

INF552: Programming Assignment 2

Yi-Li Chen

Part 1: Implementation

- All imported library

```
1. import csv
2. import numpy as np
3. from math import *
4. import pandas as pd
5. import matplotlib.pyplot as plt
6. from scipy.stats import multivariate_normal
```

- Read text file into a list

```
1. data = []
2. with open('clusters.txt') as file:
3.     reader = csv.reader(file)
4.     for line in reader:
5.         data.append([float(line[0]),float(line[1])])
```

K-means

- Initialize k centroids which is picking the points randomly among the dataset.

```
1. def initial_centroids(data, k):
2.     num_samples, dim = np.array(data).shape
3.     centroids = np.zeros((k, dim))
4.     idx = np.random.randint(len(data), size=3)
5.     for i in range(k):
6.         centroids[i] = data[idx[i]]
7.     return centroids
```

- Calculate the Euclidean Distance between two data points

```
1. def min_distance(p1, p2):
2.     return sqrt((p1[0]-p2[0])**2+(p1[1]-p2[1])**2)
```

- Find out each data point belongs to which cluster by the distance between centroids

```
1. def update2centroids(datapoint, centroids):
2.     k_num_distance = []
3.     for c in centroids:
4.         d = min_distance(c, datapoint)
5.         k_num_distance.append(d)
6.     return np.argmin(k_num_distance)
```

- Classify the data points into k cluster by the index we have found in previous function. Then, the centroids will be updated by the mean values of data in each cluster. In addition, store the difference between the old and new centroids, so that we can when it should be stopped.

```

1. def classify2cluster(data, k, centroids, cluster):
2.     cluster=[] for i in range(k)
3.     for p in data:
4.         c_idx = update2centroids(p, centroids)
5.         cluster[c_idx].append(p)
6.     # find new centroids from the clusters
7.     num_samples, dim = np.array(data).shape
8.     new_centroids = np.zeros((k, dim))
9.     for i in range(k):
10.        new_x = np.array(cluster[i][:,0]).mean()
11.        new_y = np.array(cluster[i][:,1]).mean()
12.        new_centroids[i] = [new_x,new_y]
13.        difference = new_centroids - centroids
14.    return (difference, new_centroids, cluster)

```

- Firstly, initialize the centroids. Secondly, classify the data repeatedly with terminated condition which is when there is no difference between the old and updated centroids. In addition, create a dict() to store the centroids for each cluster of every iteration.

```

1. def kmeans(data, k):
2.     mark = ['or', 'ob', 'og', 'ok', '^r', '+r', 'sr', 'dr', '<r', 'pr']
3.     centroids_dict = {'cluster1':[], 'cluster2':[], 'cluster3':[]}
4.     cluster=[] for i in range(k)
5.     centroids = initial_centroids(data, k)
6.     centroids_dict['cluster1'].append(tuple(centroids[0]))
7.     centroids_dict['cluster2'].append(tuple(centroids[1]))
8.     centroids_dict['cluster3'].append(tuple(centroids[2]))
9.     difference, new_centroids, cluster = classify2cluster(data, k, centroids, cluster)
10.    while np.any(difference != 0):
11.        cluster=[] for i in range(k)
12.        difference, new_centroids, cluster = classify2cluster(data, k, new_centroids, cluster)
13.        for i in range(3):
14.            plt.plot(np.array(cluster[i][:,0]),np.array(cluster[i][:,1]),mark[i])
15.            plt.show()
16.            centroids_dict['cluster1'].append(tuple(new_centroids[0]))

```

```

17.         centroids_dict['cluster2'].append(tuple(new_centroids[1]))
18.         centroids_dict['cluster3'].append(tuple(new_centroids[2]))
19.     return (difference, new_centroids, cluster, centroids_dict)

```

- The final result of k-means clustering:

```

1. difference, new_centroids, cluster, centroids_dict = kmeans(data, 3)

```

	Cluster1	Cluster2	Cluster3
Centroid	5.738495346032258, 5.164838081193549	-0.960652907023256, -0.652218412860465	3.2888485605151514, 1.9326883657575757

GMM

- Initialize the ric of each data points randomly, but the sum of ric for each data point needs to be 1

```

1. def initial_ric(data, k):
2.     initial_ric = np.zeros((len(data), k))
3.     random_int = np.random.randint(low=1,high=10000,size=(len(data), k))
4.     for row in range(len(random_int)):
5.         initial_ric[row] = random_int[row]/sum(random_int[row])
6.     return initial_ric

```

- Update the mean for each Gaussian distribution from the ric and data points

```

1. def calculate_new_mu(k, ric, data):
2.     new_mu = np.zeros((k, 2))
3.     # calculate the mu of each cluster
4.     for c_index in range(k):
5.         # loop all datapoints and their corresponding ric
6.         for r, p in zip(ric[:,c_index], np.array(data)):
7.             new_mu[c_index] += r*p
8.         new_mu[c_index] = new_mu[c_index] / sum(ric[:,c_index])
9.     return new_mu

```

- Update the covariance for each Gaussian distribution from the ric, data points and the mean we just got.

```

1. def calculate_new_sigma(k, ric, data, mu):
2.     new_sigma = np.array([[0, 0], [0, 0]] for i in range(0,k)], dtype='float64')
3.     for c_index in range(k):
4.         # loop all datapoints and their corresponding ric

```

```

5.         for r, p in zip(ric[:,c_index], np.array(data)):
6.             new_sigma[c_index] += r*np.outer(p-mu[c_index],p-mu[c_index])
7.             new_sigma[c_index] = new_sigma[c_index] / sum(ric[:,c_index])
8.     return new_sigma

```

- Update the pi for each Gaussian distribution from the ric, data points.

```

1. def caluculate_new_pi(k, ric, data):
2.     new_pi = np.zeros(k)
3.     for idx in range(k):
4.         new_pi[idx] = sum(ric[:,idx]) / len(data)
5.     return new_pi

```

- Calculate the new ric from these updated parameters for each data points based on Bayes rule where π is the prior weights and the likelihood is normal

```

1. def update_ric(k, data, mu, sigma, pi):
2.     pdfs = np.zeros((len(data), k))
3.     for i in range(k):
4.         pdfs[:, i] = pi[i] * multivariate_normal.pdf(data, mu[i], sigma[i])
5.     ric = pdfs / pdfs.sum(axis=1).reshape(-1, 1)
6.     return ric

```

- In sum, firstly, initialize the ric randomly. Secondly, update the μ , covariance and pi step by step for each Gaussian distribution. Then, update the ric from these updated parameters. Afterwards, go back to the previous function so that parameters can be updated again. Finally, the two steps of E-step and M-step are then repeated until convergence. We have set a condition for while loop which is it will terminate if the difference between the updated ric and old one less than 0.005

```

1. def GMM(data, k):
2.     param_dict = {'mu':[], 'covariance':[], 'pi':[]}
3.     diff = float('Inf')
4.     ric = initial_ric(data, k)
5.     while np.any(abs(diff) > 0.005):
6.         mu = caluculate_new_mu(k, ric, data)
7.         sigma = caluculate_new_sigma(k, ric, data, mu)
8.         pi = caluculate_new_pi(k, ric, data)
9.         new_ric = update_ric(k, data, mu, sigma, pi)
10.        # store all updated parameters for each iteration of GMM
11.        param_dict['mu'].append(mu)

```

```

12.     param_dict['covariance'].append(sigma)
13.     param_dict['pi'].append(pi)
14.     diff = new_ric - ric
15.     ric = new_ric
16.     # classify and store the data points into list
17.     cluster=[] for i in range(k)
18.     for idx in range(len(ric)):
19.         c_idx = np.argmax(ric[idx])
20.         cluster[c_idx].append(data[idx])
21.     # plot the cluster result
22.     for i in range(3):
23.         mark = ['or', 'ob', 'og', 'ok', '^r', '+r', 'sr', 'dr', '<r', 'pr']
24.         plt.plot(np.array(cluster[i])[:,0],np.array(cluster[i])[:,1],mark[i])
25.     plt.show()
26.     return (mu, sigma, pi, ric, param_dact)

```

- The final result of GMM clustering:

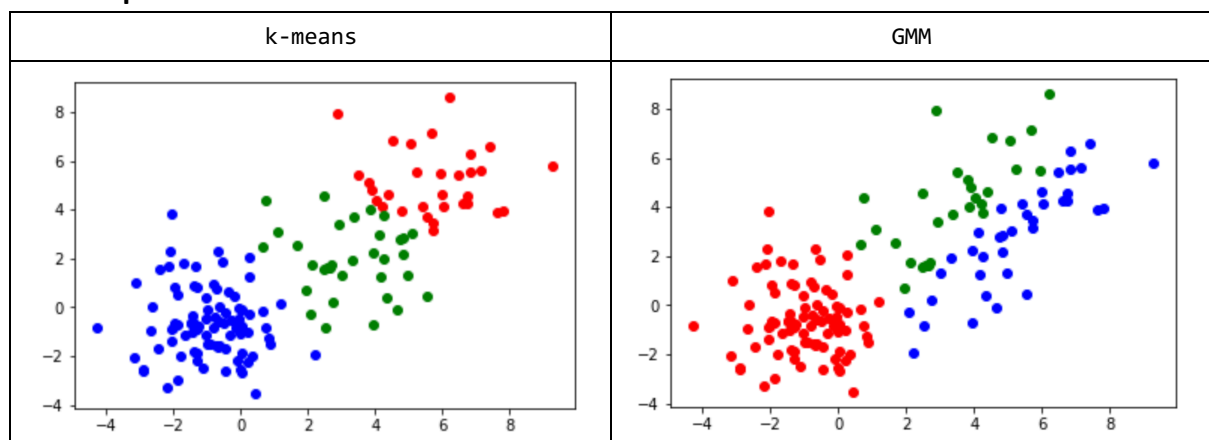
```

1. mu, sigma, pi, ric, param_dict = GMM(data, 3)

```

	Cluster1	Cluster2	Cluster3
Parameter	μ : [-0.994437 0.63388089]	μ : [4.92603572 2.94483694]	μ : [3.83434248 3.69019296]
	Σ : [[1.18193522 -0.0816692]	Σ : [[3.50297936 2.97191914]	Σ : [[3.44038384 2.08056859]
	[-0.0816692 2.02109235]]	[2.97191914 4.33618297]]	[2.08056859 5.46032776]]
	π : 0.55921906	π : 0.20585555	π : 0.23492538

The comparison of the results of k-means and GMM:



From the above figures of both results, we can easily find out that the clusters of k-means are tend to be spherical. However, the cluster of GMM are not limited to this.

In this assignment, the most challenging I encountered is how to update all parameters simultaneously. I really learn a lot of methods of data manipulation for NumPy array. Especially, when setting the stopping criteria, I transferred two for loop into one row command.