Score Method:** Removes rows where the absolute z-score is greater than 3 (outliers). Adjust this threshold based on the domain knowledge.

## Step 5: Feature Transformation

Transform features if necessary to normalize their distributions or to better represent the underlying data patterns.

```
from sklearn.preprocessing import StandardScaler
import numpy as np

Log transformation of skewed data
student_mat_cleaned['log_absences'] = student_mat_cleaned['absences'].apply(lambda x: np.log(x+1))

Standard scaling for continuous variables
scaler = StandardScaler()
scaled_features = scaler.fit_transform(student_mat_cleaned[['age', 'G1', 'G2', 'G3']])

Replace the original columns with scaled values
student_mat_cleaned[['age', 'G1', 'G2', 'G3']] = scaled_features
print(student_mat_cleaned[['age', 'G1', 'G2', 'G3']].head())
```

### Feature Transformation:

- **Log Transformation:** Reduces skewness in absences .
- **Standard Scaling:** Normalizes continuous variables for better model performance.

## Step 6: Feature Selection and Scaling

Finally, we may need to select the most relevant features for our machine learning model to avoid overfitting and improve efficiency.

```
from sklearn.feature_selection import SelectKBest, f_regression

Feature selection
X = student_mat_cleaned.drop(['G3'], axis=1) # Drop the target variable
y = student_mat_cleaned['G3']

Use SelectKBest to select top features
best_features = SelectKBest(score_func=f_regression, k=10)
fit = best_features.fit(X, y)
dfscores = pd.DataFrame(fit.scores_)
dfcolumns = pd.DataFrame(X.columns)

Concatenate dataframes for better visualization
feature_scores = pd.concat([dfcolumns, dfscores], axis=1)
feature_scores.columns = ['Feature', 'Score'] # Naming the data frame columns
print(feature_scores.nlargest(10, 'Score')) # Print 10 best features
```

### Feature Selection:

- **SelectKBest**: Chooses the top 10 features based on the highest scores from the univariate statistical test ( `f_regression` ).

## Conclusion

The data wrangling process for the **Student Performance** dataset ensures that the data is clean, well-encoded, and ready for feature engineering and model building. This in-depth approach includes handling missing values, encoding categorical variables, detecting and treating outliers, transforming features, and selecting the most relevant features for the machine learning model.

To provide a well-detailed Exploratory Data Analysis (EDA) for the **Student Performance**:

### ### Steps for Exploratory Data Analysis (EDA)

1. **Understanding Data Structure**
2. **Univariate Analysis**
3. **Bivariate Analysis**
4. **Multivariate Analysis**
5. **Correlation Analysis**
6. **Feature Distributions**
7. **Key Findings and Insights**

### ### 1. Understanding Data Structure

Before diving into the analysis, let's take a closer look at the dataset's structure, including

```
```python
import pandas as pd

# Load the dataset
student_mat = pd.read_csv('student-mat.csv')
student_por = pd.read_csv('student-por.csv')

# Display basic information
print(student_mat.info())
print(student_por.info())

# Display the first few rows
print(student_mat.head())
print(student_por.head())
```

Key Observations:

- Each dataset contains 33 columns.
- The columns include categorical (e.g., `school`, `sex`, `address`) and numerical variables (e.g., `age`, `absences`, `G1`, `G2`, `G3`).

2. Univariate Analysis

Univariate analysis focuses on analyzing each feature individually to understand their distributions, central tendencies, and variability.

Categorical Variables

Let's analyze the distribution of categorical variables such as `school`, `sex`, `address`, `famsize`, etc.

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

# Set up the matplotlib figure
fig, axes = plt.subplots(3, 3, figsize=(15, 15))
fig.suptitle('Univariate Analysis of Categorical Variables')

# Categorical variable distributions
sns.countplot(x='school', data=student_mat, ax=axes[0, 0])
sns.countplot(x='sex', data=student_mat, ax=axes[0, 1])
sns.countplot(x='address', data=student_mat, ax=axes[0, 2])
sns.countplot(x='famsize', data=student_mat, ax=axes[1, 0])
sns.countplot(x='Pstatus', data=student_mat, ax=axes[1, 1])
sns.countplot(x='Mjob', data=student_mat, ax=axes[1, 2])
sns.countplot(x='Fjob', data=student_mat, ax=axes[2, 0])
sns.countplot(x='reason', data=student_mat, ax=axes[2, 1])
sns.countplot(x='guardian', data=student_mat, ax=axes[2, 2])

plt.tight_layout()
plt.show()
```

Numerical Variables

Next, let's analyze the distributions of numerical variables like `age`, `absences`, and the three grades (`G1`, `G2`, `G3`).

```
# Set up the matplotlib figure for numerical variables
fig, axes = plt.subplots(2, 2, figsize=(12, 10))
fig.suptitle('Univariate Analysis of Numerical Variables')

# Numerical variable distributions
sns.histplot(student_mat['age'], kde=True, ax=axes[0, 0])
sns.histplot(student_mat['absences'], kde=True, ax=axes[0, 1])
sns.histplot(student_mat['G1'], kde=True, ax=axes[1, 0])
sns.histplot(student_mat['G3'], kde=True, ax=axes[1, 1])
```

```
plt.tight_layout()
plt.show()
```

Key Insights:

- **Categorical Variables:** The majority of students come from one school, with slightly more female than male students. Most students live in urban areas (address), with Mjob and Fjob categories indicating different occupations.
- **Numerical Variables:** Most students are aged between 15 and 18. Absences show a skewed distribution, suggesting a few students have high absenteeism. The grade distributions (G1 , G2 , G3) are roughly normal.

3. Bivariate Analysis

Bivariate analysis helps understand the relationships between two variables, especially how independent variables affect the dependent variable (G3 - final grade).

Categorical vs. Numerical Variables

Analyze how categorical variables impact the final grades (G3).

```
# Set up the matplotlib figure
fig, axes = plt.subplots(3, 3, figsize=(15, 15))
fig.suptitle('Bivariate Analysis: Categorical Variables vs. G3')

# Impact of categorical variables on final grade
sns.boxplot(x='school', y='G3', data=student_mat, ax=axes[0, 0])
sns.boxplot(x='sex', y='G3', data=student_mat, ax=axes[0, 1])
sns.boxplot(x='address', y='G3', data=student_mat, ax=axes[0, 2])
sns.boxplot(x='famsize', y='G3', data=student_mat, ax=axes[1, 0])
sns.boxplot(x='Pstatus', y='G3', data=student_mat, ax=axes[1, 1])
sns.boxplot(x='Mjob', y='G3', data=student_mat, ax=axes[1, 2])
sns.boxplot(x='Fjob', y='G3', data=student_mat, ax=axes[2, 0])
sns.boxplot(x='reason', y='G3', data=student_mat, ax=axes[2, 1])
sns.boxplot(x='guardian', y='G3', data=student_mat, ax=axes[2, 2])

plt.tight_layout()
plt.show()
```

Key Insights:

- **School and Address:** Students in some schools or from urban areas seem to perform slightly better.
- **Family Size (famsize) and Parental Status (Pstatus):** No significant differences observed in performance.

Numerical vs. Numerical Variables

Analyze correlations between numerical variables, such as G1, G2, and G3.

```
# Scatter plot of numerical variables
sns.pairplot(student_mat[['G1', 'G2', 'G3', 'age', 'absences']], kind='scatter', diag_kind='hist')
plt.suptitle('Bivariate Analysis: Numerical Variables', y=1.02)
plt.show()
```

Key Insights:

- **Grades (G1, G2, G3):** Strong positive correlations between G1, G2, and G3 grades, as expected.
- **Absences:** Negative correlation with G3, suggesting more absences lead to lower final grades.

4. Multivariate Analysis

Multivariate analysis helps us explore the interactions between multiple variables and how they collectively affect the target variable (G3).

```
# Visualizing the impact of multiple variables using a heatmap
plt.figure(figsize=(12, 8))
sns.heatmap(student_mat.corr(), annot=True, cmap='coolwarm', fmt=".2f")
plt.title('Correlation Heatmap')
plt.show()
```

Key Insights:

- **Correlations:** Strong correlation between G1 and G2 with G3. This suggests early performance is indicative of final performance.
- **Age and Absences:** Weak negative correlations with G3.

5. Feature Distributions

Analyzing the distribution of key features helps in understanding their spread and identifying any necessary transformations.

```
# Visualize the distribution of 'absences' after log transformation
import numpy as np

student_mat['log_absences'] = np.log1p(student_mat['absences'])
sns.histplot(student_mat['log_absences'], kde=True)
plt.title('Distribution of Log-Transformed Absences')
plt.show()
```

Key Insights:

- **Absences:** Log transformation reduces skewness, making the data more suitable for regression models.

6. Key Findings and Insights

- **Significant Predictors:** Features like `G1`, `G2`, and `absences` are significant predictors of `G3` (final grade).
- **Impact of Categorical Variables:** Certain features like `school`, `sex`, and `Mjob` show some differentiation in `G3`.
- **Feature Relationships:** Strong correlations exist between sequential grades (`G1`, `G2`, `G3`), indicating the importance of continuous performance.

Conclusion

The EDA reveals key insights into student performance in secondary schools. The grades (`G1`, `G2`, `G3`) have strong correlations, absences negatively impact performance, and several categorical variables show potential significance. This information is valuable for building predictive models and understanding student performance drivers.

Conclusions from the Exploratory Data Analysis (EDA)

Based on the comprehensive EDA performed on the **Student Performance** dataset, several significant insights and patterns were identified. These findings provide a deeper understanding of the factors influencing students' final grades (`G3`). Here are the key conclusions:

1. Strong Correlation Between Sequential Grades:

- The grades `G1` (first period) and `G2` (second period) have a very strong positive correlation with `G3` (final grade). This indicates that students' performance tends to be consistent over time. Students who perform well in earlier periods (`G1` and `G2`) are likely to perform well in the final period (`G3`). This finding suggests that early intervention based on initial grades could be beneficial in helping students improve or maintain their academic performance.

2. Impact of Absenteeism on Academic Performance:

- There is a noticeable negative correlation between the number of absences (`absences`) and the final grade (`G3`). This indicates that students who are frequently absent tend to have lower final grades. Addressing absenteeism through better student engagement, support systems, or other interventions could potentially lead to improved academic outcomes.

3. Effect of Demographic and Socioeconomic Factors:

- Certain categorical variables such as `school`, `sex`, `Mjob` (mother's job), and `Fjob` (father's job) show variations in students' final grades. For example, students from different schools or with parents in certain professions show varying academic performances. While these factors do not have as strong an impact as the sequential

grades or absences, they provide context that may be important for understanding underlying disparities and targeting specific student groups for support.

4. Weak Influence of Family and Personal Characteristics:

- Features like `famsize` (family size), `Pstatus` (parental status), and `guardian` (guardian type) do not show significant variation in the final grades (`G3`). This suggests that, within this dataset, these variables may have a limited direct impact on academic performance. However, these factors could still play a role in a broader socio-emotional context.

5. Alcohol Consumption and Academic Performance:

- Analysis of the `Da1c` (weekday alcohol consumption) and `Wa1c` (weekend alcohol consumption) variables suggests a potential negative impact on academic performance. Higher levels of alcohol consumption correlate with slightly lower grades, indicating that lifestyle factors can affect students' focus and study habits.

6. Effect of Study Time and Failures:

- The `studytime` variable (weekly study time) shows a positive association with higher grades. Conversely, the `failures` variable (number of past class failures) is negatively associated with the final grade. This indicates that consistent study time and minimizing past failures can help improve academic outcomes. Encouraging effective study habits and providing support to students with past failures could be beneficial strategies for educators.

7. Distribution of Academic Performance:

- The final grades (`G3`) are normally distributed with a slight skew toward lower grades, indicating that a majority of students achieve mid-range scores, with fewer students at the higher and lower extremes. This finding could be useful in setting targeted academic goals or interventions aimed at different segments of students based on their performance.

8. Potential Target Areas for Intervention:

- Based on the correlations and variable impacts observed, areas like improving study habits (`studytime`), reducing absenteeism (`absences`), and managing alcohol consumption (`Da1c`, `Wa1c`) are potential targets for intervention to improve student performance. Tailored support programs focusing on these aspects could result in better academic outcomes.

Final Thoughts

The EDA of the Student Performance dataset reveals a complex interplay of factors that influence academic success. While some variables have a more direct impact on the final grades (`G3`), others provide context and background that can guide more nuanced strategies for

educators and policymakers. Future modeling efforts should focus on leveraging these insights to develop predictive models and intervention strategies that are grounded in the data-driven understanding of student performance dynamics.

Start coding or generate with AI.