# Student Performance Data Analysis

Final Project by: Lim Swee Ming
Dataset from Kaggle

# Project Overview

Student Performance Dataset

Dataset Source:
https://www.kaggle.com/larsen0966/student-performance-data-set

649 students participated in this survey and a set of data values featuring 31 columns.

Across left to right, the data values consist of student's parents education and occupation, student's social life and their academic grades.

This project scope aims to reveal specific inequalities of the student's background. Teachers may see them in a better light to guide students who performed poorly and help them succeed in academics.

# Data Cleaning

- Check for missing values

- Creating new columns and assigning new values
  for correlation function at later part

  For instance, the dataset comes with three grade periods.
  Summing up altogether and categorize into low, average
  high score. Other variables from scale 1 – 5 is also
  categorize the same manner.

| family relationship | freetime | go out | Alcohol consumption | Weekday Alcohol consumption |
|---|---|---|---|---|
| 4 | 3 | 4 | 1 | 1 |
| 5 | 3 | 3 | 1 | 1 |
| 4 | 3 | 2 | 2 | 3 |
| 3 | 2 | 2 | 1 | 1 |
| 4 | 3 | 2 | 1 | 2 |
| 5 | 4 | 2 | 1 | 2 |
| 4 | 4 | 4 | 1 | 1 |
| 4 | 1 | 4 | 1 | 1 |
| 4 | 2 | 2 | 1 | 1 |
| 5 | 5 | 1 | 1 | 1 |
| 3 | 3 | 3 | 1 | 2 |

# Identify the top/bottom 20 students based on their grades

```python
1   #ranging the total scores in 3 levels - Low, Average, High Grades
2   #students who scored Lower than 25 is Low
3   #students who scored between 25 to 45 is average
4   #students who scored higher than 45 is high
5
6   def marks(total_score):
7       if(total_score<25):
8           return("low")
9       elif(total_score>=25 and total_score<45):
10          return("average")
11      elif(total_score>=45):
12          return("high")
13  student_data["grades"]=student_data["total_score"].apply(marks)
```

```python
1   counts = student_data["grades"].value_counts()
2   print(counts)
```

```
average    505
high        86
low         58
Name: grades, dtype: int64
```
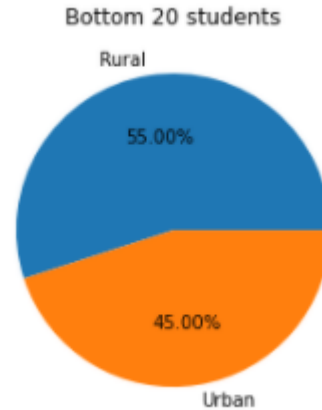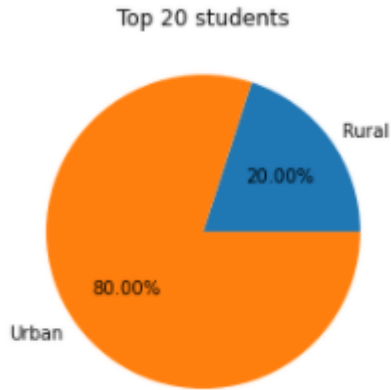


Total Score - Number of Students

# Identify the top/bottom 20 students based on their grades

```
1  top_20 = student_data.sort_values(by = ["total_score"],ascending=False).head(15)
2  top_20
```

|    | sex | age | address | famsize | Pstatus | Medu | Fedu | Mjob     | Fjob     | guardian | ... | G3 | total_score | grades  | absent | family | go_out | alcohol_weekday |
|----|-----|-----|---------|---------|---------|------|------|----------|----------|----------|-----|----|-------------|---------|--------|--------|--------|-----------------|
| 0  | F   | 18  | U       | GT3     | A       | 4    | 4    | at_home  | teacher  | mother   | ... | 11 | 22          | low     | 1      | 3      | 3      | 1               |
| 1  | F   | 17  | U       | GT3     | T       | 1    | 1    | at_home  | other    | father   | ... | 11 | 31          | average | 1      | 3      | 3      | 1               |
| 2  | F   | 15  | U       | LE3     | T       | 1    | 1    | at_home  | other    | mother   | ... | 12 | 37          | average | 2      | 3      | 2      | 2               |
| 3  | F   | 15  | U       | GT3     | T       | 4    | 2    | health   | services | mother   | ... | 14 | 42          | average | 1      | 3      | 2      | 1               |
| 4  | F   | 16  | U       | GT3     | T       | 3    | 3    | other    | other    | father   | ... | 13 | 37          | average | 1      | 3      | 2      | 1               |
| 5  | M   | 16  | U       | LE3     | T       | 4    | 3    | services | other    | mother   | ... | 13 | 37          | average | 2      | 3      | 2      | 1               |
| 6  | M   | 16  | U       | LE3     | T       | 2    | 2    | other    | other    | mother   | ... | 13 | 38          | average | 1      | 3      | 3      | 1               |
| 7  | F   | 17  | U       | GT3     | A       | 4    | 4    | other    | teacher  | mother   | ... | 13 | 36          | average | 1      | 3      | 3      | 1               |
| 8  | M   | 15  | U       | LE3     | A       | 3    | 2    | services | other    | mother   | ... | 17 | 48          | high    | 1      | 3      | 2      | 1               |
| 9  | M   | 15  | U       | GT3     | T       | 3    | 4    | other    | other    | mother   | ... | 13 | 37          | average | 1      | 3      | 1      | 1               |
| 10 | F   | 15  | U       | GT3     | T       | 4    | 4    | teacher  | health   | mother   | ... | 14 | 42          | average | 1      | 3      | 3      | 1               |
| 11 | F   | 15  | U       | GT3     | T       | 2    | 1    | services | other    | father   | ... | 13 | 35          | average | 1      | 3      | 2      | 1               |
| 12 | M   | 15  | U       | LE3     | T       | 4    | 4    | health   | services | father   | ... | 12 | 37          | average | 1      | 3      | 3      | 1               |
| 13 | M   | 15  | U       | GT3     | T       | 4    | 3    | teacher  | other    | mother   | ... | 13 | 37          | average | 1      | 3      | 3      | 1               |
| 14 | M   | 15  | U       | GT3     | A       | 2    | 2    | other    | other    | other    | ... | 15 | 43          | average | 1      | 3      | 2      | 1               |
| 15 | F   | 16  | U       | GT3     | T       | 4    | 4    | health   | other    | mother   | ... | 17 | 51          | high    | 2      | 3      | 3      | 1               |
| 16 | F   | 16  | U       | GT3     | T       | 4    | 4    | services | services | mother   | ... | 14 | 40          | average | 3      | 3      | 3      | 1               |
| 17 | F   | 16  | U       | GT3     | T       | 3    | 3    | other    | other    | mother   | ... | 14 | 41          | average | 1      | 3      | 2      | 1               |
| 18 | M   | 17  | U       | GT3     | T       | 3    | 2    | services | services | mother   | ... | 7  | 23          | low     | 1      | 3      | 3      | 2               |
| 19 | M   | 16  | U       | LE3     | T       | 4    | 3    | health   | other    | father   | ... | 12 | 36          | average | 2      | 3      | 3      | 1               |

# Analysis top/bottom 20 student's living in rural/urban

Top 20 students

Rural
20.00%

80.00%

Urban
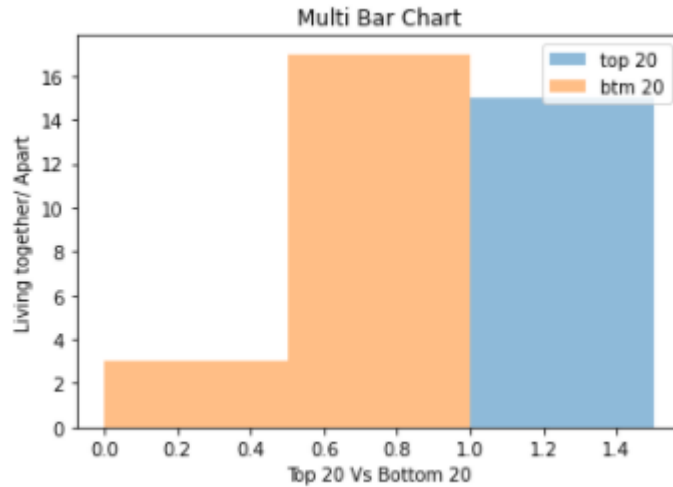
Bottom 20 students

Rural

55.00%

45.00%

Urban

The unique challenge faced by rural students is the transportation and the language they use might be affecting their scores.

# Analysis the top/bottom 20 student's living with/apart parents



Top 20 = 15 students living together with parents and 5 living apart
Btm 20 = 17 students living together with parents and 3 living apart

This graph shows no sign of disparity differences.

# Analysis the top/bottom 20 student's parents education using mean

```
1  top_20.groupby('grades')[["Medu", "Fedu"]].agg([np.mean])
```

|  | Medu mean | Fedu mean |
|---|---|---|
| grades | | |
| high | 3.133333 | 2.466667 |

```
1  btm_20.groupby('grades')[["Medu", "Fedu"]].agg([np.mean])
```

|  | Medu mean | Fedu mean |
|---|---|---|
| grades | | |
| low | 2.4 | 1.8 |

From the results, this means most students who are from top 20, both their parents education level is higher than the bottom 20.

# Performance by gender

```
]:   1  student_score_sex = student_data.groupby("sex")['total_score'].mean().reset_index()
     2  student_score_sex
```

```
]:
        sex   total_score
     0   F     35.712794
     1   M     33.669173
```

```
]:   1  student_score_sex.plot(kind = "barh")
```

`]: <matplotlib.axes._subplots.AxesSubplot at 0x21bb7bd8d30>`



Female and male academic performance is very close. However in fact the assumption of female tend to get better grades than guys might be true in this case.

# Does family support have impact towards achievements?



From this observation, family support does have great impact on the grades.
You can tell that students who have family support (3) scores higher. For student who have poor family relationship can be seen as positive skew. In another words, students are likely to score lower than 30.

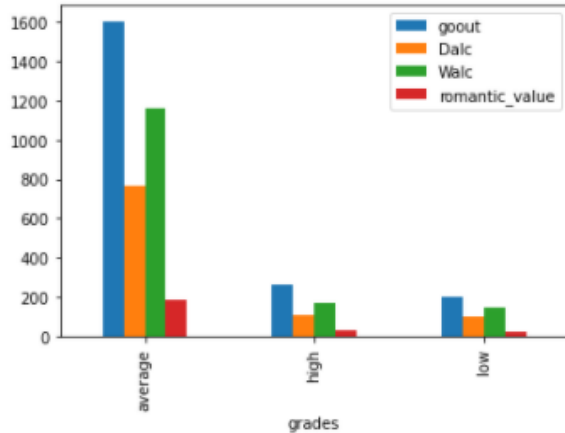# What type of job parent's hold will have impact towards student achievements?



Parents job can be seen as very impactful towards the scores. For teacher, has both the highest median and greatest variability among all. Services job has a normal distribution. While you compared at home, it has the lowest mean. It also can be seen as positive skew since median is shown closer to the line. Meaning might not be a good sign since it is towards the low scores.

# Which social factors has the most influence on academic achievement?

```
1  Social_factors = student_data.groupby('grades')[["goout", 'Dalc', 'Walc', 'romantic_value']].sum()
```

```
1  Social_factors.plot(kind = "bar")
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x21bb7dc98b0>
```
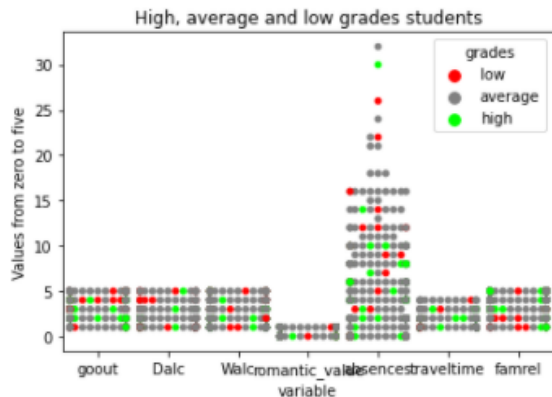


In visualizing the social factors, what we observed that students of low and high grades have quite a similar distribution.

# Which social factors has the most influence on academic achievement?

```
1  sns.swarmplot(x='variable',y='value',hue='grades' , data=melt,palette={'high':'lime','average':'grey','low':'red'})
2  plt.ylabel('Values from zero to five')
3  plt.title('High, average and low grades students')
4  plt.show
```
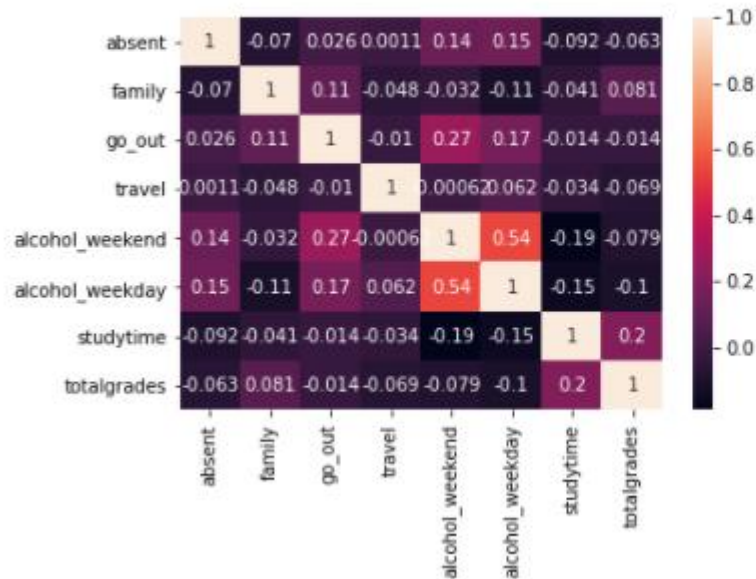
<function matplotlib.pyplot.show(*args, **kw)>



But if we were to look at this, we can construct a better perspective.

The low scorers tends to go out more often than the high scorers. Weekend, Weekday, romantic and travel time does not constitute significant difference. For family relation, you can tell that most low scorers agree they have poor relationship with family.
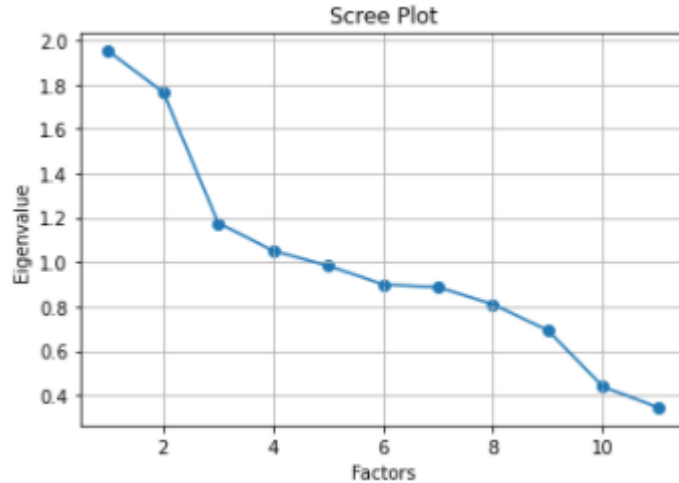
# Social factors correlation with student performance

```
1  sns.heatmap(a.corr(), annot=True)
2  plt.show()
```



This perhaps would give an in-depth insight with the correlation method where the social lifestyle are compared with their grades.

# Import Factor Analyzer – sklearn function

## To identify which factor is the most important for academics to know.



List of factors

1 = total grades
2 = study time
3 = family
4 = father education
5 = mother education
6 = travel
7 = go out
8 = absent
9 = weekday alcohol consumption
10 = weekday alcohol consumption

Notices that the total grades as the highest and dropping down to the next component. To simplify, the closest to the total score means the higher factor the academics must take note of and also take note of the factor that makes students grade drop. Retain factor 2, 3, 4 and 5 which is above the 1 criteria.

# Conclusion

Insights we have identified based on the report:

• Students who stay in urban areas are more likely to have higher score

• Students living together or apart with parents has no contribution to their grade

• Parents education level has effects on the student's grade

• Female academic performance is slightly higher than male

• Higher family support can achieve greater grades

• Parents job has an impact as well to the grades

• Study, family and parent's education factors are the most important influence for greater achievements while alcohol, absent has adverse effects. Traveling is also affecting the grades. Going out not so much of an impact.

Clearly, teachers are not able to change anything of the student's family background and can only educate on their social activities and keep an eye on students who scored poorly. The school should promote anti consumption of alcohol and absents.

# Thanks!

Any questions?