
Table of Contents

.....	1
Comments Problem 1	1
Component definitions	1
K-means	2
Generate empirical probability table	3
Comments Problem 2	4
EM	4
Generate a random vector u in d dimensions	5
Comments Problem 3	5
Comments problem 4	6
Generate d-dimensional data samples	6
d-dimensional k-means	6
Generate empirical probability table for d-dimensional data	7
d-dimentional EM	8

```
clear all;
close all;
clc;
addpath('Functions');
```

Comments Problem 1

Looking at the emprical probabilities, we see that the K-means algorithm separates the dataset into two clusters with ease. The cetainty for a specific is assigned to the designated cluster with a probability higher than 0.8. However, when we increase the number of clusters to be devided into, we observe that certain parts of the clusters in the dataset gets spread more out over the new clusters. Especially cluster 1 gets spread pretty uniformly out over two of the new clusters. The more clusters we allow for the K-means algorithm, the more uncertain is the probability for a dataset from a specific cluster to land in a new cluster. Much of the reason why this is is due to the dataset being very tight with a pretty similar mean value between the mixed gaussians.

Component definitions

```
N = 200;

% Component 1
theta_1 = 0;
m_1 = [0 0]';
pi_1 = 1/2;
lambda_1 = 2;
lambda_2 = 1;
u_1 = [cos(theta_1) sin(theta_1)]';
u_2 = [-sin(theta_1) cos(theta_1)]';
C_1 = [u_1 u_2]*diag([lambda_1,lambda_2])*inv([u_1 u_2]);

% Component 2
```

```

theta2_2 = -3*pi/4;
m_2 = [-2 1]';
pi_2 = 1/6;
lambda_a1 = 2;
lambda_a2 = 1/4;
u1_2 = [cos(theta2_2) sin(theta2_2)]';
u2_2 = [-sin(theta2_2) cos(theta2_2)]';
C_2 = [u1_2 u2_2]*diag([lambda_a1,lambda_a2])*inv([u1_2 u2_2]);

% Component 3
theta2_3 = pi/4;
m_3 = [3 2]';
pi_3 = 1/3;
lambda_b1 = 3;
lambda_b2 = 1;
u1_3 = [cos(theta2_3) sin(theta2_3)]';
u2_3 = [-sin(theta2_3) cos(theta2_3)]';
C_3 = [u1_3 u2_3]*diag([lambda_b1,lambda_b2])*inv([u1_3 u2_3]);

X = zeros(N,2);
Z = zeros(N,3);
for i = 1:N
    a = rand();
    if a < pi_1
        X(i,:) = mvnrnd(m_1,C_1);
        Z(i,:) = [1 0 0];
    elseif a < pi_1 + pi_2
        X(i,:) = mvnrnd(m_2,C_2);
        Z(i,:) = [0 1 0];
    else
        X(i,:) = mvnrnd(m_3,C_3);
        Z(i,:) = [0 0 1];
    end
end
end

```

K-means

```

N_rand_inits = 5;
K_max = 5;
m_opt = zeros(K_max,2,K_max);
C_opt = zeros(N,K_max);
for K = 2:K_max
    min_sme = inf;
    for i = 1:N_rand_inits
        C = randi(K,N,1);
        [m,C] = k_means(N,K,C,X);
        sme = SME(m,X,C);
        if sme < min_sme
            m_opt(1:K,:,K) = m;
            C_opt(:,K) = C;
            min_sme = sme;
        end
    end
end

```

```

        end
    end

    % One-hot encoding
    a = zeros(N,K_max,K_max);
    for K = 2:K_max
        for i = 1:N
            a(i,C_opt(i,K),K) = 1;
        end
    end
end

```

Generate empirical probability table

```

pk_kmeans = zeros(3,K_max,K_max);
for K = 2:K_max
    for l = 1:3
        for k = 1:K
            num_k_l = 0;
            num_l = 0;
            for i = 1:N
                if Z(i,l) == 1
                    num_l = num_l + 1;
                    if a(i,k,K) == 1
                        num_k_l = num_k_l + 1;
                    end
                end
            end
            pk_kmeans(l,k,K) = num_k_l/num_l;
        end
    end
end
plot_table(pk_kmeans,K_max,'K-means');

```

Plotting table of $P(k|i)$ generated from K-means with $K = 2$

ans =

0.1193	0.8807
0	1.0000
0.9194	0.0806

Plotting table of $P(k|i)$ generated from K-means with $K = 3$

ans =

0.0092	0.4037	0.5872
0	0.9655	0.0345
0.6935	0	0.3065

Plotting table of $P(k|i)$ generated from K-means with $K = 4$

ans =

0.1560	0.2661	0	0.5780
0	0.9655	0	0.0345
0.5323	0	0.4032	0.0645

Plotting table of $P(k|i)$ generated from K-means with $K = 5$

ans =

0.2752	0	0.1101	0.4771	0.1376
0.2069	0	0.7931	0	0
0	0.3871	0	0.0968	0.5161

Comments Problem 2

In the figure one may observe that the contours of the gaussian distributions generated with the EM algorithm fits the clusters made by the K-means algorithm very well. Also, upon expection of the average values of the probabilities, one may see that they are pretty close to the values we got from the previous excercise. Although a bit weaker probabilities, but one would be able to classify the different clusters equally based on the probabilities from the different tables generated by K-means and EM.

EM

```
m_EM = zeros(K_max,2,K_max);
C_EM = zeros(2,2,K_max, K_max);
pi_EM = zeros(K_max,K_max);
pk_EM = zeros(3,K_max,K_max);
for K = 2:K_max
    m_init = m_opt(:, :, K);
    [m_, C_, pi_, pk_] = EM(m_init, N, K, X, Z);
    m_EM(:, :, K) = m_;
    C_EM(:, :, 1:K, K) = C_;
    pi_EM(1:K, K) = pi_;
    pk_EM(:, 1:K, K) = pk_;
end
plot_table(pk_EM, K_max, 'EM');
```

Plotting table of $P(k|i)$ generated from EM with $K = 2$

ans =

0.1380	0.8620
0.0087	0.9913
0.8042	0.1958

Plotting table of $P(k|i)$ generated from EM with $K = 3$

ans =

0.0536	0.3806	0.5657
0.0108	0.9160	0.0732
0.7075	0.0123	0.2802

Plotting table of $P(k|i)$ generated from EM with $K = 4$

ans =

0.1538	0.2556	0.0026	0.5880
0.0016	0.9601	0.0000	0.0383
0.4976	0.0027	0.3965	0.1032

Plotting table of $P(k|i)$ generated from EM with $K = 5$

ans =

0.2295	0.0022	0.1213	0.5049	0.1420
0.1639	0.0000	0.7802	0.0541	0.0017
0.0161	0.3890	0.0001	0.0998	0.4949

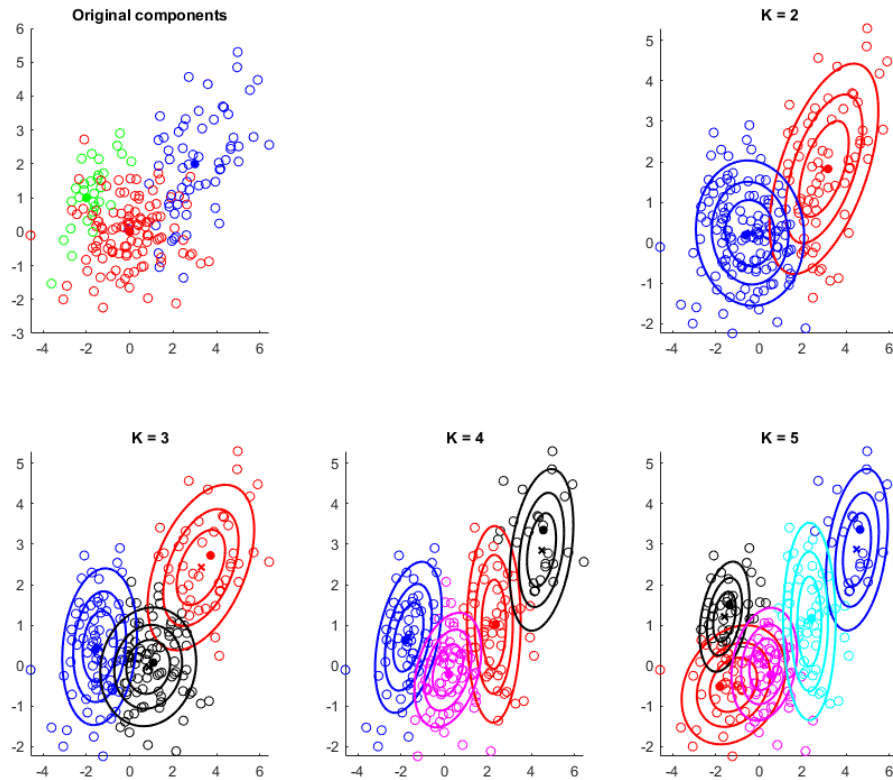
Generate a random vector u in d dimensions

```
d = 30;
u = zeros(d, 7);
for j = 1:7
    u(:,j) = generate_random_vector(d);
    while check_orthogonality(u,j) == 0
        u(:,j) = generate_random_vector(d);
    end
end
```

Comments Problem 3

One may observe in the figure that the means and the covariances for the example where we have $K=3$ clusters align pretty well. For the other numbers of clusters, the means and the covariances align well with the K-means algorithms.

```
plotData(X,Z,m_1,m_2,m_3,K_max,C_opt,m_opt,C_EM,m_EM,N);
```



Comments problem 4

A new vector is too correlated to another vector if the average absolute value of the product is greater than 0.2. (Based on experiments) The code is precented as follows.

Generate d-dimensional data samples

```
sigma2 = 0.01;
N_d = 200;
[X_d, Z_d] = generate_sample_data(u,sigma2,N_d);
```

d-dimensional k-means

```
m_opt_d = zeros(K_max,size(X_d,2),K_max);
C_opt_d = zeros(N_d,K_max);
for K = 2:K_max
    min_sme = inf;
    for i = 1:N_rand_inits
        C_d = randi(K,N_d,1);
        [m_d,C_d] = k_means(N_d,K,C_d,X_d);
        sme = SME(m_d,X_d,C_d);
        if sme < min_sme
```

```

        m_opt_d(1:K,:,K) = m_d;
        C_opt_d(:,K) = C_d;
        min_sme = sme;
    end
end
end

```

Generate empirical probability table for d-dimensional data

```

pk_kmeans_d = zeros(3,K_max,K_max);
for K = 2:K_max
    for l = 1:3
        for k = 1:K
            num_k_l = 0;
            num_l = 0;
            for i = 1:N_d
                if Z_d(i,l) == 1
                    num_l = num_l + 1;
                    if a(i,k,K) == 1
                        num_k_l = num_k_l + 1;
                    end
                end
            end
            pk_kmeans_d(l,k,K) = num_k_l/num_l;
        end
    end
end
plot_table(pk_kmeans_d,K_max,'d-dimensional K-means');

```

Plotting table of $P(k|i)$ generated from d-dimensional K-means with $K = 2$

ans =

0.3770	0.6230
0.3333	0.6667
0.3429	0.6571

Plotting table of $P(k|i)$ generated from d-dimensional K-means with $K = 3$

ans =

0.2459	0.3607	0.3934
0.2174	0.3478	0.4348
0.2000	0.3714	0.4286

Plotting table of $P(k|i)$ generated from d-dimensional K-means with $K = 4$

ans =

0.2295	0.2295	0.1475	0.3934
0.1884	0.3043	0.1739	0.3333
0.3286	0.3143	0.0571	0.3000

Plotting table of $P(k|i)$ generated from d-dimensional K-means with $K = 5$

ans =

0.2295	0.1475	0.0984	0.2951	0.2295
0.1449	0.1739	0.2319	0.2754	0.1739
0.1714	0.0429	0.1857	0.3000	0.3000

d-dimensional EM

```

m_EM_d = zeros(K_max,2,K_max);
C_EM_d = zeros(2,2,K_max, K_max);
pi_EM_d = zeros(K_max,K_max);
pk_EM_d = zeros(3,K_max,K_max);
for K = 2:K_max
    m_init = m_opt(:, :, K);
    [m_, C_, pi_, pk_] = EM(m_init, N, K, X, Z);
    m_EM_d(:, :, K) = m_;
    C_EM_d(:, :, 1:K, K) = C_;
    pi_EM_d(1:K, K) = pi_;
    pk_EM_d(:, 1:K, K) = pk_;
end
plot_table(pk_EM_d, K_max, 'EM');

```

Plotting table of $P(k|i)$ generated from EM with $K = 2$

ans =

0.1380	0.8620
0.0087	0.9913
0.8042	0.1958

Plotting table of $P(k|i)$ generated from EM with $K = 3$

ans =

0.0536	0.3806	0.5657
0.0108	0.9160	0.0732
0.7075	0.0123	0.2802

Plotting table of $P(k|i)$ generated from EM with $K = 4$

ans =

0.1538	0.2556	0.0026	0.5880
0.0016	0.9601	0.0000	0.0383

0.4976 0.0027 0.3965 0.1032

Plotting table of $P(k/i)$ generated from EM with $K = 5$

ans =

0.2295	0.0022	0.1213	0.5049	0.1420
0.1639	0.0000	0.7802	0.0541	0.0017
0.0161	0.3890	0.0001	0.0998	0.4949

Published with MATLAB® R2018a