```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
```

```
1 # Load the dataset
2 df = pd.read_csv("/content/data.csv")
```

```
1
2 # Step 1: Data Preprocessing
3 df.drop(columns=['Roll'], inplace=True)   # Drop unused column
4 df.drop_duplicates(inplace=True)          # Remove duplicate rows
5 df.dropna(inplace=True)                   # Remove rows with null values
```

```
1 # Step 2: Calculate Cumulative GPA
2 semester_cols = ['1st', '2nd', '3rd', '4th', '5th']
3 df["CGPA"] = df[semester_cols].mean(axis=1).round(2)
4
```

```
1 # Step 3: Display cleaned dataset
2 print("Cleaned Data (First 5 Rows):")
3 print(df.head())
```

```
Cleaned Data (First 5 Rows):
    1st   2nd   3rd   4th   5th  College Code  Gender  Roll no.  Subject Code  \
0  8.11  7.68  7.11  7.43  8.18           115  Female   17020.0            16
1  6.48  5.90  4.15  4.29  4.96           115    Male   17021.0            16
2  8.41  8.24  7.52  8.25  7.75           115  Female   17022.0            16
3  7.33  6.83  6.33  6.79  6.89           115    Male   17023.0            16
4  7.89  7.34  7.22  7.32  7.46           115    Male   17024.0            16

   CGPA
0  7.70
1  5.16
2  8.03
3  6.83
4  7.45
```

```
1
2 # Step 4: Basic Summary Stats
3 print("\nSummary Statistics:")
4 print(df[semester_cols + ['CGPA']].describe())
```
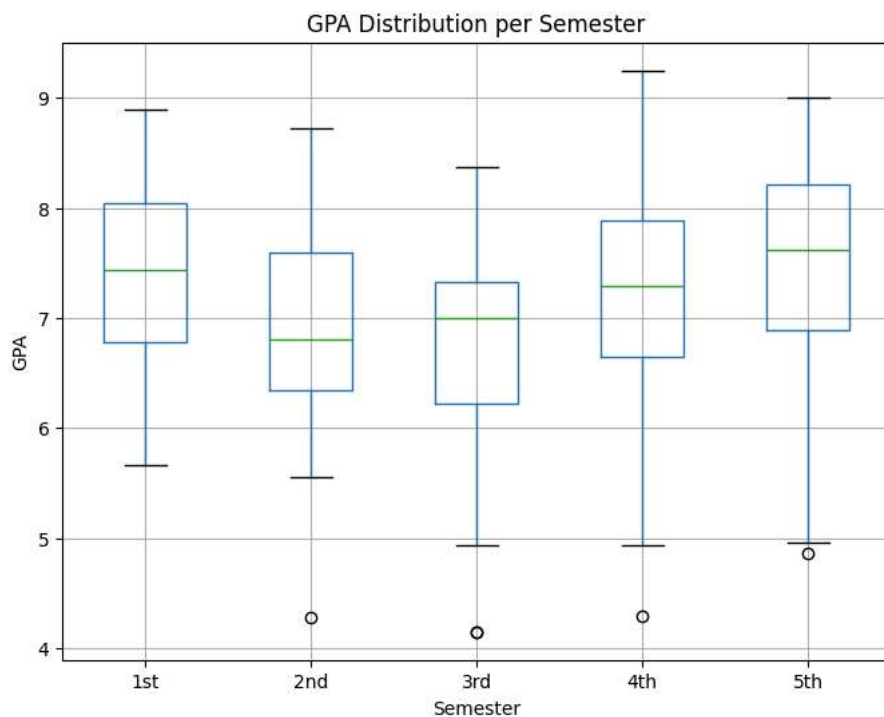
```
Summary Statistics:
             1st        2nd        3rd        4th        5th       CGPA
count  46.000000  46.000000  46.000000  46.000000  46.000000  46.000000
mean    7.397609   6.930217   6.703043   7.237826   7.527609   7.159565
std     0.798391   0.910425   0.917324   1.057981   0.967963   0.856102
min     5.670000   4.280000   4.150000   4.290000   4.860000   5.120000
25%     6.787500   6.350000   6.217500   6.650000   6.890000   6.567500
50%     7.440000   6.810000   7.000000   7.290000   7.625000   7.265000
75%     8.040000   7.590000   7.322500   7.890000   8.210000   7.817500
max     8.890000   8.720000   8.370000   9.250000   9.000000   8.560000
```

```
1
2 # Step 5: Plot GPA Trend for a Sample Student
3 sample_student = df.iloc[0]
4 plt.plot(semester_cols, sample_student[semester_cols], marker='o', color='blue')
5 plt.title(f"GPA Trend for Roll No: {int(sample_student['Roll no.'])}")
6 plt.xlabel("Semester")
7 plt.ylabel("GPA")
8 plt.ylim(0, 10)
9 plt.grid(True)
10 plt.show()
```
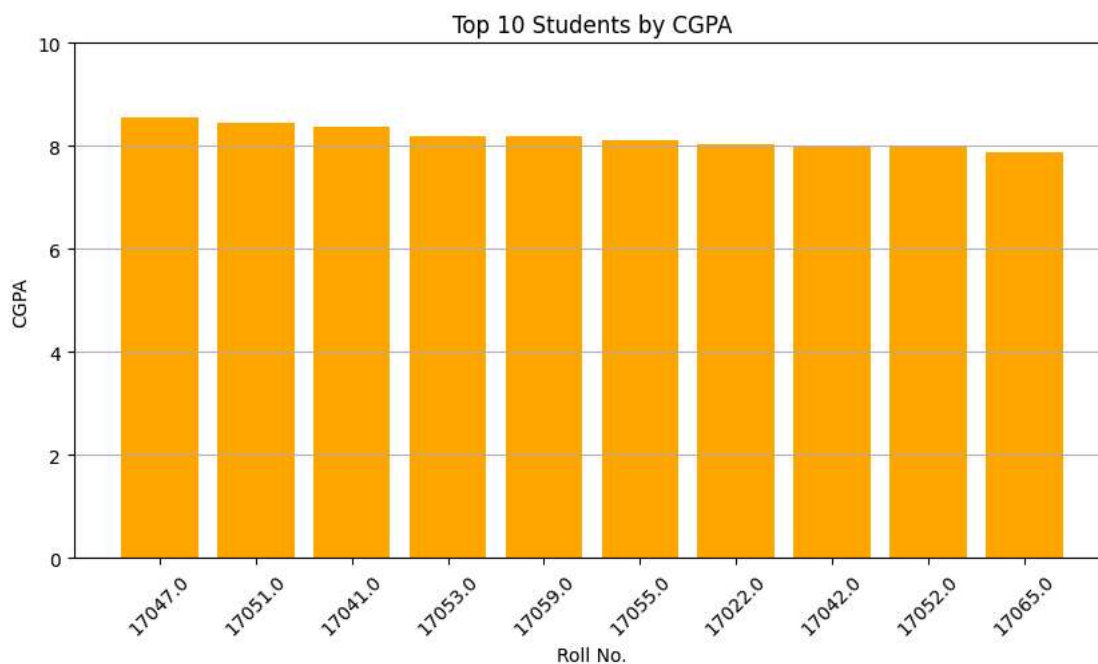
GPA Trend for Roll No: 17020

```
1 #box plot
2 plt.figure(figsize=(8, 6))
3 df[semester_cols].boxplot()
4 plt.title("GPA Distribution per Semester")
5 plt.xlabel("Semester")
6 plt.ylabel("GPA")
7 plt.grid(True)
8 plt.show()
```


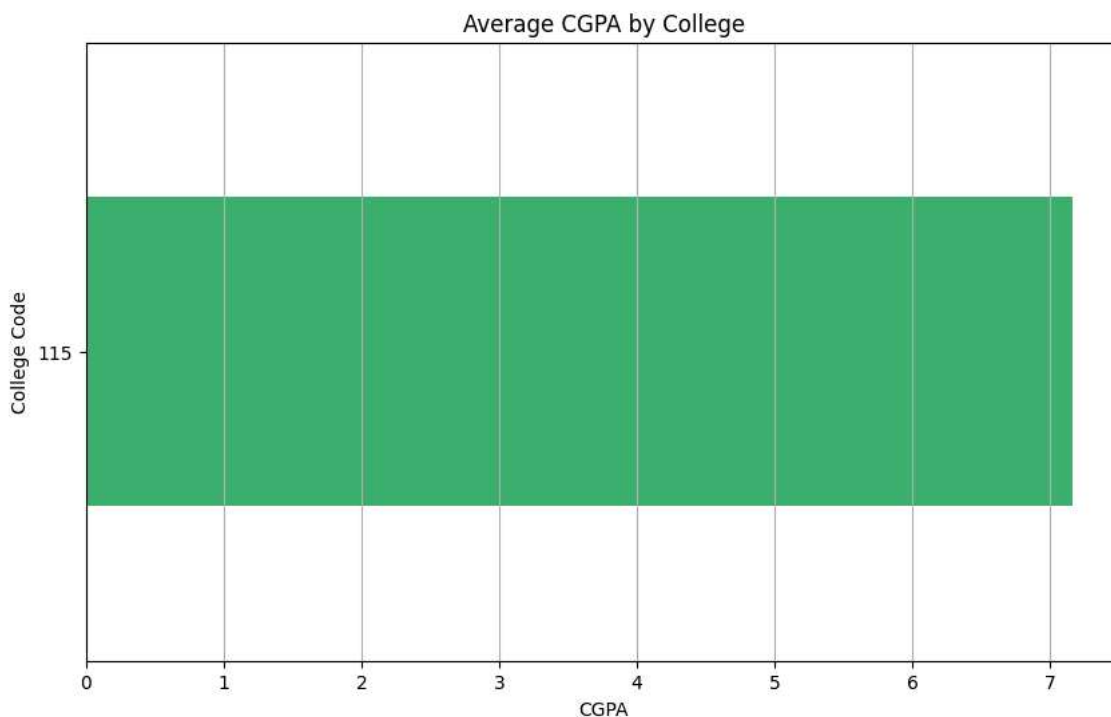
GPA Distribution per Semester

```
 1 # Top 10 Students by CGPA
 2 top_students = df.sort_values(by="CGPA", ascending=False).head(10)
 3
 4 plt.figure(figsize=(10, 5))
 5 plt.bar(top_students["Roll no."].astype(str), top_students["CGPA"], color='orange')
 6 plt.title("Top 10 Students by CGPA")
 7 plt.xlabel("Roll No.")
 8 plt.ylabel("CGPA")
 9 plt.xticks(rotation=45)
10 plt.ylim(0, 10)
11 plt.grid(axis='y')
```

```
12 plt.show()
13
```



Top 10 Students by CGPA

```
1 #College-wise Average CGPA
2 college_avg = df.groupby("College Code")["CGPA"].mean().sort_values()
3
4 plt.figure(figsize=(10, 6))
5 college_avg.plot(kind='barh', color='mediumseagreen')
6 plt.title("Average CGPA by College")
7 plt.xlabel("CGPA")
8 plt.ylabel("College Code")
9 plt.grid(axis='x')
10 plt.show()
11
```
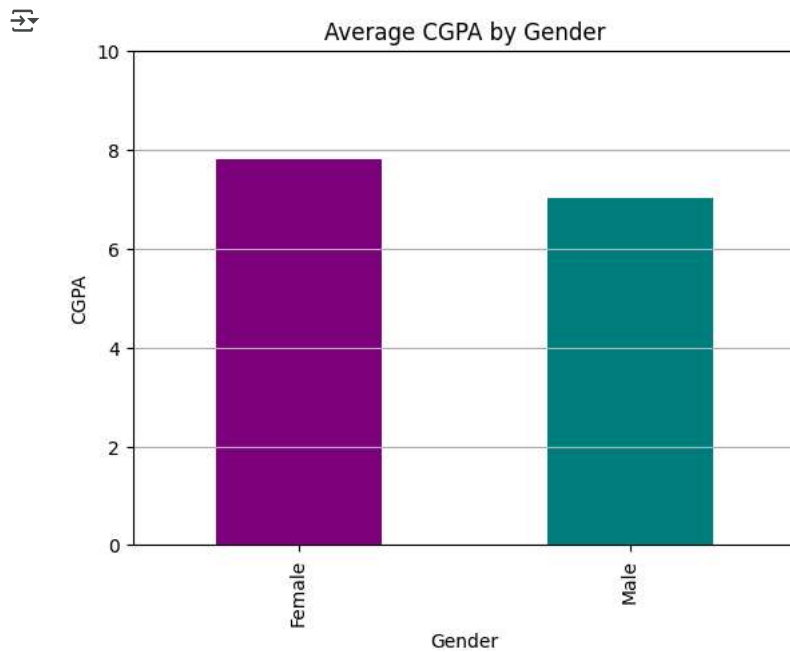


Average CGPA by College

```
1 # Step 6: Gender-wise CGPA Comparison
2 gender_avg = df.groupby("Gender")["CGPA"].mean()
3
```

```
 4 gender_avg.plot(kind='bar', color=['purple', 'teal'])
 5 plt.title("Average CGPA by Gender")
 6 plt.ylabel("CGPA")
 7 plt.xlabel("Gender")
 8 plt.ylim(0, 10)
 9 plt.grid(axis='y')
10 plt.show()
```
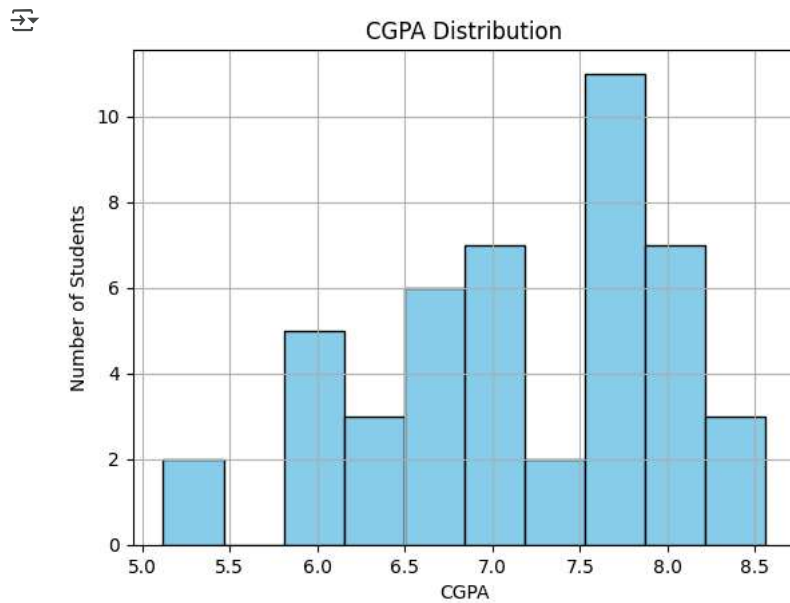

Average CGPA by Gender

```
1 # Step 7: Histogram of CGPA Distribution
2 plt.hist(df["CGPA"], bins=10, color='skyblue', edgecolor='black')
3 plt.title("CGPA Distribution")
4 plt.xlabel("CGPA")
5 plt.ylabel("Number of Students")
6 plt.grid(True)
7 plt.show()
```


CGPA Distribution

```
1 # Assign Pass/Fail (CGPA >= 5)
2 df['Status'] = np.where(df['CGPA'] >= 5, 'Pass', 'Fail')
3 print(df['Status'].value_counts())
```

```
Status
Pass    46
Name: count, dtype: int64
```
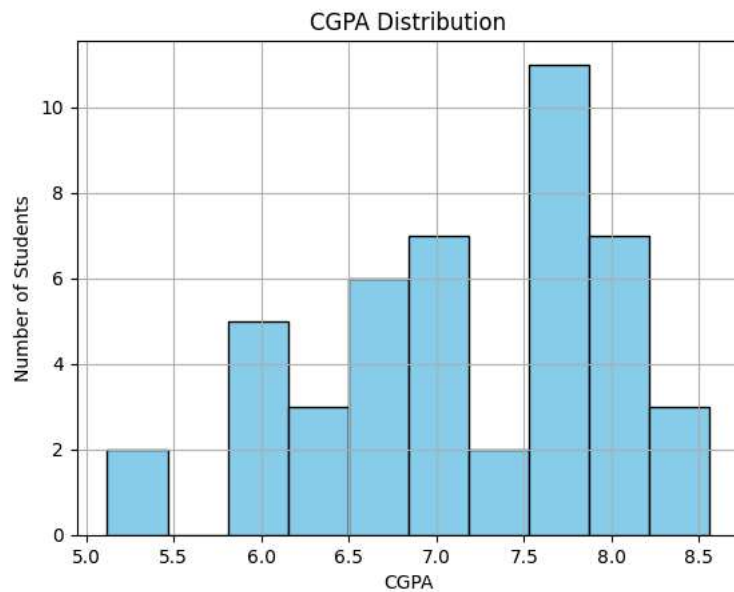
```python
1 # Show summary
2 print("\nSummary:")
3 print("Total students:", len(df))
4 print("Passed:", (df['Status'] == 'Pass').sum())
5 print("Failed:", (df['Status'] == 'Fail').sum())
6 print("Average CGPA:", df['CGPA'].mean().round(2))
```

```
Summary:
Total students: 46
Passed: 46
Failed: 0
Average CGPA: 7.16
```

```python
1 # CGPA distribution plot
2 plt.hist(df['CGPA'], bins=10, color='skyblue', edgecolor='black')
3 plt.title("CGPA Distribution")
4 plt.xlabel("CGPA")
5 plt.ylabel("Number of Students")
6 plt.grid(True)
7 plt.show()
```



```python
1 # Save final result to CSV
2 df.to_csv("/content/data.csv", index=False)
```