# NCU Introduction to Data Science Fall 2022 – HW1

**Colab link:**

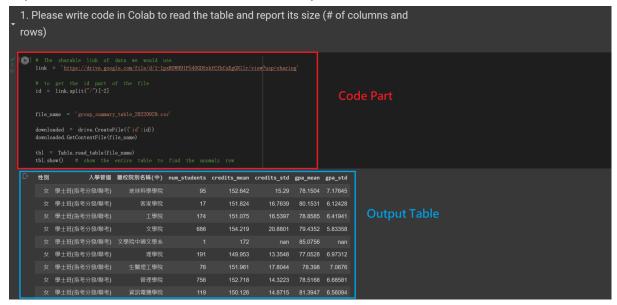**Student Info:** NCU MIS 3A 109403019 鄒翔宇

# 目錄

# Prerequisite

I use *PyDrive* to load the data from Google Drive. To use tables, I also import all of the modules called *datascience.* And, we need to do visualization, so I import *matplotlib*.



In addition, we have to enable *matplotlib* to support 繁體中文 finely, otherwise we would get the *RuntimeWarning: Glyph xxxxx missing from current font*.

# 1. Please write code in Colab to read the table and report its size (# of columns and rows)

I use the Table method *read_table* to read our imported CSV file and use the *show* method to print out the entire table to find the anomaly row.

```python
# The sharable link of data we would use
link = 'https://drive.google.com/file/d/1-1pxROWH91P54OGDSzkfCfkCuEgGNllr/view?usp=sharing'

# to get the id part of the file
id = link.split("/")[-2]

file_name = 'group_summary_table_20220929.csv'

downloaded = drive.CreateFile({'id':id})
downloaded.GetContentFile(file_name)

tbl = Table.read_table(file_name)
tbl.show()    # show the entire table to find the anomaly row
```

**Code Part**

**Output Table**

| 性別 | 入學管道 | 離校院別名稱(中) | num_students | credits_mean | credits_std | gpa_mean | gpa_std |
|---|---|---|---|---|---|---|---|
| 女 | 學士班(指考分發/聯考) | 地球科學學院 | 95 | 152.642 | 15.29 | 78.1504 | 7.17645 |
| 女 | 學士班(指考分發/聯考) | 客家學院 | 17 | 151.824 | 16.7639 | 80.1531 | 6.12428 |
| 女 | 學士班(指考分發/聯考) | 工學院 | 174 | 151.075 | 16.5397 | 78.8585 | 6.41941 |
| 女 | 學士班(指考分發/聯考) | 文學院 | 686 | 154.219 | 20.8801 | 79.4352 | 5.83358 |
| 女 | 學士班(指考分發/聯考) | 文學院中國文學系 | 1 | 172 | nan | 85.0756 | nan |
| 女 | 學士班(指考分發/聯考) | 理學院 | 191 | 149.953 | 13.3548 | 77.0528 | 6.97312 |
| 女 | 學士班(指考分發/聯考) | 生醫理工學院 | 76 | 151.961 | 17.8044 | 78.398 | 7.0676 |
| 女 | 學士班(指考分發/聯考) | 管理學院 | 756 | 152.718 | 14.3223 | 78.5168 | 6.68581 |
| 女 | 學士班(指考分發/聯考) | 資訊電機學院 | 119 | 150.126 | 14.8715 | 81.3947 | 6.56094 |

Then, I use the *num_rows* and *num_columns* methods in the Table object type to print out its size. We can see that the **# of columns is 8** and the **# of rows is 49**.
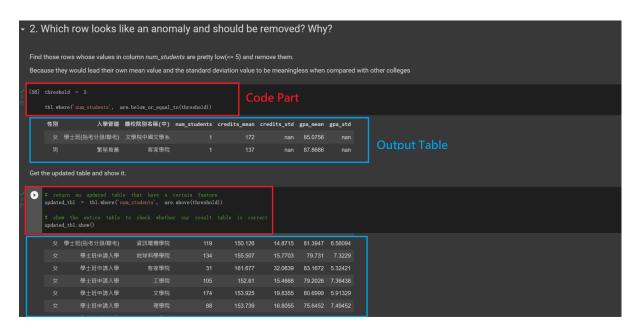
Print out the # of rows and the # of this table.

```python
[7] print("# of columns:", tbl.num_columns)
    print("# of rows:  ", tbl.num_rows)
```

**Code Part**

```
# of columns: 8
# of rows:  49
```

**Output**

# 2. Which row looks like an anomaly and should be removed? Why?

## 2.1 Please explain.

I want to find those rows whose values in column ***num_students* are pretty low(<= 5)** because they would lead their mean value and the standard deviation value to be meaningless when compared with other colleges. Also, If the number of students is 1, which would cause the denominator to become 0 when we calculate the standard deviation.

Actually, I want to use IQR to find outliers and remove them, but I found the result is meaningless.

## 2.2 Please write code in Colab to remove it and return an updated table.



Also, we can check the # of rows of the updated table and the total rows that have been removed to prove our operations are correct.
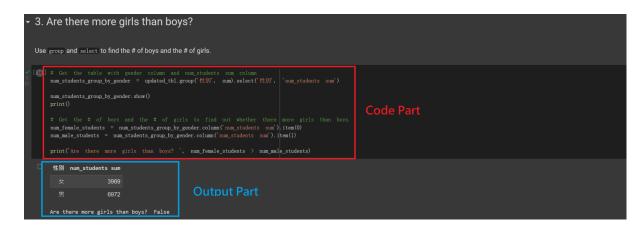


Note: I would use the updated table for the following question.

# 3. Are there more girls than boys?

No, there are boys more than girls. We can find out by the following table.

## 3.1  Please use `group` to do it.

**The # of female students** is **3969** while **the # of male students** is **6972**.
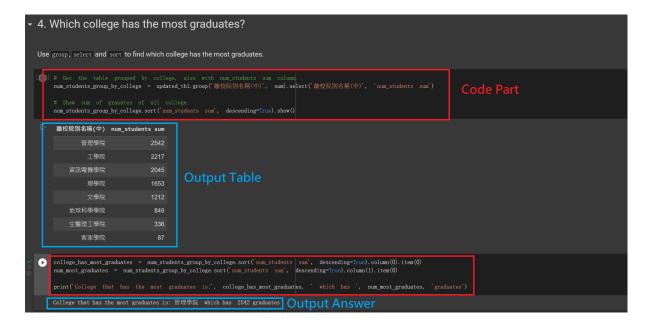


## 3.2  What % of students are female?

**36.28 %** of the students are **female**, and **63.74%** of the students are **male**.
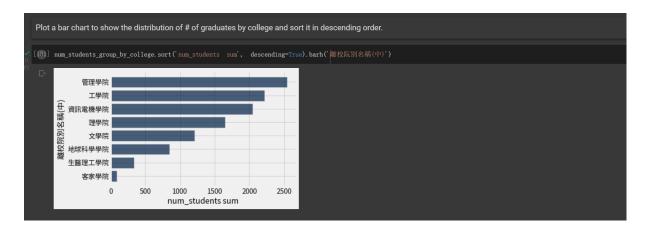
# 4. Which college has the most graduates?

The college that has the most graduates is 管理學院 which has 2542 graduates.

## 4.1  Please use `group` to do it.



## 4.2  Please plot a bar chart to show the distribution of # of graduates by college and sort it in descending order.

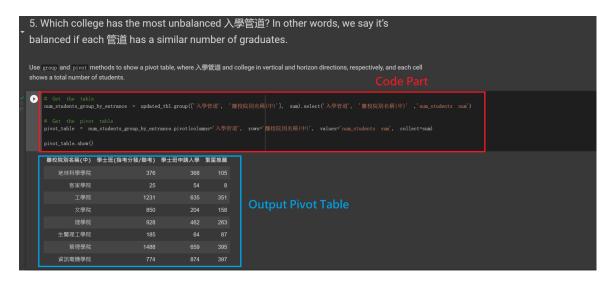I use the Table method *sort* and *barh* to plot a bar chart.

# 5. Which college has the most unbalanced 入學管道? In other words, we say it's balanced if each 管道 has a similar number of graduates.

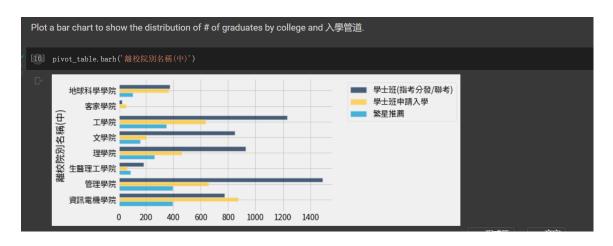管理學院 has the most unbalanced 入學管道.

## 5.1 Please use `group` to show a pivot table, where 入學管道 and college in vertical and horizon directions, respectively, and each cell shows a total number of students.

I use *group* and *pivot* methods to generate the pivot table and *show* it out.



## 5.2 Please plot a bar chart to show the distribution of # of graduates by college and 入學管道.

I use *barh* method to plot the bar chart to show the distribution of # of graduates by college and 入學管道。We can speculate 管理學院 has the most imbalanced distrubition of # of graduates by college and 入學管道 through above bar chart.

To prove our speculation, I calculate the **standard devition**(I think std could be seen as imblance level for distribution) and add it as a column to our pivot table. And then, print the *pivot_table* with sorting.



Yes! 管理學院 is on the top of table, so we can say that 管理學院 has the most unbalanced 入學管道。

# Reference

1. **Data8 docs**

2. **Ways to import CSV files in Google Colab - GeekforGeeks**

3. **Colab 進行matplotlib繪圖時顯示繁體中文**

4. **Python statistics.stdev() Method - W3school**