```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import scipy.stats as stats
```

```
1 df = pd.read_csv ("/content/student_performance_prediction.csv")
```

```
1 # configuring display options for dataframe
2 pd.set_option('display.max_columns', None) # display max columns
3 pd.set_option('display.max_rows', None) # display max rows
```

```
1 df.sample(5)
```

| | Student ID | Study Hours per Week | Attendance Rate | Previous Grades | Participation in Extracurricular Activities | Parent Education Level | Passed |
|---|---|---|---|---|---|---|---|
| **29460** | S29461 | 11.8 | 101.8 | 95.9 | No | Bachelor | No |
| **189** | S00190 | 2.4 | 82.0 | 60.4 | Yes | Associate | Yes |
| **6084** | S06085 | 9.5 | 50.4 | 28.2 | No | Bachelor | Yes |
| **25215** | S25216 | 13.4 | 82.1 | 62.3 | No | NaN | Yes |

```
1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 40000 entries, 0 to 39999
Data columns (total 7 columns):
 #   Column                                       Non-Null Count  Dtype
---  ------                                       --------------  -----
 0   Student ID                                   40000 non-null  object
 1   Study Hours per Week                         38005 non-null  float64
 2   Attendance Rate                              38008 non-null  float64
 3   Previous Grades                              38006 non-null  float64
 4   Participation in Extracurricular Activities  38000 non-null  object
 5   Parent Education Level                       38000 non-null  object
 6   Passed                                       38000 non-null  object
dtypes: float64(3), object(4)
memory usage: 2.1+ MB
```

```
1 print ("DataSet Shape:", df.shape) # shape
2 print ("Total Size:", df.size) # size
```

```
DataSet Shape: (40000, 7)
Total Size: 280000
```

```
1 df.columns # columns list
```

```
Index(['Student ID', 'Study Hours per Week', 'Attendance Rate',
       'Previous Grades', 'Participation in Extracurricular Activities',
       'Parent Education Level', 'Passed'],
      dtype='object')
```

## ⌄ Handling Missing Values

```
1 # total no. of missing values
2 total_missin_vals = df.isna().sum().sum()
3 print (f"{total_missin_vals} values are missing in the dataset.")
```

```
11981 values are missing in the dataset.
```

```
1 # total % of missing values
2 missin_vals_percent = round((total_missin_vals/df.shape[0])*100)
3 print (f"{missin_vals_percent}% of values are missing in the dataset.")
```

30% of values are missing in the dataset.

```
1 cols_wid_missin_vals = df.columns[df.isna().any()].tolist()
2 print (f"List of columns with missing values: {cols_wid_missin_vals}")
```

List of columns with missing values: ['Study Hours per Week', 'Attendance Rate', 'Previous Grades', 'Participatic

```
1 missin_vals_per_col = df [['Study Hours per Week', 'Attendance Rate', 'Previous Grades', 'Participation in Extracu
2 missin_vals_per_col
```

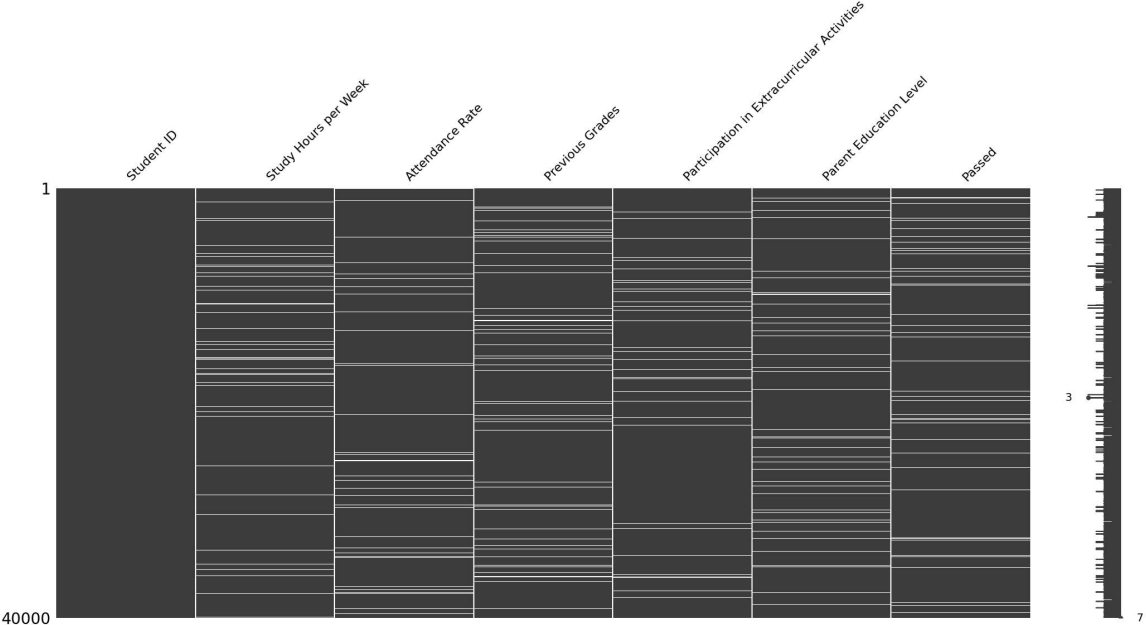|  | 0 |
|---|---|
| Study Hours per Week | 1995 |
| Attendance Rate | 1992 |
| Previous Grades | 1994 |
| Participation in Extracurricular Activities | 2000 |
| Parent Education Level | 2000 |
| Passed | 2000 |

dtype: int64

```
1 percent_missin_vals_per_col = (missin_vals_per_col/df.shape[0])*100
2 percent_missin_vals_per_col
```

|  | 0 |
|---|---|
| Study Hours per Week | 4.9875 |
| Attendance Rate | 4.9800 |
| Previous Grades | 4.9850 |
| Participation in Extracurricular Activities | 5.0000 |
| Parent Education Level | 5.0000 |
| Passed | 5.0000 |

dtype: float64

```
1 # visual for missing values in the data
2 import missingno as msno
3 msno.matrix(df)
4 plt.show()
```

```
1  df[['Study Hours per Week', 'Attendance Rate', 'Previous Grades']].sample(25)
```

| | Study Hours per Week | Attendance Rate | Previous Grades |
|---|---|---|---|
| **39166** | 10.6 | 65.3 | 76.5 |
| **23873** | 15.0 | 45.5 | 49.6 |
| **1832** | 14.8 | 95.7 | 83.0 |
| **38374** | 14.0 | 72.4 | 74.6 |
| **3806** | 10.2 | 103.4 | 43.0 |
| **26264** | 7.8 | 75.9 | 77.7 |
| **31900** | 3.2 | 45.5 | 54.5 |
| **15301** | 15.4 | 82.5 | 36.7 |
| **29081** | 12.3 | 99.4 | 72.1 |
| **12264** | 2.5 | 64.4 | 46.1 |
| **13674** | 8.9 | 94.2 | 54.6 |
| **15202** | 3.5 | 74.4 | 53.4 |
| **13680** | NaN | 74.7 | 60.5 |
| **13717** | 14.1 | 78.4 | 84.5 |
| **24313** | 17.8 | 81.5 | 52.8 |
| **27653** | 6.6 | 93.0 | 73.6 |
| **13989** | 3.4 | 120.4 | 57.7 |
| **36324** | 8.3 | 62.4 | 38.8 |
| **3214** | 13.7 | 80.8 | 69.2 |
| **877** | 11.1 | 107.6 | 51.8 |
| **27401** | 14.5 | 91.0 | 63.3 |
| **33272** | 16.9 | 81.6 | NaN |
| **19766** | 13.5 | 59.4 | 98.4 |
| **13918** | 13.9 | 77.0 | 33.4 |
| **2620** | 13.5 | 83.3 | 36.8 |

```
1 # mean interpolation for numeric columns
2
3 # Calculate the mean for each column
4 mean_values = df[['Study Hours per Week', 'Attendance Rate', 'Previous Grades']].mean()
5
6 # Fill missing values with the calculated means
7 df.fillna(mean_values, inplace=True)
```

```
1 df [['Study Hours per Week', 'Attendance Rate', 'Previous Grades']].isna().sum()
```

| | 0 |
|---|---|
| **Study Hours per Week** | 0 |
| **Attendance Rate** | 0 |
| **Previous Grades** | 0 |

dtype: int64

```
1 df [['Participation in Extracurricular Activities', 'Parent Education Level', 'Passed']].sample(25)
```

| | Participation in Extracurricular Activities | Parent Education Level | Passed |
|---|---|---|---|
| 23912 | No | Associate | Yes |
| 39156 | Yes | Master | Yes |
| 27364 | NaN | Doctorate | No |
| 18109 | NaN | High School | No |
| 1119 | No | Bachelor | No |
| 231 | Yes | High School | No |
| 10062 | Yes | Doctorate | Yes |
| 8388 | Yes | Bachelor | Yes |
| 4188 | Yes | Associate | Yes |
| 34485 | No | Master | No |
| 23753 | No | Master | No |
| 25530 | NaN | High School | No |
| 28458 | No | High School | Yes |
| 9388 | Yes | High School | No |
| 7863 | Yes | Bachelor | No |
| 8835 | No | Bachelor | No |
| 15635 | No | Doctorate | Yes |
| 37991 | Yes | NaN | Yes |
| 22483 | No | Associate | Yes |
| 35464 | Yes | Doctorate | No |
| 11579 | Yes | Bachelor | NaN |
| 2506 | Yes | Bachelor | No |
| 3186 | No | NaN | Yes |
| 39602 | No | High School | No |
| 30069 | No | NaN | Yes |

```python
1 # mode interpolation for categorical columns
2
3 # 1. Select categorical columns
4 categorical_columns = df.select_dtypes(include="object").columns
5
6 # 2. Iterate over categorical columns
7 for column in categorical_columns:
8     # 3. Calculate the mode for the current column
9     mode_value = df[column].mode()[0]
10
11     # 4. Fill missing values with the mode
12     df[column] = df[column].fillna(mode_value)
```

```python
1 df [['Participation in Extracurricular Activities', 'Parent Education Level', 'Passed']].isna().sum()
```

| | 0 |
|---|---|
| **Participation in Extracurricular Activities** | 0 |
| **Parent Education Level** | 0 |
| **Passed** | 0 |

**dtype:** int64

```python
1 df.columns
```

```
Index(['Student ID', 'Study Hours per Week', 'Attendance Rate',
       'Previous Grades', 'Participation in Extracurricular Activities',
       'Parent Education Level', 'Passed'],
      dtype='object')
```

## Handling Outliers & Skewness

## col 1: 'Student ID'

```python
1  for col in ['Study Hours per Week', 'Attendance Rate', 'Previous Grades']:
2      Q1 = df[col].quantile(0.25)
3      Q3 = df[col].quantile(0.75)
4      IQR = Q3 - Q1
5
6      # Define lower and upper bounds
7      lower_bound = Q1 - 1.5 * IQR
8      upper_bound = Q3 + 1.5 * IQR
9
10     # Identify outliers
11     outliers = df[(df[col] < lower_bound) | (df[col] > upper_bound)]
12     print(f"Outliers in {col}:")
13     print(outliers)
```

```
39516     No
39557    Yes
39571     No
39576     No
39678     No
39732    Yes
39797     No
39851    Yes
39890    Yes
```

```python
1 # Create a box plot for each numeric column
2 plt.figure(figsize=(12, 6))
3
4 plt.subplot(1, 3, 1)
5 sns.boxplot(x=df['Study Hours per Week'])
6 plt.title("Outliers in 'Study Hours per Week' Column")
7
8 plt.tight_layout()
9 plt.show()
```
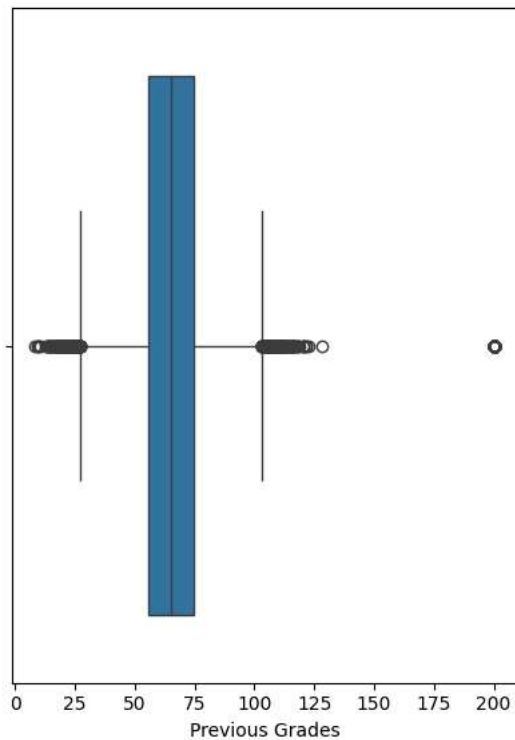


Outliers in 'Study Hours per Week' Column

```python
1 # Create a box plot for each numeric column
2 plt.figure(figsize=(12, 6))
3
4 plt.subplot(1, 3, 1)
5 sns.boxplot(x=df['Attendance Rate'])
6 plt.title("Outliers in 'Attendance Rate' Column")
7
8 plt.tight_layout()
9 plt.show()
```

Outliers in 'Attendance Rate' Column



```
1  # Create a box plot for each numeric column
2  plt.figure(figsize=(12, 6))
3
4  plt.subplot(1, 3, 1)
5  sns.boxplot(x=df['Previous Grades'])
6  plt.title("Outliers in 'Previous Grades' Column")
7
8  plt.tight_layout()
9  plt.show()
```

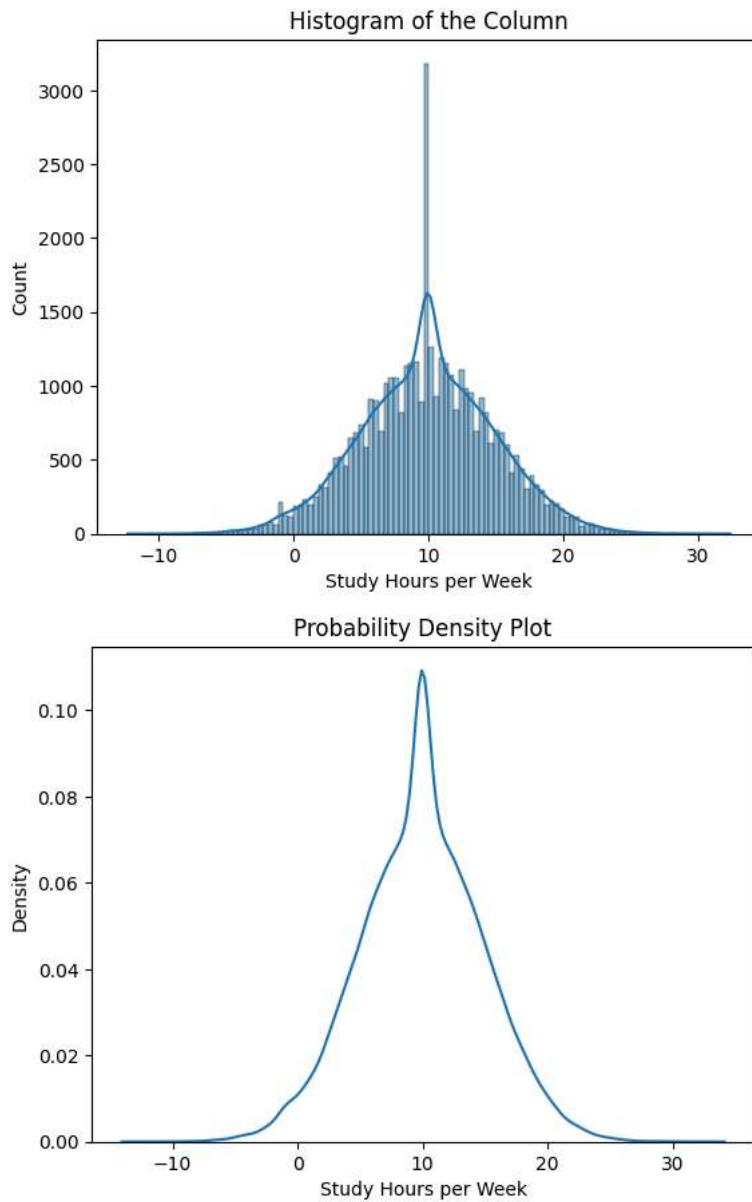Outliers in 'Previous Grades' Column



```
1 skewness = df[['Study Hours per Week', 'Attendance Rate', 'Previous Grades']].skew()
2 print("Skewness in Columns -\n", skewness)
```
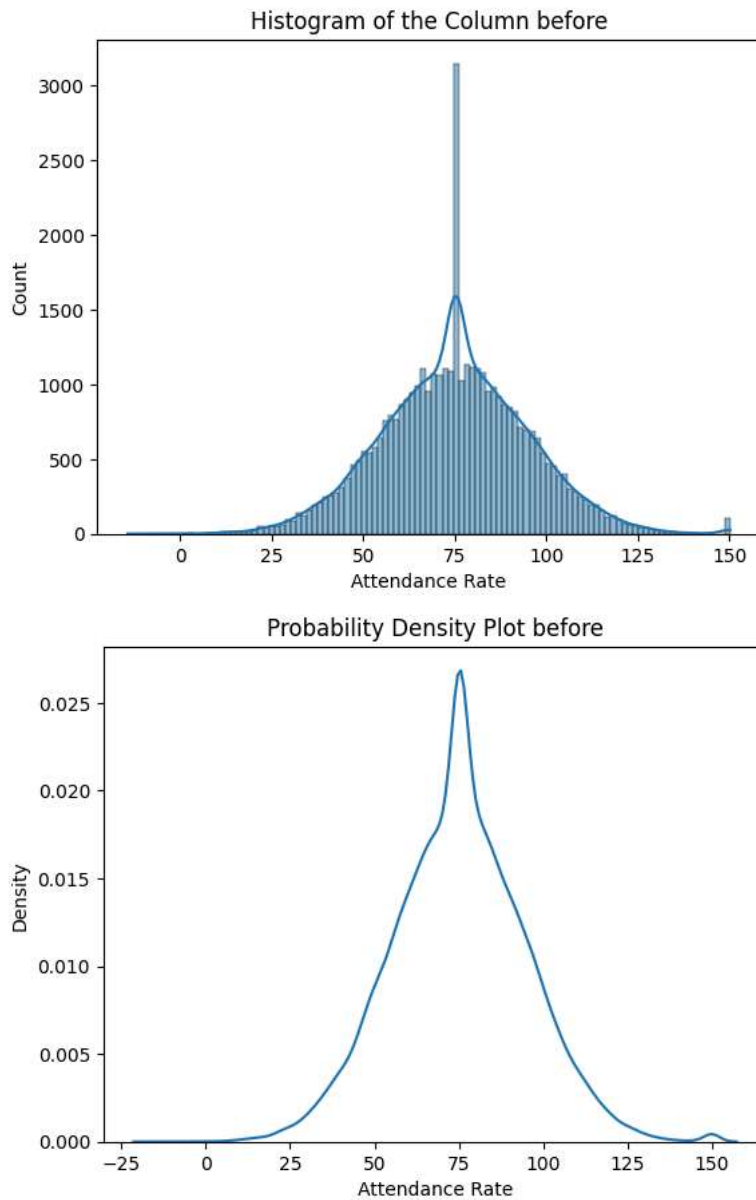
```
Skewness in Columns -
 Study Hours per Week   -0.011992
 Attendance Rate         0.091591
 Previous Grades         1.409301
 dtype: float64
```
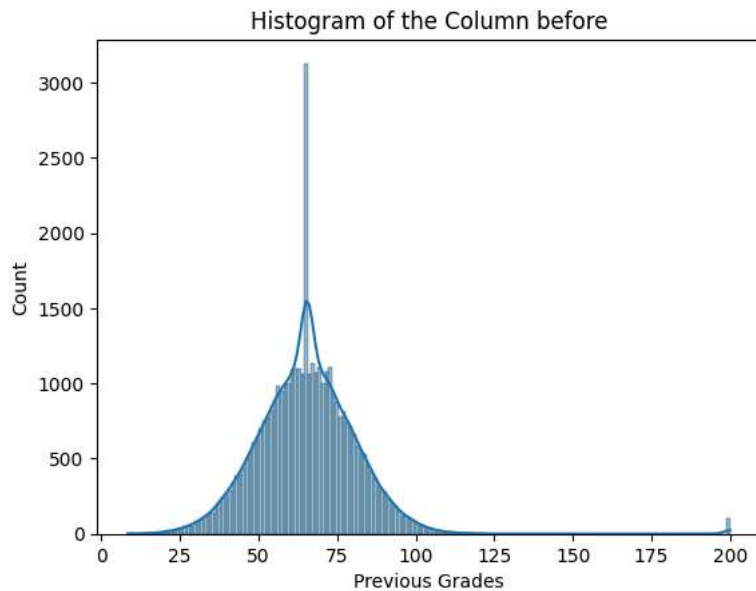
```
1 # Histogram
2 sns.histplot(df['Study Hours per Week'], kde=True)
3 plt.title('Histogram of the Column')
4 plt.show()
5
6 # Probability Density Plot
7 sns.kdeplot(df['Study Hours per Week'])
8 plt.title('Probability Density Plot')
9 plt.show()
```

## Histogram of the Column



## Probability Density Plot



```
1 # Histogram
2 sns.histplot(df['Attendance Rate'], kde=True)
3 plt.title('Histogram of the Column before')
4 plt.show()
5
6 # Probability Density Plot
7 sns.kdeplot(df['Attendance Rate'])
8 plt.title('Probability Density Plot before')
9 plt.show()
```

Histogram of the Column before



Probability Density Plot before

```
1 # Histogram
2 sns.histplot(df['Previous Grades'], kde=True)
3 plt.title('Histogram of the Column before')
4 plt.show()
5
6 # Probability Density Plot
7 sns.kdeplot(df['Previous Grades'])
8 plt.title('Probability Density Plot before')
9 plt.show()
```
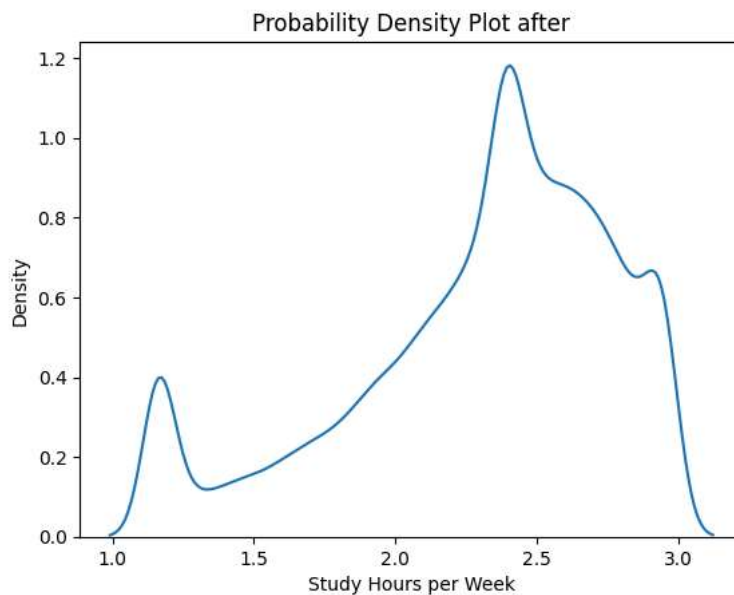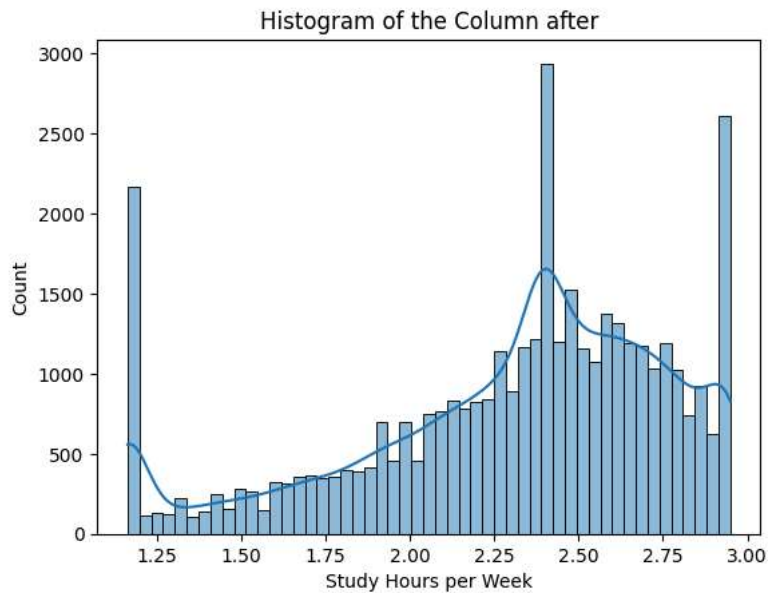
## Histogram of the Column before



## Probability Density Plot before



```
1  # Handling Negative Skewness and Outliers
2
3  # Handle negative values (if any):
4  df['Study Hours per Week'] = df['Study Hours per Week'].abs()
5
6  # Log transformation
7  df['Study Hours per Week'] = np.log1p(df['Study Hours per Week'])
8
9  # Define thresholds (adjust as needed)
10 lower_bound = df['Study Hours per Week'].quantile(0.05)
11 upper_bound = df['Study Hours per Week'].quantile(0.95)
12
13 # Cap outliers
14 df['Study Hours per Week'] = np.clip(df['Study Hours per Week'], lower_bound, upper_bound)
```
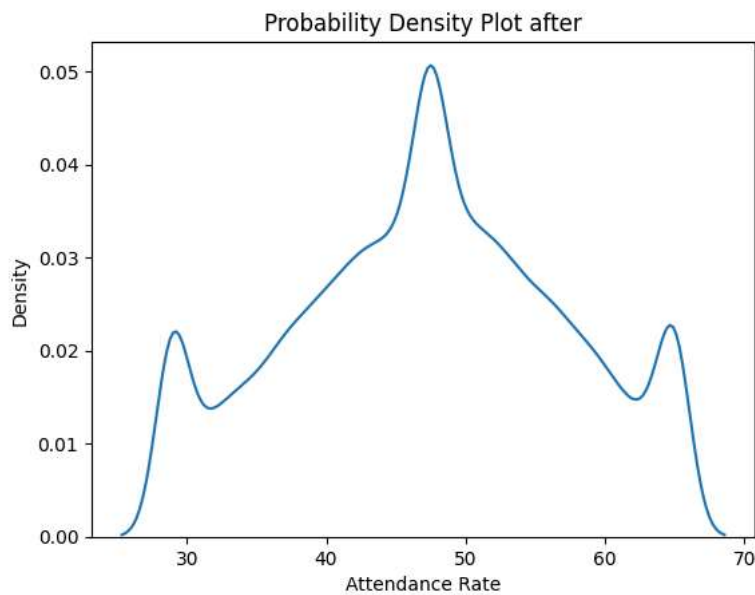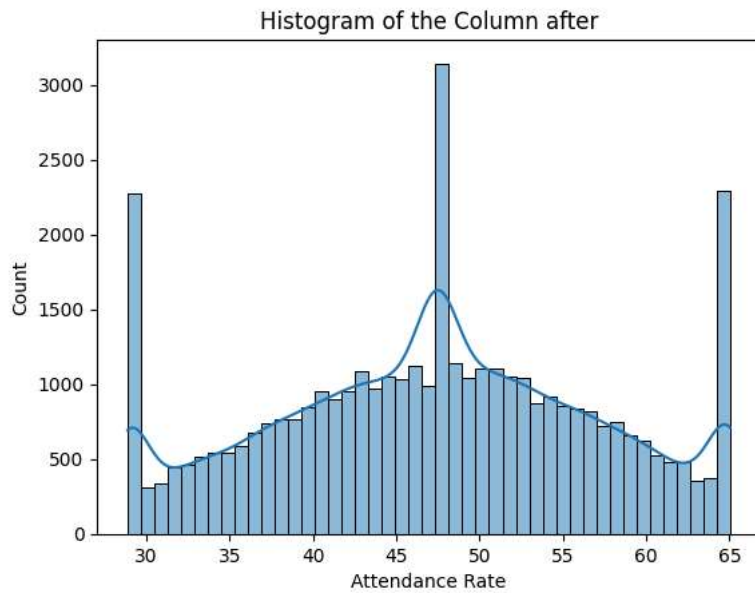
```
1 # Winsorization
2 lower_bound = df['Attendance Rate'].quantile(0.05)
3 upper_bound = df['Attendance Rate'].quantile(0.95)
4 df['Attendance Rate'] = np.clip(df['Attendance Rate'], lower_bound, upper_bound)
5
6 # Or, Box-Cox Transformation
7 from scipy.stats import boxcox
8 df['Attendance Rate'], _ = boxcox(df['Attendance Rate'] + 1)  # Add 1 to handle zero values
```

```
1 # Winsorization
2 lower_bound = df['Previous Grades'].quantile(0.05)
3 upper_bound = df['Previous Grades'].quantile(0.95)
4 df['Previous Grades'] = np.clip(df['Previous Grades'], lower_bound, upper_bound)
5
6 # Or, Box-Cox Transformation
7 from scipy.stats import boxcox
8 df['Previous Grades'], _ = boxcox(df['Previous Grades'] + 1)  # Add 1 to handle zero values
```
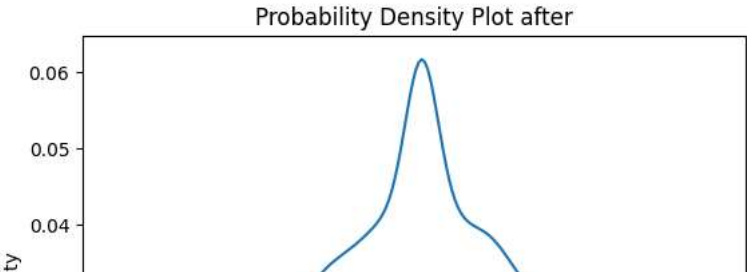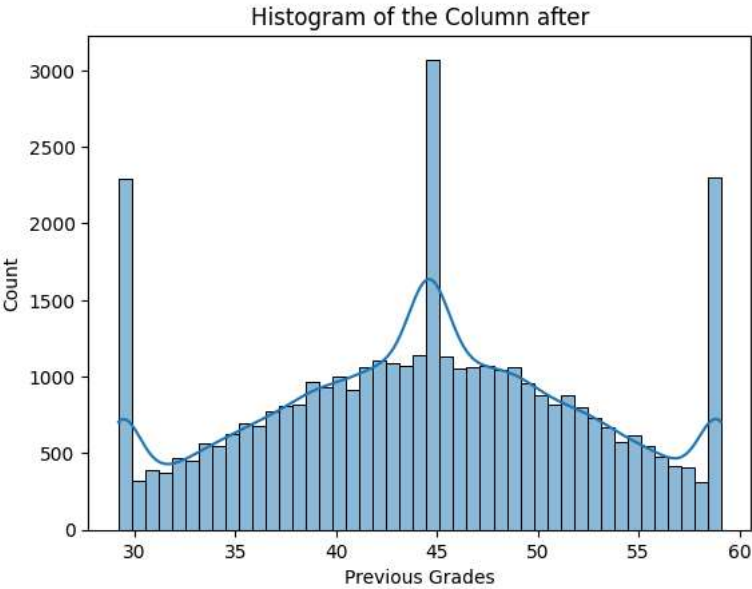
```
1 # Histogram
2 sns.histplot(df['Study Hours per Week'], kde=True)
3 plt.title('Histogram of the Column after')
4 plt.show()
5
6 # Probability Density Plot
7 sns.kdeplot(df['Study Hours per Week'])
8 plt.title('Probability Density Plot after')
9 plt.show()
```

## Histogram of the Column after



## Probability Density Plot after



```
1 # Histogram
2 sns.histplot(df['Attendance Rate'], kde=True)
3 plt.title('Histogram of the Column after')
4 plt.show()
5
6 # Probability Density Plot
7 sns.kdeplot(df['Attendance Rate'])
8 plt.title('Probability Density Plot after')
9 plt.show()
```

## Histogram of the Column after



## Probability Density Plot after



```
1 # Histogram
2 sns.histplot(df['Previous Grades'], kde=True)
3 plt.title('Histogram of the Column after')
4 plt.show()
5
6 # Probability Density Plot
7 sns.kdeplot(df['Previous Grades'])
8 plt.title('Probability Density Plot after')
9 plt.show()
```

## Histogram of the Column after



## Probability Density Plot after



```
1 df.sample(50)
```