# Project 2: K-NN, PCA, and Clustering

Submitted by: Gmon Kuzhiynaikkal
NU ID : 002724506

## Introduction:

The Dataset used for the project is seed dataset. I have collected from the UCI machine learning Repository and it has 7 attributes. The examined group comprised kernels belonging to three different varieties of wheat: Kama, Rosa and Canadian. The independent variables are area, perimeter, compactness, Kernel Length, Kernel Width, coefficient, grooveLength and the independent variable is wheat Type.

In this project, I was able to apply K-nearest neighbours on the seed dataset and I was able to further explore different types of the distance metric used by the K-nearest neighbours. I was able to interlink PCA with different clustering algorithms and I was able to find differences and variation in different clustering algorithms.

As a part of further extensions, I was able to explore one of the clustering quality metrics called Within-Cluster-Sum-of-Squared-Errors and I tried to compare with Rissannon Minimum Description Length (which is another cluster quality metric method). I have tried to plot the graph between these with different numbers of clusters(k).

# 1.Apply Nearest Neighbor classification to your data set:

## a)Write a distance metric for your data set:

I have used the euclidean distance metric for the calculation of the distance. I have created a sample data point and calculated the distance for all the data points and output of the program is given below. The equation used in the program is given below.

$$d(p,q) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \cdots + (p_n - q_n)^2}.$$

It tries to measure the distance between the two points in a N-dimensional space. The output of the distance metric is given below. I have used one point and thus the is Ex1 vector

PS: python program to run this is given as 1a.py



## B.Implement Nearest Neighbour classification for your test set:

I have implemented the nearest neighbour from the data and error terms after running the program is given is given below:



The python program to run this is 1b.py

# C.Evaluate your classifier:

I have tired to evaluate the model using the euclidean distance and the accuracy of the model is 36.5% and the output of the confusion he confusion matrix i got after running this is t [[15 0 0], [12 0 0], [14 0 0]] indicates that there are no true positive predictions made by the nearest neighbour classifier. This suggests that the classifier is not performing well and needs further investigation. The output image of the file is given below and program to get the accuracy and confusion matrix is '1b.py'



# D.Experiment with your classifier:

In order to improve the accuracy, I have tried to change the distance metric. I have tried the manhattan distance, it is giving me the same accuracy and confusion matrix as the euclidean matrix.

In order to get more accuracy, I have tried to use the more complex distance metric called the cosine similarity. I was able to get 100% accuracy and the confusion matrix found out as below:



Thus I was able to change the accuracy of the model from 36.5% to 100% by changing my distance metric and classifier with a better true positive prediction.

# 2.Experiment with clustering and PCA on two structured data sets

## A) Download and plot the example data:

The scatter plot for the first set A is shown below. I could see 6 slightly circular clusters naturally occurring from the data. The scatter plot for the Set A csv file is shown below:



The scatter plot for the Set B csv file is also 6 clusters and it more like a straight line clusters and the the scatter plot is given below for Set B:



The python file to run this is '2a.py'

## B.Cluster the example data:

**For the setA**: using the k means clustering plot diagram(k=6) and hierarchical clustering plot diagram is given below:

From the plot diagram for the data setA. I could see that both of them have produced the same shape and number of clusters using the different clustering algorithms.

**For set B:**using the k means clustering plot diagram(k=6) and hierarchical clustering plot diagram is given below:



From the plot diagram for the data setB. I could see that both of them have produced the same shape and number of clusters using the different clustering algorithms also.

PS: python program to run this is '2b.py'

## C.Compare using different numbers of clusters:

I have tried to plot the representation error with the number of clusters(k) and I could find that the error is decreasing as the number of the clusters are increasing. The plot using the setA a is given below:



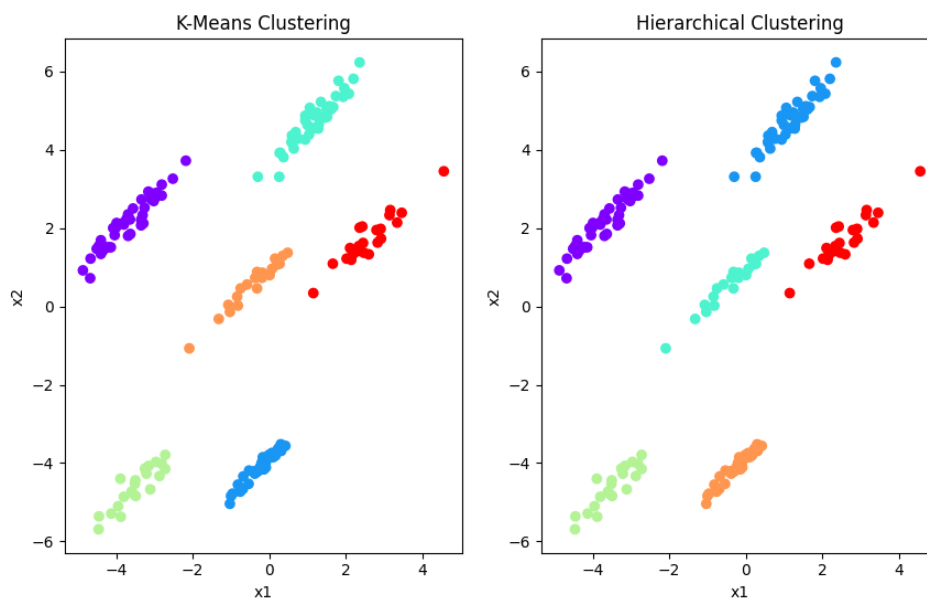Representation error almost got a stable from k=6 onwards. In order to find the cluster quality, I have used the Rissannon Minimum Description Length and the plot for RMDL to number of clusters is given below:



From the graph, I could see that k=6 has the best cluster quality as per the Rissannon Minimum Description Length. Thus for the Set A dataset, we could say that the ideal state of the data will be having 6 clusters.

PS : Python code for running this '2c.py'

# D:Apply PCA to the example data:

PCA is applied on the set A and the plot of the graph is given below the arrow head is downward. This implies that there is negative correlation with the features.



Data set A projected onto its two eigenvectors

For the set B the plot of the graph is given below and the arrow is pointing to the PC1 axis direction that there is a positive correlation with the features.



Data set B projected onto its two eigenvectors

PS: python file to run this program is '2d.py'

## E.Recluster using the projected data:

There is some difference in the Projected Data plot compared it is less scattered and the data are easier to interpret because of the lesser dimensionality. Scatter plot for Set A and Set B are given below:



PS: python code for running this is '2e.py'

# 3.Apply K-Means Clustering to your data set:

The data used for the project belongs to the seed data csv containing 7 attributes and the dependent variable is the wheat type, we are trying to classify 3 types of wheat types. The representation error vs number of cluster of the dataset is given:



If we check, we could see a variation for the cluster number 3 and plot which represent the cluster quality( used Rissannon Minimum Description Length) is given as below:



From this graph, we can conclude that, best cluster value(k value) for my dataset will be nearly equal to 3.

PS: Python program to run this program is '3.py'

## 4.Use K-Nearest Neighbour and PCA to classify activity from phone measurements:

Comparing the KNN classifier and the PCA projected data KNN classifier, it can be seen that the raw data KNN classifier has a higher accuracy (0.9016) compared to the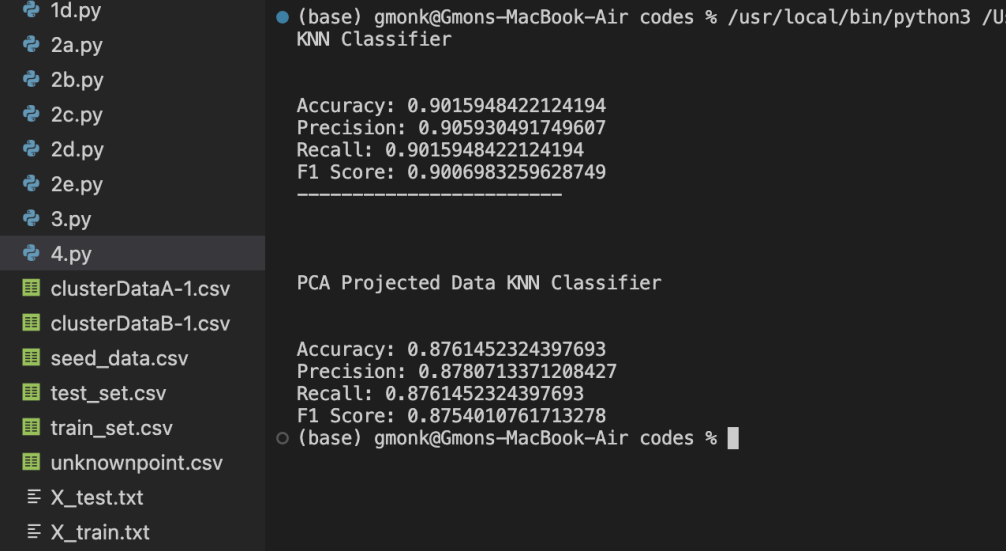 PCA projected data KNN classifier (0.8761). The precision and recall scores are also higher for the raw data KNN classifier. The F1 score, which is a weighted average of precision and recall, is also higher for the raw data KNN classifier (0.9007) compared to the PCA projected data KNN classifier (0.8754).

These results suggest that the raw data KNN classifier is performing better than the PCA projected data KNN classifier. The output of the comparison is given below on the data.



The reasons why the PCA projected data KNN classifier performed worse than the raw data KNN classifier could be that the PCA projection lost some important information in the data, leading to decreased performance. Another possibility is that the number of dimensions used in the PCA projection was not enough to explain the complex relationships in the data, leading to overfitting or underfitting

# Extensions:

I tried to implement Within-Cluster-Sum-of-Squared-Errors (WCSS) also known as inertia, this metric measures the sum of squared distances between each data point and its closest cluster centroid. A smaller WCSS score indicates a better clustering solution. I tri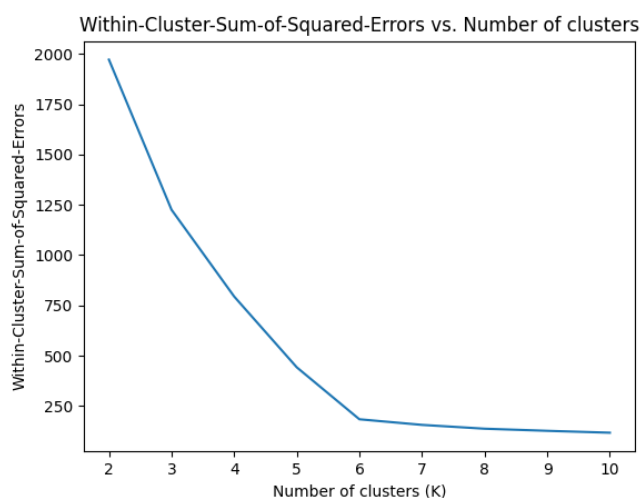ed to compare that with Rissannon Minimum Description Length on the dataset given for task 2. Tried to plot the number of clusters from (2 to 11) with the different cluster quality metric value. The plot for Rissannon Minimum Description Length is given below, for this highest value is the best and we could see that k=6 is a better value for it.



I tired to Plot WCSS, with the number of clusters and we could see that it is lower for when k=6, so it is better to take the k value for the above k mean clustering algorithm as 6

# Summary ( Reflection of what I learned from the project):

Some of the key points, I have learned from the projects are given below

- ➔ Implementing different distance metrics, which are used in the different algorithms of machine learning.
- ➔ Implementing nearest neighbour algorithm from the distance metrics that i have designed.
- ➔ Checking the accuracy of the classifies and analysing it with the confusion matrix
- ➔ Tried to improve the accuracy of the nearest neighbour classifier with different distance metrics and I was able to to improve the accuracy to 100% using cosine similarity distance metrics in the nearest neighbour.
- ➔ Tried to implement different clustering algorithms using sk-learn.
- ➔ learned to find the natural clusters in a scatter plot, before applying any cluster algorithm by plotting the data.
- ➔ Find out the variation of the representation error in the clustering algorithms with increase in the cluster algorithms.
- ➔ Find out the trend of applying principal component analysis(PCA) on the various cluster algorithms and plotting the eigenvector on scatter plot.
- ➔ I was able to find the variation of the PCA with nearest neighbour algorithms as well.
- ➔ I was able to learn different types of the cluster quality metrics, which can be used to compare different cluster quality in the algorithms.
- ➔ Learned about the representation error used in the clustering algorithms.
- ➔ As a part of the extension, I have tried to use different cluster quality metrics and tried to compare it.

# Bibliography:

Some of the importance links and materials i have used for this projects
are:

1. https://archive.ics.uci.edu/ml/datasets/seeds#

2. https://scikit-learn.org/stable/

3. Textbook by Muller and Guido

4. https://youtu.be/ukzFI9rgwfU