# University of New South Wales

## School of Electrical Engineering and Telecommunications

**INDIVIDUAL** Assessment Task: _____Part 2 Assignment 1_____

| Course Code | ELEC9782 | Course Name | Data science |
|---|---|---|---|
| Week/Session/Year | Year 4 | Lecturer | Dr. Vidhya |

| Student Number | z5036602 |
|---|---|
| Family Name | Liu |
| Given Names | Zhengyue |

### *Mark/Grade given (For official use only)*

Marker's Comments:

Since this work counts toward your formal assessment for this course, please write your name and student number where indicated above, and sign the declaration below. Attach this cover sheet to the front of your submission, so that your name and student number can be seen without any cover needing to be opened.

For further information, please see http://www.lc.unsw.edu.au/plagiarism/index.html

_____

*I declare that this assessment item is my own work, except where acknowledged, and has not been submitted for academic credit elsewhere, and acknowledge that the assessor of this item may, for the purpose of assessing this item:*

- *Reproduce this assessment item and provide a copy to another member of the University ; and/or,*
- *Communicate a copy of this assessment item to a plagiarism checking service (which may then retain a copy of the assessment item on its database for the purpose of future plagiarism checking).*

*I certify that I have read and understood the University Rules in respect of Student Academic Misconduct.*

Signature of student _____ Date: _____10/08/2019_____

# ELEC9782 Assignment 1

### z5036602 - Zhengyue LIU

### Semester 2 2019

I, Zhengyue LIU(student number z5036602), declare that the following assignment is my own work and that I have read and understood the University Rules in respect of Student Academic Misconduct.

# Question 1

## Problem (a)

The distance from multivariate normal distribution of a k-dimensional random vector $\mathbf{X} = (\mathbf{X}_1, ..., \mathbf{X}_k)^T$, $\mathbf{X_i} \sim \mathbf{N}(0, \sigma^2)$ to mean vector $\mathbf{U} = (u_1, u_2, ..., u_k)^T$ can be written as:

$$D = \sqrt{(\mathbf{X}_1 - u_1)^2 + (\mathbf{X}_2 - u_2)^2 + ... + (\mathbf{X}_k - u_k)^2}$$

$$= \sigma \sqrt{\frac{(\mathbf{X}_1 - u_1)^2}{\sigma^2} + \frac{(\mathbf{X}_2 - u_2)^2}{\sigma^2} + ... + \frac{(\mathbf{X}_k - u_k)^2}{\sigma^2}}$$

$\sum_k \left(\frac{\mathbf{X}_k - u_k}{\sigma}\right)^2$ is a standardized chi-squared distribution of degree k. Let Z be $\sum_k \left(\frac{\mathbf{X}_k - u_k}{\sigma}\right)^2$, then $D = \sigma\sqrt{Z}$. Cumulative distribution for D is:

$$\mathbf{F}_D(d) = \mathbf{P}(D \leq d) = \mathbf{P}(\sigma\sqrt{Z} < d) = \mathbf{P}(Z < \frac{d^2}{\sigma^2}) = \mathbf{F}_Z(\frac{d^2}{\sigma^2})$$

$\mathbf{F}_z$ is the cdf of chi-squared distribution of degree k(dimension).

$$f_D(d) = \mathbf{F}'_D(d) = \frac{2d}{\sigma^2}\mathbf{F}'_Z(\frac{d^2}{\sigma^2}) = \frac{2d}{\sigma^2}f_Z(\frac{d^2}{\sigma^2})$$

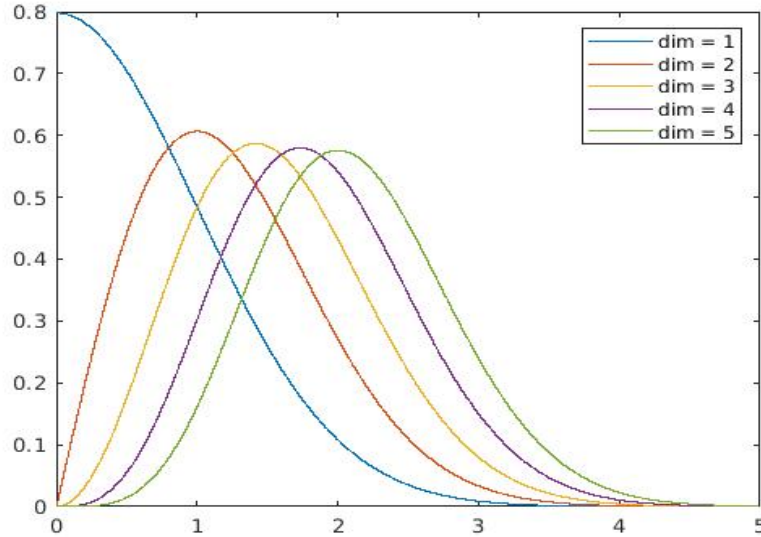$f_Z$ is the pdf of chi-squared distribution of degree k.

**PDF plots**



Figure 1: pdf of distance between data and mean for 1 dimension to 5 dimensions
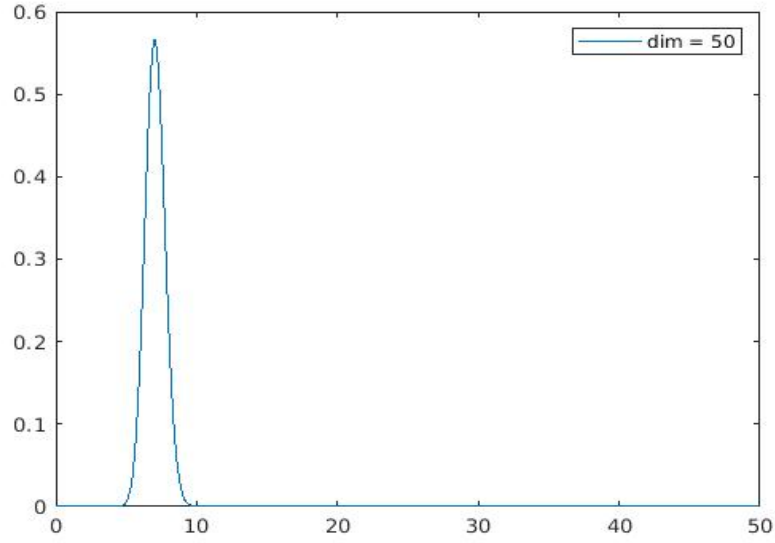
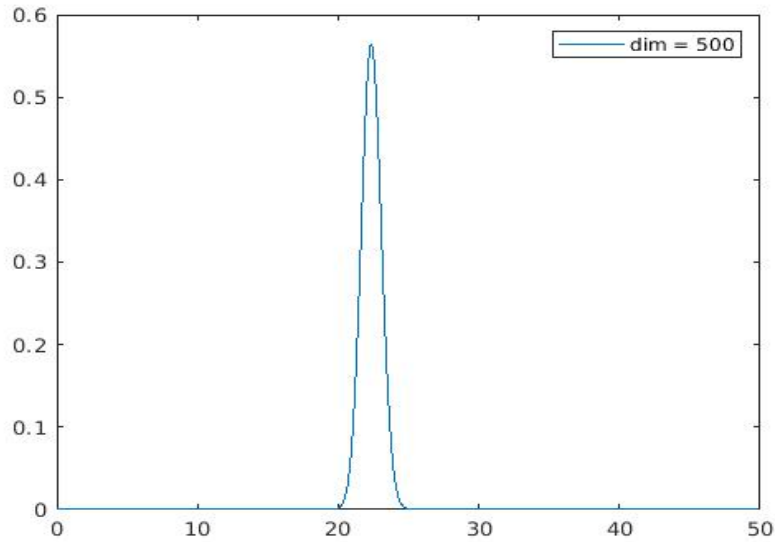Figure 2: pdf of distance between data and mean for 50 dimensions



Figure 3: pdf of distance between data and mean for 500 dimensions

Data in High dimensional space tends to be sparser than in lower dimensions. Data tends to lie further from the origin as the dimension increased. They mostly accumulate in a thin hollow sphere. As the dimension increased, the radius of the sphere increases. This phenomenon can be observed from above figures. As the dimension increases, the pdf function shifts to the right and creates a narrow probability "band".

**Matlab implementation**

The `pdf_plot.m` is used. `chi2pdf` built-in function is used to generate the chi-squared distribution. The likelihood over the different distance is calculated by using the pdf derived above. Finally, the likelihood is plotted for different dimensions.
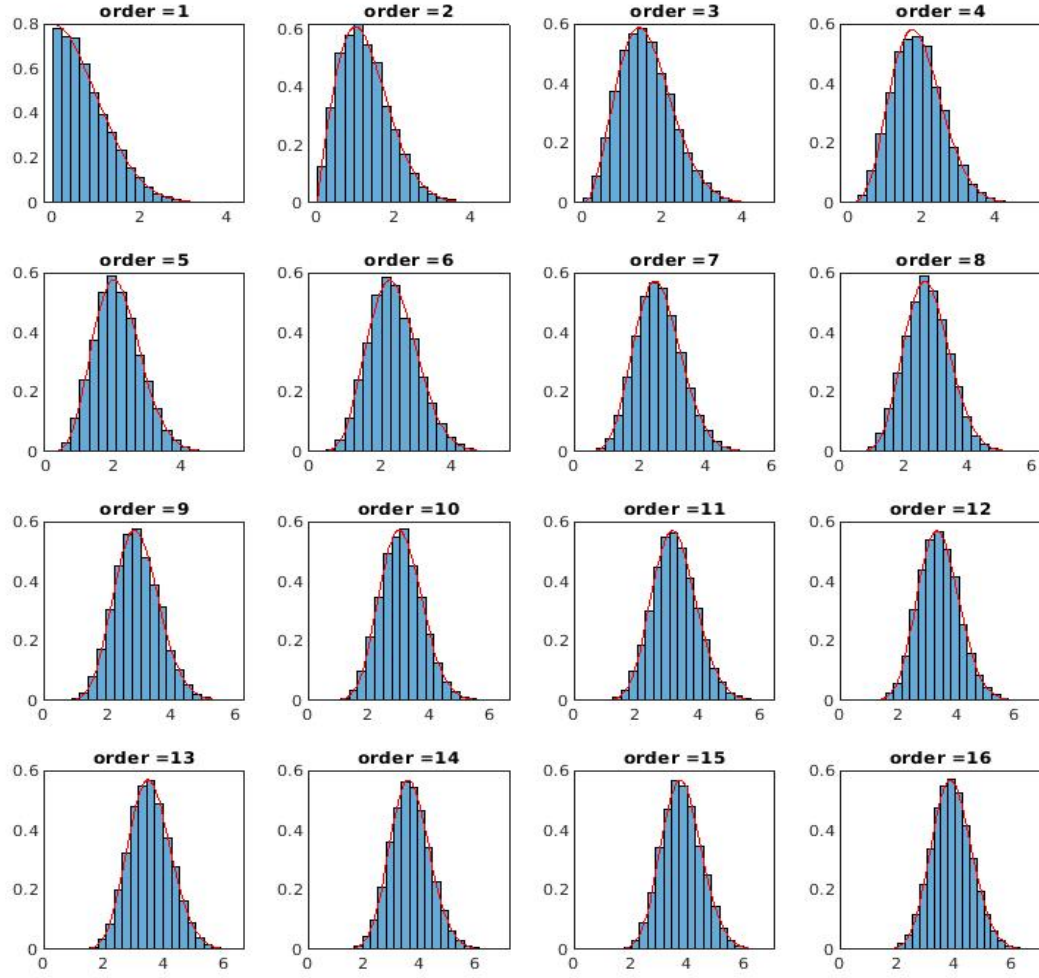
## Problem (b)



Figure 4: empirical pdf estimation by using histogram for 1 dimension to 16 dimensions, the red line is the theoretical pdf, blue bins are histogram bins. The empirical estimation agrees on theoretical values

**Matlab implementation**

The `empirical_estimates.m` is used. Firstly, multivariate guassian distributed data are generated using `randn` built-in function. The dimension of data range from 1 to 16. The distance is calculated using `norm` built-in function. Finally, the histogram of the distances is plotted for different dimensions. The theoretical pdf is also plotted on same graph for comparison. The code has been commented and attached in Q1 section in Appendix.

## Problem (c)

Cosine distance is a good measure of similarity/dissimilarity. Cosine distance calculates the angle between two data vectors, and use it as a measure of similarity. Smaller angle means two data vector are more similar to each other. In high-dimensional space, the curse of high dimensionality and sparsity will cause Euclidean distance based methods to fail. So, the angle measure between two vectors is a second best choice.

Another way to think about it is that as data mainly lies in a thin hollow sphere, and the distance between different data and the origin are similar, we can approximate such a thin hollow sphere as a spherical surface and all the data lies on this surface. The distance between different data and the origin can be considered as the radius $r$ of the spherical surface. Then, the arc length between two data points on this surface is calculated as:
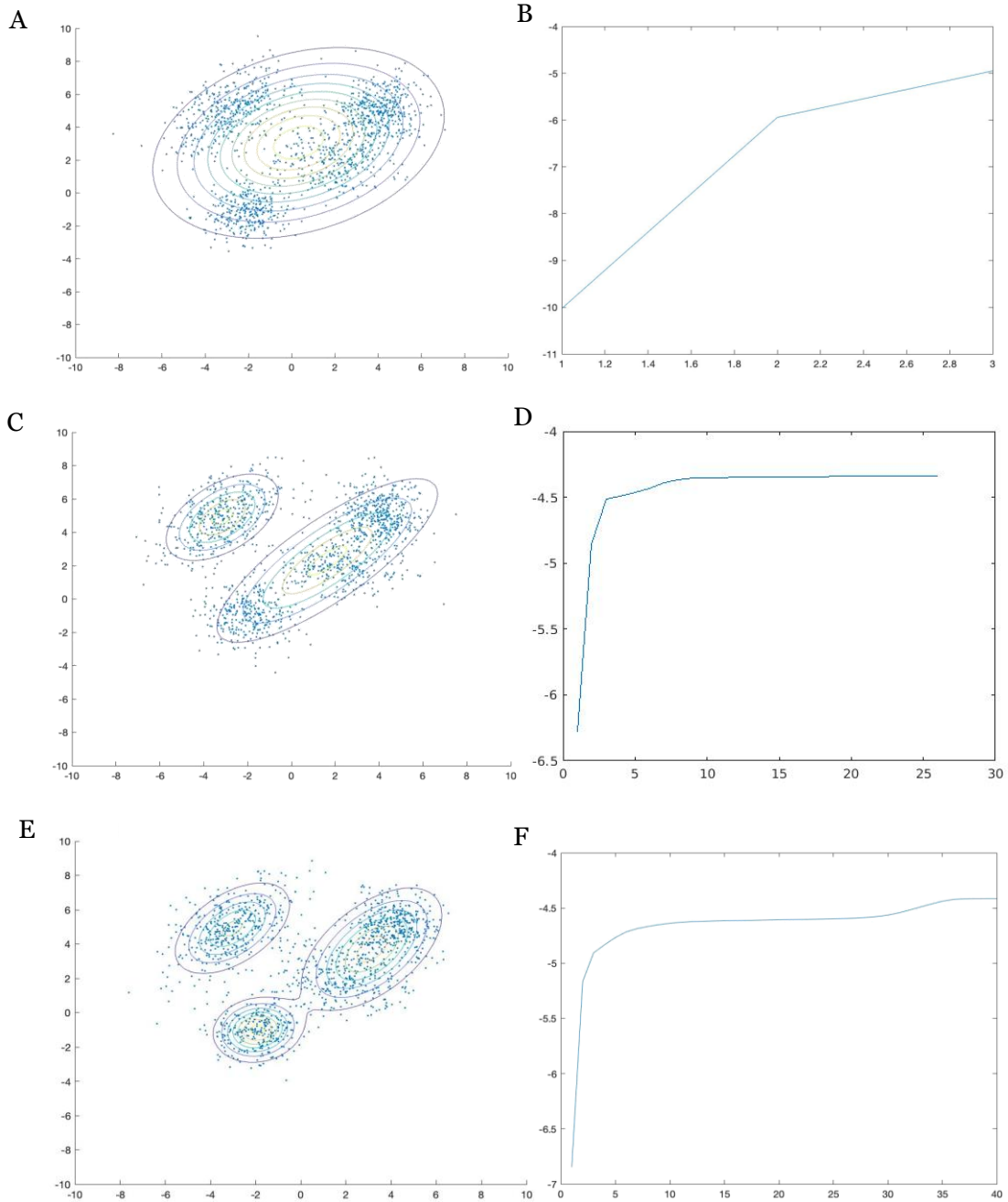
$$\text{arc length} = \theta * r$$

Hence, smaller theta means two data points are closer to each other in high-dimensional space. Whereas higher theta will give opposite result.
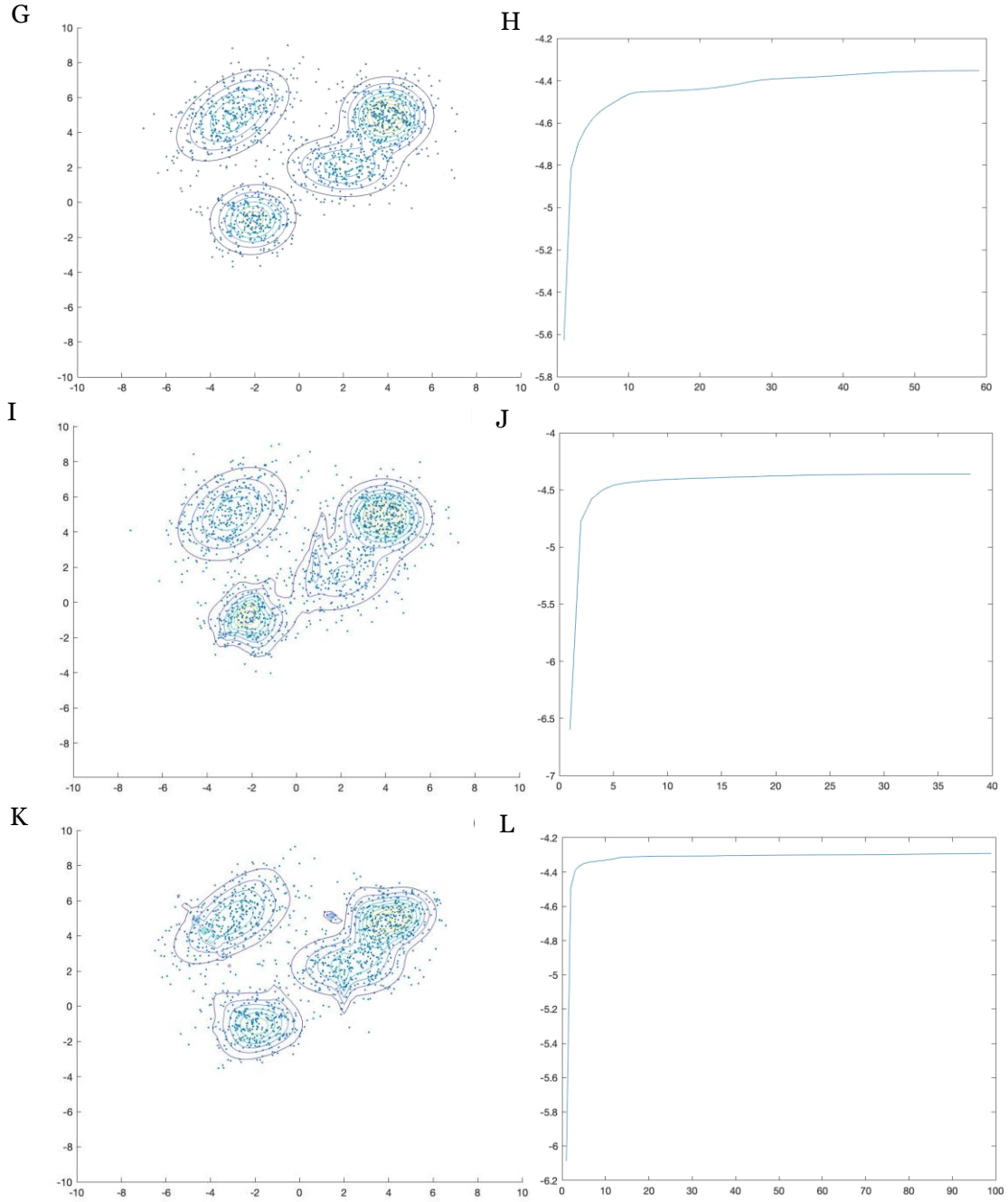
# Question 2

## Problem (a)

i.    Performance testing by visualisation (2-D)

   Testing data are generated by 2-D Gaussian mixture with 4 components:

A


B


C


D


E


F

**Figure 1**. Contour plots of fitted Gaussian mixture (left) and plots of log-likelihood of data given model v.s. EM iterations (right) when **A)** and **B)** one, **C)** and **D)** two, **E)** and **F)** three, **G)** and **H)** four, **I)** and **J)** seven, and **K)** and **L)** fourteen components were fitted.

As we can see, when we try to fit Gaussian mixture model with the number of components larger than that of the true Gaussian mixture model, the effect of extra components is negligible in contour plot. We can conclude that there is an optimal number of components for fitting a Gaussian mixture. Fitting too many components can easily introduce singularity problem and cause inaccuracy and program break-down. The way I solved the singularity problem is introduced in Matlab implementation section.

7

## ii. Performance testing Predicted parameters v.s. True parameters

Testing data generated by a 3-D GMM with 3 components:

| | $\mu_1$ | | | $\mu_2$ | | | $\mu_3$ | | |
|---|---|---|---|---|---|---|---|---|---|
| True $\mu$ | 2 | 2 | 2 | -2 | -1 | -1 | 4 | 5 | 6 |
| Predicted $\mu$ | 2.1 | 1.97 | 2.01 | -1.99 | -0.98 | -1.02 | 4.02 | 4.97 | 6.00 |
| | $\Sigma_1$ | | | $\Sigma_2$ | | | $\Sigma_3$ | | |
| True $\Sigma$ | $\begin{pmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{pmatrix}$ | | | $\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$ | | | $\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$ | | |
| Predicted $\Sigma$ | $\begin{pmatrix} 1.90 & -0.01 & 0.01 \\ -0.01 & 2.20 & 0.06 \\ 0.01 & 0.06 & 2.20 \end{pmatrix}$ | | | $\begin{pmatrix} 0.93 & -0.03 & 0.01 \\ -0.03 & 0.95 & 0.03 \\ 0.01 & 0.03 & 0.92 \end{pmatrix}$ | | | $\begin{pmatrix} 0.97 & -0.05 & 0.07 \\ -0.05 & 0.89 & 0.01 \\ 0.07 & 0.01 & 1.01 \end{pmatrix}$ | | |

Testing data generated by a 4-D GMM with 3 components:

| | $\mu_1$ | | | | $\mu_2$ | | | | $\mu_3$ | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| True $\mu$ | 3 | 2 | 1 | 5 | -2 | -1 | -1 | 0 | 4 | 5 | 6 | 7 |
| Predicted $\mu$ | 2.99 | 1.97 | 0.74 | 4.87 | -2.01 | -0.98 | -1.10 | -0.02 | 4.14 | 5.03 | 5.92 | 6.86 |
| | $\Sigma_1$ | | | | $\Sigma_2$ | | | | $\Sigma_3$ | | | |
| True $\Sigma$ | $\begin{pmatrix} 4 & 0 & 0 & 0 \\ 0 & 4 & 0 & 0 \\ 0 & 0 & 4 & 0 \\ 0 & 0 & 0 & 4 \end{pmatrix}$ | | | | $\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$ | | | | $\begin{pmatrix} 3 & 2 & 0 & 0 \\ 2 & 3 & 0 & 0 \\ 0 & 0 & 3 & 1 \\ 0 & 0 & 1 & 3 \end{pmatrix}$ | | | |
| Predicted $\Sigma$ | $\begin{pmatrix} 3.74 & -0.61 & 0.26 & -0.07 \\ -0.61 & 4.11 & 0.10 & 0.12 \\ 0.26 & 0.10 & 4.02 & -0.42 \\ -0.07 & 0.12 & -0.42 & 3.85 \end{pmatrix}$ | | | | $\begin{pmatrix} 0.89 & 0.01 & 0.03 & -0.07 \\ 0.01 & 1.01 & 0.06 & -0.02 \\ 0.03 & 0.06 & 1.11 & 0.04 \\ -0.07 & -0.02 & 0.04 & 1.01 \end{pmatrix}$ | | | | $\begin{pmatrix} 2.90 & 2.03 & 0.16 & 0.19 \\ 2.02 & 3.36 & 0.13 & 0.31 \\ 0.16 & 0.13 & 3.14 & 1.24 \\ 0.19 & 0.31 & 1.24 & 3.31 \end{pmatrix}$ | | | |

## iii. Matlab implementation:

The self-implemented EM algorithm follows the procedure taught in lecture. (https://subjects.ee.unsw.edu.au/elec9782/elec9782-B-lec4.pdf p.15). The code has been thoroughly commented and attached in section Q2 in Appendix. As the EM is a stochastic process, it sometimes will minimise the covariance to zero. So, it can maximise the likelihood. In other words, it may fit an infinite high probability Gaussian on a single data with zero covariance. This is called singularity problem. My code solves this problem by adding a little bit of noise into the covariance matrix, if the covariance is considered as zero by Matlab, this noise will be added to the covariance matrix to prevent the singularity. At the same time, my program will print the warning "Hit the singularity, noise added for calculation. Try to fit the data with less Gaussian and run multiple times."

List of build-in functions used:
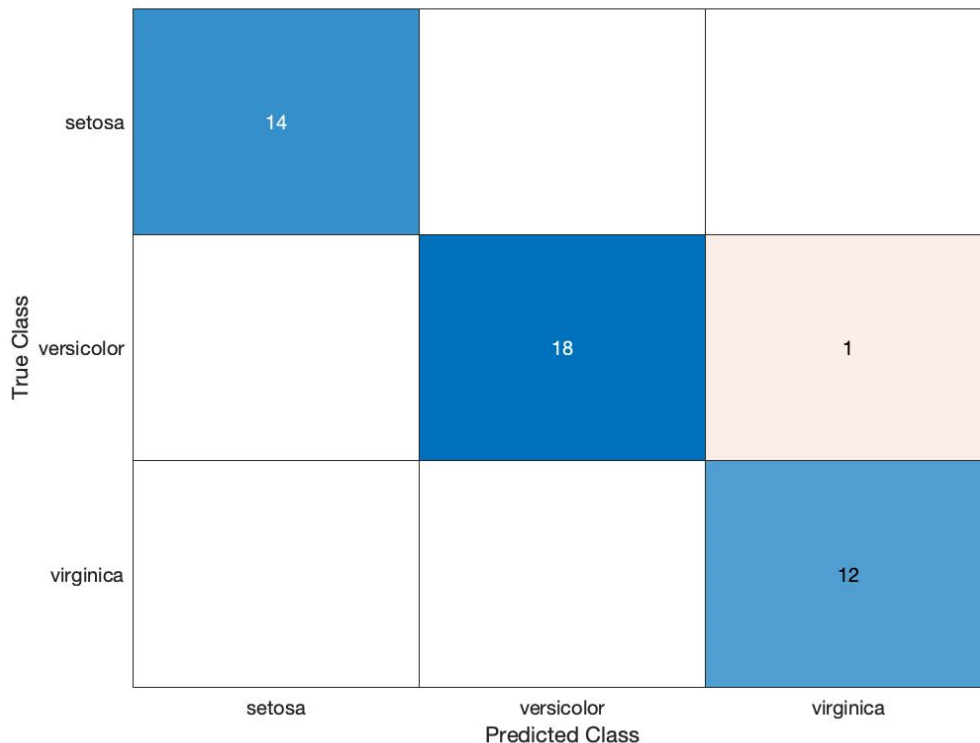
| | | |
|---|---|---|
| max() | diag() | sum() |
| min() | sqrt() | fprintf() |
| zeros() | det() | eye() |
| ones() | error() | rank() |
| length() | exp() | size() |
| plot() | isscalar() | ismatrix() |
| asset() | iscolumn() | |

# Question 3

## Problem (a)

i.    Prediction Result:



Figure 2. Confusion matrix of the prediction results.

Error rate: 2.22%

ii.    Matlab implementation:

I used logistic regression (one v.s. all) model for the classification. The data are shuffled and splitted into the training and the test set. I change the multi-class classification problem to the binary classification problem for each class. The data belongs to the chosen class will be labeled as 1(true). The data outside of the chosen class will be labeled as 0 (false) . Then, we train a logistic regression model on the chosen class by gradient descent. Each class will have their own logistic regression model.
When a new data comes in for class-label predication. We will calculate the probability score of data under each class's logistic regression model. Finally, we pick the class label which has the highest probability score as the final prediction.

The logistic regression is easy for implementation and often have promising results for supervised learning. As we have only three class labels, training three logistic regression models is not too costly.

The code has been thoroughly commented and attached in section Q3 in Appendix.

List of build-in functions used:

| | | |
|---|---|---|
| mean() | unique() | confusionchart() |
| std() | find() | |
| disp() | ones() | |
| num2str() | str2double() | |
| zeros() | sum() | |
| size() | ceil() | |
| nemul() | length() | |

# Appendix

## Q1

`pdf_plot.m`

```
1  clc;clear;
2  sigma = 1;
3  D = [0:5/1000001:5];
4  figure;
5  chi_index = (D/sigma).^2;
6  %%%caculate the pdf and plot the function
7  for order = 1:5
8      chi_distribution = chi2pdf(chi_index,order);
9      f_D = ((2*D)/sigma^2).*chi_distribution;
10     plot(D,f_D,'DisplayName', [ 'dim = ' num2str(order)]);
11     hold on
12 end
13 legend('show'); %create/show legend
14 hold off
```

`empirical_estimates.m`

```
1  %%%generate the histogram plot for distance pdf from dim 1 to k
2  n = 1e4;
3  k = 4^2;
4  figure
5  for j = 1:k
6      subplot(sqrt(k),sqrt(k),j)
7      %generate multivariate data
8      x = randn(j,n);
9      d = zeros(1,n);
10     %caculate and store the norm(distance) of the data
11     for i = 1:n
12         d(i) = norm(x(:,i));
13     end
14     %plot the histogram of the distances
15     histogram(d,20,'Normalization','pdf')
16     hold on
17
18     %%%plot theoretical pdf
19     sigma = 1;
20     pdfd = @(d) 2.*d./(sigma^2).*chi2pdf(d.^2/(sigma^2),j);
21     fplot(pdfd,[0,max(d)],'r');
22     title(['order =' num2str(j)]);
23     hold off
24 end
```

# Q2

EM.m

```matlab
1  %%%implementation of EM algorithm
2  %input:
3  %X ——> data with no restriction on dimension and observations
4  %init ——> numebr of gaussians we want
5  %output:
6  %gmm_model——> a struct stores number of componentes and paramters of model
7  function [gmm_model] = EM(X,init)
8      %%%initialisation process
9      %initilise number of componentes
10     gmm_model.num_of_components = init;
11
12     %randomly initilise gaussian means. We get maximum and minimum value
13     %along all dimensions of all value, and randomly choose value within this
14     %range as the value along each dimension of our mean.
15     maximum_X = max(X,[],'all');
16     minimum_X = min(X,[],'all');
17     r = minimum_X+(maximum_X—minimum_X).*rand(init,size(X,2));
18     gmm_model.u = r;
19     %%% initialise the covariance matrix by a constant times identity
20     %%% matrix. This constant is caculated as difference between
21     %maximum and minimum value divided by
22     %number of components.
23     sigma = zeros(size(X,2),size(X,2),init);
24     for k = 1:init
25         sigma(:,:,k) = ((maximum_X—minimum_X)/init)*eye(size(X,2));
26     end
27     gmm_model.sigma = sigma;
28     %%%initialise the weights,we assmue all weights are equal
29     weights = (1/init)*ones(1,init);
30     gmm_model.weights = weights;
31
32     %%%EM process
33     iteration = 500;                    %%maximum num of iterations
34     prev_llh_vec = [];                  %%dynamic array for storing the log—liklihood of ...
           data after each iteration
35     similarity_hitted = 0;              %%flag for checking if EM hitted singularity
36     for it = 1:iteration
37         %%%update the parameter
38         [gmm_model,prev_llh,singularity_flag] = updation(X,gmm_model);
39         %%%check if EM hitted the singularity
40         if singularity_flag == 1
41             similarity_hitted =1;
42         end
43         %%%break at log—liklihood become steady
44          if it ≠ 1 && prev_llh—prev_llh_vec(end)<1e—4
45              break
46          end
47         %%%store the log—liklihood to show it is maximised
48         prev_llh_vec = [prev_llh_vec prev_llh];
49
50     end
51
52     %%%report the similarity problem
53     if similarity_hitted ==1
54         fprintf(['Hitted the singularity,noise added for caculation \n Try to fit the ...
               ', ...
55         'data with less gaussians,and run multiple times\n']);
56
57     end
58     %%%plot log—lilihood change during the EM process.
59     plot(1:length(prev_llh_vec),prev_llh_vec);
60
61
62
63  end
```

## updation.m

```matlab
1  %%%this function update the parameters by expectation maximisation
2  %   when EM hit the singularity, similar to sklearn,small noise $1e-6 * Identity matrix$
3  %   will be added with covariance matrix to prevent singularity.
4  %   the formula used here is from ...
       https://subjects.ee.unsw.edu.au/elec9782/elec9782-B-lec4.pdf
5  %   slide 15.
6  %input:
7  %   X --- > data
8  %   gmm_model_old-->a struct stores original model parameters
9  %output:
10 %   gmm_model_new-->a struct stores new model parameters
11 %   singularity_flag --> report whether the EM hit the singularity point
12
13 function [gmm_model_new,llh,singularity_flag] = updation(X,gmm_model_old)
14     %%%extract the parameters from model struct and assert the input correctness
15     u_old = gmm_model_old.u;
16     num_of_components = gmm_model_old.num_of_components;
17     assert(size(u_old,1)==num_of_components,'number of means does not agree on num of ...
           components');
18     assert(size(u_old,2)==size(X,2),'means dimension does not agree on num of components');
19     %%%caculate the posterior probability gamma given data.
20     num_of_observations = size(X,1);
21     singularity_flag = 0;
22     [gamma_matrix,llh] = latent_pospdf(X,gmm_model_old);
23
24     %%%update the paramters
25     %%%https://subjects.ee.unsw.edu.au/elec9782/elec9782-B-lec4.pdf
26     %%%page 15.
27     N_k = sum(gamma_matrix,1);
28
29     u_new = (gamma_matrix'*X)./N_k';
30     sigma_new = zeros(size(X,2),size(X,2),num_of_components);
31     for m = 1:num_of_components
32         sigma_sum=0;
33         for n = 1:num_of_observations
34             U = (X(n,:)-u_old(m,:))'*(X(n,:)-u_old(m,:));
35             sigma_sum=sigma_sum+gamma_matrix(n,m)*U;
36
37         end
38         sigma_new(:,:,m) =sigma_sum/N_k(m);
39
40         %%%prevent the singularity,similar to sklearn approach.
41         if rank(sigma_new(:,:,m)) < size(sigma_new(:,:,m),2)
42             sigma_new(:,:,m)=sigma_new(:,:,m)+1e-6*eye(size(sigma_new(:,:,m),2));
43             singularity_flag = 1;
44         end
45
46     end
47     weights = N_k/num_of_observations;
48
49     %%%construct the new model struct with updated parameters
50     gmm_model_new.u = u_new;
51     gmm_model_new.sigma = sigma_new;
52     gmm_model_new.weights = weights;
53     gmm_model_new.num_of_components = gmm_model_old.num_of_components;
54 end
```

**latent_pospdf.m**

```matlab
1  %   Given the model and a data point,this function can caculate the posterior
2  %   pabability of p(z|x u sigma) for ***all the observations***.
3
4  %   p(z|x u sigma)= pi_k*N(x|theta_k)/ sum_{j}(pi_j*N(x|theta_j)).
5  %   it gives us for each datapoint x_i the measure of:
6  %   "prabability that a given data belongs to a certain class/probability of x_i over ...
       all classes"
7  %   this function is similar to matlab built—in function "posterior()"
8
9  %   For convenience this function also ouputs the total log—likelihood wich is
10 %   essentially the log of sum_{j}(pi_j*N(x|theta_j))(the mariginal
11 %   likelihood of the X)
12 %input: x—>data (no restrictions on observtions and dimensions)
13 %       gmm_model—>a struct stores number of components and parameters of the model
14 %
15 %output:gamma—>posteriror probability matrix, the row represents
16 %                observations,cololumn corresponds to probability belongs to
17 %                each gaussian distribution.
18 %       llh —> total log liklihood of data given current updated model
19
20
21
22 function [gamma,llh] = latent_pospdf(x,gmm_model)
23     %%%extract the value from model struct and assert the input correctness
24     u = gmm_model.u;
25     sigma = gmm_model.sigma;
26     weights = gmm_model.weights;
27     num_of_components = gmm_model.num_of_components;
28     assert(isvector(weights),'models weights input are wrong');
29     assert(num_of_components==length(weights),'components number not equal wights number');
30     assert(ismatrix(sigma(:,:,1)),'sigma dimension is wrong');
31
32     %%%caculates the pi_k*N(x|theta_k),
33     %the caculation here is matrixlised.Given the data matrix, row
34     %represents each observation, columns represents dimension.
35     %The caculated output has row represents each observation,
36     %columns represents different gaussian clusters .
37     num_of_observations = size(x,1);
38     N = zeros(num_of_observations,num_of_components);
39     for m = 1:num_of_components
40         N(:,m) = pdfmvn(x,u(m,:),sigma(:,:,m));
41     end
42     N_pi = N*diag(weights);
43
44     NUM = N_pi;
45
46     %%%caculates sum_{j}(pi_j*N(x|theta_j)).
47     DEN = sum(N_pi,2);
48
49     %%%divide each data observation(rows) by denominator.
50     gamma = NUM./DEN;
51
52     %%%caculates the log—likelihood
53     llh = sum(log(DEN))/num_of_observations;
54
55 end
```

**pdfmvn.m**

```matlab
1  %this function caculates the liklihood of data by using the pdf of 1—d gaussian
2  %or multivariate gaussian, the function is similar to built—in function
3  %mvnpdf.m
4  %input: X——>data no restricitons on rows and columns
5  %        u——>mean vector of gaussian
6  %        sigma——>covariance—matrix of gaussian
7  %output: y ——> probability
8  function y = pdfmvn(X,mu,sigma)
9      %%%1—d
10     if (iscolumn(X))
11         assert(isscalar(mu),'u is not a scalar');
12         assert(isscalar(sigma),'sigma is not a scalar');
13         constant = 1/(sigma*sqrt(2*pi));
14         X = X—mu;
15         y = constant.*exp(—X.^2/(2*sigma^2));
16     %%%2—d.For multiple data, row>1, we pick the diagonal of y as output
17     elseif(ismatrix(X))
18         dim = size(X,2);
19         assert(size(mu,2)== dim,'u dim is not right');
20         assert(size(sigma,2)== dim,'sigma dim is not right');
21         assert(size(sigma,1)== dim,'sigma dim is not right');
22         constant = 1/sqrt(det(sigma)*((2*pi)^dim));
23         X = X—mu;
24         U = X/sigma;
25         z = (—1/2)*(U*X');
26         y = diag(constant.*exp(z));
27
28     else
29
30         error('wrong input for pdfmvn');
31
32     end
33
34 end
```

**test.m**

```matlab
1  clear;clc
2  mu1 = [3 2 1 5];          % Mean of the 1st component
3  sigma1 = [4 0 0 0; 0 4 0 0; 0 0 4 0; 0 0 0 4];  % Covariance of the 1st component
4  mu2 = [—2 —1,—1, 0];       % Mean of the 2nd component
5  sigma2 = [1 0 0 0; 0 1 0 0; 0 0 1 0; 0 0 0 1];  % Covariance of the 2nd component
6  mu3 = [4 5 6 7];          % Mean of the 3rd component
7  sigma3 = [3 2 0 0; 2 3 0 0; 0 0 3 1; 0 0 1 3];  % Covariance of the 2nd component
8  mu4 = [—3 5 6 4];         % Mean of the 4th component
9  sigma4 = [2 1 1 1; 1 2 1 1; 1 1 2 1; 1 1 1 1];  % Covariance of the 2nd component
10
11 rng('shuffle')                % For reproducibility
12 r1 = mvnrnd(mu1,sigma1,300);
13 r2 = mvnrnd(mu2,sigma2,300);
14 r3 = mvnrnd(mu3,sigma3,400);
15 %r4 = mvnrnd(mu4,sigma4,400);
16 X = [r1; r2;r3];
17
18 [gmm_model] = EM(X,3);
19
20 figure;
21 scatter(X(:,1),X(:,2),10,'.')    % Scatter plot with points of size 10
22 hold on
23 gm = gmdistribution(gmm_model.u,gmm_model.sigma,gmm_model.weights);
24 gmPDF = @(x,y)reshape(pdf(gm,[x(:) y(:)]),size(x));
25 fcontour(gmPDF,[—10 10 —10 10]);
26 title('Contour plot of fitted gaussian mixture on data(14 components)')
```

# Q3

**main.m**

```matlab
1   clear;clc;
2   load fisheriris.mat
3   rng = ('default');
4   %%%split data to test data and train data
5   test_ratio = 0.3;
6
7   data_with_label = [string(species),meas];
8   categories = unique(string(species));
9   [test,train] = split_train_test(data_with_label,test_ratio);
10
11  %%%train parameter for each category(one vs all) and stored in a matrix
12  theta = zeros(length(categories),size(meas,2)+1)';
13  J = zeros(3,1);
14
15  [theta(:,1),J(1)] = model_creation(train,categories(1),0.1,10000);    %Question:
16  [theta(:,2),J(2)] = model_creation(train,categories(2),0.2,10000);    %quicker learning ...
        rate needed
17  [theta(:,3),J(3)] = model_creation(train,categories(3),0.2,10000);     %why they did ...
        not converge to same point?
18
19  %%%prediction
20  result = predict(theta,test);
21  predicted_label = categories(result);
22  true_label = test(:,1);
23
24  %%%error and confusion matrix
25  error = err_caculation (true_label,predicted_label);
26  disp(['Err is : ', num2str(error) , '%']);
27  figure
28  cm = confusionchart(true_label,predicted_label);
```

**err_calculation.m**

```matlab
1   %Caculate the error by number of unmatch between ground-truth and
2   %prediction
3   %input : ture_label --> true_label
4   %output : predicted_label -->predicted_label
5   %err --> error in percent
6   function err = err_caculation (true_label,predicted_label)
7       counter = 0;
8       for k = 1: length(true_label)
9           if true_label(k) ≠ predicted_label (k)
10              counter = counter + 1;
11          end
12
13      end
14
15      err = counter / length (true_label)*100;
16
17  end
```

## K_shuffle.m

```matlab
1  %%%The function will random shuffle positions of value in a vector
2  %For, Details see Knuth shuffle
3  %https://en.wikipedia.org/wiki/Fisher%E2%80%93Yates_shuffle
4  %input X--> numeric array
5  %output X--> shuffled numeric array
6  function [X] = K_shuffle(X)
7      n = numel(X);
8      for i = 2:n       % Knuth shuffle in forward direction:
9          w = ceil(rand * i);    % 1 ≤ w ≤ i
10         t = X(w);
11         X(w) = X(i);
12         X(i) = t;
13     end
14 end
```

## split_train_test.m

```matlab
1  %%%split the train_test data
2  %shuffle the index, split them by test_train _ratio
3  %input :
4  %test_ratio --> portion of data for test data --> data
5  %output:
6  %test --> test_data  train -- > train_data
7  function [test,train] = split_train_test(data,test_ratio)
8  %     rng = ('default');
9      num_of_observations = size(data,1);
10     shuffled_indices = K_shuffle([1:num_of_observations]);
11     test_set_size = floor(num_of_observations*test_ratio);
12     test_indices = shuffled_indices(1:test_set_size);
13     train_indices = shuffled_indices(test_set_size:end);
14     test = data(test_indices,:);
15     train = data(train_indices,:);
16
17 end
```

## model_creation.m

```matlab
1  %(critical function)
2  %learning the parameter for logistic regression
3  %input :
4  %data-->training_data str-->the category we want for logistic...
5  %regression so the data belongs to the category labeled as 1, data outside
6  %of the category labeled as 0(one vs all); alpha-->
7  %learning_rate;iter-->number of iteration for gradient decent
8  %output :
9  %theta --> trained_output J-->cost function
10 function [theta, J] = model_creation(data,str,alpha,iter)
11     [data] = split_data_by_binary_category(data,str);
12     y = str2double(data(:,1));
13     x_raw = str2double(data(:,2:end));
14     x = (x_raw - mean(x_raw))./std(x_raw);
15
16     x = [ones(size(x,1),1),x];
17     theta = zeros(size(x,2),1);
18     m = size(x,1);
19     J = 0;
20     counter = 0;
21     for k = 1:iter
22         h = sigmoid(x*theta);
23         J = (-1/m)*sum(y.*log(h)+(1-y).*log(1-h));
24          counter = counter +1;
25         theta = theta-(alpha/m)*x'*(h-y);
26     end
27
28 end
```

## split_data_by_binary_category.m

```matlab
1  %marking the data_label setted by user as 1
2  %all other data_label are marked as zeros
3  %input :
4  %data ——> data str ——> data_label that user want
5  %marked data label
6  function [data] = split_data_by_binary_category(data,str)
7      false_index = find(data(:,1) ≠str);
8      true_index = find(data(:,1) ==str);
9      data(true_index,1) = 1;
10     data(false_index,1) = 0;
11 end
```

## sigmoid.m

```matlab
1  %sigmoid function
2  function output = sigmoid(z)
3  output = 1./(1+exp(−z));
4
5
6  end
```

## predict.m

```matlab
1  %%%predict the label,given the model paramters theta
2  %input:
3  %   test——> test_data;
4  %output:
5  %   result——>the output here is the numer corresponding
6  % to the postion of label in label array.
7  function [result] = predict(theta,test)
8      test = str2double(test(:,2:end));
9      test_data = (test − mean(test))./std(test);
10     test_data = [ones(size(test,1),1),test_data];
11
12     h = sigmoid(test_data*theta);
13     [¬,col] = max(h,[],2);
14     result = col;
15
16
17 end
```