

CW3 report

Enze Zhou, Student ID:2254411

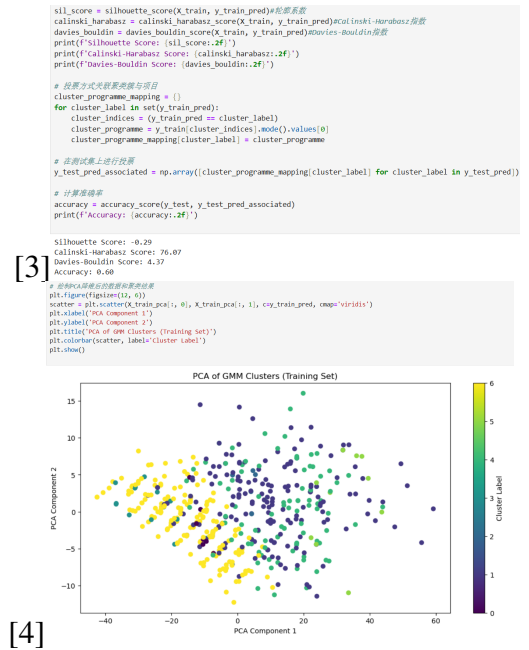
TA:Zhejun.Yang

I. GAUSSIAN MIXTURE MODEL

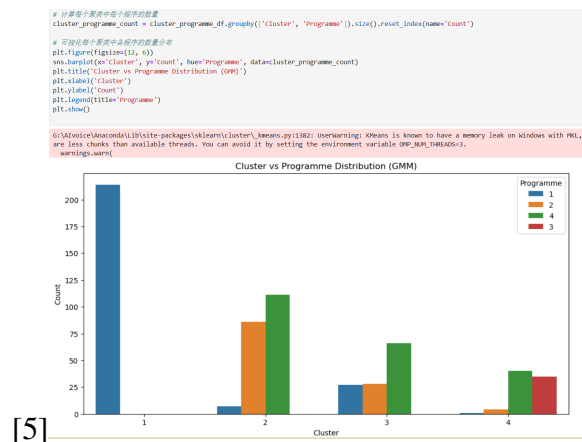
In this analysis, I first selected gender, grade, total score, MCQ, Q1, Q2, Q3, Q4, Q5 and other features, and divided the data set into training set and test set, as shown in Figure [1]. Next, I used Gaussian mixture model (GMM) to fit the training data, and selected the optimal number of model components through AIC, BIC and likelihood evaluation, as shown in Figure [2]. According to the line chart, I could judge that the optimal number of GMM components should be 7, so I used 7 components in the subsequent training of the model. In order to evaluate the clustering effect, I calculated the posterior probability of the training set sample, and used the contour coefficient, Calinski-Harabasz index and Davies-Bouldin index to measure the clustering quality.



Then, I associate each cluster in the training set with the program label, as shown in Figure [3], and use the voting method to determine the program label corresponding to each cluster, and verify it in the test set. In addition, in order to see the clustering results more directly, I conducted PCA dimensionality reduction on the data, like Figure [4], and visualized the data and clustering results after dimensionality reduction, which showed that the effect was not good.

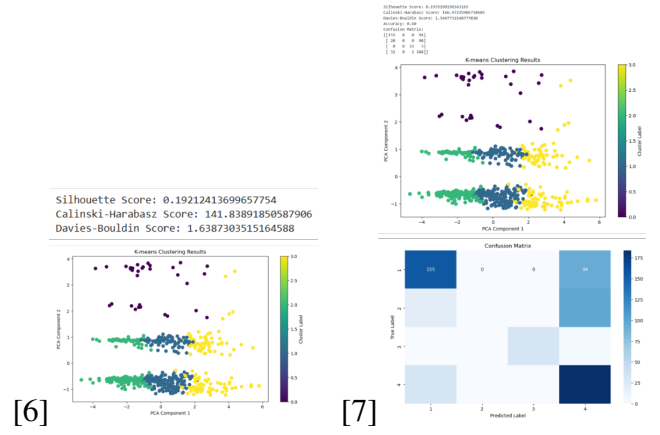


Finally, I displayed the clustering results through principal component reduction and visualization, as shown in Figure [5], and analyzed the distribution of each program label in each cluster. However, it can be seen that except cluster1, there are only programme 1, and other clusters have many programmes. This means that it has little to do with the programme, almost nothing to do with it.



I do cluster analysis mainly by using GMM,

where we are able to reveal hidden patterns and structures in the data set, especially the relationships between different features. By associating the clustering results with program labels, we can better understand the classification information in the data, helping to improve the classification model and make more accurate decisions in practical applications. This analysis method not only improves the interpretability of the data, but also provides a valuable reference for the subsequent model optimization.



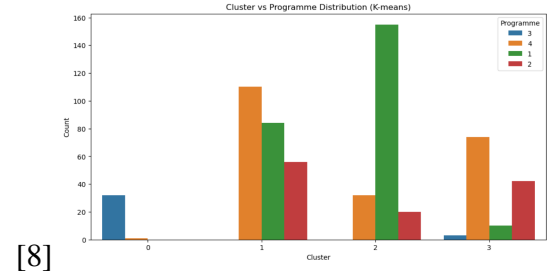
clusters are almost unrelated, and we can say that this cluster is not associated with the programme.

II. K-MEANS

I performed K-means clustering analysis on the data to explore relationships between features and map the results to actual item labels. First of all, I standardized the feature data to eliminate the influence between different dimensions. Then I use K-means algorithm to cluster the data, and since there are four programmes, I set four clusters. As shown in Figure [6], I calculated the contour coefficient, Calinski-Harabasz index and Davies-Bouldin index to evaluate the clustering quality, assess the compactness and separation of clusters, and reduce and visualize the clustering results.

I want to know whether each cluster is associated with the programme, so I associate each cluster with the actual programme label to determine the most common project label in each cluster. Finally, I categorize the data set, predict the cluster label for each sample, and map it back to the actual item label. The accuracy of classification is calculated and confusion matrix is generated to judge whether my clustering effect is good, as shown in Figure [7]. Moreover, I used PCA to reduce the dimension of the data to two dimensions, and visualized the clustering results and confusion matrix, so that the results could be seen more intuitively.

In order to find out whether each cluster has any relationship with the Programme, I analyzed the quantity distribution of each project label in each cluster and presented it visually, as shown in Figure [8]. However, it can be seen from the results that this is similar to the effect of GMM, only cluster1 may have a little relationship with the programme. Other

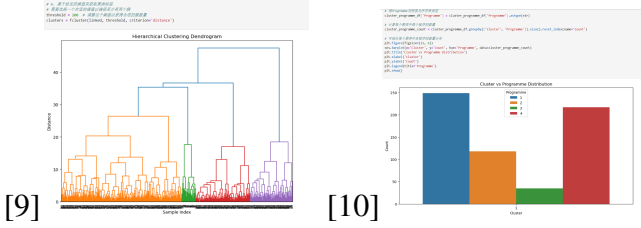


The significance of these efforts is that by using K-means clustering analysis, we are able to uncover hidden patterns and structures in the data set, revealing relationships between features. By associating the clustering results with the item labels, the classification information of the data can be better understood, helping to improve the classification model and make more accurate decisions in practical applications. In addition, the clustering quality index can be used to objectively evaluate the clustering effect and ensure the reliability of the model. The final visual display helps to intuitively understand the clustering results and item label distribution, and provides strong support for subsequent data analysis and decision making.

III. HIERARCHICAL CLUSTERING

I then performed cluster analysis on the data using hierarchical clustering methods to explore the relationships between features and map the results to actual item labels. First, I use Ward method and

Euclidean distance to perform hierarchical clustering on the data. The hierarchical clustering model is constructed by calculating the link matrix. Based on this model, a tree diagram like Figure [9] is drawn to show the hierarchical clustering results, and an appropriate threshold is selected to determine the clustering label. Then I obtain the cluster label of each sample according to the selected threshold value, and associate it with the actual item label to determine the quantity distribution of each item label in each cluster. Moreover, the quantity distribution of each item label in each cluster is analyzed and visually displayed, as shown in Figure [10]. The result shows that, like the previous two models, it has almost nothing to do with Programme.



The significance of these efforts is that by using hierarchical clustering analysis, we are able to uncover hidden patterns and structures in the data set, revealing relationships between features. Hierarchical clustering provides an intuitive way to understand the hierarchical structure and similarity relationships of data. By associating the clustering results with the item labels, the classification information of the data can be better understood, helping to improve the classification model and make more accurate decisions in practical applications.

IV. CONCLUSION

By using the Gaussian mixture model (GMM), K-Means, and Hierarchical Clustering for data clustering analysis, I delve into the relationships between features and map the results to actual item labels. In Gaussian mixture model analysis, we selected the optimal number of model components by AIC, BIC and logarithmic likelihood assessment, and found that the correlation between the samples in the training set and the item label was low. When K-means method is used for clustering, the clustering

quality is evaluated by contour coefficient, Calinski-Harabasz index and Davies-Bouldin index, and the mapping relationship with item labels is found in the clustering results, but the results show that most clusters are not strongly correlated with item labels. In the hierarchical cluster analysis, I used Ward method and Euclidean distance to build a hierarchical cluster model, visualized the clustering results through the tree graph, and found that most clusters still had a low correlation with item labels after selecting appropriate thresholds.

Although the correlation between the clustering results of each method and the item labels is not strong, through these analyses, I revealed the hidden patterns and structures in the data set, and improved the understanding of the data classification information. These analysis methods not only enhance the ability of data interpretation, but also provide valuable reference for improving classification models, which helps to make more accurate decisions in practical applications. The final visual display enables us to intuitively understand the distribution of clustering results and item labels, which provides strong support for subsequent data analysis and decision making