

# Project #7 - Expectation Maximization

LI YICHENG\*

USCID:7827077047

email: l.y.c.liyicheng@gmail.com

USC Viterbi of Engineering

## I. PROBLEM DESCRIPTION

This project details the use of EM for one popular problems of learning a mixture. We consider the problem of learning a GMM (also called EM clustering), where both the mixture weights of the Gaussians and the parameters for each component Gaussian must be learned. We illustrate using EM to learn a GMM when there are additional constraints on the GMM parameters. In section II, we talk about generating samples from a multivariate Gaussian. In section III, we discuss generate a mixture model of multivariate Gaussian. In section III, we will get into how to perform the GMM-EM estimation to learn a GMM model with 2 dimensions and 2 subpopulations. In section IV, we use the practical data from reality and perform the EM estimation to get insight of the feature of the data.

## II. MULTIVARIATE GAUSSIAN

### I. Definition

#### I.1 Bivariate Case

Suppose we have two random variable,  $N_1, N_2$ ,  $N_1 \sim (\mu_1, \sigma_1^2)$  and  $N_2 \sim (\mu_2, \sigma_2^2)$ . Consider a plane with  $point(x, y)$  on it. If we take one sample from  $N_1$  and take the value as the value of  $x$  and take one sample from  $N_2$  and take the value as  $y$ . After many sampling action, for one small area on the plane, mark as  $(x, y)$  we count the number of samples, and we may get the joint probability of  $f_{N_1, N_2}(N_1 = x, N_2 = y)$ .

In this case

$$f(x, y) = \frac{1}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} \exp\left(-\frac{1}{2(1-\rho^2)} \left[ \frac{(x-\mu_X)^2}{\sigma_X^2} + \frac{(y-\mu_Y)^2}{\sigma_Y^2} - \frac{2\rho(x-\mu_X)(y-\mu_Y)}{\sigma_X\sigma_Y} \right]\right)$$

where  $\rho$  is the correlation between  $X$  and  $Y$  and where  $\sigma_X > 0$  and  $\sigma_Y > 0$ . In this case,

$$\mu = \begin{pmatrix} \mu_X \\ \mu_Y \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \sigma_X^2 & \rho\sigma_X\sigma_Y \\ \rho\sigma_X\sigma_Y & \sigma_Y^2 \end{pmatrix}. \text{ And we can expand this model to higher dimation.}$$

---

\*github link: <https://github.com/IAMLYCHEE/EE501-PROJ7>

## I.2 High Dimensional Case

In this case the distribution has density

$$f_{\mathbf{x}}(x_1, \dots, x_k) = \frac{\exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)}{\sqrt{(2\pi)^k |\boldsymbol{\Sigma}|}}$$

$$= \frac{\exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)}{\sqrt{|2\pi \boldsymbol{\Sigma}|}},$$

where  $\mathbf{x}$  is a real  $k$ -dimensional column vector and  $|\boldsymbol{\Sigma}|$  is the determinant of  $\boldsymbol{\Sigma}$ .

## II. Construction and Sampling

let  $Z_1, Z_2, \dots, Z_n$  be i.i.d standard normal distribution, that is  $Z_i \sim N(0, 1)$ . For constant  $a_{i,j}$  and  $\mu_i$ , define:

$$X_1 = a_{11}Z_1 + a_{12}Z_2 + \dots + a_{1m}Z_m + \mu_1$$

$$X_2 = a_{21}Z_1 + a_{22}Z_2 + \dots + a_{2m}Z_m + \mu_2$$

$\vdots$

$$X_n = a_{n1}Z_1 + a_{n2}Z_2 + \dots + a_{nm}Z_m + \mu_n$$

The vector  $X = [X_1, X_2, \dots, X_n]$  is called multivariate normal

In matrix mode, this can be written as:

$$\begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{pmatrix} = \begin{pmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,m} \\ a_{2,1} & a_{2,2} & \dots & a_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n,1} & a_{n,2} & \dots & a_{n,m} \end{pmatrix} * \begin{pmatrix} Z_1 \\ Z_2 \\ \vdots \\ Z_m \end{pmatrix} + \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{pmatrix}$$

And it is easy to get that the covariance  $\Sigma = A * A^T$  where  $A$  is the matrix formed by  $a_{i,j}$ . Therefore, if we were given  $\mu$  which provides the information about the mean value of each Gaussian components and  $\Sigma$  the covariance between each two of the Gaussian components. We can decompose the  $\Sigma$  into matrix  $A$  and generate samples from standard Gaussian random variables, then we can get several samples from Gaussian multivariate.

### II.1 Implement

According to the last discussion the code is straightforward:

**filename: multiGauGenerator**

```
1 function sample = multiGauGenerator(cov,mu,sampleAmount)
2 %use: sample = multiGauGenerator(cov,mu,sampleAmount)
3 %generate multivariate gaussian random samples with given mean and
4 %covariance.
5 %Li Yicheng
6
7 %reshape the input
```

```

8 mu = reshape(mu,[length(mu),1]);
9 a = chol(cov,'lower');
10 % a * a';
11 aSize = size(cov);
12 sample = zeros(sampleAmount,aSize(1));
13 for i = 1 : sampleAmount
14     x = a * randn(aSize(2),1) + mu;
15     sample(i,:) = x';
16 end

```

## II.2 Experiment

give :  $\mu = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix}, \Sigma = \begin{pmatrix} 3 & -1 & 1 \\ -1 & 5 & 3 \\ 1 & 3 & 4 \end{pmatrix}$

we generate samples from this gaussian multivariate:

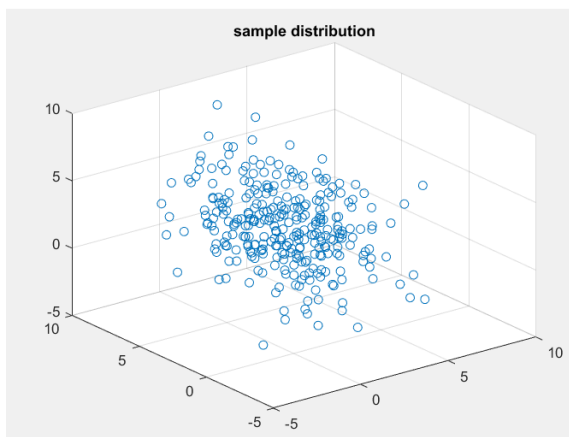
**filename: solution1.m**

```

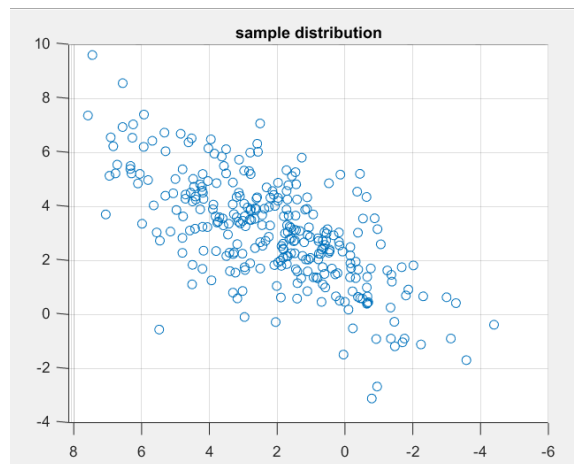
1 clear
2 sampleAmount = 300;
3 cov = [3,-1,1;-1,5,3;1,3,4];
4 mu = [1;2;3];
5 sample = multiGauGenerator(cov,mu,sampleAmount);
6 scatter3(sample(:,1),sample(:,2),sample(:,3))
7 title('sample distribution')

```

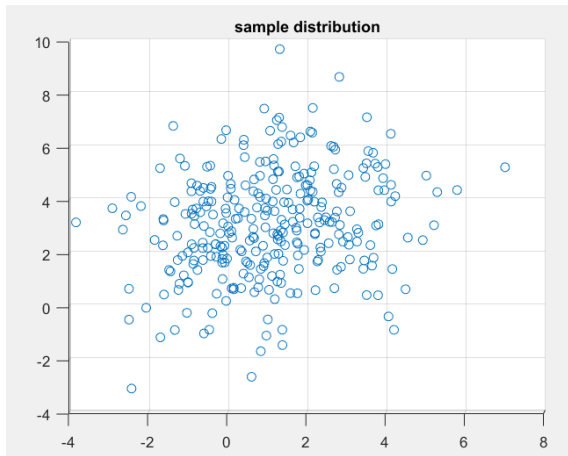
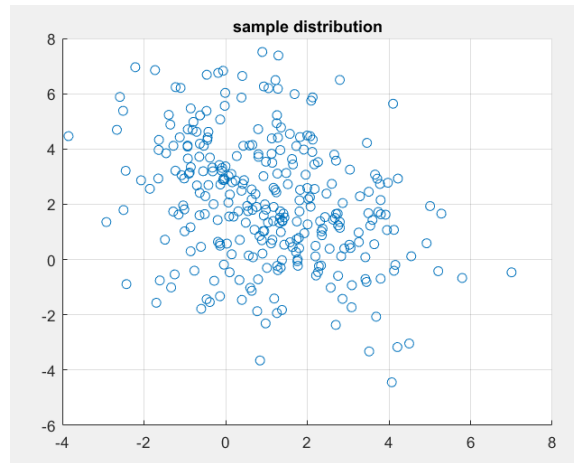
**Result:**



**Figure 1:** samples in a 3D view



**Figure 2:** samples view from  $x_1$


 Figure 3: samples view from  $x_2$ 

 Figure 4: samples view from  $x_3$ 

### III. MIXTURE MODEL

#### I. Operation

If we already have several multivariate Gaussian, we mixture them by given each of the multivariate a weight.

When sampling we can first generate a sample from standard uniform distribution, and decide which Gaussian multivariate generator should we choose and then generate samples from that generator.

#### II. Experiment

In this experiment, we implement the following mixture distribution:  $f(x) = 0.4N(-1,1) + 0.6N(1,1)$ , according to the above discussion. the code is straightforward:

**filename: solution2.m**

```

1 sampleAmount = 100000;
2 sample = zeros(sampleAmount,1);
3 for i = 1 : sampleAmount
4     if rand(1) < 0.4
5         sample(i) = randn(1) - 1;
6     else
7         sample(i) = randn(1) + 1;
8     end
9 end
10 histogram(sample, 'BinLimits', [-4,4], 'Normalization', 'probability')
11 hold on
    
```

```

12 %theoretical
13 i = 1;
14 for t = -4: 0.05 : 4
15     c = sqrt(2 * pi);
16     p1 = exp(-((t+1)^2)/2);
17     p2 = exp(-((t-1)^2)/2);
18     ft(i) = 0.4 / c * p1 + 0.6 / c * p2;
19     i = i+1;
20 end
21
22 plot(-4 : 0.05 : 4, ft/9);
23 legend('histogram','theoretical')
24 title('mixture gaussian(100000 sample)')
    
```

result:

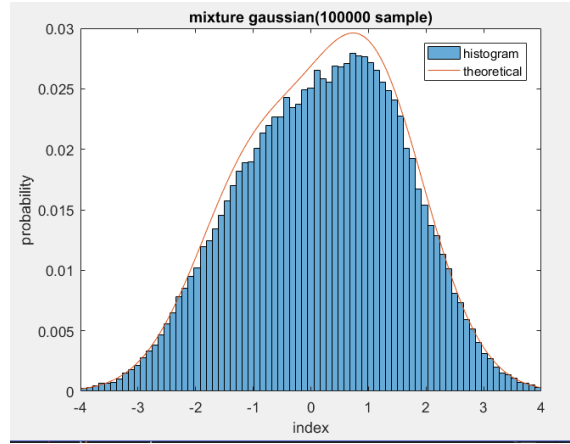


Figure 5: mixture of 2 1d GM

## IV. GMM-EM ESTIMATION

### I. Algorithm Overview

#### 1. Initialization

Choose initial estimates  $\omega_j^{(0)}, \mu_j^{(0)}, \Sigma_j^{(0)}, j = 1 \dots k$ , where  $j$  is the index of the cluster.  $k$  is the number of subpopulation.

Compute the initial log-likelihood:

$$l^{(0)} = \frac{1}{n} \sum_{i=1}^n \log \left( \sum_{j=1}^k \omega_j^{(0)} \phi(y_i | \mu_j^{(0)}, \Sigma_j^{(0)}) \right)$$

#### 2. E-step

For  $j = 1, \dots, k$  compute

$$\gamma_{ij}^{(m)} = \frac{\omega_j^{(m)} \phi(y_i | \mu_j^{(m)}, \Sigma_j^{(m)})}{\sum_{l=1}^k \omega_l^{(m)} \phi(y_i | \mu_l^{(m)}, \Sigma_l^{(m)})}$$

and

$$n_j^{(m)} = \sum_{i=1}^n \gamma_{ij}^{(m)}$$

### 3.M-step

For  $j = 1, \dots, k$  compute the new estimates

$$\omega_j^{(m+1)} = \frac{n_j^{(m)}}{n}$$

$$\mu_j^{(m+1)} = \frac{1}{n_j^{(m)}} \sum_{i=1}^n \gamma_{ij}^{(m)} y_i$$

$$\Sigma_j^{(m+1)} = \frac{1}{n_j^{(m)}} \sum_{i=1}^n \gamma_{ij}^{(m)} (y_i - \mu_j^{(m+1)})(y_i - \mu_j^{(m+1)})^T$$

**4.Convergence Check:** Compute the new likelihood:

$$l^{(m+1)} = \frac{1}{n} \sum_{i=1}^n \log \left( \sum_{j=1}^k \omega_j^{(m+1)} \phi(y_i | \mu_j^{(m+1)}, \Sigma_j^{(m+1)}) \right)$$

**Return to 2 step** if  $|l^{(m+1)} - l^{(m)}| > \delta$  for some threshold  $\delta$ .

## II. Analysis

### II.1 Initialization

First we initialize by first doing a cheaper clustering with k-means algorithm and separate the samples into two subpopulation. Then we can derive the center of the two subpopulation and take them as first estimation of  $\mu$ , set the estimation of the covariance to be identity matrices, and the first guess at the weights  $\omega_1 = \omega_2 = \dots = \omega_k$ . After we have these parameters, we can calculate the log-likelihood. Using the density function of the multivariate

**Initialize (part of GMM2d2popuEMcore.m)**

```

1 k = 2;
2 w = [1/k, 1/k];
3 [y, C] = kmeans(sample, k);
4 mu1 = C(1, :);
5 mu2 = C(2, :);
6 mu = zeros(2, 1, k);
7 mu(:, :, 1) = mu1';
8 mu(:, :, 2) = mu2';
9 covES = zeros(2, 2, k);
10 covES(:, :, 1) = eye(2);
11 covES(:, :, 2) = eye(2);
12 sumLH = 0;
13 for i = 1 : sampleAmount
14     sumLH = sumLH + log(w(1) * mvnpdf(sample(i, :), mu(:, :, 1)', covES(:, :, 1)) + ...

```

```

15         w(2) * mvnpdf(sample(i,:),mu(:,2),covES(:,2)));
16     end
17     L = sumLH/sampleAmount;

```

## II.2 E-step

Actually the fraction of  $\gamma_{i,j}$  tells us what probability  $sample[i]$  would appear when its a sample from  $GaussianMultivariate[j]$  and the nominator gives the total probability.

For example, if  $sample(i)$  have biggest probability being from subpopulation 1, then  $\gamma_{i,1} > \gamma_{i,2}$  (if the number of subpopulation is 2).

After get the  $\gamma_{i,j}$ ,  $i = 1 \cdots sampleAmount$ ,  $j = 1 \cdots k$ , we simply add the  $\gamma$  matrix by column and we can get the weight in each subpopulation.

### E-step (part of GMM2d2popuEMcore.m)

```

1  gamma = zeros(sampleAmount,k);
2  n = zeros(1,2);
3  while abs(L_new-L) > epsilon
4  %E-step
5      L = L_new;
6      for j = 1 : k
7          for i = 1 : sampleAmount
8              sumLHPoint = 0;
9              for l = 1 : k
10                 sumLHPoint = sumLHPoint + w(l) * mvnpdf(sample(i,:),mu(:,l),covES(:,l),l));
11             end
12             gamma(i,j) = w(j) * mvnpdf(sample(i,:),mu(:,j),covES(:,j)) / sumLHPoint;
13         end
14     end
15     n(1) = sum(gamma(:,1));
16     n(2) = sum(gamma(:,2));

```

## II.3 M-step

Using the  $\gamma$ , we can get the new mean and covariance, because  $\gamma$  provides information about each samples's probability of belonging to some subpopulation.

Just use that information to update mean and covariance.

### M-step (part of GMM2d2popuEMcore.m)

```

1  for j = 1 : k
2      mu(:,j) = (gamma(:,j))' * sample / n(j)';
3      sumCov = zeros(2,2);
4      for i = 1 : sampleAmount
5          sumCov = sumCov + gamma(i,j) * (sample(i,:) - mu(:,j)) * ...
6              (sample(i,:) - mu(:,j))';
7      end

```

```

8     covES(:, :, j) = sumCov / n(j);
9     end

```

## II.4 Convergence check (part of GMM2d2popuEMcore.m)

```

1     sumLH = 0;
2     for i = 1 : sampleAmount
3         sumLH = sumLH + log(w(1) * mvnpdf(sample(i, :), mu(:, :, 1)', covES(:, :, 1))) + ...
4             w(2) * mvnpdf(sample(i, :), mu(:, :, 2)', covES(:, :, 2)));
5     end
6     L_new = sumLH/sampleAmount;

```

The above implements the core algorithm of GMM-EM.

## III. Experiment

Run the following script to implement the experiment:

Experiment 1: no covariance

$$\omega = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix} \mu_1 = \begin{pmatrix} 2 \\ 5 \end{pmatrix} \mu_2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix} \Sigma_1 = \begin{pmatrix} 3 & 0 \\ 0 & 0.5 \end{pmatrix} \Sigma_2 = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$$

Experiment 2: with covariance

$$\omega = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix} \mu_1 = \begin{pmatrix} 2 \\ 5 \end{pmatrix} \mu_2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix} \Sigma_1 = \begin{pmatrix} 3 & 1 \\ 1 & 0.5 \end{pmatrix} \Sigma_2 = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 2 \end{pmatrix}$$

Experiment 3: different weight

$$\omega = \begin{pmatrix} 0.2 \\ 0.8 \end{pmatrix} \mu_1 = \begin{pmatrix} 2 \\ 5 \end{pmatrix} \mu_2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix} \Sigma_1 = \begin{pmatrix} 3 & 1 \\ 1 & 0.5 \end{pmatrix} \Sigma_2 = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 2 \end{pmatrix}$$

Experiment 4: close

$$\omega = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix} \mu_1 = \begin{pmatrix} 2 \\ 3 \end{pmatrix} \mu_2 = \begin{pmatrix} 3 \\ 2 \end{pmatrix} \Sigma_1 = \begin{pmatrix} 3 & 1 \\ 1 & 0.5 \end{pmatrix} \Sigma_2 = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 2 \end{pmatrix}$$

```

1 clear
2 close all
3
4 clear
5 %generate samples
6 %test1 : same weight well seperated spherical
7 %generate the first population
8 w = [0.5, 0.5];
9 mu1 = [2; 5];
10 mu2 = [0; 1];
11 cov1 = [3, 0; 0, 1/2];
12 cov2 = [1, 0; 0, 2];
13 GMM2d2popuEM(mu1, mu2, cov1, cov2, w, 300)

```



```
14 |
15 | clear
16 | %test2 : ellipsoidal covariance
17 | w = [0.5,0.5];
18 | mu1 = [2;5];
19 | mu2 = [0;1];
20 | cov1 = [3,1;1,1/2];
21 | cov2 = [1,1/2;1/2,2];
22 | GMM2d2popuEM(mu1,mu2,cov1,cov2,w,300)
23 |
24 | clear
25 | %test3 : different weight
26 | w = [0.2,0.8];
27 | mu1 = [2;5];
28 | mu2 = [0;1];
29 | cov1 = [3,1;1,1/2];
30 | cov2 = [1,1/2;1/2,2];
31 | GMM2d2popuEM(mu1,mu2,cov1,cov2,w,300)
32 |
33 | clear
34 | %test4 : close
35 | w = [0.5,0.5];
36 | mu1 = [2;3];
37 | mu2 = [3;2];
38 | cov1 = [3,1;1,1/2];
39 | cov2 = [1,1/2;1/2,2];
40 | GMM2d2popuEM(mu1,mu2,cov1,cov2,w,300)
```

the next page Use the second experiment to show the EM process , for the rest of the experiments,data would be provided to show the speed and equality of the algorithm

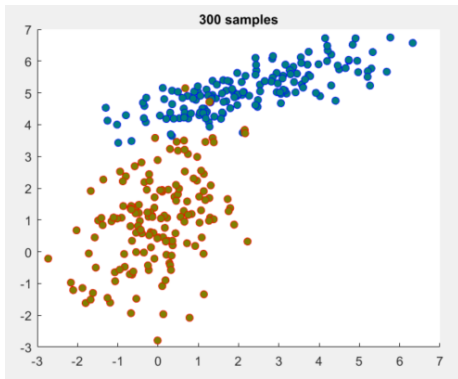


Figure 6: Samples

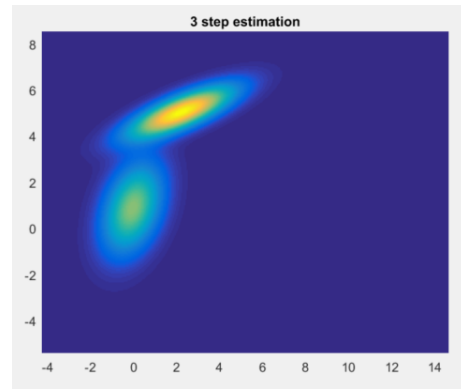


Figure 9: the 3rd step EM



Figure 7: The true expected GMM

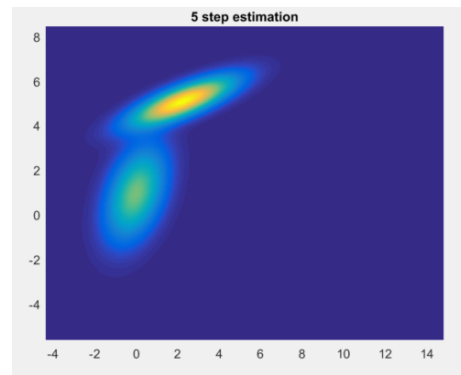


Figure 10: the 5th step EM

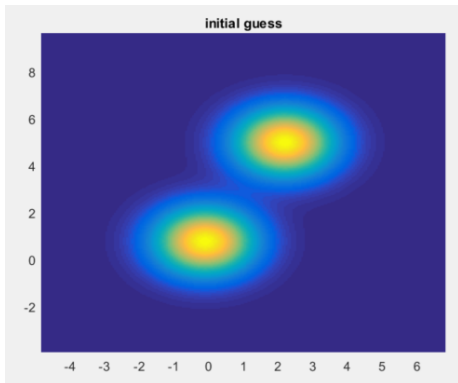


Figure 8: Initial guess

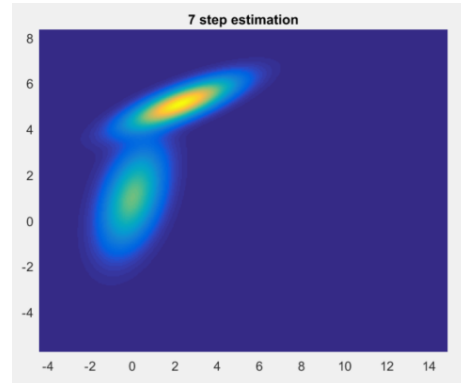


Figure 11: the 7th step EM

**result:**

**Experiment 1: true parameter**

$$\omega = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix} \mu_1 = \begin{pmatrix} 2 \\ 5 \end{pmatrix} \mu_2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix} \Sigma_1 = \begin{pmatrix} 3 & 0 \\ 0 & 0.5 \end{pmatrix} \Sigma_2 = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$$

**EM estimation**

step:5

```
mu(:,1) =
0.0599 1.0909
mu(:,2) =
2.2225 5.0133
covES(:,1) =
0.9141 -0.0927 -0.0927 1.9269
covES(:,2) =
2.9422 -0.0392 -0.0392 0.4740
w =
0.4995 0.5005
```

**Experiment 2: true parameter**

$$\omega = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix} \mu_1 = \begin{pmatrix} 2 \\ 5 \end{pmatrix} \mu_2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix} \Sigma_1 = \begin{pmatrix} 3 & 1 \\ 1 & 0.5 \end{pmatrix} \Sigma_2 = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 2 \end{pmatrix}$$

**EM estimation**

step:5

```
mu(:,1) =
2.0323 4.9962
mu(:,2) =
0.0189 1.1414
covES(:,1) =
3.0709 1.0645 1.0645 0.5731
covES(:,2) =
1.0111 0.5107 0.5107 1.9100
w =
0.4999 0.5001
```

**Experiment 3: true parameter**

$$\omega = \begin{pmatrix} 0.2 \\ 0.8 \end{pmatrix} \mu_1 = \begin{pmatrix} 2 \\ 5 \end{pmatrix} \mu_2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix} \Sigma_1 = \begin{pmatrix} 3 & 1 \\ 1 & 0.5 \end{pmatrix} \Sigma_2 = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 2 \end{pmatrix}$$

**result:**

step:12

```
mu(:,1) =
2.0810 5.0054
mu(:,2) =
0.1008 1.0737
covES(:,1) =
1.6086 0.4771 0.4771 0.3798
covES(:,2) =
0.9722 0.4338 0.4338 1.8664
w =
0.2133 0.7867
```

**Experiment 4: true parameter**

$$\omega = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix} \mu_1 = \begin{pmatrix} 2 \\ 3 \end{pmatrix} \mu_2 = \begin{pmatrix} 3 \\ 2 \end{pmatrix} \Sigma_1 = \begin{pmatrix} 3 & 1 \\ 1 & 0.5 \end{pmatrix} \Sigma_2 = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 2 \end{pmatrix}$$

**result: step:25**

```
mu(:,1) =
3.0958 2.0454
mu(:,2) =
1.9381 2.9674
covES(:,1) =
0.9031 0.6542 0.6542 2.1917
covES(:,2) =
2.5967 0.9282 0.9282 0.5100
w =
0.4805 0.5195
```

**Performance Analysis:** The weight actually has not very much influence on the speed although if the weight vary a lot from our initial guess, it may take more iterations. However, if the two Gaussian Multivariate are close, it really took a lot more iterations. And it can be seen from the data. If the threshold is set to 0.0001, the variance between the true mean and the estimation is actually lower than 5% . The covariance may vary a bit between the true parameter and the estimation, that basically is due to the amount of samples are really small and if large amount of samples are given we may get better estimation.

## V. PRACTICAL DATA EM

### I. Problem

A geyser is a hot spring characterized by an intermittent discharge of water and steam. Old Faithful is a famous cone geyser in Yellowstone National Park, Wyoming. It has a predictable geothermal discharge and since 2000, it has erupted every 44 to 125 minutes. Refer to the addendum data file that contains waiting times and the durations for 272 eruptions.

### II. Estimation using EM

**simply run the following script**

```
1 clear
2 clc
3 close all
4 duration = load('duration.mat');
5 duration = duration.duration;
6 waiting = load('waiting.mat');
7 waiting = waiting.waiting;
8 scatter(duration, waiting);
9 title('without kmeans clustering')
10 xlabel('duration')
11 ylabel('waiting')
12
13
14 %cascade data points
15 samples = [duration, waiting];
16 %shuffle
17 samples = samples(randperm(size(samples,1)),:);
18 %kmeans
19 [y,C] = kmeans(samples,2);
20
21 %figure the clusters
22 figure
23 hold on
```

```

24 scatter(samples(y == 1,1),samples(y==1,2),'x')
25 scatter(samples(y == 2,1),samples(y==2,2),'o')
26 plot(C(1,1),C(1,2),'rx','LineWidth',2)
27 plot(C(2,1),C(2,2),'ro','LineWidth',2)
28 legend('Points of cluster 1','Points of cluster 2')
29 title('Data Points with Labels by K-means Clustering')
30 xlabel('duration')
31 ylabel('waiting')
32 hold off
33
34 figure
35 GMM2d2popuEMcore(samples,0)
36 hold on
37 scatter3(samples(y == 1,1),samples(y==1,2),ones(length(samples(y==1,1)),1),'rx')
38 scatter3(samples(y == 2,1),samples(y==2,2),ones(length(samples(y==2,1)),1),'ro')
39 plot(C(1,1),C(1,2),'rx','LineWidth',2)
40 plot(C(2,1),C(2,2),'ro','LineWidth',2)
41 % legend('Points of cluster 1','Points of cluster 2')
42 title('Data Points with Labels by K-means Clustering & EM GMM estimation')
43 xlabel('duration')
44 ylabel('waiting')
45 hold off

```

### Using the EM implemented in last topic

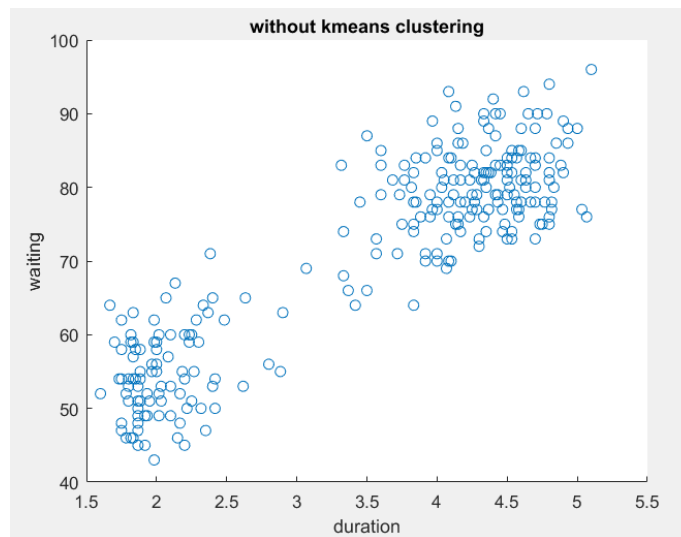


Figure 12: samples from data

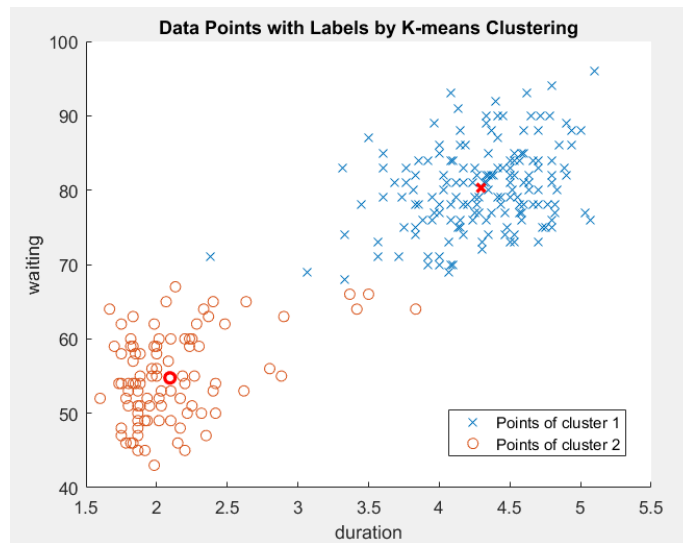


Figure 13: roughly clusering using K-mean

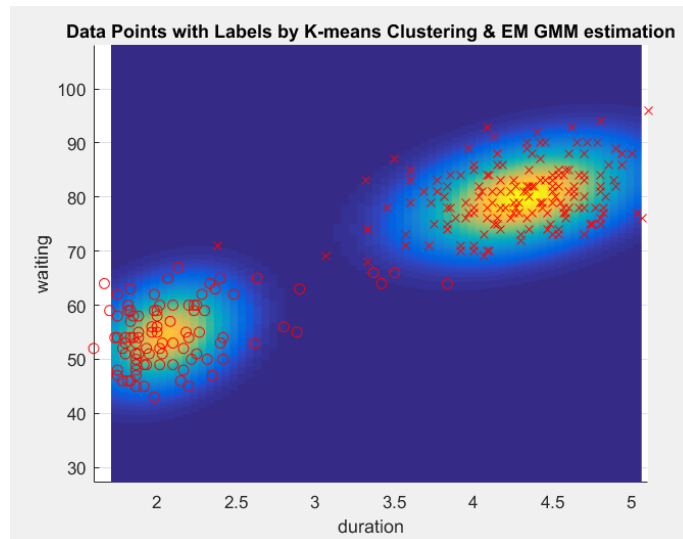


Figure 14: EM result

**data obtained**

$\mu(:,1) =$

4.2902 79.9740

$\mu(:,2) =$

2.0370 54.4843

$\text{covES}(:,1) =$

0.1693 0.9328 0.9328 35.9590

```
covES(:,2) =
0.0696 0.4400 0.4400 33.7325
```

Therefore, we can use the GMM model to estimate the behaviour of the geyser by sampling from the GMM model with parameters derived from above.

## VI. APPENDIX

filename:GMM2d2popuEMcore.m

```
1 function GMM2d2popuEMcore(sample,option)
2 sampleAmount = size(sample,1);
3 k = 2;
4 w = [1/k,1/k];
5 [y,C] = kmeans(sample,k);
6 mu1 = C(1,:);
7 mu2 = C(2,:);
8 mu = zeros(2,1,k);
9 mu(:,1) = mu1';
10 mu(:,2) = mu2';
11 covES = zeros(2,2,k);
12 covES(:,1) = eye(2);
13 covES(:,2) = eye(2);
14 %compute initial likelihood
15 sumLH = 0;
16 for i = 1 : sampleAmount
17     sumLH = sumLH + log(w(1) * mvnpdf(sample(i,:),mu(:,1)',covES(:,1)) + ...
18         w(2) * mvnpdf(sample(i,:),mu(:,2)',covES(:,2))));
19 end
20 L = sumLH/sampleAmount;
21 L_new = 0;
22 epsilon = 0.001;
23 %record data
24 disp('initial guess')
25 mu
26 covES
27 %plot the initial guess
28 plot2d2subpopuGMM(mu(:,1),mu(:,2),covES(:,1),covES(:,2),w,1,option);
29
30 %EM
31 stepCount = 1;
32 gamma = zeros(sampleAmount,k);
33 n = zeros(1,2);
34 while abs(L_new-L) > epsilon
35 %E-step
36     L = L_new;
37     for j = 1 : k
38         for i = 1 : sampleAmount
```

```

39         sumLHPoint = 0;
40         for l = 1 : k
41             sumLHPoint = sumLHPoint + w(l) * mvnpdf(sample(i,:),mu(:,l)',covES(:,l),
42                 1));
43         end
44         gamma(i,j)= w(j) * mvnpdf(sample(i,:),mu(:,j)',covES(:,j),j)/sumLHPoint;
45     end
46     n(1) = sum(gamma(:,1));
47     n(2) = sum(gamma(:,2));
48
49 %M-step
50     w = n/sampleAmount;
51     for j = 1 : k
52         mu(:,j) = (gamma(:,j)' * sample ./ n(j))';
53         sumCov = zeros(2,2);
54         for i = 1 : sampleAmount
55             sumCov = sumCov + gamma(i,j) * (sample(i,:) - mu(:,j)) * ...
56                 (sample(i,:) - mu(:,j))';
57         end
58         covES(:,j)=sumCov / n(j);
59     end
60 %the new likelihood
61
62     stepCount = stepCount + 1;
63     plot2d2subpopuGMM(mu(:,1),mu(:,2),covES(:,1),covES(:,2),w,stepCount,option);
64     disp(strcat('step:',num2str(stepCount)))
65     mu
66     covES
67
68 %Convergence check
69     sumLH = 0;
70     for i = 1 : sampleAmount
71         sumLH = sumLH + log(w(1) * mvnpdf(sample(i,:),mu(:,1)',covES(:,1)) + ...
72             w(2) * mvnpdf(sample(i,:),mu(:,2)',covES(:,2)));
73     end
74     L_new = sumLH/sampleAmount;
75 end

```

#### filename:plot2d2popuGMM.m

```

1 function plot2d2subpopuGMM(mu1,mu2,cov1,cov2,w,k,option)
2 %plot2d2subpopuGMM(mu1,mu2,cov1,cov2,w,k,option)
3 if option == 0
4     figure
5 else
6     hold on
7     pause(0.5)
8 end
9 mu1 = reshape(mu1,[2,1]);

```



```
10 mu2 = reshape(mu2,[2,1]);
11 alpha = 4.7; %scale the plot
12 if mu1(1) < mu2(1)
13     x1rangeLB = mu1(1)- alpha * sqrt(cov1(1,1));
14     x1rangeUB = mu2(1)+ alpha * sqrt(cov2(1,1));
15 else
16     x1rangeLB = mu2(1)- alpha * cov2(1,1);
17     x1rangeUB = mu1(1)+ alpha * cov1(1,1);
18 end
19 x1 = x1rangeLB:0.05:x1rangeUB;
20
21 if mu1(2) < mu2(2)
22     x2rangeLB = mu1(2)- alpha * sqrt(cov1(2,2));
23     x2rangeUB = mu2(2)+ alpha * sqrt(cov2(2,2));
24 else
25     x2rangeLB = mu2(2)- alpha * sqrt(cov2(2,2));
26     x2rangeUB = mu1(2)+ alpha * sqrt(cov1(2,2));
27 end
28 x2 = x2rangeLB: 0.05 : x2rangeUB;
29
30 [X1,X2] = meshgrid(x1,x2);
31 f = w(1) * mvnpdf([X1(:) X2(:)],mu1',cov1)+w(2) * mvnpdf([X1(:) X2(:)],mu2',cov2);
32 f = reshape(f,length(x2),length(x1));
33 h = surf(x1,x2,f);
34 if k == 0
35     title(' true GMM density ');
36 else
37     if k == 1
38         title(' initial guess ');
39     else
40         title(strcat(num2str(k),' step estimation '));
41     end
42 end
43 axis tight
44 view(0,90)
45 set(h,'LineStyle','none')
```