

# Handling missing data with expectation maximization algorithm

Loc Nguyen

Independent scholar, Vietnam

Email: ng\_phloc@yahoo.com

Homepage: www.locnguyen.net

## Abstract

Expectation maximization (EM) algorithm is a powerful mathematical tool for estimating parameter of statistical models in case of incomplete data or hidden data. EM assumes that there is a relationship between hidden data and observed data, which can be a joint distribution or a mapping function. Therefore, this implies another implicit relationship between parameter estimation and data imputation. If missing data which contains missing values is considered as hidden data, it is very natural to handle missing data by EM algorithm. Handling missing data is not a new research but this report focuses on the theoretical base with detailed mathematical proofs for fulfilling missing values with EM. Besides, multinormal distribution and multinomial distribution are the two sample statistical models which are concerned to hold missing values.

**Keywords:** expectation maximization (EM), missing data, multinormal distribution, multinomial distribution.

## 1. Introduction to expectation maximization algorithm

Literature of expectation maximization (EM) algorithm in this report is mainly extracted from the preeminent article “Maximum Likelihood from Incomplete Data via the EM Algorithm” by Arthur P. Dempster, Nan M. Laird, and Donald B. Rubin (Dempster, Laird, & Rubin, 1977). For convenience, let DLR be reference to such three authors. The preprint “Tutorial on EM algorithm” (Nguyen, 2020) by Loc Nguyen is also referred in this report.

Now we skim through an introduction of EM algorithm. Suppose there are two spaces  $\mathbf{X}$  and  $\mathbf{Y}$ , in which  $\mathbf{X}$  is *hidden space* whereas  $\mathbf{Y}$  is *observed space*. We do not know  $\mathbf{X}$  but there is a mapping from  $\mathbf{X}$  to  $\mathbf{Y}$  so that we can survey  $\mathbf{X}$  by observing  $\mathbf{Y}$ . The mapping is many-one function  $\varphi: \mathbf{X} \rightarrow \mathbf{Y}$  and we denote  $\varphi^{-1}(Y) = \{X \in \mathbf{X}: \varphi(X) = Y\}$  as all  $X \in \mathbf{X}$  such that  $\varphi(X) = Y$ . We also denote  $\mathbf{X}(Y) = \varphi^{-1}(Y)$ . Let  $f(X | \Theta)$  be the probability density function (PDF) of random variable  $X \in \mathbf{X}$  and let  $g(Y | \Theta)$  be the PDF of random variable  $Y \in \mathbf{Y}$ . Note,  $Y$  is also called observation. Equation 1.1 specifies  $g(Y | \Theta)$  as integral of  $f(X | \Theta)$  over  $\varphi^{-1}(Y)$ .

$$g(Y|\Theta) = \int_{\varphi^{-1}(Y)} f(X|\Theta) dX \quad (1.1)$$

Where  $\Theta$  is probabilistic parameter represented as a column vector,  $\Theta = (\theta_1, \theta_2, \dots, \theta_r)^T$  in which each  $\theta_i$  is a particular parameter. If  $X$  and  $Y$  are discrete, equation 1.1 is re-written as follows:

$$g(Y|\Theta) = \sum_{X \in \varphi^{-1}(Y)} f(X|\Theta)$$

According to viewpoint of Bayesian statistics,  $\Theta$  is also random variable. As a convention, let  $\Omega$  be the domain of  $\Theta$  such that  $\Theta \in \Omega$  and the dimension of  $\Omega$  is  $r$ . For example, normal distribution has two particular parameters such as mean  $\mu$  and variance  $\sigma^2$  and so we have  $\Theta = (\mu, \sigma^2)^T$ . Note that,  $\Theta$  can degrade into a scalar as  $\Theta = \theta$ . The conditional PDF of  $X$  given  $Y$ , denoted  $k(X | Y, \Theta)$ , is specified by equation 1.2.

$$k(X|Y, \Theta) = \frac{f(X|\Theta)}{g(Y|\Theta)} \quad (1.2)$$

According to DLR (Dempster, Laird, & Rubin, 1977, p. 1),  $X$  is called *complete data* and the term “incomplete data” implies existence of  $X$  and  $Y$  where  $X$  is not observed directly and  $X$  is only known by the many-one mapping  $\varphi: X \rightarrow Y$ . In general, we only know  $Y$ ,  $f(X | \Theta)$ , and  $k(X | Y, \Theta)$  and so our purpose is to estimate  $\Theta$  based on such  $Y$ ,  $f(X | \Theta)$ , and  $k(X | Y, \Theta)$ . Like MLE approach, EM algorithm also maximizes the likelihood function to estimate  $\Theta$  but the likelihood function in EM concerns  $Y$  and there are also some different aspects in EM which will be described later. Pioneers in EM algorithm firstly assumed that  $f(X | \Theta)$  belongs to exponential family with note that many popular distributions such as normal, multinomial, and Poisson belong to exponential family. Although DLR (Dempster, Laird, & Rubin, 1977) proposed a generality of EM algorithm in which  $f(X | \Theta)$  distributes arbitrarily, we should concern exponential family a little bit. Exponential family (Wikipedia, Exponential family, 2016) refers to a set of probabilistic distributions whose PDF (s) have the same exponential form according to equation 1.3 (Dempster, Laird, & Rubin, 1977, p. 3):

$$f(X|\Theta) = b(X) \exp(\Theta^T \tau(X)) / a(\Theta) \quad (1.3)$$

Where  $b(X)$  is a function of  $X$ , which is called base measure and  $\tau(X)$  is a vector function of  $X$ , which is sufficient statistic. For example, the sufficient statistic of normal distribution is  $\tau(X) = (X, XX^T)^T$ . Equation 1.3 expresses the canonical form of exponential family. Recall that  $\Omega$  is the domain of  $\Theta$  such that  $\Theta \in \Omega$ . Suppose that  $\Omega$  is a convex set. If  $\Theta$  is restricted only to  $\Omega$  then,  $f(X | \Theta)$  specifies a *regular exponential family*. If  $\Theta$  lies in a curved sub-manifold  $\Omega_0$  of  $\Omega$  then,  $f(X | \Theta)$  specifies a *curved exponential family*. The  $a(\Theta)$  is *partition function* for variable  $X$ , which is used for normalization.

$$a(\Theta) = \int_X b(X) \exp(\Theta^T \tau(X)) dX$$

As usual, a PDF is known as a popular form but its exponential family form (canonical form of exponential family) specified by equation 1.3 looks unlike popular form although they are the same. Therefore, parameter in popular form is different from parameter in exponential family form.

For example, multinormal distribution with theoretical mean  $\mu$  and covariance matrix  $\Sigma$  of random variable  $X = (x_1, x_2, \dots, x_n)^T$  has PDF in popular form is:

$$f(X|\mu, \Sigma) = (2\pi)^{-\frac{n}{2}} |\Sigma|^{-\frac{1}{2}} * \exp\left(-\frac{1}{2} (X - \mu)^T \Sigma^{-1} (X - \mu)\right)$$

Hence, parameter in popular form is  $\Theta = (\mu, \Sigma)^T$ . Exponential family form of such PDF is:

$$f(X|\theta_1, \theta_2) = (2\pi)^{-\frac{n}{2}} * \exp\left((\theta_1, \theta_2) \begin{pmatrix} X \\ XX^T \end{pmatrix}\right) / \exp\left(-\frac{1}{4} \theta_1^T \theta_2^{-1} \theta_1 - \frac{1}{2} \log|-2\theta_2|\right)$$

Where,

$$\begin{aligned} \Theta &= \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} \\ \theta_1 &= \Sigma^{-1} \mu \\ \theta_2 &= -\frac{1}{2} \Sigma^{-1} \\ b(X) &= (2\pi)^{-\frac{n}{2}} \\ \tau(X) &= \begin{pmatrix} X \\ XX^T \end{pmatrix} \\ a(\Theta) &= \exp\left(-\frac{1}{4} \theta_1^T \theta_2^{-1} \theta_1 - \frac{1}{2} \log|-2\theta_2|\right) \end{aligned}$$

The exponential family form is used to represents all distributions belonging to exponential family as canonical form. Parameter in exponential family form is called exponential family parameter. As a convention, parameter  $\Theta$  mentioned in EM algorithm is often exponential family parameter if PDF belongs to exponential family and there is no additional information.

Expectation maximization (EM) algorithm has many iterations and each iteration has two steps in which expectation step (E-step) calculates sufficient statistic of hidden data based on observed data and current parameter whereas maximization step (M-step) re-estimates parameter. When DLR proposed EM algorithm (Dempster, Laird, & Rubin, 1977), they firstly concerned that the PDF  $f(X | \Theta)$  of hidden space belongs to exponential family. E-step and M-step at the  $t^{\text{th}}$  iteration are described in table 1.1 (Dempster, Laird, & Rubin, 1977, p. 4), in which the current estimate is  $\Theta^{(t)}$ , with note that  $f(X | \Theta)$  belongs to regular exponential family.

*E-step:*

We calculate current value  $\tau^{(t)}$  of the sufficient statistic  $\tau(X)$  from observed  $Y$  and current parameter  $\Theta^{(t)}$  according to following equation:

$$\tau^{(t)} = E(\tau(X)|Y, \Theta^{(t)}) = \int_{\varphi^{-1}(Y)} k(X|Y, \Theta^{(t)})\tau(X)dX$$

*M-step:*

Basing on  $\tau^{(t)}$ , we determine the next parameter  $\Theta^{(t+1)}$  as solution of following equation:

$$E(\tau(X)|\Theta) = \int_X f(X|\Theta)\tau(X)dX = \tau^{(t)}$$

Note,  $\Theta^{(t+1)}$  will become current parameter at the next iteration ( $(t+1)^{\text{th}}$  iteration).

**Table 1.1.** E-step and M-step of EM algorithm given regular exponential PDF  $f(X|\Theta)$   
EM algorithm stops if two successive estimates are equal,  $\Theta^* = \Theta^{(t)} = \Theta^{(t+1)}$ , at some  $t^{\text{th}}$  iteration. At that time we conclude that  $\Theta^*$  is the optimal estimate of EM process. As a convention, the estimate of parameter  $\Theta$  resulted from EM process is denoted  $\Theta^*$  instead of  $\hat{\Theta}$  in order to emphasize that  $\Theta^*$  is solution of optimization problem.

For further research, DLR gave a preeminent generality of EM algorithm (Dempster, Laird, & Rubin, 1977, pp. 6-11) in which  $f(X | \Theta)$  specifies arbitrary distribution. In other words, there is no requirement of exponential family. They define the conditional expectation  $Q(\Theta' | \Theta)$  according to equation 1.4 (Dempster, Laird, & Rubin, 1977, p. 6).

$$Q(\Theta' | \Theta) = E(\log(f(X|\Theta')) | Y, \Theta) = \int_{\varphi^{-1}(Y)} k(X|Y, \Theta) \log(f(X|\Theta')) dX \quad (1.4)$$

If  $X$  and  $Y$  are discrete, equation 2.4 can be re-written as follows:

$$Q(\Theta' | \Theta) = E(\log(f(X|\Theta')) | Y, \Theta) = \sum_{X \in \varphi^{-1}(Y)} k(X|Y, \Theta) \log(f(X|\Theta'))$$

The two steps of generalized EM (GEM) algorithm aim to maximize  $Q(\Theta | \Theta^{(t)})$  at some  $t^{\text{th}}$  iteration as seen in table 1.2 (Dempster, Laird, & Rubin, 1977, p. 6).

*E-step:*

The expectation  $Q(\Theta | \Theta^{(t)})$  is determined based on current parameter  $\Theta^{(t)}$ , according to equation 1.4. Actually,  $Q(\Theta | \Theta^{(t)})$  is formulated as function of  $\Theta$ .

*M-step:*

The next parameter  $\Theta^{(t+1)}$  is a maximizer of  $Q(\Theta | \Theta^{(t)})$  with subject to  $\Theta$ . Note that  $\Theta^{(t+1)}$  will become current parameter at the next iteration (the  $(t+1)^{\text{th}}$  iteration).

**Table 1.2.** E-step and M-step of GEM algorithm

DLR proved that GEM algorithm converges at some  $t^{\text{th}}$  iteration. At that time,  $\Theta^* = \Theta^{(t+1)} = \Theta^{(t)}$  is the optimal estimate of EM process, which is an optimizer of  $L(\Theta)$ .

$$\Theta^* = \underset{\Theta}{\operatorname{argmax}} L(\Theta)$$

It is deduced from E-step and M-step that  $Q(\Theta | \Theta^{(t)})$  is increased after every iteration. How to maximize  $Q(\Theta | \Theta^{(t)})$  is the optimization problem which is dependent on applications. For example, the estimate  $\Theta^{(t+1)}$  can be solution of the equation created by setting the first-order derivative of  $Q(\Theta | \Theta^{(t)})$  regarding  $\Theta$  to be zero,  $DQ(\Theta | \Theta^{(t)}) = \mathbf{0}^T$ . If solving such equation is too

complex or impossible, some popular methods to solve optimization problem are Newton-Raphson (Burden & Faires, 2011, pp. 67-71), gradient descent (Ta, 2014), and Lagrange duality (Wikipedia, Karush–Kuhn–Tucker conditions, 2014).

In practice, if  $Y$  is observed as particular  $N$  observations  $Y_1, Y_2, \dots, Y_N$ . Let  $\mathcal{Y} = \{Y_1, Y_2, \dots, Y_N\}$  be the observed sample of size  $N$  with note that all  $Y_i$  (s) are mutually independent and identically distributed (iid). Given an observation  $Y_i$ , there is an associated random variable  $X_i$ . All  $X_i$  (s) are iid and they are not existent in fact. Each  $X_i \in \mathbf{X}$  is a random variable like  $X$ . Of course, the domain of each  $X_i$  is  $\mathbf{X}$ . Let  $\mathcal{X} = \{X_1, X_2, \dots, X_N\}$  be the set of associated random variables. Because all  $X_i$  (s) are iid, the joint PDF of  $\mathcal{X}$  is determined as follows:

$$f(\mathcal{X}|\Theta) = f(X_1, X_2, \dots, X_N|\Theta) = \prod_{i=1}^N f(X_i|\Theta)$$

Because all  $X_i$  (s) are iid and each  $Y_i$  is associated with  $X_i$ , the conditional joint PDF of  $\mathcal{X}$  given  $\mathcal{Y}$  is determined as follows:

$$k(\mathcal{X}|\mathcal{Y}, \Theta) = k(X_1, X_2, \dots, X_N|Y_1, Y_2, \dots, Y_N, \Theta) = \prod_{i=1}^N k(X_i|Y_i, \Theta) = \prod_{i=1}^N k(X_i|Y_i, \Theta)$$

The conditional expectation  $Q(\Theta' | \Theta)$  given samples  $\mathbf{X}$  and  $\mathbf{Y}$  is re-written according to equation 1.5.

$$Q(\Theta'|\Theta) = \sum_{i=1}^N \int_{\varphi^{-1}(Y_i)} k(X|Y_i, \Theta) \log(f(X|\Theta')) dX \quad (1.5)$$

Equation 1.5 is proved in (Nguyen, 2020, pp. 45-47). In case that  $f(X|\Theta)$  and  $k(X|Y_i, \Theta)$  belong to exponential family, equation 1.5 becomes equation 1.6 with an observed sample  $\mathcal{Y} = \{Y_1, Y_2, \dots, Y_N\}$ .

$$Q(\Theta'|\Theta) = \left( \sum_{i=1}^N E(\log(b(X))|Y_i, \Theta) \right) + \left( (\Theta')^T \sum_{i=1}^N \tau_{\Theta, Y_i} \right) - N \log(a(\Theta')) \quad (1.6)$$

Where,

$$E(\log(b(X))|Y_i, \Theta) = \int_{\varphi^{-1}(Y_i)} k(X|Y_i, \Theta) \log(b(X)) dX$$

$$\tau_{\Theta, Y_i} = E(\tau(X)|Y_i, \Theta) = \int_{\varphi^{-1}(Y_i)} k(X|Y_i, \Theta) \tau(X) dX$$

DLR (Dempster, Laird, & Rubin, 1977, p. 1) called  $\mathbf{X}$  as *complete data* because the mapping  $\varphi: \mathbf{X} \rightarrow \mathbf{Y}$  is many-one function. There is another case that the complete space  $\mathbf{Z}$  consists of hidden space  $\mathbf{X}$  and observed space  $\mathbf{Y}$  with note that  $\mathbf{X}$  and  $\mathbf{Y}$  are separated. There is no explicit mapping  $\varphi$  from  $\mathbf{X}$  and  $\mathbf{Y}$  but there exists a PDF of  $Z \in \mathbf{Z}$  as the joint PDF of  $X \in \mathbf{X}$  and  $Y \in \mathbf{Y}$ .

$$f(Z|\Theta) = f(X, Y|\Theta)$$

The PDF of  $Y$  becomes:

$$f(Y|\Theta) = \int_{\mathbf{X}} f(X, Y|\Theta) dX$$

The PDF  $f(Y|\Theta)$  is equivalent to the PDF  $g(Y|\Theta)$  mentioned in equation 1.1. Although there is no explicit mapping from  $\mathbf{X}$  to  $\mathbf{Y}$ , the PDF of  $Y$  above implies an implicit mapping from  $\mathbf{Z}$  to  $\mathbf{Y}$ . The conditional PDF of  $X$  given  $Z$  is specified according to Bayes' rule as follows:

$$f(Z|Y, \Theta) = f(X, Y|Y, \Theta) = f(X|Y) f(Y|Y) = f(X|Y, \Theta) = \frac{f(X, Y|\Theta)}{f(Y|\Theta)} = \frac{f(X, Y|\Theta)}{\int_{\mathbf{X}} f(X, Y|\Theta) dX}$$

The conditional PDF  $f(X|Y, \Theta)$  is equivalent to the conditional PDF  $k(X|Y, \Theta)$  mentioned in equation 1.2. Of course, given  $Y$ , we always have:

$$\int_X f(X|Y, \Theta) dX = 1$$

Equation 1.7 specifies the conditional expectation  $Q(\Theta' | \Theta)$  in case that there is no explicit mapping from  $\mathbf{X}$  to  $\mathbf{Y}$  but there exists the joint PDF of  $X$  and  $Y$ .

$$Q(\Theta' | \Theta) = \int_X f(X|Y, \Theta) \log(f(X|Y, \Theta')) dX = \int_X f(X|Y, \Theta) \log(f(X, Y | \Theta')) dX \quad (1.7)$$

Where,

$$f(X|Y, \Theta) = \frac{f(X, Y | \Theta)}{f(Y | \Theta)} = \frac{f(X, Y | \Theta)}{\int_X f(X, Y | \Theta) dX}$$

Note,  $\mathbf{X}$  is separated from  $\mathbf{Y}$  and the complete data  $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$  is composed of  $\mathbf{X}$  and  $\mathbf{Y}$ . For equation 1.7, the existence of the joint PDF  $f(X, Y | \Theta)$  can be replaced by the existence of the conditional PDF  $f(Y|X, \Theta)$  and the prior PDF  $f(X|\Theta)$  due to:

$$f(X, Y | \Theta) = f(Y|X, \Theta) f(X | \Theta)$$

In applied statistics, equation 1.4 is often replaced by equation 1.7 because specifying the joint PDF  $f(X, Y | \Theta)$  is more practical than specifying the mapping  $\varphi: \mathbf{X} \rightarrow \mathbf{Y}$ . However, equation 1.4 is more general equation 1.7 because the requirement of the joint PDF for equation 1.7 is stricter than the requirement of the explicit mapping for equation 1.4. In case that  $X$  and  $Y$  are discrete, equation 1.7 becomes:

$$Q(\Theta' | \Theta) = \sum_X P(X|Y, \Theta) \log(P(X, Y | \Theta'))$$

In practice, suppose  $Y$  is observed as a sample  $\mathcal{Y} = \{Y_1, Y_2, \dots, Y_N\}$  of size  $N$  with note that all  $Y_i$  (s) are mutually independent and identically distributed (iid). The observed sample  $\mathcal{Y}$  is associated with a hidden set (latent set)  $\mathcal{X} = \{X_1, X_2, \dots, X_N\}$  of size  $N$ . All  $X_i$  (s) are iid and they are not existent in fact. Let  $X \in \mathbf{X}$  be the random variable representing every  $X_i$ . Of course, the domain of  $X$  is  $\mathbf{X}$ . Equation 1.8 specifies the conditional expectation  $Q(\Theta' | \Theta)$  given such  $\mathcal{Y}$ .

$$Q(\Theta' | \Theta) = \sum_{i=1}^N \int_X f(X|Y_i, \Theta) \log(f(X, Y_i | \Theta')) dX \quad (1.8)$$

Equation 1.8 is a variant of equation 1.5 in case that there is no explicit mapping between  $X_i$  and  $Y_i$  but there exists the same joint PDF between  $X_i$  and  $Y_i$ . If both  $X$  and  $Y$  are discrete, equation 1.8 becomes:

$$Q(\Theta' | \Theta) = \sum_{i=1}^N \sum_X P(X|Y_i, \Theta) \log(P(X, Y_i | \Theta')) \quad (1.9)$$

If  $X$  is discrete and  $Y$  is continuous such that  $f(X, Y | \Theta) = P(X|\Theta) f(Y | X, \Theta)$  then, according to the total probability rule, we have:

$$f(Y | \Theta) = \sum_X P(X|\Theta) f(Y | X, \Theta)$$

Note, when only  $X$  is discrete, its PDF  $f(X|\Theta)$  becomes the probability  $P(X|\Theta)$ . Therefore, equation 1.10 is a variant of equation 1.8, as follows:

$$Q(\Theta' | \Theta) = \sum_{i=1}^N \sum_X P(X|Y_i, \Theta) \log(P(X|\Theta') f(Y_i | X, \Theta')) \quad (1.10)$$

Where  $P(X | Y_i, \Theta)$  is determined by Bayes' rule, as follows:

$$P(X|Y_i, \Theta) = \frac{P(X|\Theta) f(Y_i | X, \Theta)}{\sum_X P(X|\Theta) f(Y_i | X, \Theta)}$$

Equation 1.10 is the base for estimating the probabilistic mixture model by EM algorithm, which is not main subject of this report. Now we consider how to apply EM into handling missing data in which equation 1.8 is most concerned. The goal of maximum likelihood estimation (MLE), maximum a posteriori (MAP), and EM is to estimate statistical based on sample. Whereas MLE and MAP require complete data, EM accepts hidden data or incomplete data. Therefore, EM is appropriate to handle missing data which contains missing values. Indeed, estimating parameter with missing data is very natural for EM but it is necessary to have a new viewpoint in which missing data is considered as hidden data ( $X$ ). Moreover, the GEM version with joint probability (without mapping function, please see equation 1.7 and equation 1.8) is used and some changes are required. Handling missing data, which is the main subject of this report is described in next section.

## 2. Handling missing data

Let  $X = (x_1, x_2, \dots, x_n)^T$  be  $n$ -dimension random variable whose  $n$  elements are partial random variables  $x_j$  (s). Suppose  $X$  is composed of two parts such as observed part  $X_{obs}$  and missing part  $X_{mis}$  such that  $X = \{X_{obs}, X_{mis}\}$ . Note,  $X_{obs}$  and  $X_{mis}$  are considered as random variables.

$$X = \{X_{obs}, X_{mis}\} = (x_1, x_2, \dots, x_n)^T \quad (2.1)$$

When  $X$  is observed,  $X_{obs}$  and  $X_{mis}$  are determined. For example, given  $X = (x_1, x_2, x_3, x_4)^T$ , when  $X$  is observed as  $X = (x_1=1, x_2=?, x_3=4, x_4=?, x_5=9)^T$  where question mark “?” denotes missing value,  $X_{obs}$  and  $X_{mis}$  are determined as  $X_{obs} = (x_1=1, x_3=4, x_5=9)^T$  and  $X_{mis} = (x_2=?, x_4=?)^T$ . When  $X$  is observed as  $X = (x_1=?, x_2=3, x_3=4, x_4=?, x_5=?)^T$  then,  $X_{obs}$  and  $X_{mis}$  are determined as  $X_{obs} = (x_2=3, x_3=4)^T$  and  $X_{mis} = (x_1=?, x_4=?, x_5=?)^T$ . Let  $M$  be a set of indices that  $x_j$  (s) are missing when  $X$  is observed.  $M$  is called missing index set.

$$M = \{j: x_j \text{ missing}\} \text{ where } j = \overline{1, n} \quad (2.2)$$

Suppose

$$M = \{m_1, m_2, \dots, m_{|M|}\} \quad (2.3)$$

Where,

$$\begin{aligned} m_i &= \overline{1, n} \\ m_i &\neq m_j \end{aligned}$$

Let  $\bar{M}$  is complementary set of the set  $M$  given the set  $\{1, 2, \dots, n\}$ .  $\bar{M}$  is called existent index set.

$$\bar{M} = \{j: x_j \text{ existent}\} \text{ where } j = \overline{1, n} \quad (2.4)$$

$M$  or  $\bar{M}$  can be empty. They are mutual because  $\bar{M}$  can be defined based on  $M$  and vice versa.

$$M \cup \bar{M} = \{1, 2, \dots, n\}$$

$$M \cap \bar{M} = \emptyset$$

Suppose

$$\bar{M} = \{\bar{m}_1, \bar{m}_2, \dots, \bar{m}_{|\bar{M}|}\} \quad (2.5)$$

Where,

$$\begin{aligned} \bar{m}_i &= \overline{1, n} \\ \bar{m}_i &\neq \bar{m}_j \\ |M| + |\bar{M}| &= n \end{aligned}$$

We have:

$$X_{mis} = (x_j: j \in M)^T = (x_{m_1}, x_{m_2}, \dots, x_{m_{|M|}})^T \quad (2.6)$$

$$X_{obs} = (x_j: j \in \bar{M})^T = (x_{\bar{m}_1}, x_{\bar{m}_2}, \dots, x_{\bar{m}_{|\bar{M}|}})^T \quad (2.7)$$

Obviously, dimension of  $X_{mis}$  is  $|M|$  and dimension of  $X_{obs}$  is  $|\bar{M}| = n - |M|$ . Note, when composing  $X$  from  $X_{obs}$  and  $X_{mis}$  as  $X = \{X_{obs}, X_{mis}\}$ , it is required a right re-arrangement of elements in both  $X_{obs}$  and  $X_{mis}$ .

Let  $Z = (z_1, z_2, \dots, z_n)^T$  be  $n$ -dimension random variable whose each element  $z_j$  is binary random variable indicating if  $x_j$  is missing. Random variable  $Z$  is also called missingness variable.

$$z_j = \begin{cases} 1 & \text{if } x_j \text{ missing} \\ 0 & \text{if } x_j \text{ existent} \end{cases} \quad (2.8)$$

For example, given  $X = (x_1, x_2, x_3, x_4)^T$ , when  $X$  is observed as  $X = (x_1=1, x_2=?, x_3=4, x_4=?, x_5=9)^T$ , we have  $X_{obs} = (x_1=1, x_3=4, x_5=9)^T$ ,  $X_{mis} = (x_2=?, x_4=?)^T$ , and  $Z = (z_1=0, z_2=1, z_3=0, z_4=1, z_5=0)^T$ .

Generally, when  $X$  is replaced by a sample  $\mathcal{X} = \{X_1, X_2, \dots, X_N\}$  whose  $X_i$  (s) are iid, let  $Z = \{Z_1, Z_2, \dots, Z_N\}$  be a set of missingness variables associated with  $\mathcal{X}$ . All  $Z_i$  (s) are iid too.  $\mathcal{X}$  and  $Z$  can be represented as matrices. Given  $X_i$ , its associative quantities are  $Z_i$ ,  $M_i$ , and  $\bar{M}_i$ . Let  $X = \{X_{obs}, X_{mis}\}$  be random variable representing every  $X_i$ . Let  $Z$  be random variable representing every  $Z_i$ . As a convention,  $X_{obs}(i)$  and  $X_{mis}(i)$  refer to  $X_{obs}$  part and  $X_{mis}$  part of  $X_i$ . We have:

$$\begin{aligned} X_i &= \{X_{obs}(i), X_{mis}(i)\} = (x_{i1}, x_{i2}, \dots, x_{in})^T \\ X_{mis}(i) &= (x_{im_1}, x_{im_2}, \dots, x_{im_{|M|}})^T \\ X_{obs}(i) &= (x_{i\bar{m}_1}, x_{i\bar{m}_2}, \dots, x_{i\bar{m}_{|\bar{M}|}})^T \\ M_i &= \{m_{i1}, m_{i2}, \dots, m_{i|M|}\} \\ \bar{M}_i &= \{\bar{m}_{i1}, \bar{m}_{i2}, \dots, \bar{m}_{i|\bar{M}|}\} \\ Z_i &= (z_{i1}, z_{i2}, \dots, z_{in})^T \end{aligned} \quad (2.9)$$

For example, given sample of size 4,  $\mathcal{X} = \{X_1, X_2, X_3, X_4\}$  in which  $X_1 = (x_{11}=1, x_{12}=?, x_{13}=3, x_{14}=?)^T$ ,  $X_2 = (x_{21}=?, x_{22}=2, x_{23}=?, x_{24}=4)^T$ ,  $X_3 = (x_{31}=1, x_{32}=2, x_{33}=?, x_{34}=?)^T$ , and  $X_4 = (x_{41}=?, x_{42}=?, x_{43}=3, x_{44}=4)^T$  are iid. Therefore, we also have  $Z_1 = (z_{11}=0, z_{12}=1, z_{13}=0, z_{14}=1)^T$ ,  $Z_2 = (z_{21}=1, z_{22}=0, z_{23}=1, z_{24}=0)^T$ ,  $Z_3 = (z_{31}=0, z_{32}=0, z_{33}=1, z_{34}=1)^T$ , and  $Z_4 = (z_{41}=1, z_{42}=1, z_{43}=0, z_{44}=0)^T$ . All  $Z_i$  (s) are iid too.

	$x_1$	$x_2$	$x_3$	$x_4$
$X_1$	1	?	3	?
$X_2$	?	2	?	4
$X_3$	1	2	?	?
$X_4$	?	?	3	4

	$z_1$	$z_2$	$z_3$	$z_4$
$Z_1$	0	1	0	1
$Z_2$	1	0	1	0
$Z_3$	0	0	1	1
$Z_4$	1	1	0	0

Of course, we have  $X_{obs}(1) = (x_{11}=1, x_{13}=3)^T$ ,  $X_{mis}(1) = (x_{12}=?, x_{14}=?)^T$ ,  $X_{obs}(2) = (x_{22}=2, x_{24}=4)^T$ ,  $X_{mis}(2) = (x_{21}=?, x_{23}=?)^T$ ,  $X_{obs}(3) = (x_{31}=1, x_{32}=2)^T$ ,  $X_{mis}(3) = (x_{33}=?, x_{34}=?)^T$ ,  $X_{obs}(4) = (x_{43}=3, x_{44}=4)^T$ , and  $X_{mis}(4) = (x_{41}=?, x_{42}=?)^T$ . We also have  $M_1 = \{m_{11}=2, m_{12}=4\}$ ,  $\bar{M}_1 = \{\bar{m}_{11}=1, \bar{m}_{12}=3\}$ ,  $M_2 = \{m_{21}=1, m_{22}=3\}$ ,  $\bar{M}_2 = \{\bar{m}_{21}=2, \bar{m}_{22}=4\}$ ,  $M_3 = \{m_{31}=3, m_{32}=4\}$ ,  $\bar{M}_3 = \{\bar{m}_{31}=1, \bar{m}_{32}=2\}$ ,  $M_4 = \{m_{41}=1, m_{42}=2\}$ , and  $\bar{M}_4 = \{\bar{m}_{41}=3, \bar{m}_{42}=4\}$ .

Both  $X$  and  $Z$  are associated with their own PDFs, as follows:

$$\begin{aligned} f(X|\Theta) &= f(X_{obs}, X_{mis}|\Theta) \\ f(Z|X_{obs}, X_{mis}, \Phi) \end{aligned} \quad (2.10)$$

Where  $\Theta$  and  $\Phi$  are parameters of PDFs of  $X = \{X_{obs}, X_{mis}\}$  and  $Z$ , respectively. The goal of handling missing data is to estimate  $\Theta$  and  $\Phi$  given  $X$ . Sufficient statistic of  $X = \{X_{obs}, X_{mis}\}$  is composed of sufficient statistic of  $X_{obs}$  and sufficient statistic of  $X_{mis}$ .

$$\tau(X) = \tau(X_{obs}, X_{mis}) = \{\tau(X_{obs}), \tau(X_{mis})\} \quad (2.11)$$

How to compose  $\tau(X)$  from  $\tau(X_{obs})$  and  $\tau(X_{mis})$  is dependent on distribution type of the PDF  $f(X|\Theta)$ .

The joint PDF of  $X$  and  $Z$  is main object of handling missing data, which is defined as follows:

$$f(X, Z|\Theta, \Phi) = f(X_{obs}, X_{mis}, Z|\Theta, \Phi) = f(Z|X_{obs}, X_{mis}, \Phi)f(X_{obs}, X_{mis}|\Theta) \quad (2.12)$$

The PDF of  $X_{obs}$  is defined as integral of  $f(X|\Theta)$  over  $X_{mis}$ :

$$f(X_{obs}|\Theta) = \int_{X_{mis}} f(X_{obs}, X_{mis}|\Theta) dX_{mis} \quad (2.13)$$

The PDF of  $X_{mis}$  is conditional PDF of  $X_{mis}$  given  $X_{obs}$  is:

$$f(X_{mis}|X_{obs}, \Theta_M) = f(X_{mis}|X_{obs}, \Theta) = \frac{f(X|\Theta)}{f(X_{obs}|\Theta)} = \frac{f(X_{obs}, X_{mis}|\Theta)}{f(X_{obs}|\Theta)} \quad (2.14)$$

The notation  $\Theta_M$  implies that the parameter  $\Theta_M$  of the PDF  $f(X_{mis}|X_{obs}, \Theta_M)$  is derived from the parameter  $\Theta$  of the PDF  $f(X|\Theta)$ , which is function of  $\Theta$  and  $X_{obs}$  as  $\Theta_M = u(\Theta, X_{obs})$ . Thus,  $\Theta_M$  is not a new parameter and it is dependent on distribution type.

$$\Theta_M = u(\Theta, X_{obs}) \quad (2.15)$$

How to determine  $u(\Theta, X_{obs})$  is dependent on distribution type of the PDF  $f(X|\Theta)$ .

There are three types of missing data, which depends on relationship between  $X_{obs}$ ,  $X_{mis}$ , and  $Z$  (Josse, Jiang, Sportisse, & Robin, 2018):

- Missing data ( $X$  or  $\mathcal{X}$ ) is Missing Completely At Random (MCAR) if the probability of  $Z$  depends on both  $X_{obs}$  and  $X_{mis}$  such that  $f(Z|X_{obs}, X_{mis}, \Phi) = f(Z|\Phi)$ .
- Missing data ( $X$  or  $\mathcal{X}$ ) is Missing At Random (MAR) if the probability of  $Z$  depends on only  $X_{obs}$  such that  $f(Z|X_{obs}, X_{mis}, \Phi) = f(Z|X_{obs}, \Phi)$ .
- Missing data ( $X$  or  $\mathcal{X}$ ) is Missing Not At Random (MNAR) in all other cases.

There are two main approaches for handling missing data (Josse, Jiang, Sportisse, & Robin, 2018):

- Using some statistical models such as EM to estimate parameter with missing data.
- Inputting plausible values for missing values to obtain some complete samples (copies) from the missing data. Later on, every complete sample is used to produce an estimate of parameter by some estimation methods, for example, MLE and MAP. Finally, all estimates are synthesized to produce the best estimate.

Here we focus on the first approach with EM to estimate parameter with missing data. Without loss of generality, given sample  $\mathcal{X} = \{X_1, X_2, \dots, X_N\}$  in which all  $X_i$  (s) are iid, by applying equation 1.8 for GEM with the joint PDF  $f(X_{obs}, X_{mis}, Z|\Theta, \Phi)$ , we consider  $\{X_{obs}, Z\}$  as observed part and  $X_{mis}$  as hidden part. Let  $X = \{X_{obs}, X_{mis}\}$  be random variable representing all  $X_i$  (s). Let  $X_{obs}(i)$  denote observed part  $X_{obs}$  of  $X_i$  and let  $Z_i$  be missingness variable corresponding to  $X_i$ , by following equation 1.8, the expectation  $Q(\Theta', \Phi'|\Theta, \Phi)$  becomes:

$$\begin{aligned} Q(\Theta', \Phi'|\Theta, \Phi) &= \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), Z_i, \Theta, \Phi) * \log(f(X_{obs}(i), X_{mis}, Z_i|\Theta', \Phi')) dX_{mis} \\ &= \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta) * \log(f(X_{obs}(i), X_{mis}, Z_i|\Theta', \Phi')) dX_{mis} \\ &= \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) * \log(f(X_{obs}(i), X_{mis}, Z_i|\Theta', \Phi')) dX_{mis} \end{aligned}$$



$$\begin{aligned}
&= \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \\
&\quad * \log(f(X_{obs}(i), X_{mis}|\Theta', \Phi') * f(Z_i|X_{obs}(i), X_{mis}, \Theta', \Phi')) dX_{mis} \\
&= \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) * \log(f(X_{obs}(i), X_{mis}|\Theta') * f(Z_i|X_{obs}(i), X_{mis}, \Phi')) dX_{mis} \\
&= \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \\
&\quad * \left( \log(f(X_{obs}(i), X_{mis}|\Theta')) + \log(f(Z_i|X_{obs}(i), X_{mis}, \Phi')) \right) dX_{mis} \\
&= \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \log(f(X_{obs}(i), X_{mis}|\Theta')) dX_{mis} \\
&\quad + \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \log(f(Z_i|X_{obs}(i), X_{mis}, \Phi')) dX_{mis}
\end{aligned}$$

In short,  $Q(\Theta', \Phi' | \Theta, \Phi)$  is specified as follows:

$$Q(\Theta', \Phi' | \Theta, \Phi) = Q_1(\Theta' | \Theta) + Q_2(\Phi' | \Theta) \quad (2.16)$$

Where,

$$\begin{aligned}
Q_1(\Theta' | \Theta) &= \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \log(f(X_{obs}(i), X_{mis}|\Theta')) dX_{mis} \\
Q_2(\Phi' | \Theta) &= \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \log(f(Z_i|X_{obs}(i), X_{mis}, \Phi')) dX_{mis}
\end{aligned}$$

Note, unknowns of  $Q(\Theta', \Phi' | \Theta, \Phi)$  are  $\Theta'$  and  $\Phi'$ . Because it is not easy to maximize  $Q(\Theta', \Phi' | \Theta, \Phi)$  with regard to  $\Theta'$  and  $\Phi'$ , we assume that the PDF  $f(X|\Theta)$  belongs to exponential family.

$$f(X|\Theta) = f(X_{obs}, X_{mis}|\Theta) = b(X_{obs}, X_{mis}) * \exp((\Theta)^T \tau(X_{obs}, X_{mis}))/a(\Theta) \quad (2.17)$$

Note,

$$\begin{aligned}
b(X) &= b(X_{obs}, X_{mis}) \\
\tau(X) &= \tau(X_{obs}, X_{mis}) = \{\tau(X_{obs}), \tau(X_{mis})\}
\end{aligned}$$

It is easy to deduce that

$$f(X_{mis}|X_{obs}, \Theta_M) = b(X_{mis}) \exp((\Theta_M)^T \tau(X_{mis}))/a(\Theta_M) \quad (2.18)$$

Therefore,

$$f(X_{mis}|X_{obs}(i), \Theta_{M_i}) = b(X_{mis}) \exp((\Theta_{M_i})^T \tau(X_{mis}))/a(\Theta_{M_i})$$

We have:

$$\begin{aligned}
Q_1(\Theta' | \Theta) &= \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \log(f(X_{obs}(i), X_{mis}|\Theta')) dX_{mis} \\
&= \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \\
&\quad * \log(b(X_{obs}(i), X_{mis}) \exp((\Theta')^T \tau(X_{obs}(i), X_{mis}))/a(\Theta')) dX_{mis}
\end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \\
&\quad * \left( \log(b(X_{obs}(i), X_{mis})) + (\Theta')^T \tau(X_{obs}(i), X_{mis}) - \log(a(\Theta')) \right) dX_{mis} \\
&= \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \log(b(X_{obs}(i), X_{mis})) dX_{mis} \\
&\quad + \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) (\Theta')^T \tau(X_{obs}(i), X_{mis}) dX_{mis} \\
&\quad - \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \log(a(\Theta')) dX_{mis} \\
&= \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \log(b(X_{obs}(i), X_{mis})) dX_{mis} \\
&\quad + (\Theta')^T \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \tau(X_{obs}(i), X_{mis}) dX_{mis} \\
&\quad - \log(a(\Theta')) \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) dX_{mis} \\
&= \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \log(b(X_{obs}(i), X_{mis})) dX_{mis} \\
&\quad + (\Theta')^T \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \tau(X_{obs}(i), X_{mis}) dX_{mis} - N \log(a(\Theta')) \\
&= \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \log(b(X_{obs}(i), X_{mis})) dX_{mis} \\
&\quad + (\Theta')^T \sum_{i=1}^N \left\{ \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \tau(X_{obs}(i)) dX_{mis}, \right. \\
&\quad \left. \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \tau(X_{mis}) dX_{mis} \right\} - N \log(a(\Theta')) \\
&= \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \log(b(X_{obs}(i), X_{mis})) dX_{mis} \\
&\quad + (\Theta')^T \sum_{i=1}^N \left\{ \tau(X_{obs}(i)) \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) dX_{mis}, \right. \\
&\quad \left. \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \tau(X_{mis}) dX_{mis} \right\} - N \log(a(\Theta'))
\end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^N \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \log(b(X_{obs}(i), X_{mis})) dX_{mis} \\
&\quad + (\Theta')^T \sum_{i=1}^N \left\{ \tau(X_{obs}(i)), \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \tau(X_{mis}) dX_{mis} \right\} - N \log(a(\Theta'))
\end{aligned}$$

Therefore, equation 2.19 specifies  $Q_1(\Theta'|\Theta)$  given  $f(X|\Theta)$  belongs to exponential family.

$$\begin{aligned}
Q_1(\Theta'|\Theta) &= \sum_{i=1}^N E(\log(b(X_{obs}(i), X_{mis})) | \Theta_{M_i}) \\
&\quad + (\Theta')^T \sum_{i=1}^N \{ \tau(X_{obs}(i)), E(\tau(X_{mis}) | \Theta_{M_i}) \} - N \log(a(\Theta'))
\end{aligned} \tag{2.19}$$

Where,

$$\begin{aligned}
&E(\log(b(X_{obs}(i), X_{mis})) | \Theta_{M_i}) \\
&= \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \log(b(X_{obs}(i), X_{mis})) dX_{mis}
\end{aligned} \tag{2.20}$$

$$E(\tau(X_{mis}) | \Theta_{M_i}) = \int_{X_{mis}} f(X_{mis}|X_{obs}(i), \Theta_{M_i}) \tau(X_{mis}) dX_{mis} \tag{2.21}$$

At M-step of some  $t^{\text{th}}$  iteration, the next parameter  $\Theta^{(t+1)}$  is solution of the equation created by setting the first-order derivative of  $Q_1(\Theta'|\Theta)$  to be zero. The first-order derivative of  $Q_1(\Theta'|\Theta)$  is:

$$\begin{aligned}
\frac{\partial Q_1(\Theta'|\Theta)}{\partial \Theta'} &= \sum_{i=1}^N \left( E(\tau(X_{obs}(i), X_{mis}) | \Theta_{M_i}) \right)^T - N \log'(a(\Theta')) \\
&= \sum_{i=1}^N \{ \tau(X_{obs}(i)), E(\tau(X_{mis}) | \Theta_{M_i}) \}^T - N \log'(a(\Theta'))
\end{aligned}$$

By referring table 1.2, we have:

$$\log'(a(\Theta')) = (E(\tau(X)|\Theta'))^T = \int_X f(X|\Theta) (\tau(X))^T dX$$

Where,

$$\begin{aligned}
f(X|\Theta) &= f(X_{obs}, X_{mis}|\Theta) = b(X_{obs}, X_{mis}) * \exp((\Theta)^T \tau(X_{obs}, X_{mis}))/a(\Theta) \\
b(X) &= b(X_{obs}, X_{mis}) \\
\tau(X) &= \tau(X_{obs}, X_{mis}) = \{ \tau(X_{obs}), \tau(X_{mis}) \}
\end{aligned}$$

Thus, the next parameter  $\Theta^{(t+1)}$  is solution of the following equation:

$$\frac{\partial Q_1(\Theta'|\Theta)}{\partial \Theta'} = \sum_{i=1}^N \{ \tau(X_{obs}(i)), E(\tau(X_{mis}) | \Theta_{M_i}) \}^T - N (E(\tau(X)|\Theta'))^T = \mathbf{0}^T$$

This implies the next parameter  $\Theta^{(t+1)}$  is solution of the following equation:

$$E(\tau(X)|\Theta') = \frac{1}{N} \sum_{i=1}^N \{ \tau(X_{obs}(i)), E(\tau(X_{mis}) | \Theta_{M_i}) \}$$

As a result, at E-step of some  $t^{\text{th}}$  iteration, given current parameter  $\Theta^{(t)}$ , the sufficient statistic of  $X$  is calculated as follows:

$$\tau^{(t)} = \frac{1}{N} \sum_{i=1}^N \left\{ \tau(X_{obs}(i)), E\left(\tau(X_{mis}) \middle| \Theta_{M_i}^{(t)}\right) \right\} \quad (2.22)$$

Where,

$$\begin{aligned} \Theta_{M_i}^{(t)} &= u(\Theta^{(t)}, M_i) \\ E\left(\tau(X_{mis}) \middle| \Theta_{M_i}^{(t)}\right) &= \int_{X_{mis}} f(X_{mis} | X_{obs}(i), \Theta_{M_i}^{(t)}) \tau(X_{mis}) dX_{mis} \end{aligned}$$

Equation 2.22 is variant of equation 2.11 when  $f(X|\Theta)$  belongs to exponential family but how to compose  $\tau(X)$  from  $\tau(X_{obs})$  and  $\tau(X_{mis})$  is not determined exactly yet.

As a result, at M-step of some  $t^{th}$  iteration, given  $\tau^{(t)}$  and  $\Theta^{(t)}$ , the next parameter  $\Theta^{(t+1)}$  is a solution of the following equation:

$$E(\tau(X) | \Theta) = \tau^{(t)} \quad (2.23)$$

Moreover, at M-step of some  $t^{th}$  iteration, the next parameter  $\Phi^{(t+1)}$  is a maximizer of  $Q_2(\Phi | \Theta^{(t)})$  given  $\Theta^{(t)}$  as follows:

$$\Phi^{(t+1)} = \underset{\Phi}{\operatorname{argmin}} Q_2(\Phi | \Theta^{(t)}) \quad (2.24)$$

Where,

$$Q_2(\Phi | \Theta^{(t)}) = \sum_{i=1}^N \int_{X_{mis}} f(X_{mis} | X_{obs}(i), \Theta_{M_i}^{(t)}) \log(f(Z_i | X_{obs}(i), X_{mis}, \Phi)) dX_{mis} \quad (2.25)$$

How to maximize  $Q_2(\Phi | \Theta^{(t)})$  depends on distribution type of  $Z_i$  which is also formulation of the PDF  $f(Z | X_{obs}, X_{mis}, \Phi)$ . For some reasons, such as accelerating estimation speed or ignoring missingness variable  $Z$  then, the next parameter  $\Phi^{(t+1)}$  will not be estimated.

In general, the two steps of GEM algorithm for handling missing data at some  $t^{th}$  iteration are summarized in table 2.1 with assumption that the PDF of missing data  $f(X|\Theta)$  belongs to exponential family.

*E-step:*

Given current parameter  $\Theta^{(t)}$ , the sufficient statistic  $\tau^{(t)}$  is calculated according to equation 2.22.

$$\tau^{(t)} = \frac{1}{N} \sum_{i=1}^N \left\{ \tau(X_{obs}(i)), E\left(\tau(X_{mis}) \middle| \Theta_{M_i}^{(t)}\right) \right\}$$

Where,

$$\begin{aligned} \Theta_{M_i}^{(t)} &= u(\Theta^{(t)}, M_i) \\ E\left(\tau(X_{mis}) \middle| \Theta_{M_i}^{(t)}\right) &= \int_{X_{mis}} f(X_{mis} | X_{obs}(i), \Theta_{M_i}^{(t)}) \tau(X_{mis}) dX_{mis} \end{aligned}$$

*M-step:*

Given  $\tau^{(t)}$  and  $\Theta^{(t)}$ , the next parameter  $\Theta^{(t+1)}$  is a solution of equation 2.23.

$$E(\tau(X) | \Theta) = \tau^{(t)}$$

Given  $\Theta^{(t)}$ , the next parameter  $\Phi^{(t+1)}$  is a maximizer of  $Q_2(\Phi | \Theta^{(t)})$  according to equation 2.24.

$$\Phi^{(t+1)} = \underset{\Phi}{\operatorname{argmin}} Q_2(\Phi | \Theta^{(t)})$$

Where,

$$Q_2(\Phi | \Theta^{(t)}) = \sum_{i=1}^N \int_{X_{mis}} f(X_{mis} | X_{obs}(i), \Theta_{M_i}^{(t)}) \log(f(Z_i | X_{obs}(i), X_{mis}, \Phi)) dX_{mis}$$

**Table 2.1.** E-step and M-step of GEM algorithm for handling missing data given exponential PDF

GEM algorithm converges at some  $t^{\text{th}}$  iteration. At that time,  $\Theta^* = \Theta^{(t+1)} = \Theta^{(t)}$  and  $\Phi^* = \Phi^{(t+1)} = \Phi^{(t)}$  are optimal estimates. If missingness variable  $Z$  is ignored for some reasons, parameter  $\Phi$  is not estimated. Because  $X_{\text{mis}}$  is a part of  $X$  and  $f(X_{\text{mis}} | X_{\text{obs}}, \Theta_M)$  is derived directly from  $f(X|\Theta)$ , in practice we can stop GEM after its first iteration was done, which is reasonable enough to handle missing data.

An interesting application of handling missing data is to fill in or predict missing values. For instance, suppose the estimate resulted from GEM is  $\Theta^*$ , missing values represented by  $\tau(X_{\text{mis}})$  are fulfilled by expectation of  $\tau(X_{\text{mis}})$  as follows:

$$\tau(X_{\text{mis}}) = E(\tau(X_{\text{mis}}) | \Theta_M^*) \quad (2.26)$$

Where,

$$\Theta_M^* = u(\Theta^*, X_{\text{obs}})$$

Now we survey a popular case that sample  $\mathcal{X} = \{X_1, X_2, \dots, X_N\}$  whose  $X_i$  (s) are iid is MCAR data and  $f(X|\Theta)$  is multinormal PDF whereas missingness variable  $Z$  follows binomial distribution of  $n$  trials. Let  $X = \{X_{\text{obs}}, X_{\text{mis}}\}$  be random variable representing every  $X_i$ . Suppose dimension of  $X$  is  $n$ . Let  $Z$  be random variable representing every  $Z_i$ . According to equation 2.9, recall that

$$\begin{aligned} X_i &= \{X_{\text{obs}}(i), X_{\text{mis}}(i)\} = (x_{i1}, x_{i2}, \dots, x_{in})^T \\ X_{\text{mis}}(i) &= (x_{im_1}, x_{im_2}, \dots, x_{im_{|M_i|}})^T \\ X_{\text{obs}}(i) &= (x_{i\bar{m}_1}, x_{i\bar{m}_2}, \dots, x_{i\bar{m}_{|\bar{M}_i|}})^T \\ M_i &= \{m_{i1}, m_{i2}, \dots, m_{i|M_i|}\} \\ \bar{M}_i &= \{\bar{m}_{i1}, \bar{m}_{i2}, \dots, \bar{m}_{i|\bar{M}_i|}\} \\ Z_i &= (z_{i1}, z_{i2}, \dots, z_{in})^T \end{aligned}$$

The PDF of  $X$  is:

$$f(X|\Theta) = f(X_{\text{obs}}, X_{\text{mis}}|\Theta) = (2\pi)^{-\frac{n}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} (X - \mu)^T \Sigma^{-1} (X - \mu)\right) \quad (2.27)$$

Therefore,

$$f(X_i|\Theta) = f(X_{\text{obs}}(i), X_{\text{mis}}(i)|\Theta) = (2\pi)^{-\frac{n}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} (X_i - \mu)^T \Sigma^{-1} (X_i - \mu)\right)$$

The PDF of  $Z$  is:

$$f(Z|\Phi) = p^{c(Z)} (1-p)^{n-c(Z)} \quad (2.28)$$

Therefore,

$$f(Z_i|\Phi) = p^{c(Z_i)} (1-p)^{n-c(Z_i)}$$

Where  $\Theta = (\mu, \Sigma)^T$  and  $\Phi = p$ .

$$\begin{aligned} \mu &= (\mu_1, \mu_2, \dots, \mu_n)^T \\ \Sigma &= \begin{pmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1n} \\ \sigma_{21} & \sigma_{22} & \cdots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \cdots & \sigma_{nn} \end{pmatrix} \end{aligned} \quad (2.29)$$

Suppose the probability of missingness at every partial random variable  $x_j$  is  $p$  and it is independent from  $X_{\text{obs}}$  and  $X_{\text{mis}}$ . The quantity  $c(Z)$  is the number of  $z_j$  (s) in  $Z$  that equal 1. For example, if  $Z = (1, 0, 1, 0)^T$  then,  $c(Z) = 2$ . The most important task here is to define equation 2.11 and equation 2.15 in order to compose  $\tau(X)$  from  $\tau(X_{\text{obs}})$ ,  $\tau(X_{\text{mis}})$  and to extract  $\Theta_M$  from  $\Theta$  when  $f(X|\Theta)$  distributes normally.

The conditional PDF of  $X_{\text{mis}}$  given  $X_{\text{obs}}$  is also multinormal PDF.

$$\begin{aligned}
f(X_{mis}|\Theta_M) &= f(X_{mis}|X_{obs}, \Theta_M) = f(X_{mis}|X_{obs}, \Theta) \\
&= (2\pi)^{-\frac{|M|}{2}} |\Sigma_M|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(X_{mis} - \mu_M)^T \Sigma_M^{-1}(X_{mis} - \mu_M)\right)
\end{aligned} \tag{2.30}$$

Therefore,

$$\begin{aligned}
f(X_{mis}(i)|\Theta_{M_i}) &= f(X_{mis}(i)|X_{obs}(i), \Theta_{M_i}) = f(X_{mis}(i)|X_{obs}(i), \Theta) \\
&= (2\pi)^{-\frac{|M_i|}{2}} |\Sigma_{M_i}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(X_{mis}(i) - \mu_{M_i})^T \Sigma_{M_i}^{-1}(X_{mis}(i) - \mu_{M_i})\right)
\end{aligned}$$

Where  $\Theta_{M_i} = (\mu_{M_i}, \Sigma_{M_i})^T$ . We denote

$$f(X_{mis}(i)|\Theta_{M_i}) = f(X_{mis}(i)|X_{obs}(i), \Theta_{M_i})$$

Because  $f(X_{mis}(i)|X_{obs}(i), \Theta_{M_i})$  only depends on  $\Theta_{M_i}$  within normal PDF whereas  $\Theta_{M_i}$  depends on  $X_{obs}(i)$ . Determining the function  $\Theta_{M_i} = u(\Theta, X_{obs}(i))$  is now necessary to extract the parameter  $\Theta_{M_i}$  from  $\Theta$  given  $X_{obs}(i)$  when  $f(X_i|\Theta)$  is normal distribution.

Let  $\Theta_{mis} = (\mu_{mis}, \Sigma_{mis})^T$  be parameter of marginal PDF of  $X_{mis}$ , we have:

$$\begin{aligned}
f(X_{mis}|\Theta_{mis}) &= (2\pi)^{-\frac{|M|}{2}} |\Sigma_{mis}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(X_{mis} - \mu_{mis})^T (\Sigma_{mis})^{-1}(X_{mis} \right. \\
&\quad \left. - \mu_{mis})\right)
\end{aligned} \tag{2.31}$$

Therefore,

$$\begin{aligned}
f(X_{mis}(i)|\Theta_{mis}(i)) \\
&= (2\pi)^{-\frac{|M_i|}{2}} |\Sigma_{mis}(i)|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(X_{mis}(i) - \mu_{mis}(i))^T (\Sigma_{mis}(i))^{-1}(X_{mis}(i) \right. \\
&\quad \left. - \mu_{mis}(i))\right)
\end{aligned}$$

Where,

$$\begin{aligned}
\mu_{mis}(i) &= (\mu_{m_{i1}}, \mu_{m_{i2}}, \dots, \mu_{m_{i|M_i|}})^T \\
\Sigma_{mis}(i) &= \begin{pmatrix} \sigma_{m_{i1}m_{i1}} & \sigma_{m_{i1}m_{i2}} & \cdots & \sigma_{m_{i1}m_{i|M_i|}} \\ \sigma_{m_{i2}m_{i1}} & \sigma_{m_{i2}m_{i2}} & \cdots & \sigma_{m_{i2}m_{i|M_i|}} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{m_{i|M_i|}m_{i1}} & \sigma_{m_{i|M_i|}m_{i2}} & \cdots & \sigma_{m_{i|M_i|}m_{i|M_i|}} \end{pmatrix}
\end{aligned} \tag{2.32}$$

Obviously,  $\Theta_{mis}(i)$  is extracted from  $\Theta$  given indicator  $M_i$ .

Let  $\Theta_{obs} = (\mu_{obs}, \Sigma_{obs})^T$  be parameter of marginal PDF of  $X_{obs}$ , we have:

$$\begin{aligned}
f(X_{obs}|\Theta_{obs}) &= (2\pi)^{-\frac{|\bar{M}|}{2}} |\Sigma_{obs}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(X_{obs} - \mu_{obs})^T (\Sigma_{obs})^{-1}(X_{obs} \right. \\
&\quad \left. - \mu_{obs})\right)
\end{aligned} \tag{2.33}$$

Therefore,

$$\begin{aligned}
f(X_{obs}(i)|\Theta_{obs}(i)) \\
&= (2\pi)^{-\frac{|\bar{M}_i|}{2}} |\Sigma_{obs}(i)|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(X_{obs}(i) - \mu_{obs}(i))^T (\Sigma_{obs}(i))^{-1}(X_{obs}(i) \right. \\
&\quad \left. - \mu_{obs}(i))\right)
\end{aligned}$$

Where,

$$\begin{aligned}\mu_{obs}(i) &= (\mu_{\bar{m}_{i1}}, \mu_{\bar{m}_{i2}}, \dots, \mu_{\bar{m}_{i|\bar{M}_i|}})^T \\ \Sigma_{obs}(i) &= \begin{pmatrix} \sigma_{\bar{m}_{i1}\bar{m}_{i1}} & \sigma_{\bar{m}_{i1}\bar{m}_{i2}} & \cdots & \sigma_{\bar{m}_{i1}\bar{m}_{i|\bar{M}_i|}} \\ \sigma_{\bar{m}_{i2}\bar{m}_{i1}} & \sigma_{\bar{m}_{i2}\bar{m}_{i2}} & \cdots & \sigma_{\bar{m}_{i2}\bar{m}_{i|\bar{M}_i|}} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{\bar{m}_{i|\bar{M}_i|}\bar{m}_{i1}} & \sigma_{\bar{m}_{i|\bar{M}_i|}\bar{m}_{i2}} & \cdots & \sigma_{\bar{m}_{i|\bar{M}_i|}\bar{m}_{i|\bar{M}_i|}} \end{pmatrix}\end{aligned}\quad (2.34)$$

Obviously,  $\Theta_{obs}(i)$  is extracted from  $\Theta$  given indicator  $\bar{M}_i$  or  $M_i$ . We have:

$$\begin{aligned}f(X_{mis}(i)|\Theta_{mis}(i)) &= \int_{X_{obs}(i)} f(X_{obs}(i), X_{mis}(i)|\Theta) dX_{obs}(i) \\ f(X_{obs}(i)|\Theta_{obs}(i)) &= \int_{X_{mis}(i)} f(X_{obs}(i), X_{mis}(i)|\Theta) dX_{mis}(i) \\ f(X_{mis}(i)|\Theta_{M_i}) &= f(X_{mis}(i)|X_{obs}(i), \Theta) = \frac{f(X_{obs}(i), X_{mis}(i)|\Theta)}{f(X_{obs}(i)|\Theta_{obs}(i))}\end{aligned}$$

Therefore, it is easy to form the parameter  $\Theta_{M_i} = (\mu_{M_i}, \Sigma_{M_i})^T$  from  $\Theta_{mis}(i) = (\mu_{mis}(i), \Sigma_{mis}(i))^T$  and  $\Theta_{obs}(i) = (\mu_{obs}(i), \Sigma_{obs}(i))^T$  as follows (Hardle & Simar, 2013, pp. 156-157):

$$\begin{aligned}\Theta_{M_i} &= u(\Theta, X_{obs}(i)) \\ &= \begin{cases} \mu_{M_i} = \mu_{mis}(i) + (V_{obs}^{mis}(i))(\Sigma_{obs}(i))^{-1}(X_{obs}(i) - \mu_{obs}(i)) \\ \Sigma_{M_i} = \Sigma_{mis}(i) - (V_{obs}^{mis}(i))(\Sigma_{obs}(i))^{-1}(V_{obs}^{obs}) \end{cases}\end{aligned}\quad (2.35)$$

Where from  $\Theta_{mis}(i) = (\mu_{mis}(i), \Sigma_{mis}(i))^T$  and  $\Theta_{obs}(i) = (\mu_{obs}(i), \Sigma_{obs}(i))^T$  are specified by equation 2.32 and equation 2.34. Moreover the  $k \times l$  matrix  $V_{obs}^{mis}(i)$  which implies correlation between  $X_{mis}$  and  $X_{obs}$  is defined as follows:

$$V_{obs}^{mis}(i) = \begin{pmatrix} \sigma_{m_{i1}\bar{m}_{i1}} & \sigma_{m_{i1}\bar{m}_{i2}} & \cdots & \sigma_{m_{i1}\bar{m}_{i|\bar{M}_i|}} \\ \sigma_{m_{i2}\bar{m}_{i1}} & \sigma_{m_{i2}\bar{m}_{i2}} & \cdots & \sigma_{m_{i2}\bar{m}_{i|\bar{M}_i|}} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{m_{i|\bar{M}_i|}\bar{m}_{i1}} & \sigma_{m_{i|\bar{M}_i|}\bar{m}_{i2}} & \cdots & \sigma_{m_{i|\bar{M}_i|}\bar{m}_{i|\bar{M}_i|}} \end{pmatrix}\quad (2.36)$$

The  $k \times l$  matrix  $V_{mis}^{obs}(i)$  which implies correlation between  $X_{obs}$  and  $X_{mis}$  is defined as follows:

$$V_{mis}^{obs}(i) = \begin{pmatrix} \sigma_{\bar{m}_{i1}m_{i1}} & \sigma_{\bar{m}_{i1}m_{i2}} & \cdots & \sigma_{\bar{m}_{i1}m_{i|\bar{M}_i|}} \\ \sigma_{\bar{m}_{i2}m_{i1}} & \sigma_{\bar{m}_{i2}m_{i2}} & \cdots & \sigma_{\bar{m}_{i2}m_{i|\bar{M}_i|}} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{\bar{m}_{i|\bar{M}_i|}m_{i1}} & \sigma_{\bar{m}_{i|\bar{M}_i|}m_{i2}} & \cdots & \sigma_{\bar{m}_{i|\bar{M}_i|}m_{i|\bar{M}_i|}} \end{pmatrix}\quad (2.37)$$

Therefore, equation 2.35 to extract  $\Theta_{M_i}$  from  $\Theta$  given  $X_{obs}(i)$  is an instance of equation 2.15. For convenience let,

$$\begin{aligned}\mu_{M_i} &= (\mu_{M_i}(m_{i1}), \mu_{M_i}(m_{i2}), \dots, \mu_{M_i}(m_{i|\bar{M}_i|}))^T \\ \Sigma_{M_i} &= \begin{pmatrix} \Sigma_{M_i}(m_{i1}, m_{i1}) & \Sigma_{M_i}(m_{i1}, m_{i2}) & \cdots & \Sigma_{M_i}(m_{i1}, m_{i|\bar{M}_i|}) \\ \Sigma_{M_i}(m_{i2}, m_{i1}) & \Sigma_{M_i}(m_{i2}, m_{i2}) & \cdots & \Sigma_{M_i}(m_{i2}, m_{i|\bar{M}_i|}) \\ \vdots & \vdots & \ddots & \vdots \\ \Sigma_{M_i}(m_{i|\bar{M}_i|}, m_{i1}) & \Sigma_{M_i}(m_{i|\bar{M}_i|}, m_{i2}) & \cdots & \Sigma_{M_i}(m_{i|\bar{M}_i|}, m_{i|\bar{M}_i|}) \end{pmatrix}\end{aligned}\quad (2.38)$$

Equation 2.38 is result of equation 2.35. Given  $X_{mis}(i) = (x_{m_{i1}}, x_{m_{i2}}, \dots, x_{m_{i|\bar{M}_i|}})^T$  then,  $\mu_{M_i}(m_{ij})$  is estimated partial mean of  $x_{m_{ij}}$  and  $\Sigma_{M_i}(m_{iu}, m_{iv})$  is estimated partial covariance of  $x_{m_{iu}}$  and  $x_{m_{iv}}$  given the conditional PDF  $f(X_{mis} | \Theta_{M_i})$ .

At E-step of some  $t^{\text{th}}$  iteration, given current parameter  $\Theta^{(t)}$ , the sufficient statistic of  $X$  is calculated according to equation 2.22. Let,

$$\tau^{(t)} = (\tau_1^{(t)}, \tau_2^{(t)})^T = \frac{1}{N} \sum_{i=1}^N \left\{ \tau(X_{obs}(i)), E(\tau(X_{mis}) | \Theta_{M_i}^{(t)}) \right\}$$

It is necessary to calculate the sufficient with normal PDF  $f(X_i | \Theta)$ , which means that we need to define what  $\tau_1^{(t)}$  and  $\tau_2^{(t)}$  are. The sufficient statistic of  $X_{obs}(i)$  is:

$$\tau(X_{obs}(i)) = (X_{obs}(i), X_{obs}(i)(X_{obs}(i))^T)^T$$

The sufficient statistic of  $X_{mis}(i)$  is:

$$\tau(X_{mis}(i)) = (X_{mis}(i), X_{mis}(i)(X_{mis}(i))^T)^T$$

We also have:

$$E(\tau(X_{mis}) | \Theta_{M_i}^{(t)}) = \int_{X_{mis}} f(X_{mis} | \Theta_{M_i}^{(t)}) \tau(X_{mis}) dX_{mis} = \begin{pmatrix} \mu_{M_i}^{(t)} \\ \Sigma_{M_i}^{(t)} + \mu_{M_i}^{(t)}(\mu_{M_i}^{(t)})^T \end{pmatrix}$$

Due to

$$E(X_{mis}(i)(X_{mis}(i))^T | \Theta_{M_i}^{(t)}) = \Sigma_{M_i}^{(t)} + \mu_{M_i}^{(t)}(\mu_{M_i}^{(t)})^T$$

Where  $\mu_{M_i}^{(t)}$  and  $\Sigma_{M_i}^{(t)}$  are  $\mu_{M_i}$  and  $\Sigma_{M_i}$  at current iteration, respectively. By referring to equation 2.38, we have

$$\mu_{M_i}^{(t)} = (\mu_{M_i}^{(t)}(m_{i1}), \mu_{M_i}^{(t)}(m_{i2}), \dots, \mu_{M_i}^{(t)}(m_{i|M_i|}))^T$$

And

$$\Sigma_{M_i}^{(t)} + \mu_{M_i}^{(t)}(\mu_{M_i}^{(t)})^T = \begin{pmatrix} \tilde{\sigma}_{11}^{(t)}(i) & \tilde{\sigma}_{12}^{(t)}(i) & \dots & \tilde{\sigma}_{1|M_i|}^{(t)}(i) \\ \tilde{\sigma}_{21}^{(t)}(i) & \tilde{\sigma}_{22}^{(t)}(i) & \dots & \tilde{\sigma}_{2|M_i|}^{(t)}(i) \\ \vdots & \vdots & \ddots & \vdots \\ \tilde{\sigma}_{|M_i|1}^{(t)}(i) & \tilde{\sigma}_{|M_i|2}^{(t)}(i) & \dots & \tilde{\sigma}_{|M_i||M_i|}^{(t)}(i) \end{pmatrix}$$

Where,

$$\tilde{\sigma}_{uv}^{(t)}(i) = \Sigma_{M_i}^{(t)}(m_{iu}, m_{iv}) + \mu_{M_i}^{(t)}(m_{iu})\mu_{M_i}^{(t)}(m_{iv})$$

Therefore,  $\tau_1^{(t)}$  is vector and  $\tau_2^{(t)}$  is matrix and then, the sufficient statistic of  $X$  at E-step of some  $t^{\text{th}}$  iteration, given current parameter  $\Theta^{(t)}$  is defined as follows:

$$\begin{aligned} \tau^{(t)} &= (\tau_1^{(t)}, \tau_2^{(t)})^T \\ \tau_1^{(t)} &= (\bar{x}_1^{(t)}, \bar{x}_2^{(t)}, \dots, \bar{x}_n^{(t)})^T \\ \tau_2^{(t)} &= \begin{pmatrix} s_{11}^{(t)} & s_{12}^{(t)} & \dots & s_{1n}^{(t)} \\ s_{21}^{(t)} & s_{22}^{(t)} & \dots & s_{2n}^{(t)} \\ \vdots & \vdots & \ddots & \vdots \\ s_{n1}^{(t)} & s_{n2}^{(t)} & \dots & s_{nn}^{(t)} \end{pmatrix} \end{aligned} \quad (2.39)$$

Each  $\bar{x}_j^{(t)}$  is calculated as follows:

$$\bar{x}_j^{(t)} = \frac{1}{N} \sum_{i=1}^N \begin{cases} x_{ij} & \text{if } j \notin M_i \\ \mu_{M_i}^{(t)}(j) & \text{if } j \in M_i \end{cases} \quad (2.40)$$

Please see equation 2.35 and equation 2.38 to know  $\mu_{M_i}^{(t)}(j)$ . Each  $s_{uv}^{(t)}$  is calculated as follows:



$$s_{uv}^{(t)} = s_{vu}^{(t)} = \frac{1}{N} \sum_{i=1}^N \begin{cases} x_{iu}x_{iv} & \text{if } u \notin M_i \text{ and } v \notin M_i \\ x_{iu}\mu_{M_i}^{(t)}(m_{iv}) & \text{if } u \notin M_i \text{ and } v \in M_i \\ \mu_{M_i}^{(t)}(m_{iu})x_{iv} & \text{if } u \in M_i \text{ and } v \notin M_i \\ \Sigma_{M_i}^{(t)}(m_{iu}, m_{iv}) + \mu_{M_i}^{(t)}(m_{iu})\mu_{M_i}^{(t)}(m_{iv}) & \text{if } u \in M_i \text{ and } v \in M_i \end{cases} \quad (2.41)$$

Equation 2.39 is an instance of equation 2.11, which compose  $\tau(X)$  from  $\tau(X_{obs})$  and  $\tau(X_{mis})$  when  $f(X|\Theta)$  distributes normally. Following is the proof of equation 2.41.

If  $u \notin M_i$  and  $v \notin M_i$  then, the partial statistic  $x_{iu}x_{iv}$  is kept intact because  $x_{iu}$  and  $x_{iv}$  are in  $X_{obs}$  are constant with regard to  $f(X_{mis} | \Theta_{M_i}^{(t)})$ . If  $u \notin M_i$  and  $v \in M_i$  then, the partial statistic  $x_{iu}x_{iv}$  is replaced by the expectation  $E(x_{iu}x_{iv} | \Theta_{M_i}^{(t)})$  as follows:

$$\begin{aligned} E(x_{iu}x_{iv} | \Theta_{M_i}^{(t)}) &= \int_{X_{mis}} f(X_{mis} | \Theta_{M_i}^{(t)}) x_{iu}x_{iv} dX_{mis} = x_{iu} \int_{X_{mis}} f(X_{mis} | \Theta_{M_i}^{(t)}) x_{iv} dX_{mis} \\ &= x_{iu}\mu_{M_i}^{(t)}(m_{iv}) \end{aligned}$$

If  $u \in M_i$  and  $v \notin M_i$  then, the partial statistic  $x_{iu}x_{iv}$  is replaced by the expectation  $E(x_{iu}x_{iv} | \Theta_{M_i}^{(t)})$  as follows:

$$\begin{aligned} E(x_{iu}x_{iv} | \Theta_{M_i}^{(t)}) &= \int_{X_{mis}} f(X_{mis} | \Theta_{M_i}^{(t)}) x_{iu}x_{iv} dX_{mis} = x_{iv} \int_{X_{mis}} f(X_{mis} | \Theta_{M_i}^{(t)}) x_{iu} dX_{mis} \\ &= \mu_{M_i}^{(t)}(m_{iu})x_{iv} \end{aligned}$$

If  $u \in M_i$  and  $v \in M_i$  then, the partial statistic  $x_{iu}x_{iv}$  is replaced by the expectation  $E(x_{iu}x_{iv} | \Theta_{M_i}^{(t)})$  as follows:

$$E(x_{iu}x_{iv} | \Theta_{M_i}^{(t)}) = \int_{X_{mis}} f(X_{mis} | \Theta_{M_i}^{(t)}) x_{iu}x_{iv} dX_{mis} = \Sigma_{M_i}^{(t)}(m_{iu}, m_{iv}) + \mu_{M_i}^{(t)}(m_{iu})\mu_{M_i}^{(t)}(m_{iv}) \blacksquare$$

At M-step of some  $t^{\text{th}}$  iteration, given  $\tau^{(t)}$  and  $\Theta^{(t)}$ , the next parameter  $\Theta^{(t+1)} = (\mu^{(t+1)}, \Sigma^{(t+1)})^T$  is a solution of equation 2.23.

$$E(\tau(X) | \Theta) = \tau^{(t)}$$

Due to

$$E(\tau(X) | \Theta) = \begin{pmatrix} \mu \\ \Sigma \end{pmatrix}$$

Equation 2.23 becomes:

$$\begin{cases} \mu = \tau_1^{(t)} \\ \Sigma = \tau_2^{(t)} \end{cases}$$

Which means that

$$\begin{cases} \mu_j^{(t+1)} = \bar{x}_j^{(t)} \\ \sigma_{uv}^{(t+1)} = \sigma_{vu}^{(t+1)} = s_{uv}^{(t)} - \bar{x}_u^{(t)} \bar{x}_v^{(t)} \end{cases} \forall j, u, v \quad (2.42)$$

Please see equation 2.40 and equation 2.41 to know  $\bar{x}_j^{(t)}$  and  $s_{uv}^{(t)}$ .

Moreover, at M-step of some  $t^{\text{th}}$  iteration, the next parameter  $\Phi^{(t+1)} = p^{(t+1)}$  is a maximizer of  $Q_2(\Phi | \Theta^{(t)})$  given  $\Theta^{(t)}$  according to equation 2.24.

$$\Phi^{(t+1)} = \underset{\Phi}{\operatorname{argmin}} Q_2(\Phi | \Theta^{(t)})$$

Because the PDF of  $Z_i$  is:

$$f(Z_i | \Phi) = p^{c(Z_i)}(1-p)^{n-c(Z_i)}$$

The  $Q_2(\Phi | \Theta^{(t)})$  becomes:

$$\begin{aligned} Q_2(\Phi | \Theta^{(t)}) &= \sum_{i=1}^N \int_{X_{mis}} f(X_{mis} | \Theta_{M_i}^{(t)}) \log(f(Z_i | X_{obs}(i), X_{mis}, \Phi)) dX_{mis} \\ &= \sum_{i=1}^N \int_{X_{mis}} f(X_{mis} | \Theta_{M_i}^{(t)}) \log(f(Z_i | \Phi)) dX_{mis} \\ &= \sum_{i=1}^N \log(f(Z_i | \Phi)) \int_{X_{mis}} f(X_{mis} | \Theta_{M_i}^{(t)}) dX_{mis} \\ &= \sum_{i=1}^N \log(f(Z_i | \Phi)) = \sum_{i=1}^N \log(p^{c(Z_i)}(1-p)^{n-c(Z_i)}) \\ &= \sum_{i=1}^N (c(Z_i) \log(p) + (n - c(Z_i)) \log(1-p)) \end{aligned}$$

The next parameter  $\Phi^{(t+1)} = p^{(t+1)}$  is solution of the equation created by setting the first-order derivative of  $Q_2(\Phi | \Theta^{(t)})$  to be zero, which means that:

$$\frac{\partial Q_2(\Phi | \Theta^{(t)})}{\partial p} = \sum_{i=1}^N \left( \frac{c(Z_i)}{p} - \frac{n - c(Z_i)}{1-p} \right) = \frac{1}{p(1-p)} \left( \left( \sum_{i=1}^N c(Z_i) \right) - npN \right) = 0$$

It is easy to deduce that the next parameter  $p^{(t+1)}$  is:

$$p^{(t+1)} = \frac{\sum_{i=1}^N c(Z_i)}{nN} \quad (2.43)$$

In general, given sample  $\mathcal{X} = \{X_1, X_2, \dots, X_N\}$  whose  $X_i$  (s) are iid is MCAR data and  $f(X|\Theta)$  is multinormal PDF whereas missingness variable  $Z$  follows binomial distribution of  $n$  trials, GEM for handling missing data is summarized in table 2.2.

*E-step:*

Given current parameter  $\Theta^{(t)} = (\mu^{(t)}, \Sigma^{(t)})^T$ , the sufficient statistic  $\tau^{(t)}$  is calculated according to equation 2.39, equation 2.40, and equation 2.41.

$$\begin{aligned} \tau^{(t)} &= (\tau_1^{(t)}, \tau_2^{(t)})^T \\ \tau_1^{(t)} &= (\bar{x}_1^{(t)}, \bar{x}_2^{(t)}, \dots, \bar{x}_n^{(t)})^T \\ \tau_2^{(t)} &= \begin{pmatrix} s_{11}^{(t)} & s_{12}^{(t)} & \dots & s_{1n}^{(t)} \\ s_{21}^{(t)} & s_{22}^{(t)} & \dots & s_{2n}^{(t)} \\ \vdots & \vdots & \ddots & \vdots \\ s_{n1}^{(t)} & s_{n2}^{(t)} & \dots & s_{nn}^{(t)} \end{pmatrix} \end{aligned}$$

$$s_{uv}^{(t)} = s_{vu}^{(t)} = \frac{1}{N} \sum_{i=1}^N \begin{cases} x_{iu}x_{iv} & \text{if } u \notin M_i \text{ and } v \notin M_i \\ x_{iu}\mu_{M_i}^{(t)}(m_{iv}) & \text{if } u \notin M_i \text{ and } v \in M_i \\ \mu_{M_i}^{(t)}(m_{iu})x_{iv} & \text{if } u \in M_i \text{ and } v \notin M_i \\ \Sigma_{M_i}^{(t)}(m_{iu}, m_{iv}) + \mu_{M_i}^{(t)}(m_{iu})\mu_{M_i}^{(t)}(m_{iv}) & \text{if } u \in M_i \text{ and } v \in M_i \end{cases}$$

Where  $\mu_{M_i}$  and  $\Sigma_{M_i}$  are specified in equation 2.35 and equation 2.38.

*M-step:*

Given  $\tau^{(t)}$  and  $\Theta^{(t)}$ , the next parameter  $\Theta^{(t+1)} = (\mu^{(t+1)}, \Sigma^{(t+1)})^T$  is specified by equation 2.42.

$$\begin{cases} \mu_j^{(t+1)} = \bar{x}_j^{(t)} \\ \sigma_{uv}^{(t+1)} = \sigma_{vu}^{(t+1)} = s_{uv}^{(t)} - x_u^{(t)}x_v^{(t)} \quad \forall j, u, v \end{cases}$$

Given  $\Theta^{(t)}$ , the next parameter  $\Phi^{(t+1)} = p^{(t+1)}$  is specified by equation 2.43.

$$p^{(t+1)} = \frac{\sum_{i=1}^N c(Z_i)}{nN}$$

Where  $c(Z_i)$  is the number of  $z_{ij}$  (s) in  $Z_i$  that equal 1.

**Table 2.2.** E-step and M-step of GEM algorithm for handling missing data given normal PDF. As aforementioned, an interesting application of handling missing data is to fill in or predict missing values. For instance, suppose the estimate resulted from GEM is  $\Theta^* = (\mu^*, \Sigma^*)^T$ , missing part  $X_{mis} = (x_{m_1}, x_{m_2}, \dots, x_{m_{|M_i|}})^T$  is replaced by  $\mu_M^*$  as follows:

$$x_{m_j} = \mu_M^*(m_j), \forall m_j \in M \quad (2.44)$$

Note,  $\mu_M^*$  which is extracted from  $\mu^*$  is estimated mean of the conditional PDF of  $X_{mis}$  (given  $X_{obs}$ ) according to equation 2.35. Moreover,  $\mu_M^*(m_j)$  is estimated partial mean of  $x_{m_j}$  given the conditional PDF  $f(X_{mis} | \Theta_M^*)$ , please see equation 2.38 for more details about  $\mu_M^*$ . As aforementioned, in practice we can stop GEM after its first iteration was done, which is reasonable enough to handle missing data.

Now we survey another interesting case that sample  $\mathcal{X} = \{X_1, X_2, \dots, X_N\}$  whose  $X_i$  (s) are iid is MCAR data and  $f(X|\Theta)$  is multinomial PDF of  $K$  trials. We ignore missingness variable  $Z$  here because it is included in the case of multinormal PDF. Let  $X = \{X_{obs}, X_{mis}\}$  be random variable representing every  $X_i$ . Suppose dimension of  $X$  is  $n$ . According to equation 2.9, recall that

$$X_i = \{X_{obs}(i), X_{mis}(i)\} = (x_{i1}, x_{i2}, \dots, x_{in})^T$$

$$X_{mis}(i) = (x_{im_1}, x_{im_2}, \dots, x_{im_{|M_i|}})^T$$

$$X_{obs}(i) = (x_{i\bar{m}_1}, x_{i\bar{m}_2}, \dots, x_{i\bar{m}_{|\bar{M}_i|}})^T$$

$$M_i = \{m_{i1}, m_{i2}, \dots, m_{i|M_i|}\}$$

$$\bar{M}_i = \{\bar{m}_{i1}, \bar{m}_{i2}, \dots, \bar{m}_{i|\bar{M}_i|}\}$$

The PDF of  $X$  is:

$$f(X|\Theta) = f(X_{obs}, X_{mis}|\Theta) = \frac{K!}{\prod_{j=1}^n (x_j!)} \prod_{j=1}^n p_j^{x_j} \quad (2.45)$$

Where  $x_j$  are integers and  $\Theta = (p_1, p_2, \dots, p_n)^T$  is the set of probabilities such that

$$\begin{aligned} \sum_{j=1}^n p_j &= 1 \\ \sum_{j=1}^n x_j &= K \\ x_j &\in \{0, 1, \dots, K\} \end{aligned}$$

Note,  $x_j$  is the number of trials generating nominal value  $j$ . Therefore,

$$f(X_i|\Theta) = f(X_{obs}(i), X_{mis}(i)|\Theta) = \frac{K!}{\prod_{j=1}^n (x_{ij}!)} \prod_{j=1}^n p_j^{x_{ij}}$$

Where,

$$\begin{aligned} \sum_{j=1}^n x_{ij} &= K \\ x_{ij} &\in \{0, 1, \dots, K\} \end{aligned}$$

The most important task here is to define equation 2.11 and equation 2.15 in order to compose  $\tau(X)$  from  $\tau(X_{obs})$ ,  $\tau(X_{mis})$  and to extract  $\Theta_M$  from  $\Theta$  when  $f(X|\Theta)$  is multinomial PDF.

Let  $\Theta_{mis}$  be parameter of marginal PDF of  $X_{mis}$ , we have:

$$f(X_{mis}|\Theta_{mis}) = \frac{K_{mis}!}{\prod_{m_j \in M} (x_{m_j}!)} \prod_{j=1}^{|M|} \left( \frac{p_{m_j}}{P_{mis}} \right)^{x_{m_j}} \quad (2.46)$$

Therefore,

$$f(X_{mis}(i)|\Theta_{mis}(i)) = \frac{K_{mis}(i)!}{\prod_{m_j \in M_i} (x_{im_j}!)} \prod_{j=1}^{|M_i|} \left( \frac{p_{m_{ij}}}{P_{mis}(i)} \right)^{x_{im_j}}$$

Where,

$$\begin{aligned} \Theta_{mis}(i) &= \left( \frac{p_{m_{i1}}}{P_{mis}(i)}, \frac{p_{m_{i2}}}{P_{mis}(i)}, \dots, \frac{p_{m_{i|M_i|}}}{P_{mis}(i)} \right)^T \\ P_{mis}(i) &= \sum_{j=1}^{|M_i|} p_{m_{ij}} \\ K_{mis}(i) &= \sum_{j=1}^{|M_i|} x_{m_{ij}} \end{aligned} \quad (2.47)$$

Obviously,  $\Theta_{mis}(i)$  is extracted from  $\Theta$  given indicator  $M_i$ .

Let  $\Theta_{obs}$  be parameter of marginal PDF of  $X_{obs}$ , we have:

$$f(X_{obs}|\Theta_{obs}) = \frac{K_{obs}!}{\prod_{\bar{m}_j \in \bar{M}} (x_{\bar{m}_j}!)} \prod_{j=1}^{|\bar{M}|} \left( \frac{p_{\bar{m}_j}}{P_{obs}} \right)^{x_{\bar{m}_j}} \quad (2.48)$$

Therefore,

$$f(X_{obs}(i)|\Theta_{obs}(i)) = \frac{K_{obs}(i)!}{\prod_{\bar{m}_j \in \bar{M}_i} (x_{i\bar{m}_j}!) } \prod_{j=1}^{|\bar{M}_i|} \left( \frac{p_{\bar{m}_{ij}}}{P_{obs}(i)} \right)^{x_{i\bar{m}_j}}$$

Where,

$$\begin{aligned} \Theta_{obs}(i) &= \left( \frac{p_{\bar{m}_{i1}}}{P_{obs}(i)}, \frac{p_{\bar{m}_{i2}}}{P_{obs}(i)}, \dots, \frac{p_{\bar{m}_{i|\bar{M}_i|}}}{P_{obs}(i)} \right)^T \\ P_{obs}(i) &= \sum_{j=1}^{|\bar{M}_i|} p_{\bar{m}_{ij}} \\ K_{obs}(i) &= \sum_{j=1}^{|\bar{M}_i|} x_{i\bar{m}_{ij}} \end{aligned} \quad (2.49)$$

Obviously,  $\Theta_{obs}(i)$  is extracted from  $\Theta$  given indicator  $\bar{M}_i$  or  $M_i$ .

The conditional PDF of  $X_{mis}$  given  $X_{obs}$  is calculated based on the PDF of  $X$  and the marginal PDF of  $X_{obs}$  as follows:

$$\begin{aligned} f(X_{mis}|\Theta_M) &= f(X_{mis}|X_{obs}, \Theta) = \frac{f(X_{obs}, X_{mis}|\Theta)}{f(X_{obs}|\Theta_{obs})} \\ &= \frac{\frac{K!}{\prod_{j=1}^n (x_j!)} \prod_{j=1}^n p_j^{x_j}}{\frac{K_{obs}!}{\prod_{j=1}^{|\bar{M}|} x_{\bar{m}_j!}} \prod_{j=1}^{|\bar{M}|} \left( \frac{p_{\bar{m}_j}}{P_{obs}} \right)^{x_{\bar{m}_j}}} \\ &= \frac{K!}{K_{obs}!} \frac{\prod_{j=1}^{|\bar{M}|} x_{\bar{m}_j!}}{\prod_{j=1}^n (x_j!)} * \frac{\prod_{j=1}^n p_j^{x_j}}{\prod_{j=1}^{|\bar{M}|} \left( \frac{p_{\bar{m}_j}}{P_{obs}} \right)^{x_{\bar{m}_j}}} \\ &= \frac{K!}{K_{obs}!} \frac{\prod_{j=1}^{|\bar{M}|} (x_{m_j}!)}{\prod_{j=1}^{|\bar{M}|} (x_{m_j}!)} * \left( \prod_{j=1}^{|\bar{M}|} p_{m_j}^{x_{m_j}} \right) * \left( \prod_{j=1}^{|\bar{M}|} p_{\bar{m}_j}^{x_{\bar{m}_j}} \left( \frac{P_{obs}}{p_{\bar{m}_j}} \right)^{x_{\bar{m}_j}} \right) \\ &= \frac{K!}{K_{obs}!} \frac{\prod_{j=1}^{|\bar{M}|} (x_{m_j}!)}{\prod_{j=1}^{|\bar{M}|} (x_{m_j}!)} * \left( \prod_{j=1}^{|\bar{M}|} p_{m_j}^{x_{m_j}} \right) * \left( \prod_{j=1}^{|\bar{M}|} (P_{obs})^{x_{\bar{m}_j}} \right) \\ &= \frac{K!}{K_{obs}!} \frac{\prod_{j=1}^{|\bar{M}|} (x_{m_j}!)}{\prod_{j=1}^{|\bar{M}|} (x_{m_j}!)} * \left( \prod_{j=1}^{|\bar{M}|} p_{m_j}^{x_{m_j}} \right) * ((P_{obs})^{K_{obs}}) \end{aligned}$$

This implies that the conditional PDF of  $X_{mis}$  given  $X_{obs}$  is multinomial PDF of  $K$  trials.

$$\begin{aligned} f(X_{mis}|X_{obs}, \Theta_M) &= f(X_{mis}|X_{obs}, \Theta) \\ &= \frac{K!}{K_{obs}!} \frac{\prod_{j=1}^{|\bar{M}|} (x_{m_j}!)}{\prod_{j=1}^{|\bar{M}|} (x_{m_j}!)} * \left( \prod_{j=1}^{|\bar{M}|} p_{m_j}^{x_{m_j}} \right) * ((P_{obs})^{K_{obs}}) \end{aligned} \quad (2.50)$$

Therefore,

$$\begin{aligned} f(X_{mis}(i)|X_{obs}(i), \Theta_{M_i}) &= f(X_{mis}(i)|X_{obs}(i), \Theta) \\ &= \frac{K!}{K_{obs}(i)!} \frac{\prod_{j=1}^{|\bar{M}_i|} (x_{im_j}!)}{\prod_{j=1}^{|\bar{M}_i|} (x_{im_j}!)} * \left( \prod_{j=1}^{|\bar{M}_i|} p_{m_{ij}}^{x_{im_j}} \right) * ((P_{obs}(i))^{K_{obs}(i)}) \end{aligned}$$

Where

$$P_{obs}(i) = \sum_{j=1}^{|\bar{M}_i|} p_{\bar{m}_{ij}}$$

$$K_{obs}(i) = \sum_{j=1}^{|\bar{M}_i|} x_{\bar{m}_{ij}}$$

Obviously, the parameter  $\Theta_{M_i}$  of the conditional PDF  $f(X_{mis}(i)|X_{obs}(i), \Theta_{M_i})$  is:

$$\Theta_{M_i} = u(\Theta, X_{obs}(i)) = \begin{pmatrix} p_{m_1} \\ p_{m_2} \\ \vdots \\ p_{m_k} \\ P_{obs}(i) = \sum_{j=1}^{|\bar{M}_i|} p_{\bar{m}_{ij}} \end{pmatrix} \quad (2.51)$$

Therefore, equation 2.51 to extract  $\Theta_{M_i}$  from  $\Theta$  given  $X_{obs}(i)$  is an instance of equation 2.15. It is easy to check that

$$\sum_{j=1}^{|\bar{M}_i|} x_{m_{ij}} + K_{obs}(i) = K_{mis}(i) + K_{obs}(i) = K$$

$$\sum_{j=1}^{|\bar{M}_i|} p_{m_{ij}} + P_{obs}(i) = \sum_{j=1}^{|\bar{M}_i|} p_{m_{ij}} + \sum_{j=1}^{|\bar{M}_i|} p_{\bar{m}_{ij}} = \sum_{j=1}^n p_j = 1$$

At E-step of some  $t^{\text{th}}$  iteration, given current parameter  $\Theta^{(t)} = (p_1^{(t)}, p_2^{(t)}, \dots, p_n^{(t)})^T$ , the sufficient statistic of  $X$  is calculated according to equation 2.22. Let,

$$\tau^{(t)} = \frac{1}{N} \sum_{i=1}^N \left\{ \tau(X_{obs}(i)), E(\tau(X_{mis}) | \Theta_{M_i}^{(t)}) \right\}$$

The sufficient statistic of  $X_{obs}(i)$  is:

$$\tau(X_{obs}(i)) = (x_{i\bar{m}_1}, x_{i\bar{m}_2}, \dots, x_{i\bar{m}_{|\bar{M}_i|}})^T$$

The sufficient statistic of  $X_{mis}(i)$  with regard to  $f(X_{mis}(i)|X_{obs}(i), \Theta_{M_i})$  is:

$$\tau(X_{mis}(i)) = \begin{pmatrix} x_{im_1} \\ x_{im_2} \\ \vdots \\ x_{im_{|\bar{M}_i|}} \\ \sum_{j=1}^{|\bar{M}_i|} x_{\bar{m}_{ij}} \end{pmatrix}$$

We also have:

$$E(\tau(X_{mis}) | \Theta_{M_i}^{(t)}) = \int_{X_{mis}} f(X_{mis} | X_{obs}, \Theta_{M_i}^{(t)}) \tau(X_{mis}) dX_{mis} = \begin{pmatrix} Kp_{m_1} \\ Kp_{m_2} \\ \vdots \\ Kp_{m_{|\bar{M}_i|}} \\ \sum_{j=1}^{|\bar{M}_i|} Kp_{\bar{m}_{ij}} \end{pmatrix}$$

Therefore, the sufficient statistic of  $X$  at E-step of some  $t^{\text{th}}$  iteration given current parameter  $\Theta^{(t)} = (p_1^{(t)}, p_2^{(t)}, \dots, p_n^{(t)})^T$  is defined as follows:

$$\begin{aligned}\tau^{(t)} &= (\bar{x}_1^{(t)}, \bar{x}_2^{(t)}, \dots, \bar{x}_n^{(t)})^T \\ \bar{x}_j^{(t)} &= \frac{1}{N} \sum_{i=1}^N \begin{cases} x_{ij} & \text{if } j \notin M_i \\ K p_j^{(t)} & \text{if } j \in M_i \end{cases} \quad \forall j\end{aligned}\tag{2.52}$$

Equation 2.52 is an instance of equation 2.11, which compose  $\tau(X)$  from  $\tau(X_{\text{obs}})$  and  $\tau(X_{\text{mis}})$  when  $f(X|\Theta)$  is multinomial PDF.

At M-step of some  $t^{\text{th}}$  iteration, we need to maximize  $Q_1(\Theta'|\Theta)$  with following constraint

$$\sum_{j=1}^n p_j = 1$$

According to equation 2.19, we have:

$$\begin{aligned}Q_1(\Theta'|\Theta) &= \sum_{i=1}^N E(\log(b(X_{\text{obs}}(i), X_{\text{mis}})) | \Theta_{M_i}) + (\Theta')^T \sum_{i=1}^N \{\tau(X_{\text{obs}}(i)), E(\tau(X_{\text{mis}}) | \Theta_{M_i})\} \\ &\quad - N \log(a(\Theta'))\end{aligned}$$

Where quantities  $b(X_{\text{obs}}(i), X_{\text{mis}})$  and  $a(\Theta')$  belongs to the PDF  $f(X|\Theta)$  of  $X$ . Because there is the constraint  $\sum_{j=1}^n p_j = 1$ , we use Lagrange duality method to maximize  $Q_1(\Theta'|\Theta)$ . The Lagrange function  $la(\Theta', \lambda | \Theta)$  is sum of  $Q_1(\Theta'|\Theta)$  and the constraint  $\sum_{j=1}^n p_j = 1$ , as follows:

$$\begin{aligned}la(\Theta', \lambda | \Theta) &= Q_1(\Theta'|\Theta) + \lambda \left( 1 - \sum_{j=1}^n p'_j \right) \\ &= \sum_{i=1}^N E(\log(b(X_{\text{obs}}(i), X_{\text{mis}})) | \Theta_{M_i}) \\ &\quad + (\Theta')^T \sum_{i=1}^N \{\tau(X_{\text{obs}}(i)), E(\tau(X_{\text{mis}}) | \Theta_{M_i})\} - N \log(a(\Theta')) + \lambda \left( 1 - \sum_{j=1}^n p'_j \right)\end{aligned}$$

Where  $\Theta' = (p'_1, p'_2, \dots, p'_n)^T$ . Note,  $\lambda \geq 0$  is called Lagrange multiplier. Of course,  $la(\Theta', \lambda | \Theta)$  is function of  $\Theta'$  and  $\lambda$ . The next parameter  $\Theta^{(t+1)}$  that maximizes  $Q_1(\Theta'|\Theta)$  is solution of the equation formed by setting the first-order derivative of Lagrange function regarding  $\Theta'$  and  $\lambda$  to be zero.

The first-order partial derivative of  $la(\Theta', \lambda | \Theta)$  with regard to  $\Theta'$  is:

$$\begin{aligned}\frac{\partial la(\Theta', \lambda | \Theta)}{\partial \Theta'} &= \sum_{i=1}^N \left( E(\tau(X_{\text{obs}}(i), X_{\text{mis}}) | \Theta_{M_i}) \right)^T - N \log'(a(\Theta')) \\ &= \sum_{i=1}^N \{\tau(X_{\text{obs}}(i)), E(\tau(X_{\text{mis}}) | \Theta_{M_i})\}^T - N \log'(a(\Theta')) - (\lambda, \lambda, \dots, \lambda)^T\end{aligned}$$

By referring table 1.2, we have:

$$\log'(a(\Theta')) = (E(\tau(X) | \Theta'))^T = \int_X f(X|\Theta) (\tau(X))^T dX$$

Thus,

$$\frac{\partial la(\Theta', \lambda | \Theta)}{\partial \Theta'} = \sum_{i=1}^N \{\tau(X_{\text{obs}}(i)), E(\tau(X_{\text{mis}}) | \Theta_{M_i})\}^T - N (E(\tau(X) | \Theta'))^T - (\lambda, \lambda, \dots, \lambda)^T$$

The first-order partial derivative of  $la(\Theta', \lambda | \Theta)$  with regard to  $\lambda$  is:

$$\frac{\partial la(\Theta', \lambda | \Theta)}{\partial \lambda} = 1 - \sum_{j=1}^n p'_j$$

Therefore, at M-step of some  $t^{\text{th}}$  iteration, given current parameter  $\Theta^{(t)} = (p_1^{(t)}, p_2^{(t)}, \dots, p_n^{(t)})^T$ , the next parameter  $\Theta^{(t+1)}$  is solution of the following equation:

$$\begin{cases} \sum_{i=1}^N \left\{ \tau(X_{obs}(i)), E(\tau(X_{mis}) | \Theta_{M_i}^{(t)}) \right\}^T \\ -N(E(\tau(X) | \Theta))^{(t)} - (\lambda, \lambda, \dots, \lambda) = \mathbf{0}^T \\ 1 - \sum_{j=1}^n p_j = 0 \end{cases}$$

This implies:

$$\begin{cases} E(\tau(X) | \Theta) = \tau^{(t)} - \begin{pmatrix} \lambda/N \\ \lambda/N \\ \lambda/N \\ \lambda/N \end{pmatrix} \\ \sum_{j=1}^n p_j = 1 \end{cases}$$

Where,

$$\tau^{(t)} = \frac{1}{N} \sum_{i=1}^N \left\{ \tau(X_{obs}(i)), E(\tau(X_{mis}) | \Theta_{M_i}^{(t)}) \right\}$$

Due to

$$\begin{aligned} E(\tau(X) | \Theta) &= \int_X \tau(X) f(X | \Theta) dX = (Kp_1, Kp_2, \dots, Kp_n)^T \\ \tau^{(t)} &= (\bar{x}_1^{(t)}, \bar{x}_2^{(t)}, \dots, \bar{x}_n^{(t)})^T \\ \bar{x}_j^{(t)} &= \frac{1}{N} \sum_{i=1}^N \begin{cases} x_{ij} & \text{if } j \notin M_i \\ Kp_j^{(t)} & \text{if } j \in M_i \end{cases} \forall j \end{aligned}$$

We obtain  $n$  equations  $Kp_j = -\lambda/N + \bar{x}_j^{(t)}$  and 1 constraint  $\sum_{j=1}^n p_j = 1$ . Therefore, we have:

$$p_j = -\frac{\lambda}{KN} + \frac{1}{KN} \sum_{i=1}^N \begin{cases} x_{ij} & \text{if } j \notin M_i \\ Kp_j^{(t)} & \text{if } j \in M_i \end{cases} \forall j$$

Summing  $n$  equations above, we have:

$$\begin{aligned} 1 &= \sum_{j=1}^n p_j = -\frac{\lambda}{KN} + \frac{1}{KN} \sum_{j=1}^n \left( \sum_{i=1}^N \begin{cases} x_{ij} & \text{if } j \notin M_i \\ Kp_j^{(t)} & \text{if } j \in M_i \end{cases} \right) \\ &= -\frac{\lambda}{KN} + \frac{1}{KN} \sum_{i=1}^N \left( \sum_{j=1}^{|\bar{M}_i|} x_{i\bar{m}_j} + \sum_{j=1}^{|M_i|} Kp_{m_j}^{(t)} \right) \end{aligned}$$

Suppose every missing value  $x_{im_j}$  is estimated by  $Kp_{m_j}$  such that:

$$\sum_{j=1}^{|\bar{M}_i|} x_{\bar{m}_{ij}} = \sum_{j=1}^{|M_i|} Kp_{m_j}^{(t)}$$



We obtain:

$$1 = -\frac{\lambda}{KN} + \frac{1}{KN} \sum_{i=1}^N \left( \sum_{j=1}^{|\bar{M}_i|} x_{i\bar{m}_j} + \sum_{j=1}^{|M_i|} x_{im_j} \right) = -\frac{\lambda}{KN} + \frac{1}{KN} \sum_{i=1}^N K = -\frac{\lambda}{KN} + 1$$

This implies

$$\lambda = 0$$

Such that

$$p_j = \frac{1}{KN} \sum_{i=1}^N \begin{cases} x_{ij} & \text{if } j \notin M_i \\ Kp_j^{(t)} & \text{if } j \in M_i \end{cases} \forall j$$

Therefore, at M-step of some  $t^{\text{th}}$  iteration, given current parameter  $\Theta^{(t)} = (p_1^{(t)}, p_2^{(t)}, \dots, p_n^{(t)})^T$ , the next parameter  $\Theta^{(t+1)}$  is specified by following equation.

$$p_j^{(t+1)} = \frac{1}{KN} \sum_{i=1}^N \begin{cases} x_{ij} & \text{if } j \notin M_i \\ Kp_j^{(t)} & \text{if } j \in M_i \end{cases} \forall j \quad (2.53)$$

In general, given sample  $\mathcal{X} = \{X_1, X_2, \dots, X_N\}$  whose  $X_i$  (s) are iid is MCAR data and  $f(X|\Theta)$  is multinomial PDF of  $K$  trials, GEM for handling missing data is summarized in table 2.3.

*M-step:*

Given  $\tau^{(t)}$  and  $\Theta^{(t)} = (p_1^{(t)}, p_2^{(t)}, \dots, p_n^{(t)})^T$ , the next parameter  $\Theta^{(t+1)}$  is specified by equation 2.53.

$$p_j^{(t+1)} = \frac{1}{KN} \sum_{i=1}^N \begin{cases} x_{ij} & \text{if } j \notin M_i \\ Kp_j^{(t)} & \text{if } j \in M_i \end{cases} \forall j$$

**Table 2.3.** E-step and M-step of GEM algorithm for handling missing data given multinomial PDF

In table 2.3, E-step is implied in how to perform M-step. As aforementioned, in practice we can stop GEM after its first iteration was done, which is reasonable enough to handle missing data. Next section includes two examples of handling missing data with multinomial distribution and multinomial distribution.

### 3. Numerical examples

It is necessary to have an example for illustrating how to handle missing data with multinomial PDF.

**Example 3.1.** Given sample of size two,  $\mathcal{X} = \{X_1, X_2\}$  in which  $X_1 = (x_{11}=1, x_{12}=?, x_{13}=3, x_{14}=?)^T$  and  $X_2 = (x_{21}=?, x_{22}=2, x_{23}=?, x_{24}=4)^T$  are iid. Therefore, we also have  $Z_1 = (z_{11}=0, z_{12}=1, z_{13}=0, z_{14}=1)^T$  and  $Z_2 = (z_{21}=1, z_{22}=0, z_{23}=1, z_{24}=0)^T$ . All  $Z_i$  (s) are iid too.

	$x_1$	$x_2$	$x_3$	$x_4$		$z_1$	$z_2$	$z_3$	$z_4$
$X_1$	1	?	3	?	$Z_1$	0	1	0	1
$X_2$	?	2	?	4	$Z_2$	1	0	1	0

Of course, we have  $X_{obs}(1) = (x_{11}=1, x_{13}=3)^T$ ,  $X_{mis}(1) = (x_{12}=?, x_{14}=?)^T$ ,  $X_{obs}(2) = (x_{22}=2, x_{24}=4)^T$  and  $X_{mis}(2) = (x_{21}=?, x_{23}=?)^T$ . We also have  $M_1 = \{m_{11}=2, m_{12}=4\}$ ,  $\bar{M}_1 = \{\bar{m}_{11}=1, \bar{m}_{12}=3\}$ ,  $M_2 = \{m_{21}=1, m_{22}=3\}$ , and  $\bar{M}_2 = \{\bar{m}_{21}=2, \bar{m}_{22}=4\}$ . Let  $X$  and  $Z$  be random variables representing every  $X_i$  and every  $Z_i$ , respectively. Suppose  $f(X|\Theta)$  is multinomial PDF and missingness variable  $Z$  follows binomial distribution of 4 trials according to equation 2.26 and equation 2.27. Dimension of  $X$  is 4. We will estimate  $\Theta = (\mu, \Sigma)^T$  and  $\Phi = p$  based on  $\mathcal{X}$ .

$$\mu = (\mu_1, \mu_2, \mu_3, \mu_4)^T$$

$$\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \sigma_{13} & \sigma_{14} \\ \sigma_{21} & \sigma_{22} & \sigma_{23} & \sigma_{24} \\ \sigma_{31} & \sigma_{32} & \sigma_{33} & \sigma_{34} \\ \sigma_{41} & \sigma_{42} & \sigma_{43} & \sigma_{44} \end{pmatrix}$$

The parameters  $\mu$  and  $\Sigma$  are initialized arbitrarily as zero vector and identity vector whereas  $p$  is initialized 0.5 as follows:

$$\begin{aligned} \mu^{(1)} &= \left( \mu_1^{(1)} = 0, \mu_2^{(1)} = 0, \mu_3^{(1)} = 0, \mu_4^{(1)} = 0 \right)^T \\ \Sigma^{(1)} &= \begin{pmatrix} \sigma_{11}^{(1)} = 1 & \sigma_{12}^{(1)} = 0 & \sigma_{13}^{(1)} = 0 & \sigma_{14}^{(1)} = 0 \\ \sigma_{21}^{(1)} = 0 & \sigma_{22}^{(1)} = 1 & \sigma_{23}^{(1)} = 0 & \sigma_{24}^{(1)} = 0 \\ \sigma_{31}^{(1)} = 0 & \sigma_{32}^{(1)} = 0 & \sigma_{33}^{(1)} = 1 & \sigma_{34}^{(1)} = 0 \\ \sigma_{41}^{(1)} = 0 & \sigma_{42}^{(1)} = 0 & \sigma_{43}^{(1)} = 0 & \sigma_{44}^{(1)} = 1 \end{pmatrix} \\ p^{(1)} &= 0.5 \end{aligned}$$

At 1<sup>st</sup> iteration, E-step, we have:

$$\begin{aligned} X_{obs}(1) &= (x_1 = 1, x_3 = 3)^T \\ \mu_{mis}(1) &= \left( \mu_2^{(1)} = 0, \mu_4^{(1)} = 0 \right)^T \\ \Sigma_{mis}(1) &= \begin{pmatrix} \sigma_{22}^{(1)} = 1 & \sigma_{24}^{(1)} = 0 \\ \sigma_{42}^{(1)} = 0 & \sigma_{44}^{(1)} = 1 \end{pmatrix} \\ \mu_{obs}(1) &= \left( \mu_1^{(1)} = 0, \mu_3^{(1)} = 0 \right)^T \\ \Sigma_{obs}(1) &= \begin{pmatrix} \sigma_{11}^{(1)} = 1 & \sigma_{13}^{(1)} = 0 \\ \sigma_{31}^{(1)} = 0 & \sigma_{33}^{(1)} = 1 \end{pmatrix} \\ V_{obs}^{mis}(1) &= \begin{pmatrix} \sigma_{21}^{(1)} = 0 & \sigma_{23}^{(1)} = 0 \\ \sigma_{41}^{(1)} = 0 & \sigma_{43}^{(1)} = 0 \end{pmatrix} \\ V_{mis}^{obs}(1) &= \begin{pmatrix} \sigma_{12}^{(1)} = 0 & \sigma_{14}^{(1)} = 0 \\ \sigma_{32}^{(1)} = 0 & \sigma_{34}^{(1)} = 0 \end{pmatrix} \\ \mu_{M_1}^{(1)} &= \mu_{mis}(1) + \left( V_{obs}^{mis}(1) \right) \left( \Sigma_{obs}(1) \right)^{-1} \left( X_{obs}(1) - \mu_{obs}(1) \right) \\ &= \left( \mu_{M_1}^{(1)}(2) = 0, \mu_{M_1}^{(1)}(4) = 0 \right)^T \\ \Sigma_{M_1}^{(1)} &= \Sigma_{mis}(1) - \left( V_{obs}^{mis}(1) \right) \left( \Sigma_{obs}(1) \right)^{-1} \left( V_{mis}^{obs}(1) \right) = \begin{pmatrix} \Sigma_{M_1}^{(1)}(2,2) = 1 & \Sigma_{M_1}^{(1)}(2,4) = 0 \\ \Sigma_{M_1}^{(1)}(4,2) = 0 & \Sigma_{M_1}^{(1)}(4,4) = 1 \end{pmatrix} \end{aligned}$$

$$\begin{aligned} X_{obs}(2) &= (x_2 = 2, x_4 = 4)^T \\ \mu_{mis}(2) &= \left( \mu_1^{(1)} = 0, \mu_3^{(1)} = 0 \right)^T \\ \Sigma_{mis}(2) &= \begin{pmatrix} \sigma_{11}^{(1)} = 1 & \sigma_{13}^{(1)} = 0 \\ \sigma_{31}^{(1)} = 0 & \sigma_{33}^{(1)} = 1 \end{pmatrix} \\ \mu_{obs}(2) &= \left( \mu_2^{(1)} = 0, \mu_4^{(1)} = 0 \right)^T \\ \Sigma_{obs}(2) &= \begin{pmatrix} \sigma_{22}^{(1)} = 1 & \sigma_{24}^{(1)} = 0 \\ \sigma_{42}^{(1)} = 0 & \sigma_{44}^{(1)} = 1 \end{pmatrix} \\ V_{obs}^{mis}(2) &= \begin{pmatrix} \sigma_{12}^{(1)} = 0 & \sigma_{14}^{(1)} = 0 \\ \sigma_{32}^{(1)} = 0 & \sigma_{34}^{(1)} = 0 \end{pmatrix} \end{aligned}$$

$$\begin{aligned}
V_{mis}^{obs}(2) &= \begin{pmatrix} \sigma_{21}^{(1)} = 0 & \sigma_{23}^{(1)} = 0 \\ \sigma_{41}^{(1)} = 0 & \sigma_{43}^{(1)} = 0 \end{pmatrix} \\
\mu_{M_2}^{(1)} &= \mu_{mis}(2) + \left(V_{obs}^{mis}(2)\right) \left(\Sigma_{obs}(2)\right)^{-1} \left(X_{obs}(2) - \mu_{obs}(2)\right) \\
&= \left(\mu_{M_2}^{(1)}(1) = 0, \mu_{M_2}^{(1)}(3) = 0\right)^T \\
\Sigma_{M_2}^{(1)} &= \Sigma_{mis}(2) - \left(V_{obs}^{mis}(2)\right) \left(\Sigma_{obs}(2)\right)^{-1} \left(V_{mis}^{obs}\right) = \begin{pmatrix} \Sigma_{M_2}^{(1)}(1,1) = 1 & \Sigma_{M_2}^{(1)}(1,3) = 0 \\ \Sigma_{M_2}^{(1)}(3,1) = 0 & \Sigma_{M_2}^{(1)}(3,3) = 1 \end{pmatrix}
\end{aligned}$$

$$\begin{aligned}
\bar{x}_1^{(1)} &= \frac{1}{2} \left( x_{11} + \mu_{M_2}^{(1)}(1) \right) = 0.5 \\
\bar{x}_2^{(1)} &= \frac{1}{2} \left( \mu_{M_1}^{(1)}(2) + x_{22} \right) = 1 \\
\bar{x}_3^{(1)} &= \frac{1}{2} \left( x_{13} + \mu_{M_2}^{(1)}(3) \right) = 1.5 \\
\bar{x}_4^{(1)} &= \frac{1}{2} \left( \mu_{M_1}^{(1)}(4) + x_{24} \right) = 2
\end{aligned}$$

$$\begin{aligned}
s_{11}^{(1)} &= \frac{1}{2} \left( (x_{11})^2 + \left( \Sigma_{M_2}^{(1)}(1,1) + \left( \mu_{M_2}^{(1)}(1) \right)^2 \right) \right) = 1 \\
s_{12}^{(1)} &= s_{21}^{(1)} = \frac{1}{2} \left( x_{11} \mu_{M_1}^{(1)}(2) + \mu_{M_2}^{(1)}(1) x_{22} \right) = 0 \\
s_{13}^{(1)} &= s_{31}^{(1)} = \frac{1}{2} \left( x_{11} x_{13} + \left( \Sigma_{M_2}^{(1)}(1,3) + \mu_{M_2}^{(1)}(1) \mu_{M_2}^{(1)}(3) \right) \right) = 1.5 \\
s_{14}^{(1)} &= s_{41}^{(1)} = \frac{1}{2} \left( x_{11} \mu_{M_1}^{(1)}(4) + \mu_{M_2}^{(1)}(1) x_{24} \right) = 0 \\
s_{22}^{(1)} &= \frac{1}{2} \left( \left( \Sigma_{M_1}^{(1)}(2,2) + \left( \mu_{M_1}^{(1)}(2) \right)^2 \right) + (x_{22})^2 \right) = 2.5 \\
s_{23}^{(1)} &= s_{32}^{(1)} = \frac{1}{2} \left( \mu_{M_1}^{(1)}(2) x_{13} + x_{22} \mu_{M_2}^{(1)}(3) \right) = 0 \\
s_{24}^{(1)} &= s_{42}^{(1)} = \frac{1}{2} \left( \left( \Sigma_{M_1}^{(1)}(2,4) + \mu_{M_1}^{(1)}(2) \mu_{M_1}^{(1)}(4) \right) + x_{22} x_{24} \right) = 4 \\
s_{33}^{(1)} &= \frac{1}{2} \left( (x_{13})^2 + \left( \Sigma_{M_2}^{(1)}(3,3) + \left( \mu_{M_2}^{(1)}(3) \right)^2 \right) \right) = 5 \\
s_{34}^{(1)} &= s_{43}^{(1)} = \frac{1}{2} \left( x_{13} \mu_{M_1}^{(1)}(4) + \mu_{M_2}^{(1)}(3) x_{24} \right) = 0 \\
s_{44}^{(1)} &= \frac{1}{2} \left( \left( \Sigma_{M_1}^{(1)}(4,4) + \left( \mu_{M_1}^{(1)}(4) \right)^2 \right) + (x_{24})^2 \right) = 8.5
\end{aligned}$$

At 1<sup>st</sup> iteration, M-step, we have:

$$\begin{aligned}
\mu_1^{(2)} &= \bar{x}_1^{(1)} = 0.5 \\
\mu_2^{(2)} &= \bar{x}_2^{(1)} = 1 \\
\mu_3^{(2)} &= \bar{x}_3^{(1)} = 1.5 \\
\mu_4^{(2)} &= \bar{x}_4^{(1)} = 2
\end{aligned}$$

$$\begin{aligned}
\sigma_{11}^{(2)} &= s_{11}^{(1)} - \left( \bar{x}_1^{(1)} \right)^2 = 0.75 \\
\sigma_{12}^{(2)} &= \sigma_{21}^{(2)} = s_{12}^{(1)} - \bar{x}_1^{(1)} \bar{x}_2^{(1)} = -0.5
\end{aligned}$$

$$\begin{aligned}
\sigma_{13}^{(2)} &= \sigma_{31}^{(2)} = s_{13}^{(1)} - \bar{x}_1^{(1)} \bar{x}_3^{(1)} = 0.75 \\
\sigma_{14}^{(2)} &= \sigma_{41}^{(2)} = s_{14}^{(1)} - \bar{x}_1^{(1)} \bar{x}_4^{(1)} = -1 \\
\sigma_{22}^{(2)} &= s_{22}^{(1)} - \left(\bar{x}_2^{(1)}\right)^2 = 1.5 \\
\sigma_{23}^{(2)} &= \sigma_{32}^{(2)} = s_{23}^{(1)} - \bar{x}_2^{(1)} \bar{x}_3^{(1)} = -1.5 \\
\sigma_{24}^{(2)} &= \sigma_{42}^{(2)} = s_{24}^{(1)} - \bar{x}_2^{(1)} \bar{x}_4^{(1)} = 2 \\
\sigma_{33}^{(2)} &= s_{33}^{(1)} - \left(\bar{x}_3^{(1)}\right)^2 = 2.75 \\
\sigma_{34}^{(2)} &= \sigma_{43}^{(2)} = s_{34}^{(1)} - \bar{x}_3^{(1)} \bar{x}_4^{(1)} = -3 \\
\sigma_{44}^{(2)} &= s_{44}^{(1)} - \left(\bar{x}_4^{(1)}\right)^2 = 4.5
\end{aligned}$$

$$p^{(2)} = \frac{c(Z_1) + c(Z_2)}{4 * 2} = \frac{2 + 2}{4 * 2} = 0.5$$

At 2<sup>nd</sup> iteration, E-step, we have:

$$\begin{aligned}
X_{obs}(1) &= (x_1 = 1, x_3 = 3)^T \\
\mu_{mis}(1) &= \left(\mu_2^{(2)} = 1, \mu_4^{(2)} = 2\right)^T \\
\Sigma_{mis}(1) &= \begin{pmatrix} \sigma_{22}^{(2)} = 1.5 & \sigma_{24}^{(2)} = 2 \\ \sigma_{42}^{(2)} = 2 & \sigma_{44}^{(2)} = 4.5 \end{pmatrix} \\
\mu_{obs}(1) &= \left(\mu_1^{(2)} = 0.5, \mu_3^{(2)} = 1.5\right)^T \\
\Sigma_{obs}(1) &= \begin{pmatrix} \sigma_{11}^{(2)} = 0.75 & \sigma_{13}^{(2)} = 0.75 \\ \sigma_{31}^{(2)} = 0.75 & \sigma_{33}^{(2)} = 2.75 \end{pmatrix} \\
V_{obs}^{mis}(1) &= \begin{pmatrix} \sigma_{21}^{(2)} = -0.5 & \sigma_{23}^{(2)} = -1.5 \\ \sigma_{41}^{(2)} = -1 & \sigma_{43}^{(2)} = -3 \end{pmatrix} \\
V_{mis}^{obs}(1) &= \begin{pmatrix} \sigma_{12}^{(2)} = -0.5 & \sigma_{14}^{(2)} = -1 \\ \sigma_{32}^{(2)} = -1.5 & \sigma_{34}^{(2)} = -3 \end{pmatrix} \\
\mu_{M_1}^{(2)} &= \mu_{mis}(1) + \left(V_{obs}^{mis}(1)\right) \left(\Sigma_{obs}(1)\right)^{-1} (X_{obs}(1) - \mu_{obs}(1)) \\
&= \left(\mu_{M_1}^{(2)}(2) \cong 0.17, \mu_{M_1}^{(2)}(4) \cong 0.33\right)^T \\
\Sigma_{M_1}^{(2)} &= \Sigma_{mis}(1) - \left(V_{obs}^{mis}(1)\right) \left(\Sigma_{obs}(1)\right)^{-1} (V_{mis}^{obs}) = \begin{pmatrix} \Sigma_{M_1}^{(2)}(2,2) \cong 0.67 & \Sigma_{M_1}^{(2)}(2,4) \cong 0.33 \\ \Sigma_{M_1}^{(2)}(4,2) \cong 0.33 & \Sigma_{M_1}^{(2)}(4,4) \cong 1.17 \end{pmatrix}
\end{aligned}$$

$$\begin{aligned}
X_{obs}(2) &= (x_2 = 2, x_4 = 4)^T \\
\mu_{mis}(2) &= \left(\mu_1^{(2)} = 0.5, \mu_3^{(2)} = 1.5\right)^T \\
\Sigma_{mis}(2) &= \begin{pmatrix} \sigma_{11}^{(2)} = 0.75 & \sigma_{13}^{(2)} = 0.75 \\ \sigma_{31}^{(2)} = 0.75 & \sigma_{33}^{(2)} = 2.75 \end{pmatrix} \\
\mu_{obs}(2) &= \left(\mu_2^{(2)} = 1, \mu_4^{(2)} = 2\right)^T \\
\Sigma_{obs}(2) &= \begin{pmatrix} \sigma_{22}^{(2)} = 1.5 & \sigma_{24}^{(2)} = 2 \\ \sigma_{42}^{(2)} = 2 & \sigma_{44}^{(2)} = 4.5 \end{pmatrix}
\end{aligned}$$

$$\begin{aligned}
V_{obs}^{mis}(2) &= \begin{pmatrix} \sigma_{12}^{(2)} = -0.5 & \sigma_{14}^{(2)} = -1 \\ \sigma_{32}^{(2)} = -1.5 & \sigma_{34}^{(2)} = -3 \end{pmatrix} \\
V_{mis}^{obs}(2) &= \begin{pmatrix} \sigma_{21}^{(2)} = -0.5 & \sigma_{23}^{(2)} = -1.5 \\ \sigma_{41}^{(2)} = -1 & \sigma_{43}^{(2)} = -3 \end{pmatrix} \\
\mu_{M_2}^{(2)} &= \mu_{mis}(2) + \left(V_{obs}^{mis}(2)\right) \left(\Sigma_{obs}(2)\right)^{-1} (X_{obs}(2) - \mu_{obs}(2)) \\
&= \left(\mu_{M_2}^{(2)}(1) \cong 0.05, \mu_{M_2}^{(2)}(3) = 0.14\right)^T \\
\Sigma_{M_2}^{(2)} &= \Sigma_{mis}(2) - \left(V_{obs}^{mis}(2)\right) \left(\Sigma_{obs}(2)\right)^{-1} (V_{mis}^{obs}) = \begin{pmatrix} \Sigma_{M_2}^{(2)}(1,1) \cong 0.52 & \Sigma_{M_2}^{(2)}(1,3) \cong 0.07 \\ \Sigma_{M_2}^{(2)}(3,1) \cong 0.07 & \Sigma_{M_2}^{(2)}(3,3) \cong 0.7 \end{pmatrix}
\end{aligned}$$

$$\begin{aligned}
\bar{x}_1^{(2)} &= \frac{1}{2} \left( x_{11} + \mu_{M_2}^{(2)}(1) \right) \cong 0.52 \\
\bar{x}_2^{(2)} &= \frac{1}{2} \left( \mu_{M_1}^{(2)}(2) + x_{22} \right) \cong 1.1 \\
\bar{x}_3^{(2)} &= \frac{1}{2} \left( x_{13} + \mu_{M_2}^{(2)}(3) \right) \cong 1.57 \\
\bar{x}_4^{(2)} &= \frac{1}{2} \left( \mu_{M_1}^{(2)}(4) + x_{24} \right) \cong 2.17
\end{aligned}$$

$$\begin{aligned}
s_{11}^{(2)} &= \frac{1}{2} \left( (x_{11})^2 + \left( \Sigma_{M_2}^{(2)}(1,1) + \left( \mu_{M_2}^{(2)}(1) \right)^2 \right) \right) \cong 0.76 \\
s_{12}^{(2)} &= s_{21}^{(2)} = \frac{1}{2} \left( x_{11} \mu_{M_1}^{(2)}(2) + \mu_{M_2}^{(2)}(1) x_{22} \right) \cong 0.13 \\
s_{13}^{(2)} &= s_{31}^{(2)} = \frac{1}{2} \left( x_{11} x_{13} + \left( \Sigma_{M_2}^{(2)}(1,3) + \mu_{M_2}^{(2)}(1) \mu_{M_2}^{(2)}(3) \right) \right) \cong 1.54 \\
s_{14}^{(2)} &= s_{41}^{(2)} = \frac{1}{2} \left( x_{11} \mu_{M_1}^{(2)}(4) + \mu_{M_2}^{(2)}(1) x_{24} \right) \cong 0.17 \\
s_{22}^{(2)} &= \frac{1}{2} \left( \left( \Sigma_{M_1}^{(2)}(2,2) + \left( \mu_{M_1}^{(2)}(2) \right)^2 \right) + (x_{22})^2 \right) \cong 2.35 \\
s_{23}^{(2)} &= s_{32}^{(2)} = \frac{1}{2} \left( \mu_{M_1}^{(2)}(2) x_{13} + x_{22} \mu_{M_2}^{(2)}(3) \right) \cong 0.39 \\
s_{24}^{(2)} &= s_{42}^{(2)} = \frac{1}{2} \left( \left( \Sigma_{M_1}^{(2)}(2,4) + \mu_{M_1}^{(2)}(2) \mu_{M_1}^{(2)}(4) \right) + x_{22} x_{24} \right) \cong 4.19 \\
s_{33}^{(2)} &= \frac{1}{2} \left( (x_{13})^2 + \left( \Sigma_{M_2}^{(2)}(3,3) + \left( \mu_{M_2}^{(2)}(3) \right)^2 \right) \right) \cong 4.86 \\
s_{34}^{(2)} &= s_{43}^{(2)} = \frac{1}{2} \left( x_{13} \mu_{M_1}^{(2)}(4) + \mu_{M_2}^{(2)}(3) x_{24} \right) \cong 0.77 \\
s_{44}^{(2)} &= \frac{1}{2} \left( \left( \Sigma_{M_1}^{(2)}(4,4) + \left( \mu_{M_1}^{(2)}(4) \right)^2 \right) + (x_{24})^2 \right) \cong 8.64
\end{aligned}$$

At 2<sup>nd</sup> iteration, M-step, we have:

$$\begin{aligned}
\mu_1^{(3)} &= \bar{x}_1^{(2)} \cong 0.52 \\
\mu_2^{(3)} &= \bar{x}_2^{(2)} \cong 1.1 \\
\mu_3^{(3)} &= \bar{x}_3^{(2)} \cong 1.57 \\
\mu_4^{(3)} &= \bar{x}_4^{(2)} \cong 2.17
\end{aligned}$$

$$\begin{aligned}
\sigma_{11}^{(3)} &= s_{11}^{(2)} - \left(\bar{x}_1^{(2)}\right)^2 \cong 0.49 \\
\sigma_{12}^{(3)} &= \sigma_{21}^{(3)} = s_{12}^{(2)} - \bar{x}_1^{(2)} \bar{x}_2^{(2)} \cong -0.44 \\
\sigma_{13}^{(3)} &= \sigma_{31}^{(3)} = s_{13}^{(2)} - \bar{x}_1^{(2)} \bar{x}_3^{(2)} \cong 0.72 \\
\sigma_{14}^{(3)} &= \sigma_{41}^{(3)} = s_{14}^{(2)} - \bar{x}_1^{(2)} \bar{x}_4^{(2)} \cong -0.96 \\
\sigma_{22}^{(3)} &= s_{22}^{(2)} - \left(\bar{x}_2^{(2)}\right)^2 \cong 1.17 \\
\sigma_{23}^{(3)} &= \sigma_{32}^{(3)} = s_{23}^{(2)} - \bar{x}_2^{(2)} \bar{x}_3^{(2)} \cong -1.31 \\
\sigma_{24}^{(3)} &= \sigma_{42}^{(3)} = s_{24}^{(2)} - \bar{x}_2^{(2)} \bar{x}_4^{(2)} \cong 1.85 \\
\sigma_{33}^{(3)} &= s_{33}^{(2)} - \left(\bar{x}_3^{(2)}\right)^2 \cong 2.4 \\
\sigma_{34}^{(3)} &= \sigma_{43}^{(3)} = s_{34}^{(2)} - \bar{x}_3^{(2)} \bar{x}_4^{(2)} \cong -2.63 \\
\sigma_{44}^{(3)} &= s_{44}^{(2)} - \left(\bar{x}_4^{(2)}\right)^2 \cong 3.94
\end{aligned}$$

$$p^{(3)} = \frac{c(Z_1) + c(Z_2)}{4 * 2} = \frac{2 + 2}{4 * 2} = 0.5$$

Because the sample is too small for GEM to converge to an exact maximizer with small enough number of iterations, we can stop GEM at the second iteration with  $\Theta^{(3)} = \Theta^* = (\mu^*, \Sigma^*)^T$  and  $\Phi^{(3)} = \Phi^* = p^*$  when difference between  $\Theta^{(2)}$  and  $\Theta^{(3)}$  is insignificant.

$$\begin{aligned}
\mu^* &= (\mu_1^* = 0.52, \mu_2^* = 1.1, \mu_3^* = 1.57, \mu_4^* = 2.17)^T \\
\Sigma^* &= \begin{pmatrix} \sigma_{11}^* = 0.49 & \sigma_{12}^* = -0.44 & \sigma_{13}^* = 0.72 & \sigma_{14}^* = -0.96 \\ \sigma_{21}^* = -0.44 & \sigma_{22}^* = 1.17 & \sigma_{23}^* = -1.31 & \sigma_{24}^* = 1.85 \\ \sigma_{31}^* = 0.72 & \sigma_{32}^* = -1.31 & \sigma_{33}^* = 2.4 & \sigma_{34}^* = -2.63 \\ \sigma_{41}^* = -0.96 & \sigma_{42}^* = 1.85 & \sigma_{43}^* = -2.63 & \sigma_{44}^* = 3.94 \end{pmatrix} \\
p^* &= 0.5
\end{aligned}$$

As aforementioned, because  $X_{mis}$  is a part of  $X$  and  $f(X_{mis} | \Theta_M)$  is derived directly from  $f(X|\Theta)$ , in practice we can stop GEM after its first iteration was done, which is reasonable enough to handle missing data.

As aforementioned, an interesting application of handling missing data is to fill in or predict missing values. For instance, the missing part  $X_{mis}(1)$  of  $X_1 = (x_{11}=1, x_{12}=?, x_{13}=3, x_{14}=?)^T$  is fulfilled by  $\mu_{M_1}^*$  according to equation 2.44 as follows:

$$\begin{aligned}
x_{12} &= \mu_2^* = 1.1 \\
x_{14} &= \mu_4^* = 2.17
\end{aligned}$$

It is necessary to have an example for illustrating how to handle missing data with multinomial PDF.

**Example 3.2.** Given sample of size two,  $\mathcal{X} = \{X_1, X_2\}$  in which  $X_1 = (x_{11}=1, x_{12}=?, x_{13}=3, x_{14}=?)^T$  and  $X_2 = (x_{21}=?, x_{22}=2, x_{23}=?, x_{24}=4)^T$  are iid.

	$x_1$	$x_2$	$x_3$	$x_4$
$X_1$	1	?	3	?
$X_2$	?	2	?	4

Of course, we have  $X_{obs}(1) = (x_{11}=1, x_{13}=3)^T$ ,  $X_{mis}(1) = (x_{12}=?, x_{14}=?)^T$ ,  $X_{obs}(2) = (x_{22}=2, x_{24}=4)^T$  and  $X_{mis}(2) = (x_{21}=?, x_{23}=?)^T$ . We also have  $M_1 = \{m_{11}=2, m_{12}=4\}$ ,  $\bar{M}_1 = \{\bar{m}_{11}=1, \bar{m}_{12}=3\}$ ,  $M_2 = \{m_{21}=1, m_{22}=3\}$ , and  $\bar{M}_2 = \{\bar{m}_{21}=2, \bar{m}_{22}=4\}$ . Let  $X$  be random variable representing every  $X_i$ . Suppose  $f(X|\Theta)$  is multinomial PDF of 10 trials. We will estimate  $\Theta = (p_1, p_2, p_3, p_4)^T$ . The parameters  $p_1, p_2, p_3$ , and  $p_4$  are initialized arbitrarily as 0.25 as follows:

$$\Theta^{(1)} = (p_1^{(1)} = 0.25, p_2^{(1)} = 0.25, p_3^{(1)} = 0.25, p_4^{(1)} = 0.25)^T$$

At 1<sup>st</sup> iteration, M-step, we have:

$$\begin{aligned}
p_1^{(2)} &= \frac{1}{10 * 2} (1 + 10 * 0.25) = 0.175 \\
p_2^{(2)} &= \frac{1}{10 * 2} (10 * 0.25 + 2) = 0.225 \\
p_3^{(2)} &= \frac{1}{10 * 2} (3 + 10 * 0.25) = 0.275 \\
p_4^{(2)} &= \frac{1}{10 * 2} (10 * 0.25 + 4) = 0.325
\end{aligned}$$

We stop GEM after the first iteration was done, which results the estimate  $\Theta^{(2)} = \Theta^* = (p_1^*, p_2^*, p_3^*, p_4^*)^T$  as follows:

$$\begin{aligned}
p_1^* &= 0.175 \\
p_2^* &= 0.225 \\
p_3^* &= 0.275 \\
p_4^* &= 0.325
\end{aligned}$$

## 4. Conclusions

In general, GEM is a powerful tool to handle missing data, which is not so difficult except that how to extract the parameter  $\Theta_M$  of the conditional PDF  $f(X_{mis} | X_{obs}, \Theta_M)$  from the whole parameter  $\Theta$  of the PDF  $f(X|\Theta)$  is most important with note that only  $f(X|\Theta)$  is defined firstly and then  $f(X_{mis} | X_{obs}, \Theta_M)$  is derived from  $f(X|\Theta)$ . Therefore, equation 2.15 is cornerstone of this method. Note, equation 2.35 and 2.51 are instances of equation 2.15 when  $f(X|\Theta)$  is multinormal PDF or multinomial PDF.

## References

- Burden, R. L., & Faires, D. J. (2011). *Numerical Analysis* (9th Edition ed.). (M. Julet, Ed.) Brooks/Cole Cengage Learning.
- Dempster, A. P., Laird, N. M., & Rubin, D. B. (1977). Maximum Likelihood from Incomplete Data via the EM Algorithm. (M. Stone, Ed.) *Journal of the Royal Statistical Society, Series B (Methodological)*, 39(1), 1-38.
- Hardle, W., & Simar, L. (2013). *Applied Multivariate Statistical Analysis*. Berlin, Germany: Research Data Center, School of Business and Economics, Humboldt University.
- Josse, J., Jiang, W., Sportisse, A., & Robin, G. (2018). *Handling missing values*. Inria. Julie Josse. Retrieved October 12, 2020, from <http://juliejosse.com/wp-content/uploads/2018/07/LectureNotesMissing.html>
- Nguyen, L. (2020). *Tutorial on EM algorithm*. MDPI. Preprints. doi:10.20944/preprints201802.0131.v8
- Ta, P. D. (2014). *Numerical Analysis Lecture Notes*. Vietnam Institute of Mathematics, Numerical Analysis and Scientific Computing. Hanoi: Vietnam Institute of Mathematics. Retrieved 2014
- Wikipedia. (2014, August 4). *Karush–Kuhn–Tucker conditions*. (Wikimedia Foundation) Retrieved November 16, 2014, from Wikipedia website: [http://en.wikipedia.org/wiki/Karush–Kuhn–Tucker\\_conditions](http://en.wikipedia.org/wiki/Karush–Kuhn–Tucker_conditions)
- Wikipedia. (2016, March September). *Exponential family*. (Wikimedia Foundation) Retrieved 2015, from Wikipedia website: [https://en.wikipedia.org/wiki/Exponential\\_family](https://en.wikipedia.org/wiki/Exponential_family)