

# CW1 report

Enze Zhou, Student ID:2254411

TA:Zhejun.Yang

## I. IMPORT

First, import all the packages that will be used for subsequent operations. Read the file and display all the data. As shown in Figure [1].

```
[1]: import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
from sklearn.preprocessing import PCA
data = pd.read_csv("data.csv")
display(data)
```

	Index	Gender	Programme	Grade	Total	MCQ	Q1	Q2	Q3	Q4	Q5
0	1	1	3	1	45.0	20	8	4	2	10.0	0
1	2	2	3	1	40.0	20	4	2	0	10.0	0
2	3	1	4	2	26.0	24	0	0	2	10.0	0
3	4	2	1	1	30.0	24	4	0	0	10.0	2
4	5	1	2	1	20.0	20	0	2	4	10.0	0
...	...	...	...	...	...	...	...	...	...	...	...
454	455	2	1	2	15.0	10	4	4	4	10.0	0
455	456	2	4	2	40.0	20	4	4	4	10.0	0
456	457	2	1	2	15.0	10	4	4	11	10.0	0
457	458	2	1	2	15.0	10	4	4	10	10.0	3
458	459	2	1	2	20.0	10	4	4	4	2.0	0

459 rows x 12 columns

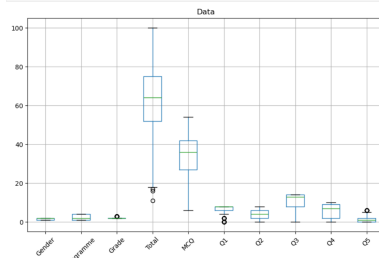
[1]

## II. TASK1:BOX PLOT DRAWING AND RELATED PROCESSING

The key is to remove the index column from the beginning. Then the box plot was drawn and the distribution of the data was observed, and it was found that total and MCQ were affected more. We can use standardization to process data to reduce the impact of scale for the raw data. As shown in Figure [2],[3].

It is important to note that because of analyzing the distribution of the feature with association of the programme that a student is enrolled, we need to remove the programme to reduce the impact of overfitting.

```
[2]: data1 = data.drop(columns=['Index'])
# 去除Index列
# 显示数据分布
plt.figure(figsize=(10, 6))
data1.boxplot()
plt.title('Data')
plt.xticks(rotation=45)
plt.show()
```



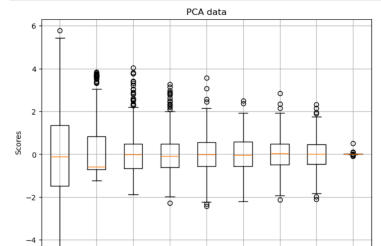
[2]

```
[3]: # 去除Programme后先标准化再画PCA
data1 = data1.drop(columns=['Programme'])

# 数据标准化
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
scaled_data = scaler.fit_transform(data1)

pca = PCA()
score = pca.fit_transform(scaled_data)

# 画PCA图
plt.figure(figsize=(8, 6))
plt.boxplot(score)
plt.title('PCA data')
plt.xlabel('Principal Component')
plt.ylabel('Scores')
plt.grid(True)
plt.show()
```



[3]

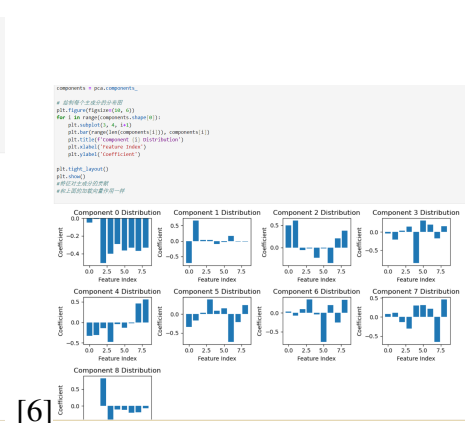
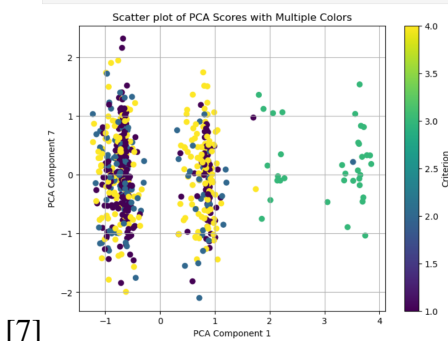
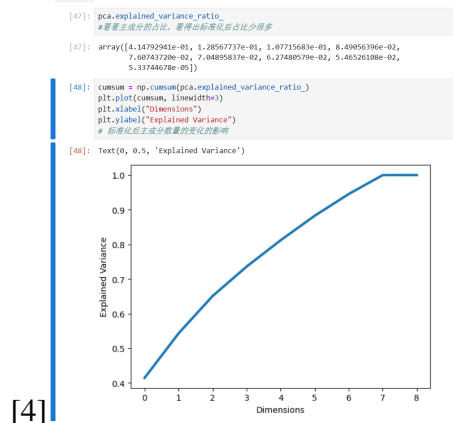
### III. TASK2:PCA AND CHOOSE A SUITABLE SET OF COMPONENT

After standardization, use scikit-learn library for PCA processing and draw the PCA diagram. Since the principal component after PCA is the linear combination with the maximum variance after projection in the original data, the distribution degree of the data can be obtained from it. In addition, the principal component can be extracted by reducing and de-correlating. As shown in Figure [4]. I also add `pca.explained-variance-ratio` to list the proportion of each component, and draw the proportion diagram of the change in the number of principal components and the diagram before and after PCA standardization, which is more convenient for me to observe the distribution of each component. It is found that the proportion of other principal components is roughly the same except the first and second principal components.

Then I drew the scatterplot of features to observe the distribution of the two features in the Programme class. As shown in Figure [5]. In addition, I made a table of principal component loading vectors and a distribution diagram of each principal component like figure [6], so as to observe the significance of contributions of each feature. I found that Q3 and Q5 contributed less, which was convenient for me to make decisions in subsequent operations.

Finally, I used different colors to represent different programmes and drew PCA scatter plots, as shown in Figure [7]. Through repeated comparisons, it was found that the second principal component and the seventh principal component could distinguish PCA results of different programmees better than other combinations. Therefore, the set of components I found is the second and the seventh.

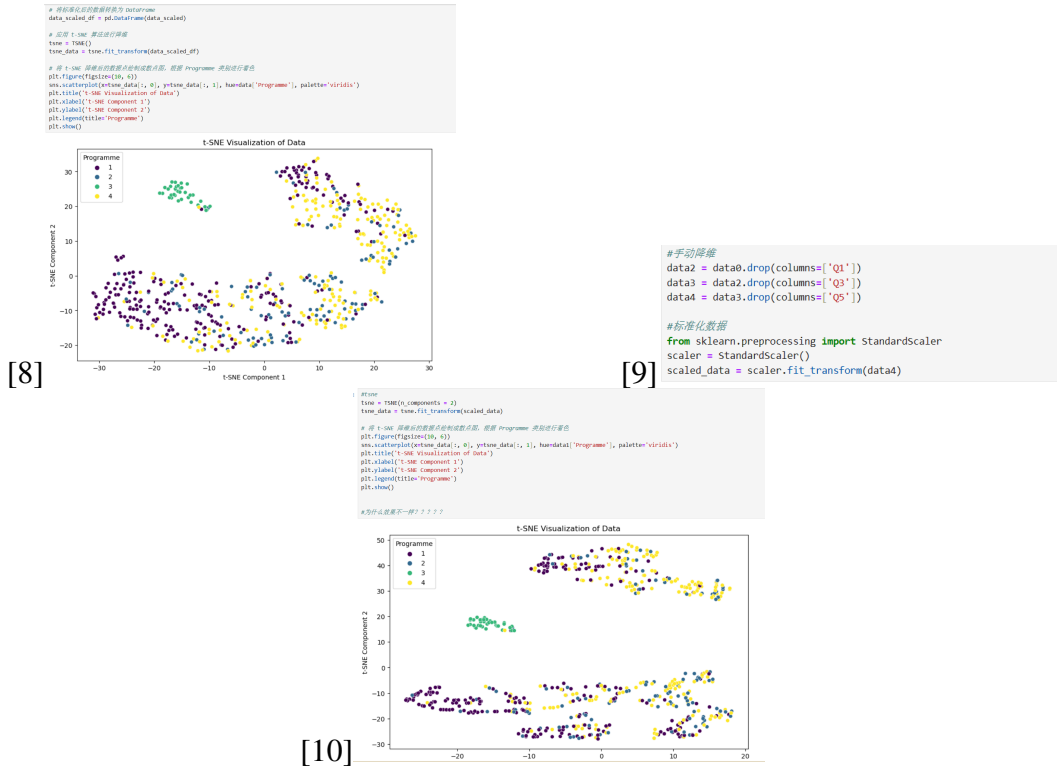
Moreover, by analyzing the variance ratio and observing the scatter plot, I found that grade and gender had the greatest effect on classification. In summary, the features I found were grade and gender.



#### IV. TASK3:T-SNE AND MANUAL DIMENSION REDUCTION

I first carried out tsne dimensionality reduction, in which I used a more complex process, I chose to re-carry out PCA and standardization, and then carried out t-SNE algorithm to reduce dimensionality, and drew the scatter plot, like figure [8], and found that only Programme3 was perfectly separated.

Not satisfied with this effect, I conducted manual dimension reduction again, removed Q1, Q3 and Q5, and conducted t-SNE algorithm dimension reduction again. The final result still only separated Programme3, and the effect was only a little better than before. The possible reason was that the feature combination selected by me in PCA was not correct. Or there are still problems in the order and analysis of my tsne, as shown in Figure [9],[10].



#### V. TASK4:VISUALISE AND COMPARE RAW FEATURES, SCALED FEATURES, PCA FEATURES

Finally, I visualized all the original features, normalized features and PCA features for comparison and observation.

