



AUSTRALIAN NATIONAL UNIVERSITY

---

COMP4670 INTRODUCTION TO MACHINE LEARNING

## Assignment 02

Edited by L<sup>A</sup>T<sub>E</sub>X

College of Engineering and Computer Science

---

© *Author*

**Jimmy Lin**

u5223173

¶ *Lecturer*

**Christfried Webers**

¶ *Tutor*

**Wen Shao**

† *Release Date*

**April. 22 2013**

‡ *Due Date*

**May. 19 2013**

$\tau$  *Time Spent*

**40 hours**

May 18, 2013

## Contents

<b>1</b>	<b>Sampling from a Multivariate Normal Distribution</b>	<b>2</b>
1.1	Explain how create vectorial data conforming to multivariate Gaussian . . . . .	2
<b>2</b>	<b>Reverse Prediction</b>	<b>6</b>
2.1	Complete optimisation problem for $\mathbf{w}^*$ . . . . .	6
2.2	Solution for $\mathbf{w}^*$ and condition of Uniqueness . . . . .	6
2.3	Optimisation problem for $\mathbf{u}$ . . . . .	7
2.4	Solution for $\mathbf{w}^*$ and condition of Uniqueness . . . . .	7
2.5	Relation between unique $\mathbf{w}^*$ and $\mathbf{u}^*$ . . . . .	8
<b>3</b>	<b>Conditional Probability and Variance via Parzen Estimator</b>	<b>9</b>
3.1	Conditional density, expectations and variance . . . . .	9
3.2	Show that the function $\sum_{n=1}^N k(x, x_n) = 1$ . . . . .	13
<b>4</b>	<b>Maximum Margin Hyperplane</b>	<b>14</b>
4.1	Two points are sufficient to determine the hyperplane . . . . .	14
4.2	Dimension of this maximum margin hyperplane . . . . .	16
4.3	Discriminant function $f(\mathbf{x})$ . . . . .	16
4.4	Discuss the values of the Lagrange parameters . . . . .	16
<b>5</b>	<b>Kernels and Feature Maps</b>	<b>17</b>
5.1	Prove $k(\mathbf{x}, \mathbf{z})$ to be Valid Kernel in 2 Dimension . . . . .	17
5.2	Generalize Input Space to Higher Dimension . . . . .	18
5.3	Generalize Kernel to Higher Power . . . . .	21
<b>6</b>	<b>EM for Trees</b>	<b>23</b>
6.1	Define latent variables and a joint generative model . . . . .	23
6.2	Derive the expected log likelihood of the joint model . . . . .	24
6.3	Derive the E-step update and show how to compute expectations . . . . .	26
6.4	Derive the M-step for EM . . . . .	27
<b>7</b>	<b>(Semi-/Un-)Supervised Learning and EM</b>	<b>32</b>
7.1	Supervised Learning . . . . .	32
7.2	Unsupervised Learning . . . . .	32
7.3	Semi-supervised Learning . . . . .	33
<b>A</b>	<b>Data Presentation</b>	<b>34</b>
A.1	Cross Validation for Supervised Learning . . . . .	34
A.2	Expectation Maximisation for Unsupervised Learning . . . . .	39

# 1 Sampling from a Multivariate Normal Distribution

## 1.1 Explain how create vectorial data conforming to multivariate Gaussian

First and foremost, we assume the **independent property** of this random number generator while producing scalar value for the vectorial data  $\mathbf{x} \in R^n$  ( $n$  is the dimensionality of the resulting data object). Formally, we describe the vectorial data  $\mathbf{x}$  and the finally resulted data matrix  $\mathbf{X}$  as follows,

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \quad \mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_m^T \end{pmatrix} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \cdots & x_{mn} \end{pmatrix} \quad (1)$$

where  $n$  is the dimensionality of each single vectorial data we aim to construct,  $m$  is the number of vectorial object we are intended to collect. Note that since the random number generator only produce scalar value, the number of times to apply this generator is  $n \times m$  if the number of such  $n$ -dimensional data objects is  $m$ .

**The way to construct the data matrix.** For each dimension, we use the random number generator with specific  $\mu$  and  $\sigma$  to produce a set of specific scalar value (specifically, the number of such scalar value is  $m$  for each dimension). And then we apply the column appending for each set of scalar values, as what we formulated in the (1).

In this way, we are able to formulate the following probability distribution for each dimension  $i = 1, \dots, n$ .

$$p(x_i) = \mathcal{N}(x_i | \mu_i, \sigma_i^2) = \frac{1}{(2\pi)^{\frac{1}{2}} \sigma_i} \exp\left(-\frac{(x_i - \mu_i)^2}{2\sigma_i^2}\right) \quad (2)$$

where  $x_i$  is the  $i$ -th element of vectorial data  $\mathbf{x}$ ,  $\mu_i$  is the mean value of the  $i$ -th dimension's univariate gaussian distribution, and  $\sigma_i$  is variance of  $i$ -th dimension's univariate gaussian distribution.

**Show  $\mathbf{x}$  conforms to Multivariate Gaussian Distribution.** The following content discuss the proof for the vectorial data  $\mathbf{x}$  generated in the way we presented above conforms to multivariate gaussian distribution.

Due to the independent assumption we have made at the beginning,

$$p(\mathbf{x}) = \prod_{i=1}^n p(x_i) \quad \text{by i.d assumption} \quad (3)$$

$$= \prod_{i=1}^n \frac{1}{(2\pi)^{\frac{1}{2}} \sigma_i} \exp\left(-\frac{(x_i - \mu_i)^2}{2\sigma_i^2}\right) \quad \text{by (2)} \quad (4)$$

$$= \frac{1}{(2\pi)^{\frac{n}{2}} \prod_{i=1}^n \sigma_i} \exp\left(-\sum_{i=1}^n \frac{(x_i - \mu_i)^2}{2\sigma_i^2}\right) \quad \text{put } \prod \text{ inside} \quad (5)$$

$$= \frac{1}{(2\pi)^{\frac{n}{2}} \prod_{i=1}^n \sigma_i} \exp\left(-\frac{1}{2} \sum_{i=1}^n \frac{(x_i - \mu_i)^2}{\sigma_i^2}\right) \quad \text{put } \frac{1}{2} \text{ outside} \quad (6)$$

Then we define mean vector  $\boldsymbol{\mu} \in R^{n \times 1}$  and covariance matrix  $\mathbf{C} \in R^{n \times n}$  of the generated vectorial object to be,

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{pmatrix} \quad \mathbf{C} = \begin{pmatrix} c_{11} & c_{12} & \cdots & c_{1n} \\ c_{21} & c_{22} & \cdots & c_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ c_{n1} & c_{n2} & \cdots & c_{nn} \end{pmatrix} \quad (7)$$

where  $\mu_i$  is the mean value of the  $i$ -th dimension's random variable  $x_i$  and  $\mathbf{C}$  is the covariance matrix of  $\mathbf{x}$ , whose entry  $c_{ij}$ , by definition, is the covariance between  $x_i$  and  $x_j$ .

Then let's focus on deriving the form of covariance matrix  $\mathbf{C}$  and the specific value of each entry  $c_{ij}$  ( $i$  is the row index and  $j$  denotes column index) based on the characteristics of data collection.

First, look at the case of  $i \neq j$  as follows.

$$c_{ij(i \neq j)} = \int \int (x_i - E(x_i))(x_j - E(x_j))P(x_i, x_j)dx_i dx_j \quad \text{definition of covariance} \quad (8)$$

$$= \int \int (x_i - E(x_i))(x_j - E(x_j))P(x_i)P(x_j)dx_i dx_j \quad \text{by i.i.d when } i \neq j \quad (9)$$

$$= \int (x_i - E(x_i))P(x_i)dx_i \int (x_j - E(x_j))P(x_j)dx_j \quad \text{independence of integration} \quad (10)$$

$$= \left( \int x_i P(x_i)dx_i - \int E(x_i)P(x_i)dx_i \right) \left( \int x_j P(x_j)dx_j - \int E(x_j)P(x_j)dx_j \right) \quad \text{linearity of integration} \quad (11)$$

$$= \left( \int x_i P(x_i)dx_i - E(x_i) \int P(x_i)dx_i \right) \left( \int x_j P(x_j)dx_j - E(x_j) \int P(x_j)dx_j \right) \quad \begin{array}{l} E_{x_i} \text{ independent on } x_i \\ E_{x_j} \text{ independent on } x_j \end{array} \quad (12)$$

where  $P(x_i, x_j)$  is the joint probability of two independent random variables from  $\mathbf{x}$ .  $P(x_i)$  is the marginal probability distribution of  $x_i$ .

By definition of the expectation, we have

$$E(x_i) = \int x_i P(x_i)dx_i \quad (13)$$

$$E(x_j) = \int x_j P(x_j)dx_j \quad (14)$$

By normalisation property of the probability distribution, we have

$$\int P(x_i)dx_i = 1 \quad (15)$$

$$\int P(x_j)dx_j = 1 \quad (16)$$

Continue the previous derivation with regard to  $\sigma_{ij}$  in (12), we have

$$c_{ij(i \neq j)} = \left( \int x_i P(x_i)dx_i - E(x_i) \int P(x_i)dx_i \right) \left( \int x_j P(x_j)dx_j - E(x_j) \int P(x_j)dx_j \right) \quad \text{from (12)} \quad (17)$$

$$= \left( E(x_i) - E(x_i) \int P(x_i)dx_i \right) \left( E(x_j) - E(x_j) \int P(x_j)dx_j \right) \quad \text{apply (13)} \quad (18)$$

$$= \left( E(x_i) - E(x_i) \right) \left( E(x_j) - E(x_j) \right) \quad \text{apply (14)} \quad (19)$$

$$= 0 \quad \text{apply (15) and (16)} \quad (20)$$

And the analysis the other case is, since the covariance of one random variable to itself is its variance (by definition). if  $i = j$ , we have

$$c_{ij(i=j)} = \sigma_i^2 \quad (21)$$

where the  $\sigma_i^2$  is the variance of the scalar random variable  $x_i$ , which is shown in the univariate gaussian distribution (2) at the very beginning.

Summarise the  $c_{ij}$ , which is the entry of covariance matrix of multivariate gaussian distribution.

$$c_{ij} = \begin{cases} 0 & \text{if } i \neq j \\ \sigma^2 & \text{if } i = j \end{cases} \quad (22)$$

That is to say,  $\mathbf{C}$  is a diagonal matrix with its diagonal entry being  $\sigma_i^2 (i = 1..n)$ .

$$\mathbf{C} = \begin{pmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n^2 \end{pmatrix} \quad (23)$$

From the form of  $\mathbf{C}$  we have derived in (23), the determinant of  $\mathbf{C}$  can be figured out

$$|\mathbf{C}| = \prod_{i=1}^n \sigma_i^2 \quad (24)$$

Obviously, all factors in (24) is non-negative, we can take square root for them

$$|\mathbf{C}|^{\frac{1}{2}} = \prod_{i=1}^n \sigma_i \quad (25)$$

Another point we can conclude from the diagonal property of covariance matrix  $\mathbf{C}$  is its inverse. Since we assume that all the variance of univariate gaussian distribution we used here is non-zero, which means the covariance matrix  $\mathbf{C}$  is full-rank and thus invertible.

$$\mathbf{C}^{-1} = \begin{pmatrix} \sigma_1^{-2} & 0 & \cdots & 0 \\ 0 & \sigma_2^{-2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n^{-2} \end{pmatrix} \quad (26)$$

Next, present the form of  $\mathbf{x} - \boldsymbol{\mu}$  before we turn to derive the scalar form of  $(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x} - \boldsymbol{\mu})$

$$\mathbf{x} - \boldsymbol{\mu} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} - \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{pmatrix} = \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \\ \vdots \\ x_n - \mu_n \end{pmatrix} \quad (27)$$

where  $n$  is dimensionality of the vectorial data we wish to construct.

According to (27) and (26), we can easily have the following deduction,

$$(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = (x_1 - \mu_1, x_2 - \mu_2, \dots, x_n - \mu_n) \begin{pmatrix} \sigma_1^{-2} & 0 & \cdots & 0 \\ 0 & \sigma_2^{-2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n^{-2} \end{pmatrix} \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \\ \vdots \\ x_n - \mu_n \end{pmatrix} \quad (28)$$

$$= \left( \frac{x_1 - \mu_1}{\sigma_1^2}, \frac{x_2 - \mu_2}{\sigma_2^2}, \dots, \frac{x_n - \mu_n}{\sigma_n^2} \right) \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \\ \vdots \\ x_n - \mu_n \end{pmatrix} \quad (29)$$

$$= \frac{(x_1 - \mu_1)^2}{\sigma_1^2} + \frac{(x_2 - \mu_2)^2}{\sigma_2^2} + \cdots + \frac{(x_n - \mu_n)^2}{\sigma_n^2} \quad (30)$$

$$= \sum_{i=1}^n \frac{(x_i - \mu_i)^2}{\sigma_i^2} \quad (31)$$

Therefore, we have

$$(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = \sum_{i=1}^n \frac{(x_i - \mu_i)^2}{\sigma_i^2} \quad (32)$$

Manipulate we have derived before for  $p(\mathbf{x})$ ,

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{n}{2}} \prod_{i=1}^n \sigma_i} \exp\left(-\frac{1}{2} \sum_{i=1}^n \frac{(x_i - \mu_i)^2}{\sigma_i^2}\right) \quad \text{from (6)} \quad (33)$$

$$= \frac{1}{(2\pi)^{\frac{n}{2}} |\mathbf{C}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} \sum_{i=1}^n \frac{(x_i - \mu_i)^2}{\sigma_i^2}\right) \quad \text{replace } |\mathbf{C}|^{\frac{1}{2}} \text{ by (25)} \quad (34)$$

$$= \frac{1}{(2\pi)^{\frac{n}{2}} |\mathbf{C}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right) \quad \text{apply (32)} \quad (35)$$

Therefore, the probability distribution of one single object  $\mathbf{x}$  collected by our approach, which is illustrated in the very beginning, can be represented as follows,

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{n}{2}} |\mathbf{C}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right) \quad (36)$$

which apprently is a multivariate Gaussian distribution with two semantic parameters, one is mean vector  $\boldsymbol{\mu}$  of  $\mathbf{x}$  and the other one covariance matrix  $\mathbf{C}$  of  $\mathbf{x}$ . Hence, it can be concluded that

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{C}) \quad (37)$$

## 2 Reverse Prediction

### 2.1 Complete optimisation problem for $\mathbf{w}^*$

Given the regression hypothesis,

$$\tilde{t}_i = \mathbf{w}^T \mathbf{x}_i \quad (38)$$

where the  $\tilde{t}_i$  is prediction by our hypothesis,  $\mathbf{w}$  is the coefficient for each dimension, and  $\mathbf{x}_i$  is the  $i$ -th input object in given set of  $n$  training data.

Since the optimal  $\mathbf{w}$  minimising the quadratic regression error of all data for a linear model, we can write down the optimisation problem as follows,

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_{i=1}^n (\tilde{t}_i - t_i)^2 \quad \text{minimise quadratic training error} \quad (39)$$

$$= \arg \min_{\mathbf{w}} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - t_i)^2 \quad \text{replace by regression hypothesis} \quad (40)$$

$$= \arg \min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - t_i)^2 \quad \text{ultimate optimisation objective} \quad (41)$$

Note that the final step is valid because the newly added coefficient  $\frac{1}{2}$  is not dependent on the  $\mathbf{w}$  and doing so is to enhance the convenience for taking derivative.

### 2.2 Solution for $\mathbf{w}^*$ and condition of Uniqueness

Define error function to be the term we wish to minimise

$$E(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - t_i)^2 \quad (42)$$

Directional derivative of  $E(\mathbf{w})$  to zero,

$$\mathcal{D}(E(\mathbf{w}))(\xi) = \langle \text{grad } E(\mathbf{w}), \xi \rangle, \quad \forall \xi \quad (43)$$

where  $\xi$  is defined in the space  $R^d$ .

Take gradient of the error function  $E(\mathbf{w})$  with regard to  $\mathbf{w}$ ,

$$\text{grad } E(\mathbf{w}) = \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - t_i) \mathbf{x}_i \quad (44)$$

From the first-order derivative (44) of error function  $E(\mathbf{w})$ , it is obvious that error function is convex function with regard to parameter  $\mathbf{w}$ . (this is because the  $\frac{\partial E(\mathbf{w})}{\partial \mathbf{w}}$  is linear to  $\mathbf{w}^T$ .) Therefore, setting the first derivative to zero would help us derive the global minima of error function.

Write the equation above (44) in the form of matrix,

$$\text{grad } E(\mathbf{w}) = \mathbf{X}^T (\mathbf{X} \mathbf{w} - \mathbf{t}) \quad (45)$$

where,  $\mathbf{X}$  is a matrix, each row of whom is  $\mathbf{x}_i$ , and  $\mathbf{t}$  is a column vector whose element is  $t_i$ .

Set  $E(\mathbf{w}) = 0$  to derive the global minimum for this convex function,

$$\mathbf{X}^T (\mathbf{X} \mathbf{w}^* - \mathbf{t}) = 0 \quad (46)$$

$$\mathbf{X}^T \mathbf{X} \mathbf{w}^* - \mathbf{X}^T \mathbf{t} = 0 \quad (47)$$

$$\mathbf{X}^T \mathbf{X} \mathbf{w}^* = \mathbf{X}^T \mathbf{t} \quad (48)$$

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} \quad (49)$$

In the derivation from (48) to (49), we have to specify condition that  $\mathbf{X}^T \mathbf{X}$  is invertible.

## 2.3 Optimisation problem for $\mathbf{u}$

The squared error for between target variable  $\mathbf{x}_i$  and  $\tilde{\mathbf{x}}_i$

$$\arg \min_{\mathbf{u}} \frac{1}{2} \sum_{i=1}^n \|t_i \mathbf{u} - \mathbf{x}_i\|^2 \quad (50)$$

which can be written as

$$\arg \min_{\mathbf{u}} \frac{1}{2} \sum_{i=1}^n (t_i \mathbf{u} - \mathbf{x}_i)^T (t_i \mathbf{u} - \mathbf{x}_i) \quad (51)$$

## 2.4 Solution for $\mathbf{w}^*$ and condition of Uniqueness

Define the error function  $E_2(\mathbf{u})$  as follows,

$$E_2(\mathbf{u}) = \frac{1}{2} \sum_{i=1}^n (t_i \mathbf{u} - \mathbf{x}_i)^T (t_i \mathbf{u} - \mathbf{x}_i) \quad (52)$$

Directional derivative of  $E_2(\mathbf{u})$  to zero,

$$\mathcal{D}(E_2(\mathbf{u}))(\xi) = \langle \text{grad } E(\mathbf{u}), \xi \rangle, \quad \forall \xi \quad (53)$$

where  $\xi$  is defined to be in the space  $R^d$ .

Take gradient of  $E_2(\mathbf{u})$ , we have

$$\text{grad } E_2(\mathbf{u}) = \sum_{i=1}^n (t_i \mathbf{u} - \mathbf{x}_i) t_i \quad (54)$$

The formula above can be written as,

$$\text{grad } E_2(\mathbf{u}) = (\mathbf{t} \mathbf{u}^T - \mathbf{X})^T \mathbf{t} \quad (55)$$

where  $\mathbf{t}$  is the column vector, each element of whom is  $t_i$ , and  $\mathbf{X}$  is matrix whose row is  $\mathbf{x}_i$ .  
Set  $\text{grad } E_2(\mathbf{u}) = 0$ , we have

$$(\mathbf{t}(\mathbf{u}^*)^T - \mathbf{X})^T \mathbf{t} = 0 \quad (56)$$

$$(\mathbf{u}^* \mathbf{t}^T - \mathbf{X}^T) \mathbf{t} = 0 \quad (57)$$

$$\mathbf{u}^* \mathbf{t}^T \mathbf{t} - \mathbf{X}^T \mathbf{t} = 0 \quad (58)$$

$$\mathbf{u}^* \mathbf{t}^T \mathbf{t} = \mathbf{X}^T \mathbf{t} \quad (59)$$

$$\mathbf{u}^* = \mathbf{X}^T \mathbf{t} (\mathbf{t}^T \mathbf{t})^{-1} \quad (60)$$

When  $\mathbf{t}^T \mathbf{t}$  is invertible, we have unique solution for  $\mathbf{u}^*$ . However, it is obvious that the  $\mathbf{t}^T \mathbf{t}$  is a scalar value, more specifically, invertibility corresponds to whether that scalar is zero. Therefore, the condition for uniqueness of  $\mathbf{u}^*$  is

$$\mathbf{t}^T \mathbf{t} \text{ is non-zero.} \quad (61)$$

which can be interpreted as

$$\sum_{i=1}^n t_i^2 \neq 0 \quad (62)$$



## 2.5 Relation between unique $\mathbf{w}^*$ and $\mathbf{u}^*$

From (49) and (60), we have

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} \quad (63)$$

$$\mathbf{u}^* = \mathbf{X}^T \mathbf{t} (\mathbf{t}^T \mathbf{t})^{-1} \quad (64)$$

Accompanied with the following conditions for unique  $\mathbf{w}^*$  and  $\mathbf{u}^*$

$$\mathbf{X}^T \mathbf{X} \text{ is invertible.} \quad (65)$$

$$\mathbf{t}^T \mathbf{t} \text{ is non-zero.} \quad (66)$$

From the expression of  $\mathbf{w}^*$  and  $\mathbf{u}^*$ , it can be concluded that

$$\mathbf{w}^* = \frac{(\mathbf{X}^T \mathbf{X})^{-1}}{(\mathbf{t}^T \mathbf{t})^{-1}} \mathbf{u}^* \quad (67)$$

which can be simplified since  $(\mathbf{t}^T \mathbf{t})^{-1}$  is scalar

$$\mathbf{w}^* = (\mathbf{t}^T \mathbf{t})(\mathbf{X}^T \mathbf{X})^{-1} \cdot \mathbf{u}^* \quad (68)$$

### 3 Conditional Probability and Variance via Parzen Estimator

#### 3.1 Conditional density, expectations and variance

As prescribed by the question,

$$p(x, t) = \frac{1}{N} \sum_{n=1}^N f(x - x_n, t - t_n) \quad (69)$$

where, the  $f(\cdot)$  is bivariate gaussian distribution centered at  $(0, 0)$  with covariance matrix  $\sigma \mathbf{I}^2$ . Therefore, we can further extend the  $p(x, t)$

$$p(x, t) = \frac{1}{N} \sum_{n=1}^N \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{x}_n - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x}_n - \boldsymbol{\mu})\right) \quad (70)$$

$$= \frac{1}{N} \sum_{n=1}^N \frac{1}{2\pi\sigma^2} \exp\left(-\frac{(x - x_n)^2 + (t - t_n)^2}{2\sigma^2}\right) \quad (71)$$

$$= \frac{1}{N} \sum_{n=1}^N \frac{1}{2\pi\sigma^2} \exp\left(-\frac{(x - x_n)^2}{2\sigma^2}\right) \exp\left(-\frac{(t - t_n)^2}{2\sigma^2}\right) \quad (72)$$

$$= \frac{1}{2\pi\sigma^2 N} \sum_{n=1}^N \exp\left(-\frac{(x - x_n)^2}{2\sigma^2}\right) \exp\left(-\frac{(t - t_n)^2}{2\sigma^2}\right) \quad (73)$$

Note that the derivation from (70) to (71) is valid in that the  $x_n$  and  $t_n$  are both scalar value and the density estimator is a bivariate gaussian distribution. In the last step of derivation, we harness the observation that  $2\pi\sigma^2$  is independent on  $n$ .

Work out the marginal probability  $p(x)$ ,

$$p(x) = \int_t p(x, t) dt \quad (74)$$

$$= \int_t \frac{1}{2\pi\sigma^2 N} \sum_{n=1}^N \exp\left(-\frac{(x - x_n)^2}{2\sigma^2}\right) \exp\left(-\frac{(t - t_n)^2}{2\sigma^2}\right) dt \quad (75)$$

$$= \frac{1}{2\pi\sigma^2 N} \sum_{n=1}^N \int_t \exp\left(-\frac{(x - x_n)^2}{2\sigma^2}\right) \exp\left(-\frac{(t - t_n)^2}{2\sigma^2}\right) dt \quad (76)$$

$$= \frac{1}{2\pi\sigma^2 N} \sum_{n=1}^N \exp\left(-\frac{(x - x_n)^2}{2\sigma^2}\right) \int_t \exp\left(-\frac{(t - t_n)^2}{2\sigma^2}\right) dt \quad (77)$$

$$= \frac{1}{2\pi\sigma^2 N} \sum_{n=1}^N \exp\left(-\frac{(x - x_n)^2}{2\sigma^2}\right) (2\pi)^{\frac{1}{2}} \sigma \quad (78)$$

$$= \frac{(2\pi)^{\frac{1}{2}} \sigma}{2\pi\sigma^2 N} \sum_{n=1}^N \exp\left(-\frac{(x - x_n)^2}{2\sigma^2}\right) \quad (79)$$

$$= \frac{1}{(2\pi)^{\frac{1}{2}} \sigma N} \sum_{n=1}^N \exp\left(-\frac{(x - x_n)^2}{2\sigma^2}\right) \quad (80)$$

Note that to derive (74), sum rule is utilised. From (74) to (75), we cite the previous statement in (73). Derivation from (75) to (76) is based on linearity of integration. The validity of derivation from (76) to (77) comes from the observation that  $\frac{1}{2\pi\sigma^2} \exp\left(-\frac{(x - x_n)^2}{2\sigma^2}\right)$  is independent on  $t$ , on which integration is applied. The support for integration exponential term (from (77) to (78)) is presented in the following derivations. And we make use of the observation that  $(2\pi)^{\frac{1}{2}} \sigma$  is independent on the summation variable  $n$ .

$$\int_t \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right) dt = \int_t (2\pi)^{\frac{1}{2}} \sigma \frac{1}{(2\pi)^{\frac{1}{2}} \sigma} \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right) dt \quad \text{basic algebra} \quad (81)$$

$$= (2\pi)^{\frac{1}{2}} \sigma \int_t \frac{1}{(2\pi)^{\frac{1}{2}} \sigma} \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right) dt \quad \sigma \text{ independent on } t \quad (82)$$

$$= (2\pi)^{\frac{1}{2}} \sigma \cdot 1 \quad \text{normalisation} \quad (83)$$

$$= (2\pi)^{\frac{1}{2}} \sigma \quad (84)$$

Note that from (82) to (83), we use the normalisation property of gaussian distribution (integration over the whole space of its argument equals to 1.)

Then we derive  $p(t|x)$ , which is the probability of target variable  $t$  conditioned on  $x$

$$p(t|x) = \frac{p(x, t)}{p(x)} \quad (85)$$

$$= \frac{\frac{1}{2\pi\sigma^2 N} \sum_{n=1}^N \exp\left(-\frac{(x-x_n)^2}{2\sigma^2}\right) \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right)}{\frac{1}{(2\pi)^{\frac{1}{2}} \sigma N} \sum_{n=1}^N \exp\left(-\frac{(x-x_n)^2}{2\sigma^2}\right)} \quad (86)$$

$$= \frac{1}{(2\pi)^{\frac{1}{2}} \sigma N} \frac{\sum_{n=1}^N \exp\left(-\frac{(x-x_n)^2}{2\sigma^2}\right) \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right)}{\sum_{n=1}^N \exp\left(-\frac{(x-x_n)^2}{2\sigma^2}\right)} \quad (87)$$

$$= \frac{1}{(2\pi)^{\frac{1}{2}} \sigma N} \sum_{n=1}^N \left( \frac{\exp\left(-\frac{(x-x_n)^2}{2\sigma^2}\right)}{\sum_{n=1}^N \exp\left(-\frac{(x-x_n)^2}{2\sigma^2}\right)} \cdot \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right) \right) \quad (88)$$

Note that first step of derivation (85) above is based on product rule. Then to derive (86) we cited (73) for and (80) for  $p(x)$  respectively. The most important step is the third one, which lies in the observation that the whole term  $\sum_{n=1}^N \exp\left(-\frac{(x-x_n)^2}{2\sigma^2}\right)$  is independent on summation variable  $n$  so that it can put into the summation on the numerator.

Then we define  $k(x, x_n)$  to be,

$$k(x, x_n) = \frac{\exp\left(-\frac{(x-x_n)^2}{2\sigma^2}\right)}{\sum_{n=1}^N \exp\left(-\frac{(x-x_n)^2}{2\sigma^2}\right)} \quad (89)$$

Therefore,  $p(t|x)$  can be represented with the help of  $k(x, x_n)$  as follows,

$$p(t|x) = \frac{1}{(2\pi)^{\frac{1}{2}} \sigma N} \sum_{n=1}^N \left( k(x, x_n) \cdot \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right) \right) \quad (90)$$

Next derive the expectation  $\mathbb{E}(t|x)$  by using probability above,

$$\mathbb{E}(t|x) = \int_t t p(t|x) dt \quad (91)$$

$$= \int_t t \frac{1}{(2\pi)^{\frac{1}{2}} \sigma N} \sum_{n=1}^N \left( k(x, x_n) \cdot \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right) \right) dt \quad (92)$$

$$= \int_t \frac{1}{(2\pi)^{\frac{1}{2}} \sigma N} \sum_{n=1}^N \left( k(x, x_n) \cdot t \cdot \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right) \right) dt \quad (93)$$

$$= \frac{1}{(2\pi)^{\frac{1}{2}} \sigma N} \sum_{n=1}^N \left( \int_t k(x, x_n) \cdot t \cdot \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right) dt \right) \quad (94)$$

$$= \frac{1}{(2\pi)^{\frac{1}{2}} \sigma N} \sum_{n=1}^N \left( k(x, x_n) \int_t t \cdot \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right) dt \right) \quad (95)$$

Note that we make use of the definition of expectation for starting up this series of derivations. For second step, we refer to (90) to substitute  $p(t|x)$  in (92). Then since  $t$  is independent on summation variable  $n$ , we place the  $t$  into the term within summation of  $n$ . Next (94) is because of the linearity of integration operator. And since  $k(x, x_n)$  independent on  $t$  we can place  $k(x, x_n)$  out of the integration over  $t$ . In the following derivation, we seek to work out the integration of  $t \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right)$  over  $t$ .

$$\int_t t \cdot \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right) dt = \int_t (2\pi)^{\frac{1}{2}} \sigma \cdot \frac{1}{(2\pi)^{\frac{1}{2}} \sigma} \cdot t \cdot \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right) dt \quad (96)$$

$$= (2\pi)^{\frac{1}{2}} \sigma \int_t t \cdot \frac{1}{(2\pi)^{\frac{1}{2}} \sigma} \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right) dt \quad (97)$$

$$= (2\pi)^{\frac{1}{2}} \sigma t_n \quad (98)$$

Thus, we have

$$\mathbb{E}(t|x) = \frac{1}{(2\pi)^{\frac{1}{2}} \sigma N} \cdot \sum_{n=1}^N \left( k(x, x_n) (2\pi)^{\frac{1}{2}} \sigma t_n \right) \quad (99)$$

$$= \frac{1}{(2\pi)^{\frac{1}{2}} \sigma N} \cdot (2\pi)^{\frac{1}{2}} \sigma \cdot \sum_{n=1}^N \left( k(x, x_n) t_n \right) \quad (100)$$

$$= \frac{1}{N} \sum_{n=1}^N \left( k(x, x_n) t_n \right) \quad (101)$$

Finally, derive the conditional variance,

$$\text{VAR}(t|x) = \int_t \left( t - \mathbb{E}(t|x) \right)^2 p(t|x) dt \quad (102)$$

$$= \int_t \left( t^2 - 2t\mathbb{E}(t|x) + \mathbb{E}^2(t|x) \right) \cdot p(t|x) dt \quad (103)$$

$$= \int_t t^2 \cdot p(t|x) dt + \int_t -2t\mathbb{E}(t|x) \cdot p(t|x) dt + \int_t \mathbb{E}^2(t|x) p(t|x) dt \quad (104)$$

$$= \int_t t^2 \cdot p(t|x) dt - 2\mathbb{E}(t|x) \int_t t \cdot p(t|x) dt + \mathbb{E}^2(t|x) \int_t p(t|x) dt \quad (105)$$

$$= \int_t t^2 \cdot p(t|x) dt - 2\mathbb{E}^2(t|x) + \mathbb{E}^2(t|x) \int_t p(t|x) dt \quad (106)$$

$$= \int_t t^2 \cdot p(t|x) dt - 2\mathbb{E}^2(t|x) + \mathbb{E}^2(t|x) \quad (107)$$

$$= \int_t t^2 \cdot p(t|x) dt - \mathbb{E}^2(t|x) \quad (108)$$

It seems that computation for first term is more involved. However, we can use some tricks to transform the integration of first term.

$$\int_t t^2 \cdot p(t|x) dt = \int_t t^2 \cdot \frac{1}{(2\pi)^{\frac{1}{2}} \sigma N} \sum_{n=1}^N \left( k(x, x_n) \cdot \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right) \right) dt \quad (109)$$

$$= \int_t \frac{1}{(2\pi)^{\frac{1}{2}} \sigma N} \sum_{n=1}^N \left( k(x, x_n) \cdot t^2 \cdot \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right) \right) dt \quad (110)$$

$$= \frac{1}{(2\pi)^{\frac{1}{2}} \sigma N} \sum_{n=1}^N \left( \int_t k(x, x_n) \cdot t^2 \cdot \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right) dt \right) \quad (111)$$

$$= \frac{1}{(2\pi)^{\frac{1}{2}} \sigma N} \sum_{n=1}^N \left( k(x, x_n) \int_t t^2 \cdot \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right) dt \right) \quad (112)$$

$$= \frac{1}{(2\pi)^{\frac{1}{2}} \sigma N} \sum_{n=1}^N \left( k(x, x_n) \int_t \left( (t-t_n)^2 + 2tt_n - t_n^2 \right) \cdot \exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right) dt \right) \quad (113)$$

Then we focus on manipulating the integration in (113)

$$\int_t \left( (t - t_n)^2 + 2tt_n - t_n^2 \right) \cdot \exp\left(-\frac{(t - t_n)^2}{2\sigma^2}\right) dt \quad (114)$$

$$= \int_t (t - t_n)^2 \exp\left(-\frac{(t - t_n)^2}{2\sigma^2}\right) dt + \int_t 2tt_n \cdot \exp\left(-\frac{(t - t_n)^2}{2\sigma^2}\right) dt + \int_t -t_n^2 \cdot \exp\left(-\frac{(t - t_n)^2}{2\sigma^2}\right) dt \quad (115)$$

$$= \int_t (t - t_n)^2 \exp\left(-\frac{(t - t_n)^2}{2\sigma^2}\right) dt + 2t_n \int_t t \cdot \exp\left(-\frac{(t - t_n)^2}{2\sigma^2}\right) dt - t_n^2 \int_t \exp\left(-\frac{(t - t_n)^2}{2\sigma^2}\right) dt \quad (116)$$

$$= \int_t (t - t_n)^2 \exp\left(-\frac{(t - t_n)^2}{2\sigma^2}\right) dt + 2t_n \cdot ((2\pi)^{\frac{1}{2}} \sigma t_n) - t_n^2 ((2\pi)^{\frac{1}{2}} \sigma) \quad (117)$$

$$= \int_t (t - t_n)^2 \exp\left(-\frac{(t - t_n)^2}{2\sigma^2}\right) dt + 2t_n^2 (2\pi)^{\frac{1}{2}} \sigma - t_n^2 (2\pi)^{\frac{1}{2}} \sigma \quad (118)$$

$$= (2\pi)^{\frac{1}{2}} \sigma^3 + t_n^2 (2\pi)^{\frac{1}{2}} \sigma \quad (119)$$

$$= (2\pi)^{\frac{1}{2}} \sigma (\sigma^2 + t_n^2) \quad (120)$$

Explain the derivation for  $\exp\left(-\frac{(t-t_n)^2}{2\sigma^2}\right)$  in above manipulation,

$$\int_t (t - t_n)^2 \exp\left(-\frac{(t - t_n)^2}{2\sigma^2}\right) dt = \int_t (2\pi)^{\frac{1}{2}} \sigma \frac{1}{(2\pi)^{\frac{1}{2}} \sigma} (t - t_n)^2 \exp\left(-\frac{(t - t_n)^2}{2\sigma^2}\right) dt \quad (121)$$

$$= (2\pi)^{\frac{1}{2}} \sigma \int_t (t - t_n)^2 \cdot \frac{1}{(2\pi)^{\frac{1}{2}} \sigma} \exp\left(-\frac{(t - t_n)^2}{2\sigma^2}\right) dt \quad (122)$$

$$= (2\pi)^{\frac{1}{2}} \sigma \cdot \sigma^2 \quad (123)$$

$$= (2\pi)^{\frac{1}{2}} \sigma^3 \quad (124)$$

Finally, turn back to conditional variance (108),

$$\text{VAR}(t|x) = \int_t t^2 \cdot p(t|x) dt - \mathbb{E}^2(t|x) \quad (125)$$

$$= \frac{1}{(2\pi)^{\frac{1}{2}} \sigma N} \sum_{n=1}^N \left( k(x, x_n) \int_t \left( (t - t_n)^2 + 2tt_n - t_n^2 \right) \cdot \exp\left(-\frac{(t - t_n)^2}{2\sigma^2}\right) dt \right) - \mathbb{E}^2(t|x) \quad (126)$$

$$= \frac{1}{(2\pi)^{\frac{1}{2}} \sigma N} \sum_{n=1}^N \left( k(x, x_n) \cdot (2\pi)^{\frac{1}{2}} \sigma (\sigma^2 + t_n^2) \right) - \mathbb{E}^2(t|x) \quad (127)$$

$$= \frac{(2\pi)^{\frac{1}{2}} \sigma}{(2\pi)^{\frac{1}{2}} \sigma N} \sum_{n=1}^N \left( k(x, x_n) (\sigma^2 + t_n^2) \right) - \mathbb{E}^2(t|x) \quad (128)$$

$$= \frac{1}{N} \sum_{n=1}^N \left( k(x, x_n) (\sigma^2 + t_n^2) \right) - \left( \frac{1}{N} \sum_{n=1}^N (k(x, x_n) t_n) \right)^2 \quad (129)$$

$$= \frac{1}{N} \sum_{n=1}^N k(x, x_n) \sigma^2 + \frac{1}{N} \sum_{n=1}^N k(x, x_n) t_n^2 - \frac{1}{N^2} \left( \sum_{n=1}^N k(x, x_n) t_n \right)^2 \quad (130)$$

As we will derive in the next question,  $\sum_{n=1}^N k(x, x_n) = 1$  will hold. But we use it here first,

$$\text{VAR}(t|x) = \frac{1}{N} \sum_{n=1}^N k(x, x_n) \sigma^2 + \frac{1}{N} \sum_{n=1}^N k(x, x_n) t_n^2 - \frac{1}{N^2} \left( \sum_{n=1}^N k(x, x_n) t_n \right)^2 \quad (131)$$

### 3.2 Show that the function $\sum_{n=1}^N k(x, x_n) = 1$

In last section, we have defined the  $k(x, x_n)$  in (89) as

$$k(x, x_n) = \frac{\exp\left(-\frac{(x-x_n)^2}{2\sigma^2}\right)}{\sum_{n=1}^N \exp\left(-\frac{(x-x_n)^2}{2\sigma^2}\right)}$$

Now we start to prove that

$$\sum_{n=1}^N k(x, x_n) = 1. \quad (132)$$

Based on the definition, we add a summation operator over variable  $n$

$$\sum_{n=1}^N k(x, x_n) = \sum_{n=1}^N \left( \frac{\exp\left(-\frac{(x-x_n)^2}{2\sigma^2}\right)}{\sum_{m=1}^N \exp\left(-\frac{(x-x_m)^2}{2\sigma^2}\right)} \right) \quad (133)$$

Since the denominator is one term that has executed summation over variable  $n$ , it is obvious that the denominator  $\sum_{m=1}^N \exp\left(-\frac{(x-x_m)^2}{2\sigma^2}\right)$  is independent on  $n$ . Thus, it is valid to pull the denominator out of the newly added summation operator as follows.

$$\sum_{n=1}^N k(x, x_n) = \frac{1}{\sum_{m=1}^N \exp\left(-\frac{(x-x_m)^2}{2\sigma^2}\right)} \cdot \sum_{n=1}^N \exp\left(-\frac{(x-x_n)^2}{2\sigma^2}\right) \quad (134)$$

which can be written as

$$\sum_{n=1}^N k(x, x_n) = \frac{\sum_{n=1}^N \exp\left(-\frac{(x-x_n)^2}{2\sigma^2}\right)}{\sum_{m=1}^N \exp\left(-\frac{(x-x_m)^2}{2\sigma^2}\right)} \quad (135)$$

It is obvious that the numerator and denominator are of the same value. Thus

$$\sum_{n=1}^N k(x, x_n) = 1 \quad (136)$$

## 4 Maximum Margin Hyperplane

### 4.1 Two points are sufficient to determine the hyperplane

We start from the optimisation objective if the maximum margin hyperplane is desired.

$$\arg \max \text{ margin} \quad (137)$$

As widely known, margin is defined to be the minimal distance a set of data objects of the same class to the decision surface (hyperplane). Therefore, we have

$$\arg \max \min \text{ distance} \quad (138)$$

Now introduce the expression for the distance of each data objects to one hyperplane from lecture notes.

$$\text{distance} = \frac{t_n(\mathbf{w}^T \mathbf{x}_n + b)}{\|\mathbf{w}\|} \quad (139)$$

As derivation provided above, we have the following complete optimisation objective for the problem of maximum margin hyperplane.

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \min \left( \frac{t_n(\mathbf{w}^T \mathbf{x}_n + b)}{\|\mathbf{w}\|} \right) \quad (140)$$

Then we need to transform the acquired optimisation objective. Since the  $\|\mathbf{w}\|$  is independent on  $n$ , which is the index of data objects, we can extract the norm  $\|\mathbf{w}\|$  outward, and the resulted optimisation objective is to be

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \frac{1}{\|\mathbf{w}\|} \min (t_n(\mathbf{w}^T \mathbf{x}_n + b)) \quad (141)$$

Then we simplify the problem by extracting a common factor from  $\mathbf{w}$  and  $b$  to make the following happen,

$$\forall n, t_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \quad (142)$$

Note that the above trick is available since the case here only involves 2 points and thus the space is linearly separable for given data objects. And we can further extend the consequence of the trick as follows,

$$\min (t_n(\mathbf{w}^T \mathbf{x}_n + b)) = 1 \quad (143)$$

In this manner, the optimisation objective is converted to be

$$\begin{aligned} \mathbf{w}^* &= \arg \max_{\mathbf{w}} \frac{1}{\|\mathbf{w}\|} \\ \text{s.t. } &\forall n, t_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \end{aligned} \quad (144)$$

Further simplify the optimisation objective, we have

$$\begin{aligned} \mathbf{w}^* &= \arg \min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{s.t. } &\forall n, t_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \end{aligned} \quad (145)$$

Construct Lagrangian as follows,

$$\mathcal{L} = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{n=1}^N a_n (t_n(\mathbf{w}^T \mathbf{x}_n + b) - 1) \quad (146)$$

$$\begin{aligned} \text{s.t. } &a_n \geq 0 \\ &t_n(\mathbf{w}^T \mathbf{x}_n + b) - 1 \geq 0 \\ &a_n(1 - t_n(\mathbf{w}^T \mathbf{x}_n + b)) = 0 \end{aligned} \quad (147)$$

Given two points  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , which are of different class, it is necessary to treat those two data objects to have minimal distance to the hyperplane (margin) for each class. Besides, two-point-case for SVM must be linearly separable (two points cannot be equal to each other).

$$t_1(\mathbf{w}^T \mathbf{x}_1 + b) = 1 \quad (148)$$

$$t_2(\mathbf{w}^T \mathbf{x}_2 + b) = 1 \quad (149)$$

And the corresponding Lagrangian would be

$$\mathcal{L} = \frac{1}{2} \mathbf{w}^T \mathbf{w} + a_1(1 - t_1(\mathbf{w}^T \mathbf{x}_1 + b)) + a_2(1 - t_2(\mathbf{w}^T \mathbf{x}_2 + b)) \quad (150)$$

$$= \frac{1}{2} \mathbf{w}^T \mathbf{w} + (a_1 + a_2) - a_1 t_1(\mathbf{w}^T \mathbf{x}_1 + b) - a_2 t_2(\mathbf{w}^T \mathbf{x}_2 + b) \quad (151)$$

To simplify the problem, we specify the  $t_1$  and  $t_2$  to be +1, -1 respectively, that is

$$t_1 = 1 \quad (152)$$

$$t_2 = -1 \quad (153)$$

Therefore, the derived formula can be represented as

$$\mathbf{w}^T \mathbf{x}_1 + b = 1 \quad (154)$$

$$\mathbf{w}^T \mathbf{x}_2 + b = -1 \quad (155)$$

And the lagrangian can be simplified to be

$$\mathcal{L} = \frac{1}{2} \mathbf{w}^T \mathbf{w} + (a_1 + a_2) - a_1(\mathbf{w}^T \mathbf{x}_1 + b) + a_2(\mathbf{w}^T \mathbf{x}_2 + b) \quad (156)$$

Take derivative with regard to  $\mathbf{w}^T$ , and  $b$ , we have

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{w} + (a_2 \mathbf{x}_2 - a_1 \mathbf{x}_1) \quad (157)$$

$$\frac{\partial \mathcal{L}}{\partial b} = a_2 - a_1 \quad (158)$$

$$\frac{\partial \mathcal{L}}{\partial a_1} = 1 - \mathbf{w}^T \mathbf{x}_1 - b \quad (159)$$

$$\frac{\partial \mathcal{L}}{\partial a_2} = 1 + \mathbf{w}^T \mathbf{x}_2 + b \quad (160)$$

Set partial derivative of lagrangian term with regard to  $\mathbf{w}$ ,  $b$ ,  $a_1$  and  $a_2$ , we have

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 0, \quad \frac{\partial \mathcal{L}}{\partial b} = 0, \quad \frac{\partial \mathcal{L}}{\partial a_1} = 0, \quad \frac{\partial \mathcal{L}}{\partial a_2} = 0 \quad (161)$$

Solve the set of equation, we have

$$a_1 = a_2 = \frac{2}{(\mathbf{x}_2 - \mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1)} = \frac{2}{\|\mathbf{x}_2 - \mathbf{x}_1\|^2} \quad (162)$$

$$\mathbf{w} = \frac{2(\mathbf{x}_1 - \mathbf{x}_2)}{(\mathbf{x}_2 - \mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1)} = \frac{2(\mathbf{x}_1 - \mathbf{x}_2)}{\|\mathbf{x}_2 - \mathbf{x}_1\|^2} \quad (163)$$

$$b = \frac{(\mathbf{x}_2 - \mathbf{x}_1)^T (\mathbf{x}_1 + \mathbf{x}_2)}{(\mathbf{x}_2 - \mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1)} = \frac{(\mathbf{x}_2 - \mathbf{x}_1)^T (\mathbf{x}_1 + \mathbf{x}_2)}{\|\mathbf{x}_2 - \mathbf{x}_1\|^2} \quad (164)$$

Note that  $\|\cdot\|$  is l2-norm.

As we can see, whatever the dimension of given data objects  $\mathbf{x}_1$  and  $\mathbf{x}_2$  is, two points are sufficient to determine the hyperplane (the  $\mathbf{w}$  and  $b$  can be figured out).



## 4.2 Dimension of this maximum margin hyperplane

Obviously, the dimension of this maximum margin hyperplane is  $D - 1$ , where the  $D$  is dimensionality of two given points  $\mathbf{x}_1$  and  $\mathbf{x}_2$ .

Formally, this is because the hyperplane we found is perpendicular to the  $\mathbf{x}_1 - \mathbf{x}_2$  (as shown in the above equation, direction vector  $\mathbf{w}$  of the hyperplane is the same as  $\mathbf{x}_1 - \mathbf{x}_2$ ), and thus all points in this hyperplane lose the freedom in the direction  $\mathbf{x}_1 - \mathbf{x}_2$ . And since the given point  $\mathbf{x}_1$  and  $\mathbf{x}_2$  is of dimension  $D$ , the dimension of hyperplane should be  $D - 1$ .

## 4.3 Discriminant function $f(\mathbf{x})$

The discriminant function  $f(\mathbf{x})$  can be

$$f(\mathbf{x}) = \mathbf{w}^T \cdot \mathbf{x} + b \quad (165)$$

$$= \frac{2(\mathbf{x}_1 - \mathbf{x}_2)^T}{\|\mathbf{x}_2 - \mathbf{x}_1\|^2} \cdot \mathbf{x} + \frac{(\mathbf{x}_2 - \mathbf{x}_1)^T(\mathbf{x}_1 + \mathbf{x}_2)}{\|\mathbf{x}_2 - \mathbf{x}_1\|^2} \quad (166)$$

Since we only care about the sign of  $f(\mathbf{x})$ , the  $f(\mathbf{x})$  can be further simplified by varying the scale.

$$f(\mathbf{x}) = 2(\mathbf{x}_1 - \mathbf{x}_2)^T \cdot \mathbf{x} + (\mathbf{x}_2 - \mathbf{x}_1)^T(\mathbf{x}_1 + \mathbf{x}_2) \quad (167)$$

$$= (\mathbf{x}_2 - \mathbf{x}_1)^T(\mathbf{x}_1 + \mathbf{x}_2 - 2\mathbf{x}) \quad (168)$$

Classification result:

$$y(\mathbf{x}) = \begin{cases} 1 & \text{if } (\mathbf{x}_2 - \mathbf{x}_1)^T(\mathbf{x}_1 + \mathbf{x}_2 - 2\mathbf{x}) > 0 \\ \text{reject} & \text{if } (\mathbf{x}_2 - \mathbf{x}_1)^T(\mathbf{x}_1 + \mathbf{x}_2 - 2\mathbf{x}) = 0 \\ -1 & \text{if } (\mathbf{x}_2 - \mathbf{x}_1)^T(\mathbf{x}_1 + \mathbf{x}_2 - 2\mathbf{x}) < 0 \end{cases} \quad (169)$$

Note that when  $f(\mathbf{x}) = 0$ , this maximum margin classifier cannot determine which class the given point  $\mathbf{x}$  is in. Therefore, this classifier reject the decision of class of  $\mathbf{x}$ .

## 4.4 Discuss the values of the Lagrange parameters

From (146), we have general version of lagrangian  $\mathcal{L}$

$$\begin{aligned} \mathcal{L} &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{n=1}^N a_n (1 - t_n(\mathbf{w}^T \mathbf{x}_n + b)) \\ \text{s.t.} \quad &a_n \geq 0 \\ &1 - t_n(\mathbf{w}^T \mathbf{x}_n + b) \leq 0 \\ &a_n (1 - t_n(\mathbf{w}^T \mathbf{x}_n + b)) = 0 \end{aligned} \quad (170)$$

Based on the KKT condition, we have

$$\begin{cases} a_n = 0 & \text{if } 1 - t_n(\mathbf{w}^T \mathbf{x}_n + b) < 0 \\ a_n > 0 & \text{if } 1 - t_n(\mathbf{w}^T \mathbf{x}_n + b) = 0 \end{cases} \quad (171)$$

Since we have known that when condition  $1 - t_n(\mathbf{w}^T \mathbf{x}_n + b) = 0$ , the point  $\mathbf{x}_n$  has the margin (the ones with smallest distance to hyperplane in separate space). That is,  $\mathbf{x}_n$  is the smallest margin that affects the position and orientation of the hyperplane,  $a_n$  can be of any value that is bigger than zero. Conversely, if  $\mathbf{x}_n$  is not the margin,  $a_n$  must be zero to ignore the value of the condition term, since the optimisation objective of our model is to find a hyperplane so that the margin is maximum (non-margin objects are not taken into consideration).

Therefore, to some extent, the functionality of lagrangian parameter  $a_n$  is to ignore points that are unnecessary to be considered.

In our problem, since there is only two points and they must be respective margin in its located space, both of these two points satisfy  $1 - t_n(\mathbf{w}^T \mathbf{x}_n + b) = 0$ . And hence, these two points participate in determining the position and orientation of hyperplane.

## 5 Kernels and Feature Maps

### 5.1 Prove $k(\mathbf{x}, \mathbf{z})$ to be Valid Kernel in 2 Dimension

To prove  $k(\mathbf{x}, \mathbf{z})$  is a valid kernel in 2 dimension, we should figure out the specific feature mapping that corresponds to the  $k(\mathbf{x}, \mathbf{z})$ . But first we declare

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \quad \mathbf{z} = \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} \quad (172)$$

Then we start to figure out the feature map,

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^3 \quad (173)$$

$$= (x_1 z_1 + x_2 z_2 + 1)^3 \quad (174)$$

$$= (x_1 z_1 + x_2 z_2)^3 + 3(x_1 z_1 + x_2 z_2)^2 + 3(x_1 z_1 + x_2 z_2) + 1 \quad (175)$$

$$= x_1^3 z_1^3 + 3x_1^2 z_1^2 x_2 z_2 + 3x_1^1 z_1^1 x_2^2 z_2^2 + x_2^3 z_2^3 \\ + 3x_1^2 z_1^2 + 6x_1^1 z_1^1 x_2^1 z_2^1 + 3x_2^2 z_2^2 + 3x_1^1 z_1^1 + 3x_2^1 z_2^1 + 1 \quad (176)$$

$$= x_1^3 \cdot z_1^3 + \sqrt{3}x_1^2 x_2 \cdot \sqrt{3}z_1^2 z_2 + \sqrt{3}x_1 x_2^2 \cdot \sqrt{3}z_1 z_2^2 + x_2^3 \cdot z_2^3 \\ + \sqrt{3}x_1^2 \cdot \sqrt{3}z_1^2 + \sqrt{6}x_1 x_2 \cdot \sqrt{6}z_1 z_2 + \sqrt{3}x_2^2 \cdot \sqrt{3}z_2^2 \\ + \sqrt{3}x_1 \cdot \sqrt{3}z_1 + \sqrt{3}x_2 \cdot \sqrt{3}z_2 + 1 \cdot 1 \quad (177)$$

Next, define feature map  $\Phi(\mathbf{x})$  to be,

$$\Phi(\mathbf{x}) = \begin{pmatrix} x_1^3 \\ \sqrt{3}x_1^2 x_2 \\ \sqrt{3}x_1 x_2^2 \\ x_2^3 \\ \sqrt{3}x_1^2 \\ \sqrt{6}x_1 x_2 \\ \sqrt{3}x_2^2 \\ \sqrt{3}x_1 \\ \sqrt{3}x_2 \\ 1 \end{pmatrix} \quad (178)$$

Obviously,  $k(\mathbf{x}, \mathbf{z})$  can be decomposed as inner product of  $\Phi(\mathbf{x})$  and  $\Phi(\mathbf{z})$

$$k(\mathbf{x}, \mathbf{z}) = \langle \Phi(\mathbf{x}), \Phi(\mathbf{z}) \rangle \quad (179)$$

where the inner product is defined to be,

$$\langle A, B \rangle = \text{tr}(A^T B) \quad (180)$$

Therefore, we can conclude that

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^3 \text{ is a valid kernel when input space has only two dimensions.}$$

## 5.2 Generalize Input Space to Higher Dimension

**Yes!** We can use the same method in section 5.1 to validate the extrapolation on dimension  $n$ . The specific technique is to utilise the Mathematical Induction.

**Proof Objective:**

$$\forall n \geq 2, \dim(\mathbf{x}) = \dim(\mathbf{z}) = n, k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^3 \text{ is valid kernel.} \quad (181)$$

where the  $\dim()$  operator returns the dimensionality of its argument.

**Base case:**

the statement that kernel  $k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^3$  is valid for  $n = 2$  has been proven in (179).

**Inductive cases:** Assume that the

$$k_M(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^3, \dim(\mathbf{x}) = \dim(\mathbf{z}) = M \text{ is valid kernel function.} \quad (182)$$

Prove that

$$k_{M+1}(\mathbf{x}^*, \mathbf{z}^*) = ((\mathbf{x}^*)^T \mathbf{z}^* + 1)^3, \text{ is also valid kernel function.} \quad (183)$$

where  $\mathbf{x}^* = (\mathbf{x}, x_{M+1})$ ,  $\mathbf{z}^* = (\mathbf{z}, z_{M+1})$ .

**Proof for Inductive Cases:**

First of all, extend  $k_M(\mathbf{x}, \mathbf{z})$ ,

$$k_M(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^3 \quad (184)$$

$$= \left( \left( \sum_{i=1}^M x_i z_i \right) + 1 \right)^3 \quad (185)$$

$$= \left( \sum_{i=1}^M x_i z_i \right)^3 + 3 \left( \sum_{i=1}^M x_i z_i \right)^2 + 3 \left( \sum_{i=1}^M x_i z_i \right) + 1 \quad (186)$$

And from the assumption of inductive cases (182), we have,

$$\exists \Phi(\cdot), \text{ such that } \Phi_M^T(\mathbf{x}) \Phi_M(\mathbf{z}) = k_M(\mathbf{x}, \mathbf{z}) \quad (187)$$

Then use  $\Phi_M(\cdot)$  to denote the  $\Phi(\cdot)$  above in the extension of  $k_M(\mathbf{x}, \mathbf{z})$

$$k_M(\mathbf{x}, \mathbf{z}) = \left( \sum_{i=1}^M x_i z_i \right)^3 + 3 \left( \sum_{i=1}^M x_i z_i \right)^2 + 3 \left( \sum_{i=1}^M x_i z_i \right) + 1 = \Phi_M^T(\mathbf{x}) \Phi_M(\mathbf{z}) \quad (188)$$

Now, start to manipulate  $k_{M+1}(\mathbf{x}^*, \mathbf{z}^*)$  as follows,

$$k_{M+1}(\mathbf{x}^*, \mathbf{z}^*) = \left( \sum_{i=1}^{M+1} x_i z_i \right)^3 \quad (189)$$

$$= \left( \left( \sum_{i=1}^M x_i z_i \right) + (x_{M+1} z_{M+1} + 1) \right)^3 \quad (190)$$

$$\begin{aligned} &= \left( \sum_{i=1}^M x_i z_i \right)^3 + 3 \left( \sum_{i=1}^M x_i z_i \right)^2 (x_{M+1} z_{M+1} + 1) \\ &\quad + 3 \left( \sum_{i=1}^M x_i z_i \right) (x_{M+1} z_{M+1} + 1)^2 + (x_{M+1} z_{M+1} + 1)^3 \end{aligned} \quad (191)$$

$$\begin{aligned} &= \left( \sum_{i=1}^M x_i z_i \right)^3 + 3 \left( \sum_{i=1}^M x_i z_i \right)^2 x_{M+1} z_{M+1} + 3 \left( \sum_{i=1}^M x_i z_i \right)^2 \\ &\quad + 3 \left( \sum_{i=1}^M x_i z_i \right) (x_{M+1} z_{M+1})^2 + 6 \left( \sum_{i=1}^M x_i z_i \right) (x_{M+1} z_{M+1}) + 3 \left( \sum_{i=1}^M x_i z_i \right) \\ &\quad + (x_{M+1} z_{M+1})^3 + 3 (x_{M+1} z_{M+1})^2 + 3 x_{M+1} z_{M+1} + 1 \end{aligned} \quad (192)$$

Next, replace the sum of boxed terms in (192) by using (188), we have the followings,

$$\begin{aligned}
 k_{M+1}(\mathbf{x}^*, \mathbf{z}^*) = & \underbrace{\Phi_M^T(\mathbf{x})\Phi_M(\mathbf{z})}_{\sqrt{6}\Phi_{(M+1)}^{***}(\mathbf{x}^*) \cdot \sqrt{6}\Phi_{(M+1)}^{***}(\mathbf{z}^*)} + \underbrace{3(\sum_{i=1}^M x_i z_i)^2 x_{M+1} z_{M+1}}_{x_{M+1}^3 \cdot z_{M+1}^3} + \underbrace{3(\sum_{i=1}^M x_i z_i)(x_{M+1} z_{M+1})^2}_{\sqrt{3}x_{M+1}^2 \cdot \sqrt{3}z_{M+1}^2} \\
 & + \underbrace{6(\sum_{i=1}^M x_i z_i)(x_{M+1} z_{M+1})}_{\sqrt{3}x_{M+1} \cdot \sqrt{3}z_{M+1}} + \underbrace{(x_{M+1} z_{M+1})^3}_{x_{M+1}^3 \cdot z_{M+1}^3} + \underbrace{3(x_{M+1} z_{M+1})^2}_{\sqrt{3}x_{M+1}^2 \cdot \sqrt{3}z_{M+1}^2} + \underbrace{3x_{M+1} z_{M+1}}_{\sqrt{3}x_{M+1} \cdot \sqrt{3}z_{M+1}}
 \end{aligned} \tag{193}$$

In above formula, the decomposition of each term is presented. For safety, we also write down the formal representation as follows,

$$\Phi_M^T(\mathbf{x})\Phi_M(\mathbf{z}) = \Phi_M(\mathbf{x}) \cdot \Phi_M(\mathbf{z}) \tag{194}$$

$$3(\sum_{i=1}^M x_i z_i)^2 x_{M+1} z_{M+1} = \sqrt{3}\Phi_{(M+1)}^*(\mathbf{x}^*) \cdot \sqrt{3}\Phi_{(M+1)}^*(\mathbf{z}^*) \tag{195}$$

$$3(\sum_{i=1}^M x_i z_i)(x_{M+1} z_{M+1})^2 = \sqrt{3}\Phi_{(M+1)}^{**}(\mathbf{x}^*) \cdot \sqrt{3}\Phi_{(M+1)}^{**}(\mathbf{z}^*) \tag{196}$$

$$6(\sum_{i=1}^M x_i z_i)(x_{M+1} z_{M+1}) = \sqrt{6}\Phi_{(M+1)}^{***}(\mathbf{x}^*) \cdot \sqrt{6}\Phi_{(M+1)}^{***}(\mathbf{z}^*) \tag{197}$$

$$(x_{M+1} z_{M+1})^3 = x_{M+1}^3 \cdot z_{M+1}^3 \tag{198}$$

$$3(x_{M+1} z_{M+1})^2 = \sqrt{3}x_{M+1}^2 \cdot \sqrt{3}z_{M+1}^2 \tag{199}$$

$$3x_{M+1} z_{M+1} = \sqrt{3}x_{M+1} \cdot \sqrt{3}z_{M+1} \tag{200}$$

where  $\Phi_{(M+1)}^*$ ,  $\Phi_{(M+1)}^{**}$ ,  $\Phi_{(M+1)}^{***}$  are defined as follows,

$$\Phi_{(M+1)}^*(\mathbf{x}^*) = \begin{pmatrix} x_1 x_1 x_{M+1} \\ x_1 x_2 x_{M+1} \\ \vdots \\ x_1 x_j x_{M+1} \\ x_1 x_{j+1} x_{M+1} \\ \vdots \\ x_1 x_M x_{M+1} \\ x_2 x_1 x_{M+1} \\ \vdots \\ x_i x_j x_{M+1} \\ x_i x_{j+1} x_{M+1} \\ \vdots \\ x_i x_M x_{M+1} \\ x_{i+1} x_1 x_{M+1} \\ \vdots \\ x_M x_M x_{M+1} \end{pmatrix} \tag{201}$$

$$\Phi_{(M+1)}^{**}(\mathbf{x}^*) = \begin{pmatrix} x_1 x_{M+1}^2 \\ x_2 x_{M+1}^2 \\ \vdots \\ x_i x_{M+1}^2 \\ x_{i+1} x_{M+1}^2 \\ \vdots \\ x_M x_{M+1}^2 \end{pmatrix} \tag{202}$$

$$\Phi_{(M+1)}^{***}(\mathbf{x}^*) = \begin{pmatrix} x_1 x_{M+1} \\ x_2 x_{M+1} \\ \vdots \\ x_i x_{M+1} \\ x_{i+1} x_{M+1} \\ \vdots \\ x_M x_{M+1} \end{pmatrix} \tag{203}$$

Therefore, we can figure out the feature mapping  $\Phi_{M+1}$ , such that the following is satisfied,

$$k_{M+1}(\mathbf{x}^*, \mathbf{z}^*) = \Phi_{M+1}^T(\mathbf{x}^*)\Phi_{M+1}(\mathbf{z}^*) \tag{204}$$

where the specific feature mapping is of the form as follows,

$$\Phi_{M+1}(\mathbf{x}^*) = \begin{pmatrix} \Phi_M(\mathbf{x}) \\ \sqrt{3}\Phi_{(M+1)}^*(\mathbf{x}^*) \\ \sqrt{3}\Phi_{(M+1)}^{**}(\mathbf{x}^*) \\ \sqrt{6}\Phi_{(M+1)}^{***}(\mathbf{x}^*) \\ x_{M+1}^3 \\ \sqrt{3}x_{M+1}^2 \\ \sqrt{3}x_{M+1} \end{pmatrix} \quad (205)$$

Hence,  $k_{M+1}(\mathbf{x}^*, \mathbf{z}^*)$  is also a valid function since we have proven that

$$\exists \Phi_{M+1}(\mathbf{x}^*), \text{ such that } k_{M+1}(\mathbf{x}^*, \mathbf{z}^*) = \langle \Phi_{M+1}(\mathbf{x}^*), \Phi_{M+1}(\mathbf{z}^*) \rangle \quad (206)$$

where the definition of inner product follows (180).

Hence, we successfully prove the inductive part (183). Based on the theory of Mathematical Induction, it is concluded that

$$\forall n \geq 2, \dim(\mathbf{x}) = \dim(\mathbf{z}) = n, k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^3 \text{ is valid kernel.} \quad (207)$$

### 5.3 Generalize Kernel to Higher Power

**Still Yes.** We can still prove the validity of generalized kernel by finding feature map. The technique for proving generalization on  $p$  is rather similiar to section 5.2.

**Proof Objective:**

$$\forall p \geq 3, k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^p \text{ is a valid kernel.} \quad (208)$$

**Base Case:**  $p = 3$

For the case of dimension  $n = 1$ , it is obvious that the  $k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}\mathbf{z} + 1)^3$  is a valid kernel. We provide a quick proof in the followings.

$$k(x, z) = (xz + 1)^3 = x^3 z^3 + 3x^2 z^2 + 3xz + 1 \quad (209)$$

We can define feature mapping  $\Phi(x)$  to be

$$\Phi(x) = \begin{pmatrix} x^3 \\ \sqrt{3}x^2 \\ \sqrt{3}x \\ 1 \end{pmatrix} \quad (210)$$

In this manner, we can represent  $k(x, z)$  as inner product of two feature mapping.

$$k(x, z) = \langle \Phi(x), \Phi(z) \rangle \quad (211)$$

where the definition of inner product follows (180).

Thus, the statement that  $k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}\mathbf{z} + 1)^3$  is a valid kernel for dimension  $n = 1$  holds.

And the statement that  $k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^3$  is valid kernel for dimensions  $n \geq 2$  has been proven in last section 5.2.

Therefore, it can be conclude that for any positive integer dimension of  $\mathbf{x}$  and  $\mathbf{z}$ , where  $\dim(\mathbf{x}) = \dim(\mathbf{z}) = n \geq 0$ , we have  $k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^3$  is valid kernel.

**Inductive Cases:**

Assume that

$$k_Q(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^Q \text{ is valid kernel function} \quad (212)$$

Prove that

$$k_{Q+1}(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^{Q+1} \text{ is also valid kernel function} \quad (213)$$

**Proof for Inductive Cases:**

Start from the assumption (212) and the definition of valid kernel,

$$\exists \Phi(\mathbf{x}), \text{ such that } k_Q(\mathbf{x}, \mathbf{z}) = \langle \Phi(\mathbf{x}), \Phi(\mathbf{z}) \rangle = \Phi^T(\mathbf{x})\Phi(\mathbf{z}) \quad (214)$$

Since the  $k_M(\mathbf{x}, \mathbf{z})$  we talk about has the form,

$$k_Q(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^Q = \left( \sum_{i=1}^n x_i z_i + 1 \right)^Q \quad (215)$$

where  $n$  is the dimensionality of vector  $\mathbf{x}$  and  $\mathbf{z}$ .

We can derive the feature map for  $k_Q(\mathbf{x}, \mathbf{z})$ , denoted by  $\Phi_Q$  as follows,

$$(\Phi_Q(\mathbf{x}))^T \Phi_Q(\mathbf{z}) = \left( \sum_{i=1}^n x_i z_i + 1 \right)^Q \quad (216)$$

Then start to manipulate the  $k_{Q+1}(\mathbf{x}, \mathbf{z})$ ,

$$k_{Q+1}(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^{Q+1} \quad (217)$$

$$= (\mathbf{x}^T \mathbf{z} + 1)^Q (\mathbf{x}^T \mathbf{z} + 1) \quad (218)$$

$$= \underbrace{(\Phi_Q(\mathbf{x}))^T \Phi_Q(\mathbf{z}) \left( \sum_{i=1}^n x_i z_i \right)}_{x_i \Phi_Q(\mathbf{x}) \cdot z_i \Phi_Q(\mathbf{z})} + \underbrace{(\Phi_Q(\mathbf{x}))^T \Phi_Q(\mathbf{z})}_{\Phi_Q(\mathbf{x}) \cdot \Phi_Q(\mathbf{z})} \quad (219)$$

Note that the derivation from (218) to (219) make use of what we have already derived formula (216).

Obviously, we can derive the feature map  $\Phi_{Q+1}(\mathbf{x})$  for  $k_{Q+1}(\mathbf{x}, \mathbf{z})$ , that is,

$$\Phi_{Q+1}(\mathbf{x}) = \begin{pmatrix} x_1 \Phi_Q(\mathbf{x}) \\ \vdots \\ x_i \Phi_Q(\mathbf{x}) \\ \vdots \\ x_n \Phi_Q(\mathbf{x}) \\ \Phi_Q(\mathbf{x}) \end{pmatrix} \quad (220)$$

Such that,

$$k_{Q+1}(\mathbf{x}, \mathbf{z}) = (\Phi_{Q+1}(\mathbf{x}))^T \Phi_{Q+1}(\mathbf{z}) = \langle \Phi_{Q+1}(\mathbf{x}), \Phi_{Q+1}(\mathbf{z}) \rangle \Rightarrow \text{valid kernel} \quad (221)$$

where the definition of inner product follows (180).

Note that each term  $x_i \Phi_Q(\mathbf{x})$  ( $i = 1..n$ ) in (220) represents one vector whose element is multiplication of  $x_i$  and one corresponding element of original feature mapping  $\Phi_Q(\mathbf{x})$ . And the resulted feature mapping  $\Phi_{Q+1}(\mathbf{x})$  can be regarded as a huge vector formed by stacking the shown set of vectors in (220).

Therefore, we successfully prove the inductive cases (213). Based on the theory of Mathematical Induction, since base case and inductive cases are proven, it is concluded that

$$\forall p \geq 3, k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^p \text{ is a valid kernel.} \quad (222)$$

## 6 EM for Trees

### 6.1 Define latent variables and a joint generative model

In this case, the definition of latent variable  $\mathbf{z} = \{0, 1\}^k, k = \{1, 2, \dots, K\}$ , represents the real tree from which one measurement  $\mathbf{x}_n$  was generated (put it another way,  $z_k$  means the measurement corresponds to the tree whose index is  $k$ ). We model the tree one laser comes from as latent variable because this is not directly observed but are rather inferred.

Note that here we use the 1-of-k coding scheme for  $\mathbf{z}$ , that is

$$z_k = 1 \text{ or } 0 \quad (223)$$

$$\forall \mathbf{z}, \sum_{k=1}^K z_k = 1 \quad (224)$$

The joint generative model of the observations and latent variables given the parameters  $p(X, Z)$  means under current distribution of trees in this forest, the likelihood of collecting set of measurements  $X$  and such measurements to one certain tree.

Then we provide the derivation of joint distribution  $p(\mathbf{x}_n, \mathbf{z}_n)$ .

Based on the probability distribution of measurement  $\mathbf{x}_n$  conditioned on  $z_{nk}$ ,

$$p(\mathbf{x}_n | z_{nk} = 1) = \mathcal{N}(\boldsymbol{\mu}_k | \sigma^2 \mathbf{I}) \quad (225)$$

Since we use 1-of-k coding scheme, at all states of  $\mathbf{z}_n$ , only one element  $z_{nk}$  can be one while others are zeros. That is, no measurement is from two trees, which conforms to our assumption in the question that the trees have a small diameter and thus the offset of measurement from the true center is negligible. Therefore, we can represent  $p(\mathbf{x}_n | \mathbf{z}_n)$  in the way as follows,

$$p(\mathbf{x}_n | \mathbf{z}_n) = \prod_{k=1}^K \left( \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I}) \right)^{z_{nk}} \quad (226)$$

Similarly according to the 1-of-k coding scheme, set the  $p(z_k = 1) = \pi_k$ , in order to derive  $p(\mathbf{z})$ , we have

$$p(\mathbf{z}_n) = \prod_{k=1}^K \pi_k^{z_{nk}} \quad (227)$$

Consider the joint distribution of laser-range measurement  $p(\mathbf{x}_n, \mathbf{z}_n)$

$$p(\mathbf{x}_n, \mathbf{z}_n) = p(\mathbf{z}_n) p(\mathbf{x}_n | \mathbf{z}_n) \quad \text{product rule} \quad (228)$$

$$= \prod_{k=1}^K \pi_k^{z_{nk}} \prod_{k=1}^K \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})^{z_{nk}} \quad \text{cite (226) and (227)} \quad (229)$$

$$= \prod_{k=1}^K \left( \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I}) \right)^{z_{nk}} \quad \text{abstract the multiplication} \quad (230)$$

Again, based on the 1-of-k coding scheme, for each measurement  $\mathbf{x}_n$ , the corresponding latent variable has only one element  $z_{nk} = 1$  and others  $z_{nk} = 0$ . Therefore, all terms except one that  $z_{nk} = 1$  is one because the power is zero. We can easily work out the multiplication over  $k$  and derive the following representation of  $p(\mathbf{x}_n, \mathbf{z}_n)$ .

$$p(\mathbf{x}_n, \mathbf{z}_n) = 1 \times 1 \cdots \times \left( \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I}) \right) \times \cdots \times 1 \quad (231)$$

$$= \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I}) \quad (232)$$

where we use  $k$  to denote always the non-zero element index of  $\mathbf{z}$ .

Based on the sum rule, we derive the marginal probability of  $\mathbf{x}_n$ ,

$$p(\mathbf{x}_n) = \sum_{\mathbf{z}_n} p(\mathbf{x}_n, \mathbf{z}_n) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I}) \quad (233)$$



## 6.2 Derive the expected log likelihood of the joint model

To derive the expected log likelihood, we start from the log likelihood of the whole collected data sets  $\mathbf{X}$ , each row of whom is  $\mathbf{x}_n$ .

$$\ln p(\mathbf{X}|\theta) = \ln \prod_{n=1}^N p(\mathbf{x}_n|\theta) \quad \text{by i.i.d} \quad (234)$$

$$= \sum_{n=1}^N \ln p(\mathbf{x}_n|\theta) \quad \text{put log function in} \quad (235)$$

Next, to further extend the derivation shown above, we have to apply the following series of tricks.

$$p(\mathbf{x}_n, \mathbf{z}_n|\theta) = p(\mathbf{x}_n|\theta) \cdot p(\mathbf{z}_n|\mathbf{x}_n, \theta) \quad (236)$$

$$\ln p(\mathbf{x}_n, \mathbf{z}_n|\theta) = \ln p(\mathbf{x}_n|\theta) + \ln p(\mathbf{z}_n|\mathbf{x}_n, \theta) \quad (237)$$

$$\ln p(\mathbf{x}_n|\theta) = \ln p(\mathbf{x}_n, \mathbf{z}_n|\theta) - \ln p(\mathbf{z}_n|\mathbf{x}_n, \theta) \quad (238)$$

Multiply for both side a distribution  $q(\mathbf{z}_n)$ , for which  $\sum_{\mathbf{z}_n} q(\mathbf{z}_n) = 1$ .

$$q(\mathbf{z}_n) \cdot \ln p(\mathbf{x}_n|\theta) = q(\mathbf{z}_n) \cdot \ln p(\mathbf{x}_n, \mathbf{z}_n|\theta) - q(\mathbf{z}_n) \cdot \ln p(\mathbf{z}_n|\mathbf{x}_n, \theta) \quad (239)$$

And summing over  $\mathbf{z}_n$  for both side.

$$\sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln p(\mathbf{x}_n|\theta) \right) = \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln p(\mathbf{x}_n, \mathbf{z}_n|\theta) \right) - \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln p(\mathbf{z}_n|\mathbf{x}_n, \theta) \right) \quad (240)$$

Since  $\ln p(\mathbf{x}_n|\theta)$  is independent on  $\mathbf{z}_n$ , the left-hand-side term in above equation can be simplified as follows,

$$\sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln p(\mathbf{x}_n|\theta) \right) = \left( \sum_{\mathbf{z}_n} q(\mathbf{z}_n) \right) \cdot \ln p(\mathbf{x}_n|\theta) = \ln p(\mathbf{x}_n|\theta) \quad (241)$$

Note that in above derivation, we use the normalisation property of  $\sum_{\mathbf{z}_n} q(\mathbf{z}_n) = 1$ .

Take the simplified term back to equation (240), we have

$$\ln p(\mathbf{x}_n|\theta) = \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln p(\mathbf{x}_n, \mathbf{z}_n|\theta) \right) - \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln p(\mathbf{z}_n|\mathbf{x}_n, \theta) \right) \quad (242)$$

$$\begin{aligned} &= \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln p(\mathbf{x}_n, \mathbf{z}_n|\theta) \right) - \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln p(\mathbf{z}_n|\mathbf{x}_n, \theta) \right) \\ &\quad + \sum_{\mathbf{z}_n} q(\mathbf{z}_n) \ln q(\mathbf{z}_n) - \sum_{\mathbf{z}_n} q(\mathbf{z}_n) \ln q(\mathbf{z}_n) \end{aligned} \quad (243)$$

$$\begin{aligned} &= \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln p(\mathbf{x}_n, \mathbf{z}_n|\theta) - q(\mathbf{z}_n) \ln q(\mathbf{z}_n) \right) \\ &\quad - \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln p(\mathbf{z}_n|\mathbf{x}_n, \theta) - q(\mathbf{z}_n) \ln q(\mathbf{z}_n) \right) \end{aligned} \quad (244)$$

$$= \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln \frac{p(\mathbf{x}_n, \mathbf{z}_n|\theta)}{q(\mathbf{z}_n)} \right) - \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln \frac{p(\mathbf{z}_n|\mathbf{x}_n, \theta)}{q(\mathbf{z}_n)} \right) \quad (245)$$

Then put the subtraction into the term within the summation, we have

$$\ln p(\mathbf{x}_n|\theta) = \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln \frac{p(\mathbf{x}_n, \mathbf{z}_n|\theta)}{q(\mathbf{z}_n)} \right) + \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln \frac{q(\mathbf{z}_n)}{p(\mathbf{z}_n|\mathbf{x}_n, \theta)} \right) \quad (246)$$

Take above equation (246) into what we have previously derived (235),

$$\ln p(\mathbf{X}|\theta) = \sum_{n=1}^N \left( \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln \frac{p(\mathbf{x}_n, \mathbf{z}_n|\theta)}{q(\mathbf{z}_n)} \right) + \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln \frac{q(\mathbf{z}_n)}{p(\mathbf{z}_n|\mathbf{x}_n, \theta)} \right) \right) \quad (247)$$

$$= \sum_{n=1}^N \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln \frac{p(\mathbf{x}_n, \mathbf{z}_n|\theta)}{q(\mathbf{z}_n)} \right) + \sum_{n=1}^N \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln \frac{q(\mathbf{z}_n)}{p(\mathbf{z}_n|\mathbf{x}_n, \theta)} \right) \quad (248)$$

Obviously, the second term is KL divergence between  $q(\mathbf{z}_n)$  and  $p(\mathbf{z}_n|\mathbf{x}_n, \theta)$ . As we all know, KL divergence is always greater than zero, that is

$$KL\left(q(\mathbf{z}_n) \parallel p(\mathbf{z}_n|\mathbf{x}_n, \theta)\right) = \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln \frac{q(\mathbf{z}_n)}{p(\mathbf{z}_n|\mathbf{x}_n, \theta)} \right) \geq 0 \quad (249)$$

That is to say, the first term is always lower bound of  $\ln p(\mathbf{X}|\theta)$ . Formally, we have

$$\ln p(\mathbf{X}|\theta) \geq \sum_{n=1}^N \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln \frac{p(\mathbf{x}_n, \mathbf{z}_n|\theta)}{q(\mathbf{z}_n)} \right) \quad (250)$$

Now we apply the Expectation-Maximisation algorithm to derive the expected complete log likelihood. We first fix  $\theta$  and find smallest KL divergence in E-step to maximise the lower bound, that is, based on the property of KL divergence, to find a distribution  $q(\mathbf{z}_n)$  equaling to distribution  $p(\mathbf{z}_n|\mathbf{x}_n, \theta)$ , in which case the KL divergence is zero. And then in the M-step, we fix found  $q(\mathbf{z}_n)$  and maximise the lower bound with regard to parameter  $\theta$ .

That is to say, in each M-step, we have

$$q(\mathbf{z}_n) = p(\mathbf{z}_n|\mathbf{x}_n, \theta^{old}) \quad (251)$$

Now we turn to the optimisation objective in M-step. The optimisation of lower bound with regard to  $\theta$  can be written as,

$$\arg \max_{\theta} \sum_{n=1}^N \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln \frac{p(\mathbf{x}_n, \mathbf{z}_n|\theta)}{q(\mathbf{z}_n)} \right) \quad (252)$$

That is to optimise the following objective

$$\arg \max_{\theta} \left( \sum_{n=1}^N \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln p(\mathbf{x}_n, \mathbf{z}_n|\theta) \right) - \sum_{n=1}^N \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln q(\mathbf{z}_n) \right) \right) \quad (253)$$

Since  $q(\mathbf{z}_n)$  is not dependent on  $\theta$ , the optimisation objective above can be rewritten as

$$\arg \max_{\theta} \sum_{n=1}^N \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln p(\mathbf{x}_n, \mathbf{z}_n|\theta) \right) \quad (254)$$

Based on equation (251), we derived optimisation objective for M-step as follows.

$$\arg \max_{\theta} \sum_{n=1}^N \sum_{\mathbf{z}_n} \left( p(\mathbf{z}_n|\mathbf{x}_n, \theta^{old}) \cdot \ln p(\mathbf{x}_n, \mathbf{z}_n|\theta) \right) \quad (255)$$

As we can see, based on the definition of expectation, the term to be optimised above is exactly the expected log likelihood of this joint model.

$$Q(\theta, \theta^{old}) = \sum_{n=1}^N \sum_{\mathbf{z}_n} \left( p(\mathbf{z}_n|\mathbf{x}_n, \theta^{old}) \cdot \ln p(\mathbf{x}_n, \mathbf{z}_n|\theta) \right) \quad (256)$$

### 6.3 Derive the E-step update and show how to compute expectations

First from the (248), we have

$$\ln p(\mathbf{X}|\theta) = \sum_{n=1}^N \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln \frac{p(\mathbf{x}_n, \mathbf{z}_n|\theta)}{q(\mathbf{z}_n)} \right) + \sum_{n=1}^N \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln \frac{q(\mathbf{z}_n)}{p(\mathbf{z}_n|\mathbf{x}_n, \theta)} \right) \quad (257)$$

In the Expectation-step, we are gonna to minimise the KL divergence as follows. As widely known, the lowest value of KL divergence is zero if and only if the  $q(\mathbf{z}_n) = p(\mathbf{z}_n|\mathbf{x}_n)$ .

$$KL(q(\mathbf{z}_n) \parallel p(\mathbf{z}_n|\mathbf{x}_n, \theta)) = \sum_{\mathbf{z}_n} \left( q(\mathbf{z}_n) \cdot \ln \frac{q(\mathbf{z}_n)}{p(\mathbf{z}_n|\mathbf{x}_n, \theta)} \right) \geq 0 \quad (258)$$

Next we are going to derive the  $q(\mathbf{z}_n)$  by working out the form of  $p(\mathbf{z}_n|\mathbf{x}_n)$ . But first we should derive the expression of marginal distribution  $p(\mathbf{x}_n)$  and

To derive the Expectation step, then we manipulate the  $p(\mathbf{z}|\mathbf{x})$ , which is the expected probability of one data object  $\mathbf{x}$  being assigned to one cluster (tree in this case), given the observation pattern (laser measurement in this case).

$$p(\mathbf{z}_n|\mathbf{x}_n) = \frac{p(\mathbf{z}_n, \mathbf{x}_n)}{p(\mathbf{x}_n)} \quad \text{by product rule} \quad (259)$$

$$= \frac{\pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \sigma^2 \mathbf{I})}{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \sigma^2 \mathbf{I})} \quad \text{cite (232) and (233)} \quad (260)$$

where the  $\pi_k$ , which is  $p(z_{nk} = 1)$ , can be treated as the coefficient (significance) of  $k$ -th mixture of gaussian component. And the specific form of  $\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \sigma^2 \mathbf{I})$  is as follows,

$$\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \sigma^2 \mathbf{I}) = \frac{1}{(2\pi)^{\frac{D}{2}} \sigma^D} \exp\left(-\frac{(\mathbf{x}_n - \boldsymbol{\mu}_k)^T (\mathbf{x}_n - \boldsymbol{\mu}_k)}{2\sigma^2}\right) \quad (261)$$

Now we figure out the form of  $p(\mathbf{z}_n|\mathbf{x}_n)$  and thus the  $q(\mathbf{z}_n)$  that maximise the lower bound for each data object  $n$  can be computed.

By convention, we utilise the notation  $\gamma_{nk}$  to denote the probability of one unlabel measurement  $\mathbf{x}_n$  being attributed to certain tree  $\mathbf{z}_n$  given the value of measurement  $\mathbf{x}_n$ , that is

$$\gamma_{nk} \stackrel{\text{def}}{=} p(\mathbf{z}_n|\mathbf{x}_n) \quad (262)$$

By (260) to replace the  $p(\mathbf{z}_n|\mathbf{x}_n)$ , we have

$$\gamma_{nk} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \sigma^2 \mathbf{I})}{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \sigma^2 \mathbf{I})} \quad (263)$$

The procedure of computing expectation of responsibilities.

- For start-up., if no given knowledge of approximate location of each tree (no preknowledge of parameters of gaussian components), randomly initialise the parameters of all  $k$  gaussian components presented above.
- E-step: Then for each unlabelled measurement  $\mathbf{x}_n$ , we should figure out its probability of coming from all  $K$  tree components (all possible state of  $\mathbf{z}$ ), that is, the responsibility of each component to explain currently treated measurement  $\mathbf{x}_n$ . In this step, we ought to have  $N \times K$  times of computation for this responsibility.

## 6.4 Derive the M-step for EM

In the solution to 6.2, we have figured out in advance the optimisation objective of the M-step, which is to maximise the following term

$$Q(\theta, \theta^{old}) = \sum_{n=1}^N \sum_{\mathbf{z}_n} \left( p(\mathbf{z}_n | \mathbf{x}_n, \theta^{old}) \cdot \ln p(\mathbf{x}_n, \mathbf{z}_n | \theta) \right) \quad (264)$$

This can be rewritten as the follows based on (232)

$$Q(\theta, \theta^{old}) = \sum_{n=1}^N \sum_{k=1}^K \left( \gamma_{nk} \cdot \ln (\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})) \right) \quad (265)$$

This is valid since we have figured out the expression of  $p(\mathbf{z}_n | \mathbf{x}_n, \theta^{old})$  in 6.3 and that probability term has been computed in E-step, the  $p(\mathbf{z}_n | \mathbf{x}_n, \theta^{old})$  is a constant in M-step. Hence, we just use  $\gamma_{nk}$  to represent  $p(\mathbf{z}_n | \mathbf{x}_n, \theta^{old})$ .

Manipulate the partial derivative of  $Q$  with regard to  $\boldsymbol{\mu}_k$

$$\frac{\partial \sum_{n=1}^N \sum_{k=1}^K \left( \gamma_{nk} \cdot \ln (\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})) \right)}{\partial \boldsymbol{\mu}_k} \quad (266)$$

$$= \sum_{n=1}^N \frac{\partial \sum_{k=1}^K \left( \gamma_{nk} \cdot \ln (\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})) \right)}{\partial \boldsymbol{\mu}_k} \quad (267)$$

$$= \sum_{n=1}^N \frac{\partial \left( \gamma_{nk} \cdot \ln (\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})) \right)}{\partial \boldsymbol{\mu}_k} \quad (268)$$

$$= \sum_{n=1}^N \gamma_{nk} \cdot \frac{\partial \ln (\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I}))}{\partial \boldsymbol{\mu}_k} \quad (269)$$

$$= \sum_{n=1}^N \gamma_{nk} \cdot \frac{1}{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})} \cdot \frac{\partial (\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I}))}{\partial \boldsymbol{\mu}_k} \quad (270)$$

$$= \sum_{n=1}^N \gamma_{nk} \cdot \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})}{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})} \cdot \frac{\partial \left( - \frac{(\mathbf{x}_n - \boldsymbol{\mu}_k)^T (\mathbf{x}_n - \boldsymbol{\mu}_k)}{2\sigma^2} \right)}{\partial \boldsymbol{\mu}_k} \quad (271)$$

$$= \sum_{n=1}^N \gamma_{nk} \cdot \frac{\partial \left( - \frac{(\mathbf{x}_n - \boldsymbol{\mu}_k)^T (\mathbf{x}_n - \boldsymbol{\mu}_k)}{2\sigma^2} \right)}{\partial \boldsymbol{\mu}_k} \quad (272)$$

$$= \sum_{n=1}^N \gamma_{nk} \cdot \frac{(\mathbf{x}_n - \boldsymbol{\mu}_k)}{\sigma^2} \quad (273)$$

Since the  $\sigma$  is independent on the  $n$ , the final result of derivative of likelihood w.r.t  $\boldsymbol{\mu}_k$  can be represented as

$$\frac{\partial Q}{\partial \boldsymbol{\mu}_k} = \frac{1}{\sigma^2} \sum_{n=1}^N \gamma_{nk} \cdot (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (274)$$

Take derivative to the likelihood w.r.t  $\boldsymbol{\mu}_k$  to get the extremum,

$$\frac{\partial Q}{\partial \boldsymbol{\mu}_k} = 0 \quad (275)$$

Solve the equation above to acquire the extremum,

$$\sum_{n=1}^N \gamma_{nk} \cdot \frac{(\mathbf{x}_n - \boldsymbol{\mu}_k)}{\sigma^2} = 0 \quad (276)$$

$$\frac{1}{\sigma^2} \sum_{n=1}^N \gamma_{nk} \cdot (\mathbf{x}_n - \boldsymbol{\mu}_k) = 0 \quad (277)$$

$$\sum_{n=1}^N \gamma_{nk} \cdot (\mathbf{x}_n - \boldsymbol{\mu}_k) = 0 \quad (278)$$

$$\left( \sum_{n=1}^N \gamma_{nk} \mathbf{x}_n - \sum_{n=1}^N \gamma_{nk} \boldsymbol{\mu}_k \right) = 0 \quad (279)$$

$$\sum_{n=1}^N \gamma_{nk} \boldsymbol{\mu}_k = \sum_{n=1}^N \gamma_{nk} \mathbf{x}_n \quad (280)$$

$$\boldsymbol{\mu}_k \sum_{n=1}^N \gamma_{nk} = \sum_{n=1}^N \gamma_{nk} \mathbf{x}_n \quad (281)$$

$$\boldsymbol{\mu}_k = \frac{\sum_{n=1}^N \gamma_{nk} \mathbf{x}_n}{\sum_{n=1}^N \gamma_{nk}} \quad (282)$$

Note that from (277) to (278), we assume that  $\sigma^2 \neq 0$ . The derivation from (280) to (281) is valid since  $\boldsymbol{\mu}_k$  is independent on summation variable  $n$ .

Then we define  $N_k$  to be effective number of measurements  $\mathbf{x}_n$  being explained by  $k$ -th tree component. That is

$$N_k = \sum_{n=1}^N \gamma_{nk} \quad (283)$$

Therefore, we have

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} \mathbf{x}_n \quad (284)$$

Similarly, take derivative to  $Q$  w.r.t  $\sigma$ ,

$$\frac{\partial \sum_{n=1}^N \sum_{k=1}^K \left( \gamma_{nk} \cdot \ln (\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})) \right)}{\partial \sigma} \quad (285)$$

$$= \sum_{n=1}^N \frac{\partial \sum_{k=1}^K \left( \gamma_{nk} \cdot \ln (\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})) \right)}{\partial \sigma} \quad (286)$$

$$= \sum_{n=1}^N \sum_{k=1}^K \frac{\partial \left( \gamma_{nk} \cdot \ln (\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})) \right)}{\partial \sigma} \quad (287)$$

$$= \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \cdot \frac{\partial \left( \ln (\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})) \right)}{\partial \sigma} \quad (288)$$

$$= \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \cdot \frac{1}{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})} \cdot \frac{\partial \left( \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I}) \right)}{\partial \sigma} \quad (289)$$

$$= \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \cdot \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})}{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})} \cdot \left( -\frac{2}{\sigma} + \frac{(\mathbf{x}_n - \boldsymbol{\mu}_k)^T (\mathbf{x}_n - \boldsymbol{\mu}_k)}{\sigma^3} \right) \quad (290)$$

$$= \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \cdot \left( -\frac{2}{\sigma} + \frac{(\mathbf{x}_n - \boldsymbol{\mu}_k)^T (\mathbf{x}_n - \boldsymbol{\mu}_k)}{\sigma^3} \right) \quad (291)$$

Set the derivative to zero, we will obtain the extremum for  $\sigma$

$$\sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \cdot \left( -\frac{2}{\sigma} + \frac{(\mathbf{x}_n - \boldsymbol{\mu}_k)^T (\mathbf{x}_n - \boldsymbol{\mu}_k)}{\sigma^3} \right) = 0 \quad (292)$$

Since  $\sigma$  is independent on  $n$  and  $k$ , we have

$$\sigma^3 \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \cdot \left( -2\sigma^2 + (\mathbf{x}_n - \boldsymbol{\mu}_k)^T (\mathbf{x}_n - \boldsymbol{\mu}_k) \right) = 0 \quad (293)$$

$$\sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \cdot \left( -2\sigma^2 + (\mathbf{x}_n - \boldsymbol{\mu}_k)^T (\mathbf{x}_n - \boldsymbol{\mu}_k) \right) = 0 \quad (294)$$

Further manipulation as follows,

$$\sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} 2\sigma^2 = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (295)$$

Since both  $\sigma$  is independent on  $n$  and  $k$ , we have

$$2\sigma^2 \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (296)$$

It is evident that

$$\sigma^2 = \frac{\sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T (\mathbf{x}_n - \boldsymbol{\mu}_k)}{2 \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk}} \quad (297)$$

Based on the definition of  $\gamma_{nk} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})}{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})}$  in (263), we have

$$\sum_{k=1}^K \gamma_{nk} = \sum_{k=1}^K \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})}{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})} \quad (298)$$

$$= \frac{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})}{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})} \quad (299)$$

$$= 1 \quad (300)$$

Note that the derivation from (298) to (299) is valid because the whole term denominator  $\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})$  is independent on  $k$ .

Hence, the extremum of  $\sigma^2$  can be simplified as follows.

$$\sigma^2 = \frac{\sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T (\mathbf{x}_n - \boldsymbol{\mu}_k)}{2 \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk}} \quad (301)$$

$$= \frac{\sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T (\mathbf{x}_n - \boldsymbol{\mu}_k)}{2 \sum_{n=1}^N 1} \quad (302)$$

$$= \frac{\sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T (\mathbf{x}_n - \boldsymbol{\mu}_k)}{2N} \quad (303)$$

$$= \frac{1}{2N} \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (304)$$

$$= \frac{1}{2N} \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (305)$$

Finally, let's derive the extremum of  $Q$  with regard to  $\pi_k$ , but since  $\pi_k$  is constrained by condition  $\sum_{k=1}^K \pi_k = 1$ , the derivation would be a bit more complicated.

We construct Lagrangian  $\mathcal{L}'(\boldsymbol{\pi}, \lambda)$  with lagrangian multiplier  $\lambda$

$$\mathcal{L}'(\boldsymbol{\pi}, \lambda) = \sum_{n=1}^N \sum_{k=1}^K \left( \gamma_{nk} \cdot \ln(\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})) \right) + \lambda \left( 1 - \sum_{k=1}^K \pi_k \right) \quad (306)$$

Then take derivative with regard to  $\pi_k$ , we have

$$\frac{\partial \mathcal{L}'(\boldsymbol{\pi}, \lambda)}{\partial \pi_k} = \sum_{n=1}^N \left( \gamma_{nk} \cdot \frac{\mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})}{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})} \right) - \lambda \quad (307)$$

$$= \sum_{n=1}^N \left( \gamma_{nk} \cdot \frac{1}{\pi_k} \right) - \lambda \quad (308)$$

$$= \left( \frac{1}{\pi_k} \sum_{n=1}^N \gamma_{nk} \right) - \lambda \quad (309)$$

$$= \frac{N_k}{\pi_k} - \lambda \quad (310)$$

Set the derivative towards langrangian to zero, we have

$$\frac{N_k}{\pi_k} - \lambda = 0 \quad (311)$$

$$\pi_k = \frac{N_k}{\lambda} \quad (312)$$

Take the expression  $\pi_k$  back to  $\mathcal{L}'(\boldsymbol{\pi}, \lambda)$ , we have

$$\mathcal{L}(\lambda) = \sum_{n=1}^N \sum_{k=1}^K \left( \gamma_{nk} \cdot \ln(\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})) \right) + \lambda \left( 1 - \sum_{k=1}^K \pi_k \right) \quad (313)$$

$$= \sum_{n=1}^N \sum_{k=1}^K \left( \gamma_{nk} \cdot \ln \left( \frac{N_k}{\lambda} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I}) \right) \right) + \lambda - \lambda \cdot \sum_{k=1}^K \frac{N_k}{\lambda} \quad (314)$$

$$= \sum_{n=1}^N \sum_{k=1}^K \left( \gamma_{nk} \cdot \ln(N_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})) - \gamma_{nk} \cdot \ln \lambda \right) + \lambda - \sum_{k=1}^K N_k \quad (315)$$

$$= \sum_{n=1}^N \sum_{k=1}^K \left( \gamma_{nk} \cdot \ln(N_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \sigma^2 \mathbf{I})) - \gamma_{nk} \cdot \ln \lambda \right) + \lambda - N \quad (316)$$

Take the derivative to  $\mathcal{L}(\lambda)$  with regard to  $\lambda$ ,

$$\frac{\partial \mathcal{L}(\lambda)}{\partial \lambda} = \sum_{n=1}^N \sum_{k=1}^K \left( -\frac{\gamma_{nk}}{\lambda} \right) + 1 \quad \text{eliminate irrelevant terms} \quad (317)$$

$$= -\frac{1}{\lambda} \sum_{n=1}^N \sum_{k=1}^K (\gamma_{nk}) + 1 \quad \lambda \text{ independent on } n \text{ and } k \quad (318)$$

$$= -\frac{1}{\lambda} \left( \sum_{n=1}^N 1 \right) + 1 \quad \text{cite (300)} \quad (319)$$

$$= -\frac{N}{\lambda} + 1 \quad (320)$$

Set derivative to zero to get minima of  $\mathcal{L}(\lambda)$

$$-\frac{N}{\lambda} + 1 = 0 \quad (321)$$

$$\lambda = N \quad (322)$$

Take the final result back to what we have previously derived in (312).

$$\pi_k = \frac{N_k}{N} \quad (323)$$

which is the update formula for  $\pi_k$  in M-step.

In summary, we acquire all the M-step update as follows.

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} \mathbf{x}_n \quad (324)$$

$$\sigma^2 = \frac{1}{2N} \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (325)$$

$$\pi_k = \frac{N_k}{N} \quad (326)$$

where  $N_k$  is defined to be effective number of measurements  $\mathbf{x}_n$  being explained by  $k$ -th tree component.



## 7 (Semi-/Un-)Supervised Learning and EM

### 7.1 Supervised Learning

We apply the supervised learning on Iris Data and use 10-fold cross validation to figure out the error rate by using the classification approach above.

**Data Presentation.** Specific values of mean vector and covariance matrix of each gaussian distribution in each cross validation training set are shown in the appendix [Cross Validation for Supervised Learning](#). But here we just present the overall error rate in the required 10-fold cross validation and specific error rates of each validation set.

**Average ErrorRate:** 2.6667%

**Specific Error in each set:**

Set 1	Set 2	Set 3	Set 4	Set 5
0.000%	6.667%	0.000%	0.000%	0.000%
Set 6	Set 7	Set 8	Set 9	Set 10
6.667%	0.000%	6.667%	6.667%	0.000%

Note that my implementation for obtaining cross validation set in the above 10-fold cross validation is called "repeatedly random subsampling validation", which is widely used. (Refer to [http://en.wikipedia.org/wiki/Cross-validation\\_\(statistics\)](http://en.wikipedia.org/wiki/Cross-validation_(statistics)))

To run my program, please type the following in the command.

```
python ./Supervised.py
```

### 7.2 Unsupervised Learning

We apply Expectation-Maximisation technique to implement such a unsupervised learning task. For specific results of this task, please see the appendix [Expectation Maximisation for Unsupervised Learning](#).

item	Case 1	Case 2	Case 3	Case 4	Case 5
Converged log likelihood	-185.276241	-196.953446	-180.185510	-180.185529	-180.185537

As illustrated in the above table, which indicates the log likelihood while algorithm converges (likelihood has little change in adjacent iteration, threshold all set to  $10e^{-6}$ .) It is obvious that each time we run EM algorithm to get three 4D gaussian distributions, the results are different from each others. There are many possible reasons as follows.

- MultiModal Distribution
- Sensitivity to Initilisation

Note that here we count a pair of E-step and M-step as one iteration.

To run my program, please type the following in the command.

```
python ./Unsupervised.py
```

Note that all the parameters have default values, for details of command line argument, please see help documentation by typing

```
python ./Unsupervised.py -h
```

### 7.3 Semi-supervised Learning

**Approach Description:** Use 5% of labeled data to construct initial gaussian components, that is, to fit the parameters of initial gaussian components by using MLE, and then utilise 95% of unlabeled data to adjust (perhaps with less weight comparing to the labeled data) the parameters of components by Expectation-Maximisation (EM) algorithm.

One strength of this approach lies in unnecessary of subjectively determining number of mixture components  $k$  for unsupervised learning. It is quite tough for human to determine one "good"  $k$  when human has no preknowledge of specific problem. Even though sometime large  $k$  may lead to small error or large likelihood, which is desired mathematically, the result may not be interpretable in practice, which means the  $k$  make no sense in reality. If we apply the model above, we just need to label the data that has been observed as objects of different classes, and thus decrease the risk of bad  $k$  (number of components).

On top of that, since only 5% of data objects are provided for supervised learning in the first stage, the fitted parameters may be good qualitatively (it captures the approximate location of one class of objects) but not precise quantitatively (fitted parameters in a measure is far from the real one because the 5% data happens to be located in the brink of the whole class's distribution). But if we after that apply the EM algorithm to adjust the acquired parameters with a relatively enormous size of dataset, the concerns about precision of parameters fitted by small size of dataset (5% of the whole) can be removed.

Another significant reason of doing so comes from the weakness of random initialisation for startup in EM algorithm. If bad random initialisation is encountered, the parameters of model such as  $\mu$ ,  $\Sigma$  and  $\pi$  in mixture of gaussians would be likely to converge to bad local optima rather than global optima or relatively satisfying optima. Therefore, training task has to be run for several times to derive a number of results, then compare the acquired results and pick up the one with smallest error or largest likelihood, which obviously is unexpectedly time-consuming. On the contrast, if a set of labeled data is introduced, even though in a limited quantity, we can make use of those data to fit initial gaussian model (perhaps maximum likelihood estimation) to avoid running one machine learning task repeatedly (Sometimes tremendous amount data is involved in a task, we may not wish to run them for several times because they would be either impossible or too costly.).

## A Data Presentation

### A.1 Cross Validation for Supervised Learning

Cross Validation Set 1:

Error of this Testing Set: 0.00%

**For class Iris-setosa:**

$$\boldsymbol{\mu}_0 = (5.01777778 \quad 3.44666667 \quad 1.46444444 \quad 0.25111111) \quad (327)$$

$$\boldsymbol{\Sigma}_0 = \begin{pmatrix} 0.12104040 & 0.10028788 & 0.01951010 & 0.01088889 \\ 0.10028788 & 0.15163636 & 0.01351515 & 0.00937879 \\ 0.01951010 & 0.01351515 & 0.03143434 & 0.00663131 \\ 0.01088889 & 0.00937879 & 0.00663131 & 0.01210101 \end{pmatrix} \quad (328)$$

**For class Iris-versicolor:**

$$\boldsymbol{\mu}_1 = (5.92444444 \quad 2.76666667 \quad 4.26000000 \quad 1.33333333) \quad (329)$$

$$\boldsymbol{\Sigma}_1 = \begin{pmatrix} 0.27416162 & 0.08174242 & 0.18645455 & 0.05553030 \\ 0.08174242 & 0.09909091 & 0.08000000 & 0.03977273 \\ 0.18645455 & 0.08000000 & 0.22927273 & 0.07409091 \\ 0.05553030 & 0.03977273 & 0.07409091 & 0.03909091 \end{pmatrix} \quad (330)$$

**For class Iris-virginica:**

$$\boldsymbol{\mu}_2 = (6.63777778 \quad 2.96444444 \quad 5.58666667 \quad 2.01555556) \quad (331)$$

$$\boldsymbol{\Sigma}_2 = \begin{pmatrix} 0.41967677 & 0.10637374 & 0.31574242 & 0.06098990 \\ 0.10637374 & 0.10961616 & 0.08133333 & 0.05011111 \\ 0.31574242 & 0.08133333 & 0.31845455 & 0.05748485 \\ 0.06098990 & 0.05011111 & 0.05748485 & 0.07497980 \end{pmatrix} \quad (332)$$

Cross Validation Set 2:

Error of this Testing Set: 6.67%

**For class Iris-setosa:**

$$\boldsymbol{\mu}_0 = (5.01333333 \quad 3.43555556 \quad 1.47555556 \quad 0.24666667) \quad (333)$$

$$\boldsymbol{\Sigma}_0 = \begin{pmatrix} 0.12345455 & 0.10451515 & 0.01737879 & 0.01163636 \\ 0.10451515 & 0.15370707 & 0.01225253 & 0.01080303 \\ 0.01737879 & 0.01225253 & 0.03143434 & 0.00662121 \\ 0.01163636 & 0.01080303 & 0.00662121 & 0.01209091 \end{pmatrix} \quad (334)$$

**For class Iris-versicolor:**

$$\boldsymbol{\mu}_1 = (5.97777778 \quad 2.76888889 \quad 4.27111111 \quad 1.32222222) \quad (335)$$

$$\boldsymbol{\Sigma}_1 = \begin{pmatrix} 0.26494949 & 0.08724747 & 0.18161616 & 0.05573232 \\ 0.08724747 & 0.10037374 & 0.07635354 & 0.03707071 \\ 0.18161616 & 0.07635354 & 0.21528283 & 0.06679293 \\ 0.05573232 & 0.03707071 & 0.06679293 & 0.03494949 \end{pmatrix} \quad (336)$$

**For class Iris-virginica:**

$$\boldsymbol{\mu}_2 = (6.60888889 \quad 2.96444444 \quad 5.55333333 \quad 2.01111111) \quad (337)$$

$$\boldsymbol{\Sigma}_2 = \begin{pmatrix} 0.41946465 & 0.10305051 & 0.32156061 & 0.05694444 \\ 0.10305051 & 0.11143434 & 0.07557576 & 0.04926768 \\ 0.32156061 & 0.07557576 & 0.31618182 & 0.04734848 \\ 0.05694444 & 0.04926768 & 0.04734848 & 0.07782828 \end{pmatrix} \quad (338)$$

Cross Validation Set 3:

Error of this Testing Set: 0.00%

**For class Iris-setosa:**

$$\mu_0 = (5.01041667 \quad 3.41875000 \quad 1.46250000 \quad 0.25000000) \quad (339)$$

$$\Sigma_0 = \begin{pmatrix} 0.12520833 & 0.09873670 & 0.01635638 & 0.01095745 \\ 0.09873670 & 0.13900266 & 0.01135638 & 0.01159574 \\ 0.01635638 & 0.01135638 & 0.03132979 & 0.00638298 \\ 0.01095745 & 0.01159574 & 0.00638298 & 0.01106383 \end{pmatrix} \quad (340)$$

**For class Iris-versicolor:**

$$\mu_1 = (5.91428571 \quad 2.75952381 \quad 4.25238095 \quad 1.32857143) \quad (341)$$

$$\Sigma_1 = \begin{pmatrix} 0.29344948 & 0.08937282 & 0.19679443 & 0.05519164 \\ 0.08937282 & 0.10783391 & 0.08851336 & 0.04191638 \\ 0.19679443 & 0.08851336 & 0.23718931 & 0.07700348 \\ 0.05519164 & 0.04191638 & 0.07700348 & 0.03721254 \end{pmatrix} \quad (342)$$

**For class Iris-virginica:**

$$\mu_2 = (6.60666667 \quad 2.98000000 \quad 5.57333333 \quad 2.03777778) \quad (343)$$

$$\Sigma_2 = \begin{pmatrix} 0.39972727 & 0.09786364 & 0.29290909 & 0.04542424 \\ 0.09786364 & 0.11118182 & 0.07445455 & 0.05100000 \\ 0.29290909 & 0.07445455 & 0.29018182 & 0.04512121 \\ 0.04542424 & 0.05100000 & 0.04512121 & 0.08104040 \end{pmatrix} \quad (344)$$

Cross Validation Set 4:

Error of this Testing Set: 0.00%

**For class Iris-setosa:**

$$\mu_0 = (5.04047619 \quad 3.47619048 \quad 1.45952381 \quad 0.24761905) \quad (345)$$

$$\Sigma_0 = \begin{pmatrix} 0.12929733 & 0.09586527 & 0.01704413 & 0.01168409 \\ 0.09586527 & 0.12771196 & 0.00828107 & 0.01042973 \\ 0.01704413 & 0.00828107 & 0.02978513 & 0.00782811 \\ 0.01168409 & 0.01042973 & 0.00782811 & 0.01182346 \end{pmatrix} \quad (346)$$

**For class Iris-versicolor:**

$$\mu_1 = (5.92978723 \quad 2.76170213 \quad 4.25744681 \quad 1.32978723) \quad (347)$$

$$\Sigma_1 = \begin{pmatrix} 0.27909343 & 0.08551341 & 0.18586031 & 0.05604995 \\ 0.08551341 & 0.09806660 & 0.07985661 & 0.03964385 \\ 0.18586031 & 0.07985661 & 0.21423682 & 0.07107771 \\ 0.05604995 & 0.03964385 & 0.07107771 & 0.03735430 \end{pmatrix} \quad (348)$$

**For class Iris-virginica:**

$$\mu_2 = (6.59565217 \quad 2.95000000 \quad 5.56086957 \quad 2.01739130) \quad (349)$$

$$\Sigma_2 = \begin{pmatrix} 0.37953623 & 0.07911111 & 0.28915942 & 0.06074396 \\ 0.07911111 & 0.09277778 & 0.06244444 & 0.05111111 \\ 0.28915942 & 0.06244444 & 0.30110145 & 0.05514010 \\ 0.06074396 & 0.05111111 & 0.05514010 & 0.07613527 \end{pmatrix} \quad (350)$$

Cross Validation Set 5:

Error of this Testing Set: 0.00%

**For class Iris-setosa:**

$$\boldsymbol{\mu}_0 = (4.99565217 \quad 3.43043478 \quad 1.45652174 \quad 0.25000000) \quad (351)$$

$$\boldsymbol{\Sigma}_0 = \begin{pmatrix} 0.11642512 & 0.09502415 & 0.01380676 & 0.00977778 \\ 0.09502415 & 0.14527536 & 0.01024155 & 0.00822222 \\ 0.01380676 & 0.01024155 & 0.03140097 & 0.00622222 \\ 0.00977778 & 0.00822222 & 0.00622222 & 0.01144444 \end{pmatrix} \quad (352)$$

**For class Iris-versicolor:**

$$\boldsymbol{\mu}_1 = (5.91739130 \quad 2.76956522 \quad 4.24565217 \quad 1.31521739) \quad (353)$$

$$\boldsymbol{\Sigma}_1 = \begin{pmatrix} 0.27257971 & 0.09298551 & 0.18429952 & 0.05439614 \\ 0.09298551 & 0.09994203 & 0.08808696 & 0.04180676 \\ 0.18429952 & 0.08808696 & 0.22698068 & 0.07351208 \\ 0.05439614 & 0.04180676 & 0.07351208 & 0.03865217 \end{pmatrix} \quad (354)$$

**For class Iris-virginica:**

$$\boldsymbol{\mu}_2 = (6.61162791 \quad 2.96744186 \quad 5.56511628 \quad 2.01627907) \quad (355)$$

$$\boldsymbol{\Sigma}_2 = \begin{pmatrix} 0.41629014 & 0.10205426 & 0.32232004 & 0.05933001 \\ 0.10205426 & 0.11320044 & 0.07598007 & 0.04911406 \\ 0.32232004 & 0.07598007 & 0.31375415 & 0.05677187 \\ 0.05933001 & 0.04911406 & 0.05677187 & 0.07901440 \end{pmatrix} \quad (356)$$

Cross Validation Set 6:

Error of this Testing Set: 6.67%

**For class Iris-setosa:**

$$\boldsymbol{\mu}_0 = (5.03043478 \quad 3.42826087 \quad 1.46956522 \quad 0.25217391) \quad (357)$$

$$\boldsymbol{\Sigma}_0 = \begin{pmatrix} 0.12616425 & 0.10689855 & 0.01428019 & 0.00971014 \\ 0.10689855 & 0.15273913 & 0.01510145 & 0.01049275 \\ 0.01428019 & 0.01510145 & 0.02749758 & 0.00606763 \\ 0.00971014 & 0.01049275 & 0.00606763 & 0.01143961 \end{pmatrix} \quad (358)$$

**For class Iris-versicolor:**

$$\boldsymbol{\mu}_1 = (5.94000000 \quad 2.76666667 \quad 4.25777778 \quad 1.32222222) \quad (359)$$

$$\boldsymbol{\Sigma}_1 = \begin{pmatrix} 0.26972727 & 0.09227273 & 0.19331818 & 0.06113636 \\ 0.09227273 & 0.10909091 & 0.09106061 & 0.04598485 \\ 0.19331818 & 0.09106061 & 0.23840404 & 0.08073232 \\ 0.06113636 & 0.04598485 & 0.08073232 & 0.04267677 \end{pmatrix} \quad (360)$$

**For class Iris-virginica:**

$$\boldsymbol{\mu}_2 = (6.57045455 \quad 2.95000000 \quad 5.53863636 \quad 2.05227273) \quad (361)$$

$$\boldsymbol{\Sigma}_2 = \begin{pmatrix} 0.33561839 & 0.05872093 & 0.26767970 & 0.05251057 \\ 0.05872093 & 0.09279070 & 0.04965116 & 0.05337209 \\ 0.26767970 & 0.04965116 & 0.29777484 & 0.05281712 \\ 0.05251057 & 0.05337209 & 0.05281712 & 0.07604123 \end{pmatrix} \quad (362)$$

Cross Validation Set 7:

Error of this Testing Set: 0.00%

**For class Iris-setosa:**

$$\boldsymbol{\mu}_0 = (5.03777778 \quad 3.45111111 \quad 1.46888889 \quad 0.25111111) \quad (363)$$

$$\boldsymbol{\Sigma}_0 = \begin{pmatrix} 0.12194949 & 0.09393434 & 0.01438384 & 0.01075253 \\ 0.09393434 & 0.12801010 & 0.00844444 & 0.01210101 \\ 0.01438384 & 0.00844444 & 0.03173737 & 0.00662626 \\ 0.01075253 & 0.01210101 & 0.00662626 & 0.01164646 \end{pmatrix} \quad (364)$$

**For class Iris-versicolor:**

$$\boldsymbol{\mu}_1 = (5.94666667 \quad 2.77777778 \quad 4.26888889 \quad 1.32666667) \quad (365)$$

$$\boldsymbol{\Sigma}_1 = \begin{pmatrix} 0.26072727 & 0.07310606 & 0.17921212 & 0.05077273 \\ 0.07310606 & 0.09176768 & 0.07611111 & 0.03833333 \\ 0.17921212 & 0.07611111 & 0.22810101 & 0.07334848 \\ 0.05077273 & 0.03833333 & 0.07334848 & 0.03972727 \end{pmatrix} \quad (366)$$

**For class Iris-virginica:**

$$\boldsymbol{\mu}_2 = (6.58666667 \quad 2.95777778 \quad 5.55777778 \quad 2.02888889) \quad (367)$$

$$\boldsymbol{\Sigma}_2 = \begin{pmatrix} 0.40709091 & 0.08328788 & 0.29192424 & 0.04471212 \\ 0.08328788 & 0.09931313 & 0.05908586 & 0.05033838 \\ 0.29192424 & 0.05908586 & 0.28658586 & 0.04261111 \\ 0.04471212 & 0.05033838 & 0.04261111 & 0.08073737 \end{pmatrix} \quad (368)$$

Cross Validation Set 8:

Error of this Testing Set: 6.67%

**For class Iris-setosa:**

$$\boldsymbol{\mu}_0 = (4.99782609 \quad 3.42173913 \quad 1.46086957 \quad 0.24565217) \quad (369)$$

$$\boldsymbol{\Sigma}_0 = \begin{pmatrix} 0.13266184 & 0.10471498 & 0.01769082 & 0.01121256 \\ 0.10471498 & 0.15062802 & 0.01264734 & 0.00920773 \\ 0.01769082 & 0.01264734 & 0.03265700 & 0.00649275 \\ 0.01121256 & 0.00920773 & 0.00649275 & 0.01142512 \end{pmatrix} \quad (370)$$

**For class Iris-versicolor:**

$$\boldsymbol{\mu}_1 = (5.94090909 \quad 2.75454545 \quad 4.23409091 \quad 1.31136364) \quad (371)$$

$$\boldsymbol{\Sigma}_1 = \begin{pmatrix} 0.28154334 & 0.08934461 & 0.20089852 & 0.05905920 \\ 0.08934461 & 0.10393235 & 0.08484144 & 0.03959831 \\ 0.20089852 & 0.08484144 & 0.23718288 & 0.07448732 \\ 0.05905920 & 0.03959831 & 0.07448732 & 0.03730973 \end{pmatrix} \quad (372)$$

**For class Iris-virginica:**

$$\boldsymbol{\mu}_2 = (6.58000000 \quad 2.96000000 \quad 5.57111111 \quad 2.00888889) \quad (373)$$

$$\boldsymbol{\Sigma}_2 = \begin{pmatrix} 0.43254545 & 0.10100000 & 0.33395455 & 0.04563636 \\ 0.10100000 & 0.11336364 & 0.08086364 & 0.04945455 \\ 0.33395455 & 0.08086364 & 0.32846465 & 0.05412626 \\ 0.04563636 & 0.04945455 & 0.05412626 & 0.07582828 \end{pmatrix} \quad (374)$$

Cross Validation Set 9:

Error of this Testing Set: 6.67%

**For class Iris-setosa:**

$$\boldsymbol{\mu}_0 = (5.00444444 \quad 3.43111111 \quad 1.44666667 \quad 0.23777778) \quad (375)$$

$$\boldsymbol{\Sigma}_0 = \begin{pmatrix} 0.13816162 & 0.10963131 & 0.01728788 & 0.01119192 \\ 0.10963131 & 0.15219192 & 0.01033333 & 0.00788889 \\ 0.01728788 & 0.01033333 & 0.02800000 & 0.00410606 \\ 0.01119192 & 0.00788889 & 0.00410606 & 0.00876768 \end{pmatrix} \quad (376)$$

**For class Iris-versicolor:**

$$\boldsymbol{\mu}_1 = (5.94772727 \quad 2.78409091 \quad 4.26363636 \quad 1.33636364) \quad (377)$$

$$\boldsymbol{\Sigma}_1 = \begin{pmatrix} 0.25371564 & 0.08403277 & 0.17014799 & 0.05078224 \\ 0.08403277 & 0.10695032 & 0.08429175 & 0.04268499 \\ 0.17014799 & 0.08429175 & 0.19306554 & 0.06437632 \\ 0.05078224 & 0.04268499 & 0.06437632 & 0.03725159 \end{pmatrix} \quad (378)$$

**For class Iris-virginica:**

$$\boldsymbol{\mu}_2 = (6.56956522 \quad 2.97826087 \quad 5.54347826 \quad 2.03043478) \quad (379)$$

$$\boldsymbol{\Sigma}_2 = \begin{pmatrix} 0.42216425 & 0.10376812 & 0.32135266 & 0.05516908 \\ 0.10376812 & 0.11151691 & 0.08052174 & 0.05156522 \\ 0.32135266 & 0.08052174 & 0.32428986 & 0.05464734 \\ 0.05516908 & 0.05156522 & 0.05464734 & 0.08038647 \end{pmatrix} \quad (380)$$

Cross Validation Set 10:

Error of this Testing Set: 0.00%

**For class Iris-setosa:**

$$\boldsymbol{\mu}_0 = (4.98888889 \quad 3.39777778 \quad 1.46000000 \quad 0.24222222) \quad (381)$$

$$\boldsymbol{\Sigma}_0 = \begin{pmatrix} 0.12146465 & 0.09747475 & 0.01568182 & 0.00934343 \\ 0.09747475 & 0.13976768 & 0.01172727 & 0.01009596 \\ 0.01568182 & 0.01172727 & 0.03200000 & 0.00559091 \\ 0.00934343 & 0.01009596 & 0.00559091 & 0.01067677 \end{pmatrix} \quad (382)$$

**For class Iris-versicolor:**

$$\boldsymbol{\mu}_1 = (5.89534884 \quad 2.76279070 \quad 4.22790698 \quad 1.30465116) \quad (383)$$

$$\boldsymbol{\Sigma}_1 = \begin{pmatrix} 0.28188261 & 0.09410853 & 0.19156146 & 0.05430786 \\ 0.09410853 & 0.10239203 & 0.09225360 & 0.04589147 \\ 0.19156146 & 0.09225360 & 0.23253599 & 0.07415282 \\ 0.05430786 & 0.04589147 & 0.07415282 & 0.03902547 \end{pmatrix} \quad (384)$$

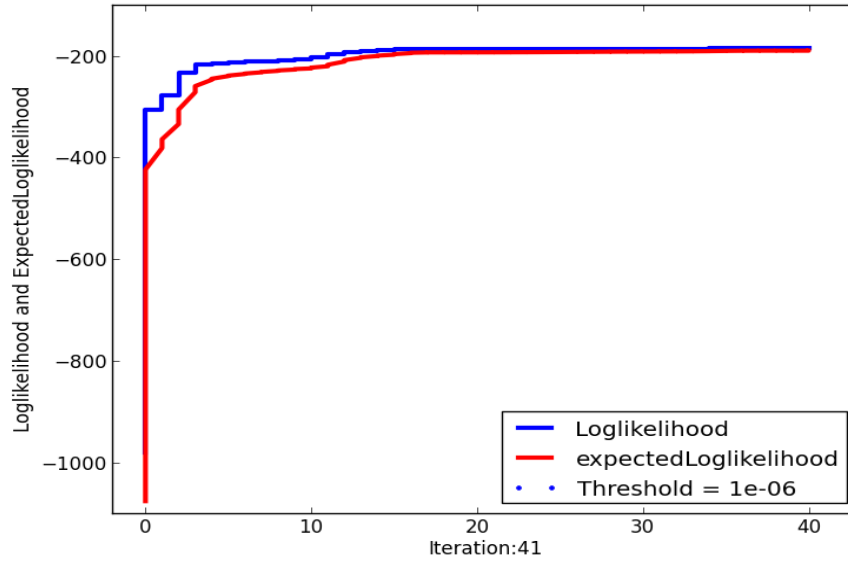
**For class Iris-virginica:**

$$\boldsymbol{\mu}_2 = (6.59574468 \quad 2.96595745 \quad 5.57021277 \quad 2.03191489) \quad (385)$$

$$\boldsymbol{\Sigma}_2 = \begin{pmatrix} 0.42954672 & 0.10093895 & 0.32117484 & 0.05122572 \\ 0.10093895 & 0.10925069 & 0.07874653 & 0.05154487 \\ 0.32117484 & 0.07874653 & 0.31691952 & 0.05118871 \\ 0.05122572 & 0.05154487 & 0.05118871 & 0.07917669 \end{pmatrix} \quad (386)$$

## A.2 Expectation Maximisation for Unsupervised Learning

### Specific data in first random initialisation



Initilisation of parameters:

$$\mu_0 = (6.94825182 \quad 3.43028943 \quad 5.13165042 \quad 1.73043353) \quad (387)$$

$$\mu_1 = (6.03430983 \quad 2.44263481 \quad 4.19416609 \quad 1.34434648) \quad (388)$$

$$\mu_2 = (7.41659032 \quad 2.74307646 \quad 3.54490011 \quad 0.37515563) \quad (389)$$

$$\Sigma_0 = \Sigma_1 = \Sigma_2 = I \quad (390)$$

$$\pi = (0.63777939 \quad 0.17825037 \quad 0.18397024) \quad (391)$$

Iteration Quantity: 41

Parameters of resulted Guassians:

$$\mu_0 = (6.30229683 \quad 2.90022864 \quad 5.06049687 \quad 1.76245416) \quad (392)$$

$$\Sigma_0 = \begin{pmatrix} 0.42223868 & 0.09703482 & 0.43821003 & 0.15608046 \\ 0.09703482 & 0.10551541 & 0.11569370 & 0.07015344 \\ 0.43821003 & 0.11569370 & 0.62024483 & 0.24770913 \\ 0.15608046 & 0.07015344 & 0.24770913 & 0.16391035 \end{pmatrix} \quad (393)$$

$$\mu_1 = (6.06077537 \quad 2.73103951 \quad 4.13451491 \quad 1.24428925) \quad (394)$$

$$\Sigma_1 = \begin{pmatrix} 0.44986364 & 0.20624275 & 0.31552694 & 0.10820035 \\ 0.20624275 & 0.10624364 & 0.13904845 & 0.05150205 \\ 0.31552694 & 0.13904845 & 0.23264190 & 0.07652083 \\ 0.10820035 & 0.05150205 & 0.07652083 & 0.02840178 \end{pmatrix} \quad (395)$$

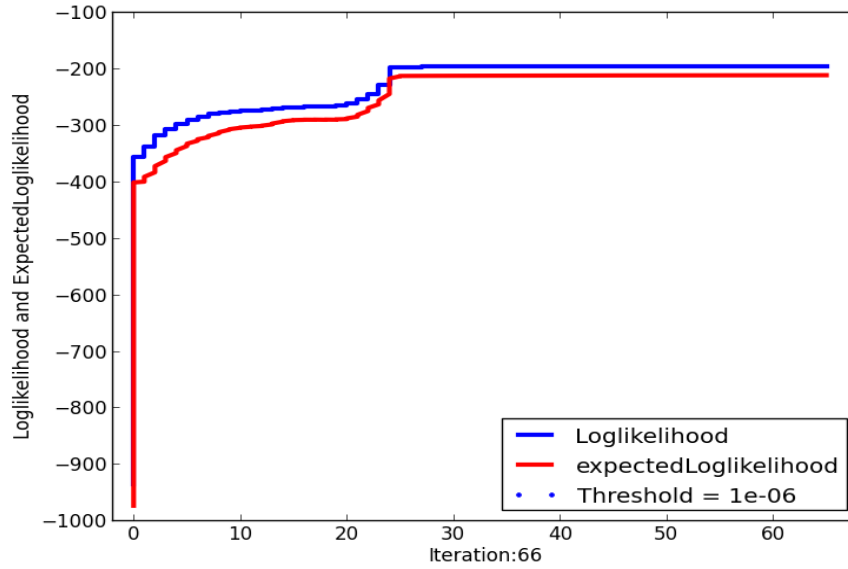
$$\mu_2 = (5.00600023 \quad 3.42800051 \quad 1.46200007 \quad 0.24599998) \quad (396)$$

$$\Sigma_2 = \begin{pmatrix} 0.12176394 & 0.09723179 & 0.01602797 & 0.01012402 \\ 0.09723179 & 0.14081549 & 0.01146392 & 0.00911203 \\ 0.01602797 & 0.01146392 & 0.02955600 & 0.00594801 \\ 0.01012402 & 0.00911203 & 0.00594801 & 0.01088400 \end{pmatrix} \quad (397)$$

$$\pi = (0.55543518 \quad 0.11123164 \quad 0.33333318) \quad (398)$$



## Specific data in second random initialisation



Initilisation of parameters:

$$\mu_0 = (5.47252304 \quad 3.04349496 \quad 3.20982127 \quad 1.74800758) \quad (399)$$

$$\mu_1 = (5.37077336 \quad 3.19556116 \quad 5.46395648 \quad -0.03621684) \quad (400)$$

$$\mu_2 = (6.08856967 \quad 3.14829166 \quad 5.74406265 \quad 1.40878059) \quad (401)$$

$$\Sigma_0 = \Sigma_1 = \Sigma_2 = I \quad (402)$$

$$\pi = (0.92942981 \quad 0.06546349 \quad 0.0051067) \quad (403)$$

Iteration Quantity: 66

Parameters of resulted Guassians:

$$\mu_0 = (5.00600011 \quad 3.42800025 \quad 1.46200004 \quad 0.24599999) \quad (404)$$

$$\Sigma_0 = \begin{pmatrix} 0.12176397 & 0.09723189 & 0.01602799 & 0.01012401 \\ 0.09723189 & 0.14081575 & 0.01146396 & 0.00911202 \\ 0.01602799 & 0.01146396 & 0.02955600 & 0.00594800 \\ 0.01012401 & 0.00911202 & 0.00594800 & 0.01088400 \end{pmatrix} \quad (405)$$

$$\mu_1 = (6.11764829 \quad 2.84264967 \quad 4.67930452 \quad 1.59705880) \quad (406)$$

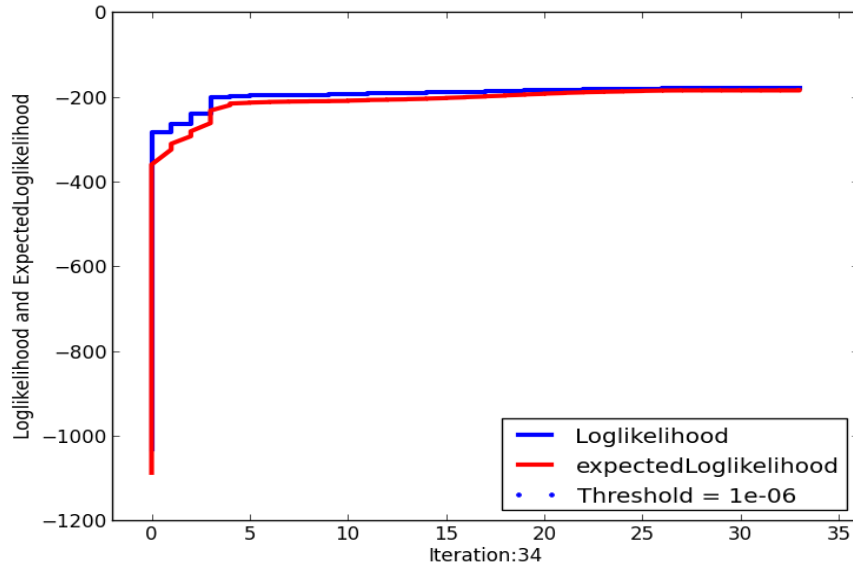
$$\Sigma_1 = \begin{pmatrix} 0.31117193 & 0.10177389 & 0.28584863 & 0.14245983 \\ 0.10177389 & 0.09664535 & 0.12171721 & 0.07841352 \\ 0.28584863 & 0.12171721 & 0.47102262 & 0.24838159 \\ 0.14245983 & 0.07841352 & 0.24838159 & 0.17198206 \end{pmatrix} \quad (407)$$

$$\mu_2 = (6.90609957 \quad 3.00296152 \quad 5.91751825 \quad 2.02823662) \quad (408)$$

$$\Sigma_2 = \begin{pmatrix} 0.47944373 & 0.10318151 & 0.37851617 & -0.00947992 \\ 0.10318151 & 0.14649679 & 0.06689303 & 0.02639494 \\ 0.37851617 & 0.06689303 & 0.33139405 & 0.01686062 \\ -0.00947992 & 0.02639494 & 0.01686062 & 0.05638487 \end{pmatrix} \quad (409)$$

$$\pi = (0.33333326 \quad 0.54461183 \quad 0.12205491) \quad (410)$$

## Specific data in third random initialisation



Initilisation of parameters:

$$\mu_0 = (5.58180881 \quad 4.84821904 \quad 4.19996763 \quad 1.66923484) \quad (411)$$

$$\mu_1 = (5.98138886 \quad 3.26599148 \quad 2.48486145 \quad -0.38105446) \quad (412)$$

$$\mu_2 = (5.81470396 \quad 4.52844168 \quad 4.91642107 \quad 3.99937301) \quad (413)$$

$$\Sigma_0 = \Sigma_1 = \Sigma_2 = I \quad (414)$$

$$\pi = (0.5694223 \quad 0.22911455 \quad 0.20146316) \quad (415)$$

Iteration Quantity: 34

Parameters of resulted Guassians:

$$\mu_0 = (5.91506499 \quad 2.77785242 \quad 4.20175063 \quad 1.29704350) \quad (416)$$

$$\Sigma_0 = \begin{pmatrix} 0.27532206 & 0.09691983 & 0.18469154 & 0.05440596 \\ 0.09691983 & 0.09264056 & 0.09113461 & 0.04299639 \\ 0.18469154 & 0.09113461 & 0.20070652 & 0.06101118 \\ 0.05440596 & 0.04299639 & 0.06101118 & 0.03201270 \end{pmatrix} \quad (417)$$

$$\mu_1 = (5.00600000 \quad 3.42800000 \quad 1.46200000 \quad 0.24600000) \quad (418)$$

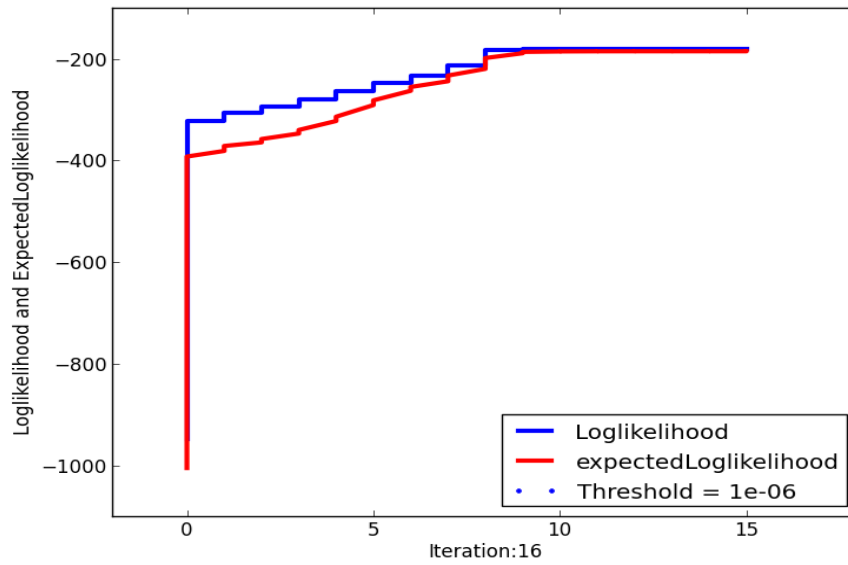
$$\Sigma_1 = \begin{pmatrix} 0.12176400 & 0.09723200 & 0.01602800 & 0.01012400 \\ 0.09723200 & 0.14081600 & 0.01146400 & 0.00911200 \\ 0.01602800 & 0.01146400 & 0.02955600 & 0.00594800 \\ 0.01012400 & 0.00911200 & 0.00594800 & 0.01088400 \end{pmatrix} \quad (419)$$

$$\mu_2 = (6.54467220 \quad 2.94870862 \quad 5.47980119 \quad 1.98476235) \quad (420)$$

$$\Sigma_2 = \begin{pmatrix} 0.38704847 & 0.09220736 & 0.30276907 & 0.06159530 \\ 0.09220736 & 0.11034094 & 0.08425917 & 0.05599446 \\ 0.30276907 & 0.08425917 & 0.32766484 & 0.07441712 \\ 0.06159530 & 0.05599446 & 0.07441712 & 0.08573678 \end{pmatrix} \quad (421)$$

$$\pi = (0.29931064 \quad 0.33333333 \quad 0.36735603) \quad (422)$$

## Specific data in fourth random initialisation



Initilisation of parameters:

$$\mu_0 = (4.93395858 \quad 3.72289358 \quad 2.88809351 \quad 1.81106115) \quad (423)$$

$$\mu_1 = (5.74870051 \quad 1.17972676 \quad 3.84560644 \quad 0.78416837) \quad (424)$$

$$\mu_2 = (7.63689574 \quad 2.95945543 \quad 5.19796747 \quad 1.71749004) \quad (425)$$

$$\Sigma_0 = \Sigma_1 = \Sigma_2 = I \quad (426)$$

$$\pi = (0.86693085 \quad 0.06884158 \quad 0.06422758) \quad (427)$$

Iteration Quantity: 16

Parameters of resulted Guassians:

$$\mu_0 = (5.00600000 \quad 3.42800000 \quad 1.46200000 \quad 0.24600000) \quad (428)$$

$$\Sigma_0 = \begin{pmatrix} 0.12176400 & 0.09723200 & 0.01602800 & 0.01012400 \\ 0.09723200 & 0.14081600 & 0.01146400 & 0.00911200 \\ 0.01602800 & 0.01146400 & 0.02955600 & 0.00594800 \\ 0.01012400 & 0.00911200 & 0.00594800 & 0.01088400 \end{pmatrix} \quad (429)$$

$$\mu_1 = (5.91510204 \quad 2.77785066 \quad 4.20180128 \quad 1.29706329) \quad (430)$$

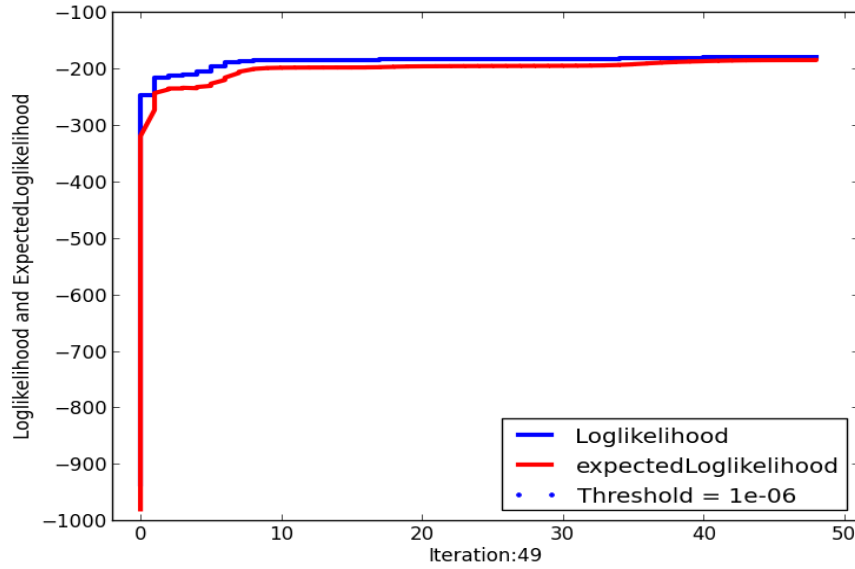
$$\Sigma_1 = \begin{pmatrix} 0.27532744 & 0.09691654 & 0.18470680 & 0.05441350 \\ 0.09691654 & 0.09264040 & 0.09113279 & 0.04299613 \\ 0.18470680 & 0.09113279 & 0.20073024 & 0.06102147 \\ 0.05441350 & 0.04299613 & 0.06102147 & 0.03201724 \end{pmatrix} \quad (431)$$

$$\mu_2 = (6.54468870 \quad 2.94872272 \quad 5.47985468 \quad 1.98479721) \quad (432)$$

$$\Sigma_2 = \begin{pmatrix} 0.38706100 & 0.09220597 & 0.30277029 & 0.06158850 \\ 0.09220597 & 0.11033973 & 0.08424916 & 0.05598870 \\ 0.30277029 & 0.08424916 & 0.32763930 & 0.07439355 \\ 0.06158850 & 0.05598870 & 0.07439355 & 0.08572418 \end{pmatrix} \quad (433)$$

$$\pi = (0.33333333 \quad 0.29933787 \quad 0.36732879) \quad (434)$$

## Specific data in fifth random initialisation



Initilisation of parameters:

$$\mu_0 = (4.74122040 \quad 3.15279417 \quad 1.77909908 \quad 0.80293032) \quad (435)$$

$$\mu_1 = (6.36558270 \quad 4.20423977 \quad 4.12724521 \quad 1.24996769) \quad (436)$$

$$\mu_2 = (7.34388669 \quad 3.54991993 \quad 4.39371314 \quad 2.05949636) \quad (437)$$

$$\Sigma_0 = \Sigma_1 = \Sigma_2 = I \quad (438)$$

$$\pi = (0.02053414 \quad 0.0965958 \quad 0.88287006) \quad (439)$$

Iteration Quantity: 49

Parameters of resulted Guassians:

$$\mu_0 = (5.00600000 \quad 3.42800000 \quad 1.46200000 \quad 0.24600000) \quad (440)$$

$$\Sigma_0 = \begin{pmatrix} 0.12176400 & 0.09723200 & 0.01602800 & 0.01012400 \\ 0.09723200 & 0.14081600 & 0.01146400 & 0.00911200 \\ 0.01602800 & 0.01146400 & 0.02955600 & 0.00594800 \\ 0.01012400 & 0.00911200 & 0.00594800 & 0.01088400 \end{pmatrix} \quad (441)$$

$$\mu_1 = (5.91484076 \quad 2.77783583 \quad 4.20129188 \quad 1.29686598) \quad (442)$$

$$\Sigma_1 = \begin{pmatrix} 0.27531479 & 0.09697069 & 0.18462390 & 0.05437052 \\ 0.09697069 & 0.09265277 & 0.09115580 & 0.04299921 \\ 0.18462390 & 0.09115580 & 0.20053150 & 0.06093600 \\ 0.05437052 & 0.04299921 & 0.06093600 & 0.03197655 \end{pmatrix} \quad (443)$$

$$\mu_2 = (6.54438345 \quad 2.94859425 \quad 5.47921795 \quad 1.98439211) \quad (444)$$

$$\Sigma_2 = \begin{pmatrix} 0.38703762 & 0.09221031 & 0.30286958 & 0.06172628 \\ 0.09221031 & 0.11033498 & 0.08432941 & 0.05603638 \\ 0.30286958 & 0.08432941 & 0.32798018 & 0.07468452 \\ 0.06172628 & 0.05603638 & 0.07468452 & 0.08588114 \end{pmatrix} \quad (445)$$

$$\pi = (0.33333333 \quad 0.29903553 \quad 0.36763114) \quad (446)$$