# EE511 Assignment 7

Zechen Ha

April 6, 2020

# Question 1

## Experiment

We can generate the multivariate normal random variables by transformation. Construct a random vector Y based on standard normal random vector Z.

$$Y = AZ + b$$

There is always:

$$\mu_Y = \mu_Z + b$$
$$\Sigma_Y = A\Sigma_Z A^T$$

where

$$\Sigma_Y = AA^T$$
$$\mu_Z = [0]$$

So we can write script to generate multivariate normal random vector Y with these formula.

## Code

```
1  import numpy as np
2
3  X = np.random.randn(3,1)
4  mu = np.array([[1],[2],[3]])
5  sigma = np.array([[3,-1,1],[-1,5,3],[1,3,4]])
6  value, vector = np.linalg.eig(sigma)
7  A = vector.dot(np.diag(np.sqrt(value)))
8  Y = A.dot(X)+mu
9  print(Y)
```

Listing 1: Q1 code

# Question 2

## Experiment

We can complete the assignment by just generating two Gaussian random variable and get the mixture with given algebra. Simulate several times and get the histogram.

## Code

```
1  import numpy as np
2  import matplotlib.pyplot as plt
3  from scipy.stats import norm
4
5  N = 2000
```

```
6  X = []
7  for i in range(N):
8      x = 0.4*np.random.normal(-1,1)+0.6*np.random.normal(1,1) #mu and std
9      X.append(x)
10 X = np.array(X)
11 mu = 0.2
12 variance = 1
13 t = np.linspace(norm.ppf(0.05,mu),norm.ppf(0.95,mu),100)
14 plt.hist(X,alpha=0.7, edgecolor='black',density=True,bins=20)
15 plt.plot(t,norm.pdf(t,mu),'r')
16 plt.title("histogram and pdf of mixture",fontsize=15)
17 plt.show()
```
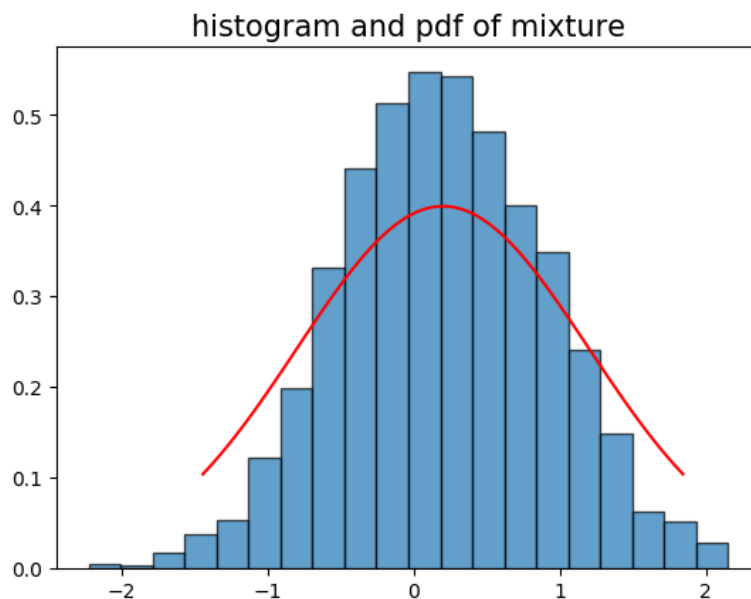
Listing 2: Q2 code

**Result**



Figure 1: Histogram and theoretical pdf of the mixture

## Question 3

### Code

```
1  import numpy as np
2  from scipy.stats import multivariate_normal
3  import matplotlib.pyplot as plt
4  from sklearn.mixture import GaussianMixture
```

```
5  import time
6
7  mu1 = [0,0]
8  mu2 = [-1,3]
9  cov1 = [[0.5, 0], [0, 0.5]]
10 cov2 =  [[1, 0], [0, 1]]
11 p = 0.6
12 sample = 300
13 data = []
14 prelabel = []
15
16 for i in range(sample):
17     if np.random.rand()<p:
18         data.append(multivariate_normal.rvs(mu1,cov1))
19         prelabel.append(0)
20     else:
21         data.append(multivariate_normal.rvs(mu2,cov2))
22         prelabel.append(1)
23
24 data = np.array(data);prelabel = np.array(prelabel)
25 start_time = time.time()
26 gmm = GaussianMixture(n_components=2, covariance_type='full',random_state
       =1).fit(data)
27 end_time = time.time()
28 label = gmm.predict(data)
29 err_num = sum(prelabel!=label)
30 err_rate = err_num/sample
31 print(err_rate)
32 print(end_time-start_time)
33 print(gmm.means_,'\n',gmm.covariances_)
34 print(gmm.weights_)
```

Listing 3: Q3 code

## Result

To make the quality of EM algorithm more clear, I define the quality as the accuracy rate. The speed of the algorithm is the time spent to calculate the expectation. The comparison is shown as follow.

When the co-variance is spherical, I use the example

$$\mu_1 = \begin{vmatrix} 0 \\ 0 \end{vmatrix} \Sigma_1 = \begin{vmatrix} 0.5 & 0 \\ 0 & 0.5 \end{vmatrix}$$

$$\mu_2 = \begin{vmatrix} -1 \\ 3 \end{vmatrix} \Sigma_2 = \begin{vmatrix} 1 & 0 \\ 0 & 1 \end{vmatrix}$$

The accuracy rate is 0.97, and the time spent is 0.004s.

When the co-variance is ellipsoidal, I use the example

$$\mu_1 = \begin{vmatrix} 0 \\ 0 \end{vmatrix} \Sigma_1 = \begin{vmatrix} 0.5 & 0 \\ 0 & 1 \end{vmatrix}$$

$$\mu_2 = \begin{vmatrix} -1 \\ 3 \end{vmatrix} \Sigma_2 = \begin{vmatrix} 1 & 0 \\ 0 & 5 \end{vmatrix}$$

The accuracy rate is 0.87, and the time spent is 0.009s.

The above examples are all well-seperated, when the co-variance is closed and ellipsoidal:

$$\mu_1 = \begin{vmatrix} 0 \\ 0 \end{vmatrix} \Sigma_1 = \begin{vmatrix} 0.5 & 0 \\ 0 & 1 \end{vmatrix}$$

$$\mu_2 = \begin{vmatrix} 0 \\ -0.5 \end{vmatrix} \Sigma_2 = \begin{vmatrix} 1 & 0 \\ 0 & 5 \end{vmatrix}$$

The accuracy rate is 0.71, and the time spent is 0.007s.

# Question 4

## Code

```
1  import numpy as np
2  from scipy.stats import multivariate_normal
3  import matplotlib.pyplot as plt
4  from sklearn.mixture import GaussianMixture
5  from sklearn.cluster import KMeans
6
7  data = np.loadtxt(open("faithful.dat.txt",'rb'),skiprows=26)
8  data = data[:,1:3]
9  kmeans = KMeans(n_clusters=2, random_state=1).fit(data)
10 label = kmeans.predict(data)
11 centers = kmeans.cluster_centers_
12 fig1 = plt.figure()
13 ax1 = fig1.add_subplot(111)
14 ax1.scatter(data[:,0],data[:,1],c=label,s=40)
15 ax1.scatter(centers[:,0],centers[:,1],c='red',s=80)
16 ax1.set_title("Kmeans Clusters",fontsize=15)
17 ax1.set_xlabel("eruptions")
18 ax1.set_ylabel("waiting")
19 plt.show()
20
21 gmm = GaussianMixture(n_components=2,covariance_type='full',random_state
       =1).fit(data)
22 labelgm = gmm.predict(data)
23 mu = gmm.means_
24 covariance = gmm.covariances_
25 t = gmm.weights_
26 x,y = np.mgrid[1:5.5:0.01,40:100:0.1]
27 pos = np.empty(x.shape+(2,))
28 pos[:,:,0] = x; pos[:,:,1] = y
29 fig2 = plt.figure()
```

**EE 511 Assignment 7**

```python
30 ax2 = fig2.add_subplot(111)
31 ax2.contourf(x,y,t[0]*multivariate_normal.pdf(pos,mu[0,:],covariance
     [0,:,:])+t[1]*multivariate_normal.pdf(pos,mu[1,:],covariance[1,:,:]))
32 ax2.scatter(data[:,0],data[:,1],c=labelgm,s=20,cmap='winter')
33 ax2.set_title("Gaussian mixture",fontsize=15)
34 ax2.set_xlabel("eruptions")
35 ax2.set_ylabel("waiting")
36 plt.show()
```

Listing 4: Q4 Code

## Result

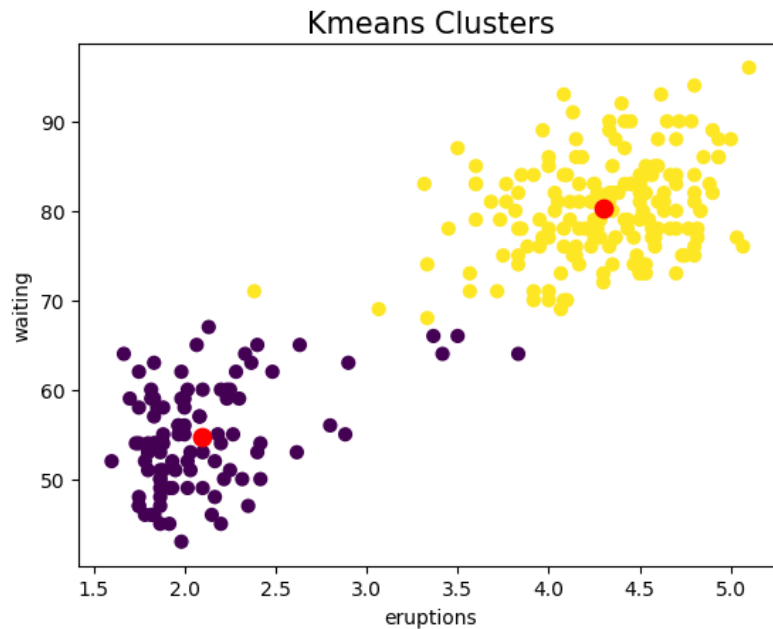The k-means scatter graph is shown as follow:



Figure 2: k-means clusters

The EM algorithm estimates the parameters $\mu$ and $\Sigma$, and get the probability density function based on the final parameters.
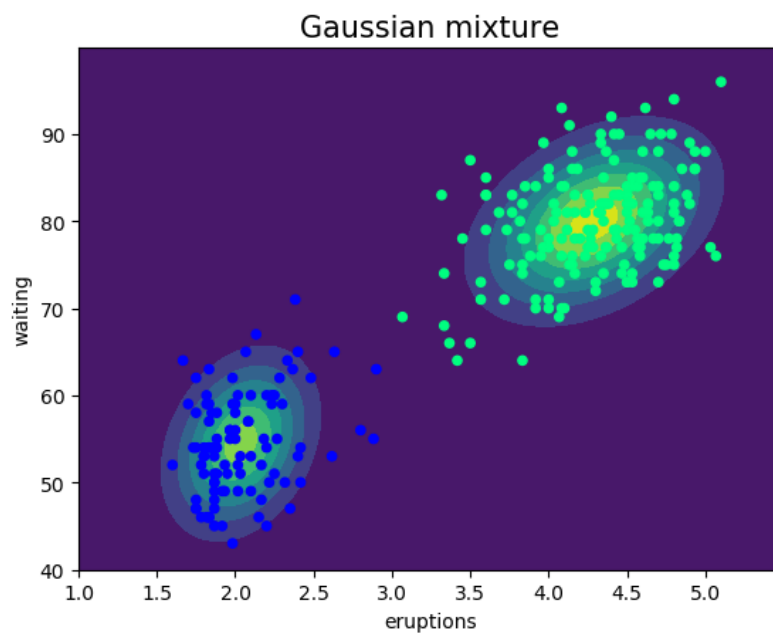
# EE 511 Assignment 7

Figure 3: GMM-EM

Different from the k-means algorithm, the GMM-EM is a soft classifier. We classify the data points on the graph based on the estimated probability of belonging to which cluster. When the probability of belonging to the cluster A is larger than B, we classify that it belongs to A. So the probability of the point belonging to B also exists.