

FIT1043 A1 KAN ROU YI 34661093

```

In [173]: import pandas as pd
          from matplotlib import pyplot as plt
          import matplotlib.pyplot as plt

In [174]: df = pd.read_csv('Student_List.csv')
df

Out[174]:
   StudentID  Age  ParentalEducation  StudyTimeWeekly  Absences  Tutoring  ParentalSupport  Extracurricular  Sports  Music  Volunteering  GPA  GradeClass
0         1540   18      Some College      10.318918      5      No      2      No  Yes  No  No  No  2.65594  C
1         2939   16      Bachelor's      6.517803      2     Yes      2      No  Yes  No  No  No  3.474562  B
2         2877   15      Some College      0.815700      1     No      1      No  No  Yes  Yes  No  2.808878  C
3         1628   16      High School      6.504335      8     No      3      No  No  No  No  No  2.195546  D
4         2052   15      Some College      2.516047     14     Yes      3      Yes  No  No  No  No  2.253871  D
...
1495      1944   15      No Education     10.586878     12     No      1      No  No  No  No  Yes  1.621912  F
1496      1506   16      Some College      3.278634      4     Yes      4      Yes  No  No  No  No  3.244862  B
1497      3208   18      Some College      7.680100      0     Yes      3      No  No  No  No  No  3.040730  F
1498      3237   17      Some College     15.078754     24     No      4      No  No  Yes  No  No  1.245091  F
1499      1149   15      Higher          1.360205     22     No      1      Yes  Yes  No  Yes  No  1.007226  F

1500 rows × 13 columns

```

A1. Exploring the dataset

1. How many students are there in this dataset?

```

In [172]: num_students = df.shape[0]
          print(f"Total number of students: {num_students}")
          Total number of students: 1500

2. What is the age range of students in this dataset?

In [179]: df.describe()

Out[179]:
   StudentID      Age  StudyTimeWeekly  Absences  ParentalSupport      GPA
count  1500.000000  1500.000000      1500.000000      1500.000000  1500.000000
mean    2209.142667   16.464000      9.625516   14.540000      2.131333   1.901045
std     693.068389   1.124372      5.718608    8.459667   1.116654   0.923638
min     1002.000000   15.000000      0.001057    0.000000    0.000000   0.000000
25%    1599.750000   15.000000      4.717155    7.000000    0.000000   1.171647
50%    2211.500000   16.000000      9.362838   15.000000    2.000000   1.681249
75%    2604.750000   17.000000     14.374170   22.000000    3.000000   2.623266
max     3089.000000   18.000000     19.876084   29.000000    4.000000   4.000000

```

```

In [169]: min = df['Age'].min()
          min = df['Age'].min()
          range = max - min
          print('Range:', range)
          Range: 3

3. Identify and discuss the data types of each column in the dataset.

In [172]: df.dtypes

Out[172]:
StudentID      int64
Age            int64
ParentalEducation  object
StudyTimeWeekly  float64
Absences        int64
Tutoring        object
ParentalSupport  int64
Extracurricular  object
Sports          object
Music          object
Volunteering    object
GPA            float64
GradeClass      object
dtype: object

```

4. GradeClass corresponds to different categories of grades, where, for example, 'A' corresponds to the highest grade (GPA ≥ 3.5) and 'F' corresponds to the lowest grade (GPA < 2.0). Calculate the percentage of students in each grade category.

```

In [174]: grade_counts = df['GradeClass'].value_counts()
          total_students = len(df)
          grade_percentages = (grade_counts / total_students) * 100
          grade_percentages

Out[174]:
GradeClass
F    8.000000
D    6.333333
C    16.266667
B    12.200000
A     4.800000
Name: count, dtype: float64

```

5. Create a pie chart to show the proportion of students in each GradeClass and discuss your observations.

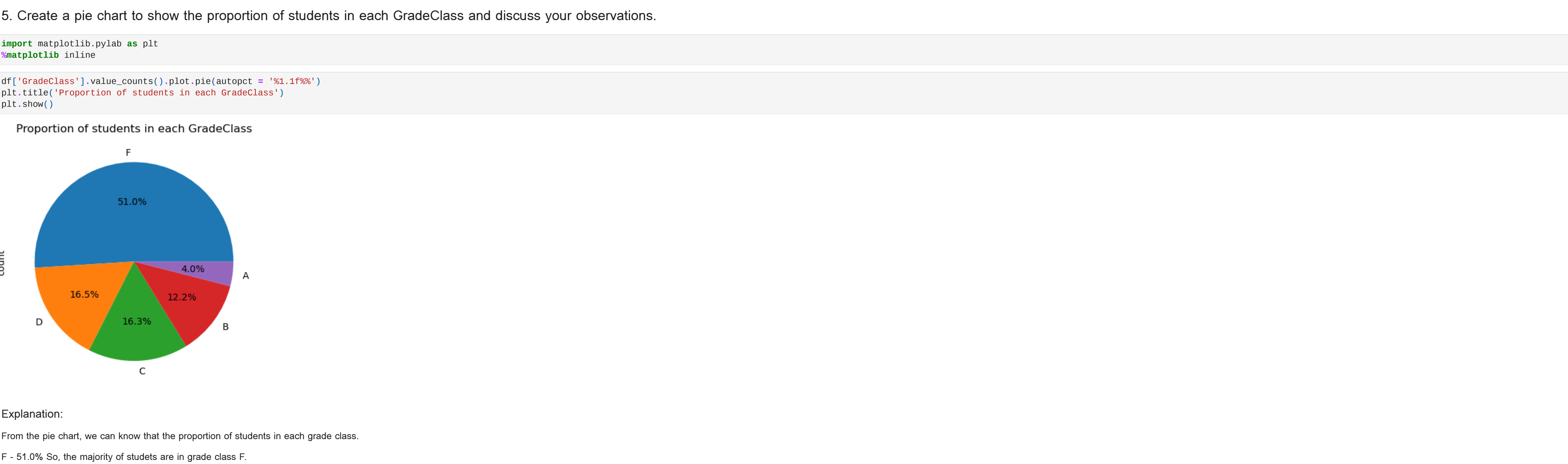
```

In [178]: import matplotlib.pyplot as plt
          %matplotlib inline

In [177]: df['GradeClass'].value_counts().plot.pie(autopct = '%1.1f%%')
          plt.title('Proportion of students in each GradeClass')
          plt.show()

Proportion of students in each GradeClass

```



A2. Exploring Parental Education

1. Determine the frequency of different levels of parental education.

- How many parents have the highest education level (denoted by Higher) in this dataset?
- How many have no education?
- What is the most common level of parental education?

```

In [176]: from IPython.display import HTML
          df = pd.read_csv('Student_List.csv')
          df.ParentalEducation

Out[176]:
   StudentID  Age  ParentalEducation  StudyTimeWeekly  Absences  Tutoring  ParentalSupport  Extracurricular  Sports  Music  Volunteering  GPA  GradeClass
0         1540   18      Some College      10.318918      5     No      2      No  Yes  No  No  No  2.65594  C
1         2939   16      Bachelor's      6.517803      2     Yes      2      No  Yes  No  No  No  3.474562  B
2         2877   15      Some College      0.815700      1     No      1      No  No  Yes  Yes  No  2.808878  C
3         1628   16      High School      6.504335      8     No      3      No  No  No  No  No  2.195546  D
4         2052   15      Some College      2.516047     14     Yes      3      Yes  No  No  No  No  2.253871  D
...
1495      1944   15      No Education     10.586878     12     No      1      No  No  No  No  Yes  1.621912  F
1496      1506   16      Some College      3.278634      4     Yes      4      Yes  No  No  No  No  3.244862  B
1497      3208   18      Some College      7.680100      0     Yes      3      No  No  No  No  No  3.040730  F
1498      3237   17      Some College     15.078754     24     No      4      No  No  Yes  No  No  1.245091  F
1499      1149   15      Higher          1.360205     22     No      1      Yes  Yes  No  Yes  No  1.007226  F

1500 rows × 13 columns

```

```

In [170]: # Create a function to count
          highest_education_level = df['ParentalEducation'].value_counts().get('Higher')
          print('Highest education level:', highest_education_level)
          Highest education level: 77

In [181]: no_education = df['ParentalEducation'].value_counts().get('No Education')
          print('No Education:', no_education)
          No Education: None

In [174]: # Create a function
          most_common_education_level = df['ParentalEducation'].mode()
          print('Most common education level:', most_common_education_level)
          Most common education level: 0      Some College
          Name: ParentalEducation, dtype: object

```

2. Replace the values in the 'ParentalEducation' column according to following table:

```

In [184]: # Create a dictionary
          column = {
              "No Education": 0,
              "High School": 1,
              "Some College": 2,
              "Bachelor's": 3,
              "Higher": 4
          }

          # Create a function
          df['ParentalEducation'] = df['ParentalEducation'].replace(column)
          df

```

```

Out[184]:
   StudentID  Age  ParentalEducation  StudyTimeWeekly  Absences  Tutoring  ParentalSupport  Extracurricular  Sports  Music  Volunteering  GPA  GradeClass
0         1540   18      2              10.318918      5     No      2      No  Yes  No  No  No  2.65594  C
1         2939   16      3              6.517803      2     Yes      2      No  Yes  No  No  No  3.474562  B
2         2877   15      2              0.815700      1     No      1      No  No  Yes  Yes  No  2.808878  C
3         1628   16      1              6.504335      8     No      3      No  No  No  No  No  2.195546  D
4         2052   15      2              2.516047     14     Yes      3      Yes  No  No  No  No  2.253871  D
...
1495      1944   15      0              10.586878     12     No      1      No  No  No  No  Yes  1.621912  F
1496      1506   16      2              3.278634      4     Yes      4      Yes  No  No  No  No  3.244862  B
1497      3208   18      2              7.680100      0     Yes      3      No  No  No  No  No  3.040730  F
1498      3237   17      2             15.078754     24     No      4      No  No  Yes  No  No  1.245091  F
1499      1149   15      4              1.360205     22     No      1      Yes  Yes  No  Yes  No  1.007226  F

1500 rows × 13 columns

```

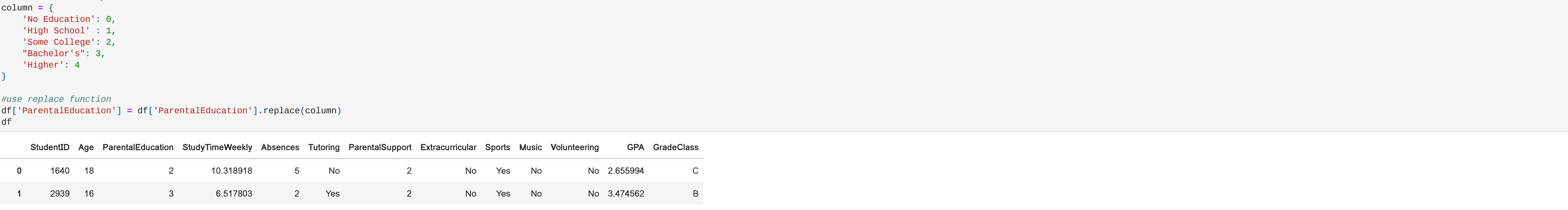
3. Use a boxplot to visualise the distribution of GPA for each level of ParentalEducation.

```

In [178]: df.boxplot(column = 'GPA', by = 'ParentalEducation')
          plt.xlabel('Parental Education')
          plt.ylabel('Student's GPA')
          plt.title('Relationship between parental education and student's GPA')
          plt.show()

Boxplot grouped by ParentalEducation
Relationship between parental education and student's GPA

```



A3. GPA distribution and Correlation Analysis

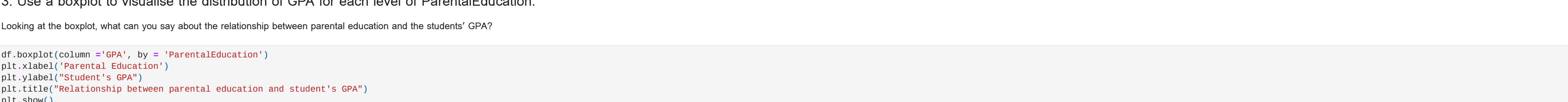
1. Create a histogram to visualise the distribution of GPA. Summarise your observations.

```

In [143]: df.GPA.hist(bins=5, color = 'turquoise')
          plt.xlabel('GPA')
          plt.ylabel('Number of Students')
          plt.title('Distribution of GPA')
          plt.show()

Distribution of GPA

```



2. Does a higher weekly study time correlate with better GPA? Use a scatter plot to visualise this relationship and calculate the correlation coefficient.

```

In [146]: plt.scatter(df['StudyTimeWeekly'], df['GPA'], color = 'lightblue')
          plt.xlabel('Study Time Weekly')
          plt.ylabel('GPA')
          plt.title('Relationship between Weekly Study Time and GPA')
          plt.show()

Relationship between Weekly Study Time and GPA

```



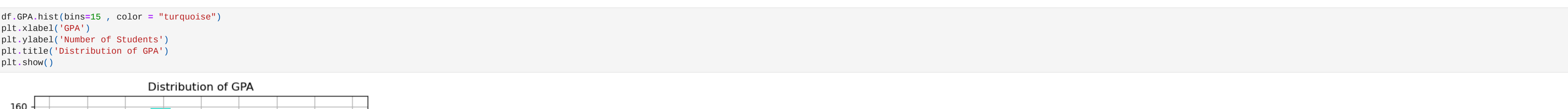
3. How about any correlation between absences and GPA? Use a scatter plot to visualise this relationship and calculate the correlation coefficient.

```

In [149]: plt.scatter(df['Absences'], df['GPA'], color = 'darkwood')
          plt.xlabel('Absences')
          plt.ylabel('GPA')
          plt.title('Relationship between Absences and GPA')
          plt.show()

Relationship between Absences and GPA

```



A4.Extracurricular Activities

1. Select a list of students who are involved in Sports, Music, Volunteering and Extracurricular activities. How many students are there? Assume this is called group A

```

In [167]: # Use 'A' to make sure all the criteria is selected
          # Create a function to select students involved in the activity
          group_A = df[(df['Sports'] == 'Yes') & (df['Music'] == 'Yes') & (df['Volunteering'] == 'Yes') & (df['Extracurricular'] == 'Yes')]
          # Create a function to count the number of students in group A
          num_students = len(group_A)
          print('Number of students in group A:', num_students)
          Number of students in group A: 6

```

2. Now Select students who did not involve in any of the above activities. How many students are there? Assume this is called group B.

```

In [168]: # Use 'B' to make sure the student does not involve in the activity
          group_B = df[(df['Sports'] == 'No') & (df['Music'] == 'No') & (df['Volunteering'] == 'No') & (df['Extracurricular'] == 'No')]
          print('Number of students in group B:', num_students)
          Number of students in group B: 432

```

3. Compare the mean GPA of students in group A versus group B. What can you say?

```

In [163]: # Create a function
          mean_GPA_A = group_A['GPA'].mean()
          mean_GPA_B = group_B['GPA'].mean()
          print('Mean GPA of students in group A:', mean_GPA_A)
          print('Mean GPA of students in group B:', mean_GPA_B)
          GPA of students in group A: 2.4475262217999996
          GPA of students in group B: 1.7312226865324607

```

Explanation:
   
 The mean of GPA of students in group A is 2.45 while mean of GPA of students in group B is 1.73. This means that students of group A are performing better than students of group B.

A5. Exploring Parental Support and Tutoring

1. Aggregate the data by ParentalSupport and find the mean and median GPA for each group. Also, find the number of students who are 18 years old in each group.

```

In [166]: # Groupby the data frame
          select_GPA = df.groupby('ParentalSupport')
          # Use agg() function
          new_data = df.groupby('ParentalSupport')[['GPA']].agg(['mean', 'median'])
          print(new_data)

ParentalSupport      mean      median
0      1.521602  1.471872
1      1.735855  1.748655
2      1.845914  1.817007
3      2.068174  2.096869
4      2.227639  2.215516

```

```

In [172]: # Create a function
          new_data = df.groupby('ParentalSupport')
          # Create a dictionary to map the old value to new labels
          new_data = {
              'None': 0,
              'Low': 1,
              'Moderate': 2,
              'High': 3,
              'Very High': 4
          }
          new_data = new_data.replace(new_data)
          new_data

Out[172]:
ParentalSupport      mean      median
None      1.521602  1.471872
Low       1.735855  1.748655
Moderate  1.845914  1.817007
High      2.068174  2.096869
Very High 2.227639  2.215516

```

```

In [174]: # Create a function to count the number of students who are 18 in each group
          students_18 = df[(df['Age'] == 18)]
          # Create a function to count the number of students in group B
          students_18_count = students_18.groupby('ParentalSupport').size()
          print(students_18_count)

ParentalSupport      size
0              33
1              88
2             118
3             118
dtype: int64

```

```

In [176]: # Create a function to count the number of students who are 18 in each group
          grouped_data = df.groupby('ParentalSupport').agg({
              'GPA': ['mean', 'median'],
              'Age': lambda x: (x == 18).sum()
          })
          grouped_data.columns = ['mean GPA', 'median GPA', 'Number of 18-Year-Old Students']
          print(grouped_data)

ParentalSupport      mean GPA      median GPA      Number of 18-Year-Old Students
0      1.521602  1.471872              33
1      1.735855  1.748655              88
2      1.845914  1.817007              98
3      2.068174  2.096869             118
4      2.227639  2.215516              95

```

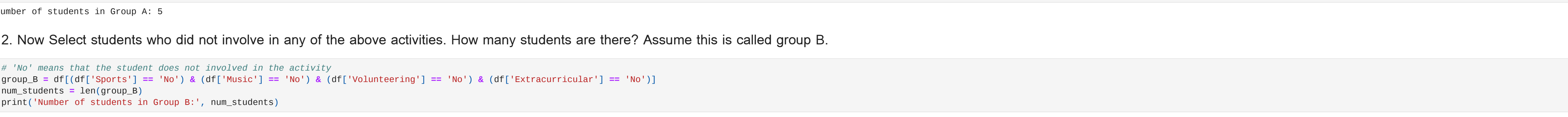
2. Plot a side by side bar chart to visualise the mean and median GPA in each level of parental support and discuss your observations. Your bar chart should illustrate each level of parental support as below.

```

In [179]: # Create a function to plot the mean and median GPA for each level of parental support
          new_data.plot(kind='bar', figsize=(18, 6), width=0.8, color = ('mean', 'lightblue', 'median': 'lightpink'))
          plt.xlabel('Parental Support', fontsize=16)
          plt.ylabel('GPA', fontsize=16)
          plt.title('Analysis of GPA and Age Distribution by Parental Support', fontsize=18)
          plt.xticks(rotation=45)
          plt.legend(['GPA Type', 'Age'], loc='upper left')
          plt.grid(True, axis='y')
          plt.show()

Analysis of GPA and Age Distribution by Parental Support

```



3. Calculate the average GPA for students who receive Tutoring versus those who don't.

```

In [185]: average_gpa = df.groupby('Tutoring')['GPA'].mean()
          average_gpa

Out[185]:
Tutoring
No      1.848887
Yes     2.322588
Name: GPA, dtype: float64

```

```

In [183]: average_gpa.plot(kind='bar', x = 'Tutoring', y = 'GPA', color = ('lightblue', 'red'))
          plt.title('Average GPA for students who receive Tutoring versus those who don't', fontsize=16)
          plt.xlabel('Tutoring', fontsize=16)
          plt.ylabel('GPA', fontsize=16)
          plt.show()

Average GPA for students who receive Tutoring versus those who don't

```

