

```

1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 import random
6 import pandas as pd
7 import time

```

Double-click (or enter) to edit

```

1 from google.colab import files
2 uploaded = files.upload()

```

Choose Files A2Q1.csv

- **A2Q1.csv**(text/csv) - 40000 bytes, last modified: 3/13/2024 - 100% done
- Saving A2Q1.csv to A2Q1 (1).csv

## Question 1

- (1) You are given a data-set with 400 data points in  $\{0, 1\}^{50}$  generated from a mixture of some distribution in the file A2Q1.csv. (Hint: Each datapoint is a flattened version of a  $\{0, 1\}^{10 \times 5}$  matrix.)
- (i) Determine which probabilistic *mixture* could have generated this data (It is not a Gaussian mixture). Derive the EM algorithm for your choice of mixture and show your calculations. Write a piece of code to implement the algorithm you derived by setting the number of mixtures  $K = 4$ . Plot the log-likelihood (averaged over 100 random initializations) as a function of iterations.
  - (ii) Assume that the same data was infact generated from a mixture of Gaussians with 4 mixtures. Implement the EM algorithm and plot the log-likelihood (averaged over 100 random initializations of the parameters) as a function of iterations. How does the plot compare with the plot from part (i)? Provide insights that you draw from this experiment.
  - iii. Run the K-means algorithm with  $K = 4$  on the same data. Plot the objective of  $K - means$  as a function of iterations.

## 1. Expectation Maximisation Algorithms

### Bernoulli Mixture Models

```

1 def BernoulliIterator(D,K,N,mu_bernoulli,Pi_bernoulli,df):
2     Epsilon = 1e-10
3     likelihoods_bernoulli = np.zeros((N, K))
4     for i in range(K):
5         likelihoods_bernoulli[:, i] = np.sum(df * np.log(mu_bernoulli[i] + Epsilon) + (1 - df) * np.log(1 - mu_bernoulli[i] + Epsilon), axis=1)
6     weighted_likelihoods_bernoulli = likelihoods_bernoulli + np.log(Pi_bernoulli)
7     lambd_bernoulli = np.exp(weighted_likelihoods_bernoulli - np.max(weighted_likelihoods_bernoulli, axis=1)[:, np.newaxis])
8     lambd_bernoulli /= np.sum(lambd_bernoulli, axis=1)[:, np.newaxis]
9     log_likelihood_bernoulli = np.sum(weighted_likelihoods_bernoulli)
10    return Pi_bernoulli,mu_bernoulli,lambd_bernoulli,log_likelihood_bernoulli

```

```

1 def BernoulliIterator2(D,K,N,lambd_bernoulli,df):
2     Pi_bernoulli = np.sum(lambd_bernoulli, axis=0) / N
3     mu_bernoulli = np.dot(lambd_bernoulli.T,df) / np.sum(lambd_bernoulli, axis = 0)[:, np.newaxis]
4     return Pi_bernoulli,mu_bernoulli,lambd_bernoulli

```

```

1 def RandomBernoulliVariables(N,D,K):
2     Pi_bernoulli = np.random.rand(K)
3     Pi_bernoulli = Pi_bernoulli / np.sum(Pi_bernoulli)
4     mu_bernoulli = np.random.rand(K, D)
5     lambd_bernoulli = np.random.rand(N, K)
6     row_sums = lambd_bernoulli.sum(axis=1)
7     lambd_bernoulli = lambd_bernoulli / row_sums[:, np.newaxis]
8     return Pi_bernoulli, mu_bernoulli, lambd_bernoulli

```

```

1 def Bernoulli_Algorithm(N,D,K,no_of_epochs,df,tol,no_of_initialisations):
2     log_like_bernoulli = np.empty((no_of_initialisations,no_of_epochs))
3     for t in range(no_of_initialisations):
4         prev_log_likelihood_bernoulli = -np.inf
5         Pi_bernoulli, mu_bernoulli, lambd_bernoulli = RandomBernoulliVariables(N,D,K)
6         for epoch in range(no_of_epochs):
7             Pi_bernoulli,mu_bernoulli,lambd_bernoulli,log_likelihood = BernoulliIterator(D,K,N,mu_bernoulli,Pi_bernoulli,df)
8             Pi_bernoulli,mu_bernoulli,lambd_bernoulli = BernoulliIterator2(D,K,N,lambd_bernoulli,df)
9             if log_likelihood >= prev_log_likelihood_bernoulli:
10                 log_like_bernoulli[t][epoch] = log_likelihood
11             else :
12                 break
13 mean_log_like_bernoulli = np.nanmean(log_like_bernoulli, axis = 0)
14 print(mean_log_like_bernoulli)
15 return mean_log_like_bernoulli

```

```

1 df = pd.read_csv('A2Q1.csv')
2 N = df.shape[0]
3 D = df.shape[1]
4 K = 4
5 no_of_epochs = 8
6 tolerance = 1e-6          # Setting max error
7 no_of_initialisations = 100
8 N,D

```

(399, 50)

```

1 mean_log_like_bernoulli = Bernoulli_Algorithm(N,D,K,no_of_epochs,df,tolerance,no_of_initialisations)

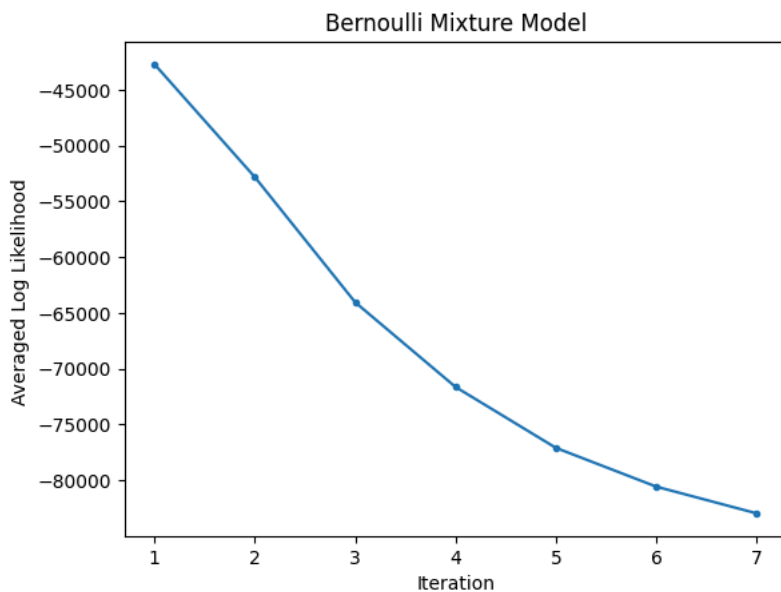
[-82067.80222419 -42709.65311479 -52831.47534691 -64050.25099231
 -71641.08654335 -77099.63528571 -80575.23756963 -82970.84655691]

```

```

1 plt.plot(range(1, len(mean_log_like_bernoulli) ), mean_log_like_bernoulli[1:len(mean_log_like_bernoulli)], marker='.')
2 plt.xlabel('Iteration')
3 plt.ylabel('Averaged Log Likelihood ')
4 plt.title('Bernoulli Mixture Model')
5 plt.show()

```



## ▼ Gaussian Mixture Models

```

1 def Initialisation_Gaussian(N,D,K,df):
2     Pi_gaus = np.random.rand(K)
3     Pi_gaus = Pi_gaus / np.sum(Pi_gaus)
4     mu_gaus = np.random.rand(K, D)
5     sigma_gaus = np.random.rand(K, D)
6     lambd_gaus = np.random.rand(N, K)
7     row_sums = lambd_gaus.sum(axis=1)
8     lambd_gaus = lambd_gaus / row_sums[:, np.newaxis]
9     return Pi_gaus, mu_gaus, sigma_gaus, lambd_gaus

```

```

1 def expectation_gaussians(D, K, N, mu_gaus, sigma_gaus, Pi_gaus, lambd_gaus, df):
2     Epsilon = 1e-10
3     likelihoods = np.zeros((N, K))
4     for i in range(K):
5         for j in range(N):
6             mahalnobis_dist = 0
7             for d in range(D):
8                 mahalnobis_dist += (df.iloc[j,d] - mu_gaus[i][d])** 2 / (2 * sigma_gaus[i][d] + Epsilon)
9
10
11         # The log likelihood using the Mahalanobis distance
12         likelihoods[j, i] = -0.5 * (D * np.log(2 * np.pi * np.prod(sigma_gaus[i]) + Epsilon)) + mahalnobis_dist
13     weighted_likelihoods = likelihoods + np.log(Pi_gaus + Epsilon)
14     log_likelihood = np.sum(weighted_likelihoods)
15     lambd_gaus = np.exp(weighted_likelihoods - np.max(weighted_likelihoods, axis=1)[: , np.newaxis])
16     lambd_gaus /= np.sum(lambd_gaus, axis=1)[: , np.newaxis]
17     lambd_gaus = np.nan_to_num(lambd_gaus, nan = 0)
18
19     return Pi_gaus, mu_gaus, lambd_gaus, log_likelihood

```

```

1 def maximization_gaussian(D, K, N, mu_gaus, sigma_gaus, lambd_gaus, df):
2     Epsilon = 1e-10
3     Nk = np.sum(lambd_gaus, axis=0)
4     Pi_gaus = Nk / N
5     kaka = np.dot(lambd_gaus.T, df)
6     mu_gaus = np.dot(lambd_gaus.T, df) / (Nk[: , np.newaxis] + Epsilon)
7     for k in range(K):
8         sigma_gaus[k] = np.sum(lambd_gaus[:, k] * np.sum((df - mu_gaus[k])**2, axis=1)) / ( Nk[k] + Epsilon) # Changed calculation of sigma_gaus
9     return Pi_gaus, mu_gaus, sigma_gaus

```

```

1 def EM_Gaussian_Algorithm(N,D,K,no_of_epochs,df,tol,no_of_initialisations):
2     prev_log_likelihood = -np.inf
3     log_like = np.empty((no_of_initialisations,no_of_epochs))
4     for t in range(no_of_initialisations):
5         Pi_gaus, mu_gaus, sigma_gaus, lambd_gaus = Initialisation_Gaussian(N,D,K,df)
6         for epoch in range(no_of_epochs):
7             Pi_gaus,mu_gaus,lambd_gaus,log_likelihood = expectation_gaussians(D, K, N, mu_gaus, sigma_gaus, Pi_gaus, lambd_gaus, df)
8             Pi_gaus,mu_gaus,lambd_gaus = maximization_gaussian(D, K, N, mu_gaus, sigma_gaus, lambd_gaus, df)
9             log_like[t][epoch] = log_likelihood
10        prev_log_likelihood = log_likelihood
11    mean_log_like_gaus = np.nanmean(log_like, axis = 0)
12    return mean_log_like_gaus,mu_gaus, sigma_gaus, lambd_gaus, Pi_gaus

```

```

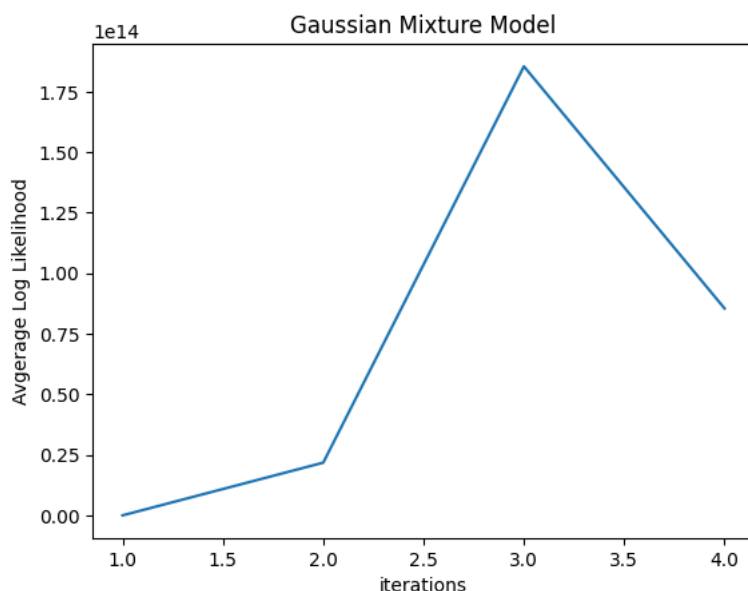
1 df = pd.read_csv('A2Q1.csv')
2 N = df.shape[0]
3 D = df.shape[1]
4 K = 4
5 no_of_epochs = 4
6 tolerance = 1e-6
7 no_of_initialisations = 100
8 mean_log_like,mu_gaus,sigma_gaus,lambd_gaus,Pi_gaus = EM_Gaussian_Algorithm(N,D,K,no_of_epochs,df,tolerance,no_of_initialisations)

```

```

1 plt.plot(range(1, len(mean_log_like) + 1), mean_log_like)
2 plt.xlabel('iterations')
3 plt.ylabel('Avgerage Log Likelihood')
4 plt.title('Gaussian Mixture Model')
5 plt.show()

```



## ▼ K Means Clustering

```
1 def euclidean_distance(x1, x2):
2     distance = np.sqrt(np.sum((x1 - x2) ** 2))
3     return distance
```

```
1 #Initializing with Random values
2 def initialize_centroids(df, K):
3     random_indices = np.random.choice(df.shape[0], size=K, replace=False)
4     centroids_initial = df[random_indices]
5     return centroids_initial
```

```
1 # Assign each data point to the nearest centroid
2 def clusterIdentifier(df, centroids):
3     clusters = []
4
5     for data_point in df:
6         distances = [euclidean_distance(data_point, centroid) for centroid in centroids]
7         cluster_distances = np.argmin(distances)
8         clusters.append(cluster_distances)
9
10    return np.array(clusters)
```

```
1 # new means for the clusters
2 def newMeans(df, clusters, K):
3     centroids = np.zeros((K, df.shape[1]))
4
5     for cluster in range(K):
6         cluster_points = df[clusters == cluster]
7
8         if len(cluster_points) > 0:
9             centroids[cluster] = np.mean(cluster_points, axis=0)
10
11    return centroids
```

```
1 # Objective Function
2 def error(data, centroids, clusters):
3     objective = 0
4
5     for i, centroid in enumerate(centroids):
6         cluster_points = data[clusters == i]
7
8         if len(cluster_points) > 0:
9             objective += np.sum((cluster_points - centroid) ** 2)
10
11    return objective
```

```
1 def Kmeans(data, k, max_iter):
2
3     centroids = initialize_centroids(data, k)
4     errs = [] # Objectives for each iteration
5     centroid_history = [centroids]
6
7     for iter in range(max_iter):
8         clusters = clusterIdentifier(data, centroids)
9         new_centroids = newMeans(data, clusters, k)
10        objective = error(data, new_centroids, clusters)
11        errs.append(objective)
12
13        if np.allclose(new_centroids, centroids):
14            break
15
16        centroids = new_centroids
17        centroid_history.append(centroids)
18
19    return clusters, centroids, errs, centroid_history
```

```
1 K = 4
2 df = pd.read_csv('A2Q1.csv')
3 N = df.shape[0]
4 D = df.shape[1]
5 data = df.values
6 max_iter = 500
7
8 clusters, centroids, errs, centroid_history = Kmeans(data, K, max_iter)
9 print(errs)
```

[2218.376663784021, 2178.6300101471334, 1983.401461492841, 1754.3724009010446, 1742.8341061377648, 1739.2729630782853, 1736.179526977088, 1736.179

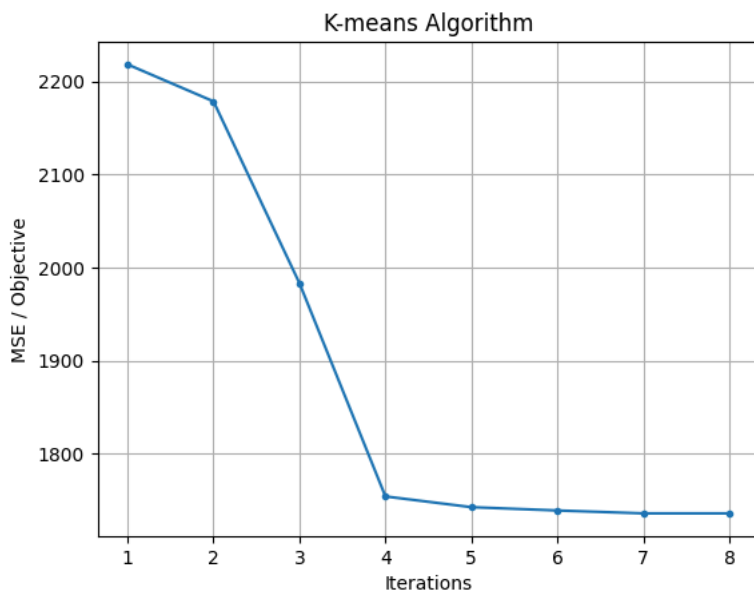


```
1 plt.plot(range(1, len(errs) + 1), errs, marker = '.')
```

```

2 plt.xlabel('Iterations')
3 plt.ylabel('MSE / Objective')
4 plt.title('K-means Algorithm')
5 plt.grid(True)
6 plt.show()

```



```

1 def plot_data_points(data,n):
2     num_rows = int(np.ceil(n / 5))
3     fig, axes = plt.subplots(nrows=2, ncols=5, figsize=(10, 2*num_rows))
4     fig.suptitle('Binary image depiction of some 10*5 groups')
5     for i, ax in enumerate(axes.flatten()):
6         if i < 10:
7             ax.imshow(data[i].reshape(10, 5), cmap='hot_r', vmin=0, vmax=1)
8             ax.set_xticks([])
9             ax.set_yticks([])
10            ax.set_xticklabels([])
11            ax.set_yticklabels([])
12
13     plt.tight_layout()
14     plt.show()
15
16 plot_data_points(data,10)
17

```

Binary image depiction of some 10\*5 groups

