

EM Algorithm Clear Explained

Hanyi Mao

March 2025

This vignette is inspired by an excellent textbook [1]. Section **2** and **3** are based on section **I.8.7** of [1], explained in a hopefully more understandable way.

1 Notation

Here are some symbols and conventions we will use in this vignette.

Notation	Description and Conventions
$\boldsymbol{\theta}$	The full parameter set of a given model (e.g., $\boldsymbol{\theta} = (\boldsymbol{\mu}, \boldsymbol{\Sigma})$).
D	Domain of definition.
t	Timestep.
n	Index of datapoints, $n = 1, 2, \dots, N$.
$\boldsymbol{\theta}^t$	Estimation of $\boldsymbol{\theta}$ at timestep t .
\mathbf{x}	The full observed data, $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$; not using a matrix notation as different datapoints may have different length
\mathbf{x}_n	A data point, treated as column vectors.
\mathbf{x}_n^T	The transpose of a vector or a matrix.
\mathbf{z}	Latent variable(s).
π_k, γ_{nk}	Some scalar parameters.
$\int_{\mathbf{z}} p(\mathbf{z}) d\mathbf{z}$	Can also mean $\sum_{\mathbf{z}} p(\mathbf{z})$ if \mathbf{z} is a discrete variable.

Table 1: Notations and conventions used in this vignette.

2 Bound Optimization

The EM algorithm utilizes the idea of **bound optimization** (also known as **MM algorithms**), so we first look at this to have some intuitions in mind. Assume we want to *maximize* a function:

$$l(\boldsymbol{\theta}) : D \rightarrow \mathbb{R}$$

In many cases, we can't do this maximization directly because l is probably intractable (meaning it's hard to compute). Now suppose we get some function

$$g(\boldsymbol{\theta}; \boldsymbol{\theta}^t) : D \rightarrow \mathbb{R}, \quad t = 0, 1, 2, \dots$$

in which the first $\boldsymbol{\theta} \in D$ is the independent variable and the second $\boldsymbol{\theta}^t \in D$ is the parameter of this function. (Notice that in many other literature, this function is denoted as $Q(\boldsymbol{\theta}^t; \boldsymbol{\theta}^t)$. We choose g rather than Q here in case people confuse this Q with the variational distribution q which we will see later.) It has following properties:

$$\forall \boldsymbol{\theta} \quad g(\boldsymbol{\theta}; \boldsymbol{\theta}^t) \leq l(\boldsymbol{\theta}) \quad (1)$$

$$g(\boldsymbol{\theta}^t; \boldsymbol{\theta}^t) = l(\boldsymbol{\theta}^t) \quad (2)$$

which means $g(\boldsymbol{\theta}; \boldsymbol{\theta}^t)$ is a lower bound of $l(\boldsymbol{\theta})$ and it's tight at $\boldsymbol{\theta}^t$. Now we can do some interesting things with this.

To formalize, we start at a random $\boldsymbol{\theta}^0$. We update the parameter at each timestep t :

$$\boldsymbol{\theta}^{t+1} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} g(\boldsymbol{\theta}; \boldsymbol{\theta}^t) \quad (3)$$

We can prove that this new $\boldsymbol{\theta}^{t+1}$ **improves** the value of l , because

$$l(\boldsymbol{\theta}^{t+1}) \geq g(\boldsymbol{\theta}^{t+1}; \boldsymbol{\theta}^t) \geq g(\boldsymbol{\theta}^t; \boldsymbol{\theta}^t) = l(\boldsymbol{\theta}^t) \quad (4)$$

The first inequality follows Equation (1), the second inequality follows Equation (3), and the third inequality follows Equation (2). Then we can use $g(\boldsymbol{\theta}; \boldsymbol{\theta}^{t+1})$ to find the next $\boldsymbol{\theta}^{t+2}$. Repeating this until convergence, we can maximize $l(\boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$. This process is visualized in figure 1. An animated version of visualization can be found in this project repository https://github.com/echeloni11/EM_vignette/, filename **BoundOptimization.mp4**.

3 The EM Algorithm

3.1 Derive a lower bound for log likelihood

In most cases of statistics/machine learning, our goal is to provide a maximum likelihood estimation of $\boldsymbol{\theta}$ given data \boldsymbol{x} .

$$\hat{\boldsymbol{\theta}}_{MLE} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\boldsymbol{x}; \boldsymbol{\theta}) \quad (5)$$

Sometimes we can't directly compute the likelihood as there are some unobserved variables in our modeling. The assumption here is that if we observed the latent variable \boldsymbol{z} , we can compute $p(\boldsymbol{x}, \boldsymbol{z}; \boldsymbol{\theta})$.

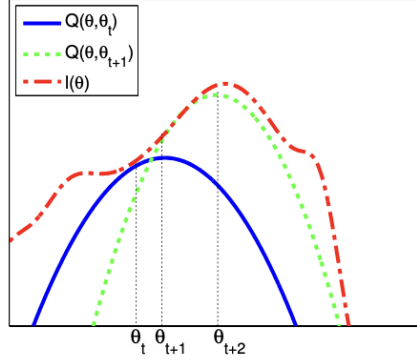


Figure 1: Visualization of Bound optimization. Adapted from [1]

With that in mind, let's see how to derive a tractable lower bound for log likelihood, so that we can apply the idea of bound optimization we saw (notice that maximizing likelihood is equal to maximizing log likelihood).

Given observed data $\mathbf{x} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$,

$$\log p(\mathbf{x}; \boldsymbol{\theta}) = \sum_{n=1}^N \log p(\mathbf{x}_n; \boldsymbol{\theta}) \quad (6)$$

$$= \sum_{n=1}^N \log \int_{\mathbf{z}_n} p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta}) d\mathbf{z}_n \quad (7)$$

The log operation outside of integral is hard to track. We use a trick called **Jensen's Inequality Trick** (which is the core of **Variational Inference (VI)** methods). The trick works by first plugging in an arbitrary auxiliary distribution (or **variational distribution**) q_n for each n , then use Jensen's Inequality (that is, $f(\mathbb{E}[x]) \geq \mathbb{E}[f(x)]$ for concave function f , and notice that log is concave)

to push the log operation into the integral.

$$\log p(\mathbf{x}; \boldsymbol{\theta}) = \sum_{n=1}^N \log p(\mathbf{x}_n; \boldsymbol{\theta}) \quad (8)$$

$$= \sum_{n=1}^N \log \int_{\mathbf{z}_n} p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta}) d\mathbf{z}_n \quad (9)$$

$$= \sum_{n=1}^N \log \int_{\mathbf{z}_n} q_n(\mathbf{z}_n) \frac{p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta})}{q_n(\mathbf{z}_n)} d\mathbf{z}_n \quad (10)$$

$$= \sum_{n=1}^N \log \mathbb{E}_{q_n(\mathbf{z}_n)} \left[\frac{p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta})}{q_n(\mathbf{z}_n)} \right] \quad (11)$$

$$\geq \sum_{n=1}^N \mathbb{E}_{q_n(\mathbf{z}_n)} \left[\log \frac{p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta})}{q_n(\mathbf{z}_n)} \right] \quad (12)$$

$$=: g(\boldsymbol{\theta}; \{q_1, q_2, \dots, q_N\}) \quad (13)$$

Now $g(\boldsymbol{\theta}; \{q_1, q_2, \dots, q_N\})$ is a lower bound of log likelihood. This is also called the **Evidence Lower Bound** or **ELBO**, as the log marginal likelihood $p(\mathbf{x}; \boldsymbol{\theta})$ is also known as the evidence. Notice that ELBO is a function with independent variable $\boldsymbol{\theta}$.

3.2 Find the q to make the bound tight

Recall that in bound optimization, we need the lower bound $g(\boldsymbol{\theta}; \boldsymbol{\theta}^t)$ to take $\boldsymbol{\theta}^t$ as parameter and be tight at $\boldsymbol{\theta}^t$. So what q should we choose such that $g(\boldsymbol{\theta}; \{q_1, q_2, \dots, q_N\})$ has such properties? With some algebraic manipulation, we can see that

$$g(\boldsymbol{\theta}; \{q_1, q_2, \dots, q_N\}) := \sum_{n=1}^N \mathbb{E}_{q_n(\mathbf{z}_n)} \left[\log \frac{p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta})}{q_n(\mathbf{z}_n)} \right] \quad (14)$$

$$= \sum_{n=1}^N \mathbb{E}_{q_n(\mathbf{z}_n)} \left[\log \frac{p(\mathbf{z}_n | \mathbf{x}_n; \boldsymbol{\theta}) p(\mathbf{x}_n; \boldsymbol{\theta})}{q_n(\mathbf{z}_n)} \right] \quad (15)$$

$$= \sum_{n=1}^N \mathbb{E}_{q_n(\mathbf{z}_n)} \left[\log \frac{p(\mathbf{z}_n | \mathbf{x}_n; \boldsymbol{\theta})}{q_n(\mathbf{z}_n)} + \log p(\mathbf{x}_n; \boldsymbol{\theta}) \right] \quad (16)$$

$$= \sum_{n=1}^N \left[\mathbb{E}_{q_n(\mathbf{z}_n)} \left[\log \frac{p(\mathbf{z}_n | \mathbf{x}_n; \boldsymbol{\theta})}{q_n(\mathbf{z}_n)} \right] + \log p(\mathbf{x}_n; \boldsymbol{\theta}) \right] \quad (17)$$

$$= \sum_{n=1}^N \left[-D_{\text{KL}}(q_n(\mathbf{z}_n) || p(\mathbf{z}_n | \mathbf{x}_n; \boldsymbol{\theta})) \right] + \log p(\mathbf{x}; \boldsymbol{\theta}) \quad (18)$$

where $D_{\text{KL}}(p(\mathbf{z})||q(\mathbf{z})) := \mathbb{E}_q[\log \frac{q(\mathbf{z})}{p(\mathbf{z})}]$ is the Kullback-Leibler divergence which is a metric that measures the distance between two distributions. The key takeaway is that $D_{\text{KL}}(q(\mathbf{z})||p(\mathbf{z})) \geq 0$, and $D_{\text{KL}}(q(\mathbf{z})||p(\mathbf{z})) = 0$ if and only if $q = p$.

Thus, from Equation (18) we can see that $g(\boldsymbol{\theta}; \{q_1, q_2, \dots, q_N\}) = \log p(\mathbf{x}; \boldsymbol{\theta})$ when $q_n(\mathbf{z}_n) = p(\mathbf{z}_n|\mathbf{x}_n; \boldsymbol{\theta})$ for all n . This means that the **ELBO** is a tight lower bound when the variational distribution q is the posterior distribution of \mathbf{z} given \mathbf{x} .

Now if we set $q_n(\mathbf{z}_n) = p(\mathbf{z}_n|\mathbf{x}_n; \boldsymbol{\theta}^t)$, we have that

$$g(\boldsymbol{\theta}; \{p(\mathbf{z}_n|\mathbf{x}_n; \boldsymbol{\theta}^t)|n = 1, 2, \dots, N\}) \leq \log p(\mathbf{x}; \boldsymbol{\theta}) \quad (19)$$

and

$$g(\boldsymbol{\theta}^t; \{p(\mathbf{z}_n|\mathbf{x}_n; \boldsymbol{\theta}^t)|n = 1, 2, \dots, N\}) = \log p(\mathbf{x}; \boldsymbol{\theta}^t) \quad (20)$$

3.3 EM is Bound Optimization

We are almost there! Now let's put everything together and see how to use the idea of bound optimization to derive our EM algorithm. Given all above, it's natural to define that

$$l(\boldsymbol{\theta}) := \log p(\mathbf{x}; \boldsymbol{\theta}) \quad (21)$$

$$g(\boldsymbol{\theta}; \boldsymbol{\theta}^t) := g(\boldsymbol{\theta}; \{p(\mathbf{z}_n|\mathbf{x}_n; \boldsymbol{\theta}^t)|n = 1, 2, \dots, N\}) \quad (22)$$

$$= \sum_{n=1}^N \mathbb{E}_{p(\mathbf{z}_n|\mathbf{x}_n; \boldsymbol{\theta}^t)} [\log \frac{p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta})}{p(\mathbf{z}_n|\mathbf{x}_n; \boldsymbol{\theta}^t)}] \quad (23)$$

Following the idea of bound optimization, all we need to do is to randomly pick a starting $\boldsymbol{\theta}^0$ and use Equation (3) to recursively update $\boldsymbol{\theta}^t$! Now we can see the process of EM algorithm and why it is named this way.

3.3.1 E (Expectation) step

During this step, all we need is to compute $g(\boldsymbol{\theta}; \boldsymbol{\theta}^t)$. From above we have

$$g(\boldsymbol{\theta}; \boldsymbol{\theta}^t) = \sum_{n=1}^N \mathbb{E}_{p(\mathbf{z}_n|\mathbf{x}_n; \boldsymbol{\theta}^t)} [\log \frac{p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta})}{p(\mathbf{z}_n|\mathbf{x}_n; \boldsymbol{\theta}^t)}] \quad (24)$$

$$= \sum_{n=1}^N [\mathbb{E}_{p(\mathbf{z}_n|\mathbf{x}_n; \boldsymbol{\theta}^t)} [\log p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta})] - \mathbb{E}_{p(\mathbf{z}_n|\mathbf{x}_n; \boldsymbol{\theta}^t)} p(\mathbf{z}_n|\mathbf{x}_n; \boldsymbol{\theta}^t)] \quad (25)$$

from which we can see that the second term is a constant term with respect to θ . Thus, we throw the second term and only need

$$\sum_{n=1}^N \mathbb{E}_{p(z_n|x_n;\theta^t)}[\log p(\mathbf{x}_n, \mathbf{z}_n; \theta)] \quad (26)$$

which is the sum of N **expectations** of log joint probabilities (or **complete-data likelihood**) over the posterior distributions.

3.3.2 M (Maximization) step

During this step, we want to compute

$$\theta^{t+1} = \underset{\theta}{\operatorname{argmax}} g(\theta; \theta^t) = \underset{\theta}{\operatorname{argmax}} \mathbb{E}_{p(z_n|x_n;\theta^t)}[\log p(\mathbf{x}_n, \mathbf{z}_n; \theta)] \quad (27)$$

This computation can be non-trivial. Luckily, for exponential family, we can solve this maximization by matching the moments of sufficient statistics. We will see more details in Example 4.3 and Appendix B.

4 Example in Biology: scRNA-Seq

4.1 Background Description and Problem Formalization

Single-cell RNA sequencing (scRNA-seq) allows us to profile gene expression at the resolution of individual cells. In our modeling framework, we consider a dataset $\mathbf{x} \in \mathbb{R}^{N \times d}$ where each row $\mathbf{x}_n \in \mathbb{R}^d$ corresponds to the gene expression profile of a single cell (with d genes measured). Although the true distribution of gene expression is complex, we assume for simplicity that each cell is drawn from a multivariate Gaussian distribution. This assumption, while a simplification, provides a tractable setting to illustrate the EM algorithm.

4.2 Clustering

In this subsection, we consider the problem of clustering cells into distinct groups. We assume a simplified situation in which the data are generated from a mixture of K multivariate Gaussian distributions with covariance matrices equal to the identity, i.e.,

$$\Sigma_k = I \quad \text{for all } k = 1, \dots, K.$$

Let $z_n \in \{1, 2, \dots, K\}$ denote the latent variable indicating the cluster assignment for cell n , and let π_k be the prior probability of $z_n = k$. Under this model,

the complete-data log likelihood for cell n is given by

$$\log p(\mathbf{x}_n, z_n = k; \boldsymbol{\theta}) = \log \pi_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, I) \quad (28)$$

$$= \log \pi_k - \frac{1}{2} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 + \text{const}, \quad (29)$$

where $\boldsymbol{\theta} = \{\pi_k, \boldsymbol{\mu}_k \mid k = 1, 2, \dots, K\}$.

EM Algorithm Derivation:

1. **E-step:** We first compute the posterior probabilities (also called *responsibilities*), which indicate how much the cell \mathbf{x}_n is "responsible" for cluster k :

$$\gamma_{nk}^t := p(z_n = k \mid \mathbf{x}_n; \boldsymbol{\theta}^t) = \frac{\pi_k^t \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k^t, I)}{\sum_{j=1}^K \pi_j^t \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_j^t, I)}.$$

Then, we compute the expectation of the complete-data log likelihood (ignoring constant terms):

$$\begin{aligned} g(\boldsymbol{\theta}; \boldsymbol{\theta}^t) &= \sum_{n=1}^N \mathbb{E}_{p(z_n \mid \mathbf{x}_n; \boldsymbol{\theta}^t)} [\log p(\mathbf{x}_n, z_n; \boldsymbol{\theta})] \\ &= \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk}^t \left[\log \pi_k - \frac{1}{2} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 \right]. \end{aligned} \quad (30)$$

2. **M-step:** Differentiating $g(\boldsymbol{\theta}; \boldsymbol{\theta}^t)$ with respect to $\boldsymbol{\mu}_k$ and setting the derivative to zero yields

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

and similarly, the update for the mixing coefficients is

$$\pi_k^{(t+1)} = \frac{1}{N} \sum_{n=1}^N \gamma_{nk}^t.$$

At the end of the EM algorithm, we obtain estimates $\{\hat{\boldsymbol{\mu}}_k, \hat{\pi}_k\}$ along with the posterior distribution indicating the probability of each data point belonging to each cluster:

$$p(z_n = k \mid \mathbf{x}_n; \boldsymbol{\theta}) = \gamma_{nk} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, I)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_j, I)}. \quad (31)$$

Assigning a Definite Class: What if we want to assign an exact class to each cell? We can do so by taking the class with the highest posterior probability.

In other words, we assign

$$\hat{z}_n = \arg \max_k p(z_n = k \mid \mathbf{x}_n; \boldsymbol{\theta}) = \arg \max_k \gamma_{nk}.$$

In fact, we can incorporate this hard assignment directly into the EM algorithm. That is, instead of computing the expectation in each E-step, we approximate it using the mode. This method is known as **Hard EM**. In our example, the hard EM process is to set

$$z_n^t = \arg \max_k \gamma_{nk}^t,$$

so that the expected complete-data log likelihood is approximated by

$$g(\boldsymbol{\theta}; \boldsymbol{\theta}^t) \approx \sum_{n=1}^N \sum_{k=1}^K \mathbb{I}_{\{z_n^t=k\}} \left[\log \pi_k - \frac{1}{2} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 \right], \quad (32)$$

where $\mathbb{I}_{\{z_n^t=k\}}$ is the indicator function that equals 1 if $z_n^t = k$ and 0 otherwise. In fact, if we assume further that the mixing coefficients are uniform, i.e., $\pi_k = \frac{1}{K}$ for all k , what we get is surprisingly the famous **K-Means** algorithm. The proof of this is left as practice; you can find the answer at Appendix A.

4.3 Missing Data

Now consider a different scenario where each cell is assumed to be generated from the same multivariate Gaussian distribution with parameters $\boldsymbol{\theta} = (\boldsymbol{\mu}, \Sigma)$, but only a subset of the entries in \mathbf{x}_n is observed. In other words, some entries in the data matrix \mathbf{x} are missing at random (MAR). We denote by $\mathbf{x}_n^{\text{obs}}$ and $\mathbf{x}_n^{\text{mis}}$ the observed and missing parts of \mathbf{x}_n , respectively. The complete-data likelihood is given by

$$p(\mathbf{x}_n^{\text{obs}}, \mathbf{x}_n^{\text{mis}}; \boldsymbol{\theta}) = \mathcal{N} \left(\begin{pmatrix} \mathbf{x}_n^{\text{obs}} \\ \mathbf{x}_n^{\text{mis}} \end{pmatrix} \middle| \begin{pmatrix} \boldsymbol{\mu}^{\text{obs}} \\ \boldsymbol{\mu}^{\text{mis}} \end{pmatrix}, \Sigma \right).$$

As before, our goal is to maximize the expected complete-data log-likelihood. We define

$$g(\boldsymbol{\theta}; \boldsymbol{\theta}^t) = \sum_{n=1}^N \mathbb{E}_{p(\mathbf{x}_n^{\text{mis}} \mid \mathbf{x}_n^{\text{obs}}; \boldsymbol{\theta}^t)} [\log p(\mathbf{x}_n^{\text{obs}}, \mathbf{x}_n^{\text{mis}}; \boldsymbol{\theta})]. \quad (33)$$

This computation can be non-trivial. However, if the complete-data probability is in the exponential family (e.g., Gaussian and many other common distributions), then we can exploit its structure to simplify the M-step. Unfortunately, the computation is still not easy. The basic idea is to use the result of conditional probability of multivariate Gaussians to impute the missing data $\mathbf{x}_n^{\text{mis}}$ as well as the expectation of second moment $(\mathbf{x}_n^{\text{mis}}, \mathbf{x}_n^{\text{obs}})(\mathbf{x}_n^{\text{mis}}, \mathbf{x}_n^{\text{obs}})^T$ at each

step. A detailed introduction to exponential family and how do we utilize its property for EM can be found at Appendix B.

References

- [1] Kevin P. Murphy. *Probabilistic Machine Learning: An introduction*. MIT Press, 2022. URL: <http://probml.github.io/book1>.

A Proof that Hard EM Reduces to K-Means

Now the mixing coefficients are uniform, i.e., $\pi_k = \frac{1}{K}$ for all k . Under this condition, $\log \pi_k$ is constant across all clusters and can be dropped from the maximization. The denominator can also be dropped as it is shared across different k . Then, the responsibility in the E-step becomes

$$z_n^t = \arg \max_k \gamma_{nk}^t = \arg \max_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k^t, I).$$

Recalling that the Gaussian density with identity covariance is given by

$$\mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k^t, I) \propto \exp \left(-\frac{1}{2} \|\mathbf{x}_n - \boldsymbol{\mu}_k^t\|^2 \right),$$

maximizing $\mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k^t, I)$ is equivalent to minimizing the squared Euclidean distance:

$$z_n^t = \arg \min_k \|\mathbf{x}_n - \boldsymbol{\mu}_k^t\|^2.$$

This is exactly the assignment rule in K-Means.

In the M-step, with hard assignments, the update for $\boldsymbol{\mu}_k$ becomes

$$\boldsymbol{\mu}_k^{(t+1)} = \frac{\sum_{n=1}^N \mathbb{I}_{\{z_n^t=k\}} \mathbf{x}_n}{\sum_{n=1}^N \mathbb{I}_{\{z_n^t=k\}}}.$$

This is the standard update in K-Means, where each cluster mean is recalculated as the average of all points assigned to that cluster.

Thus, under the assumptions of identity covariance and uniform mixing coefficients, the Hard EM procedure for mixture of Gaussians is equivalent to the standard K-Means algorithm:

- **E-step:** Assign each data point to the cluster with the nearest mean (minimum Euclidean distance).
- **M-step:** Update each cluster mean as the average of the data points assigned to that cluster.

B Matching the sufficient statistics for Exponential Family

In general, a distribution in the exponential family can be written as

$$p(\mathbf{x}; \boldsymbol{\theta}) = h(\mathbf{x}) \exp \left(\boldsymbol{\theta}^\top \mathcal{T}(\mathbf{x}) - A(\boldsymbol{\theta}) \right), \quad (34)$$

where:

- $\boldsymbol{\theta}$ is the *natural parameter*,
- $\mathcal{T}(\mathbf{x})$ is the *sufficient statistic*,
- $A(\boldsymbol{\theta})$ is the *log-partition function* (or cumulant generating function),
- $h(\mathbf{x})$ is the base measure.

A key property of the log-partition function is that its first derivative with respect to $\boldsymbol{\theta}$ equals the expectation of the sufficient statistic:

$$\frac{\partial A(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbb{E}_{p(\mathbf{x}; \boldsymbol{\theta})} [\mathcal{T}(\mathbf{x})].$$

Exploiting this property, we can rewrite the expected complete-data log-likelihood as

$$g(\boldsymbol{\theta}; \boldsymbol{\theta}^t) = \sum_{n=1}^N \left(\mathbb{E}_{q_n(\mathbf{z}_n; \boldsymbol{\theta}^t)} [\mathcal{T}(\mathbf{x}_n, \mathbf{z}_n)]^\top \boldsymbol{\theta} - A(\boldsymbol{\theta}) \right), \quad (35)$$

where $q_n(\mathbf{z}_n; \boldsymbol{\theta}^t) = p(\mathbf{z}_n \mid \mathbf{x}_n^{\text{obs}}; \boldsymbol{\theta}^t)$ is the posterior of the latent variables given the observed data under the current parameter estimates.

Taking the derivative of $g(\boldsymbol{\theta}; \boldsymbol{\theta}^t)$ with respect to $\boldsymbol{\theta}$ and setting it to zero yields

$$\frac{\partial g(\boldsymbol{\theta}; \boldsymbol{\theta}^t)}{\partial \boldsymbol{\theta}} = \sum_{n=1}^N \left(\mathbb{E}_{q_n(\mathbf{z}_n; \boldsymbol{\theta}^t)} [\mathcal{T}(\mathbf{x}_n, \mathbf{z}_n)] - \frac{\partial A(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right) = 0.$$

Using the property of $A(\boldsymbol{\theta})$, this condition is equivalent to

$$\frac{1}{N} \sum_{n=1}^N \mathbb{E}_{q_n(\mathbf{z}_n; \boldsymbol{\theta}^t)} [\mathcal{T}(\mathbf{x}_n, \mathbf{z}_n)] = \mathbb{E}_{p(\mathbf{x}; \boldsymbol{\theta})} [\mathcal{T}(\mathbf{x})]. \quad (36)$$

For a multivariate Gaussian distribution, the sufficient statistics are

$$\mathcal{T}(\mathbf{x}) = (\mathbf{x}, \mathbf{x}\mathbf{x}^\top)$$

Plugging this to the right side Equation (36), we have that

$$\mathbb{E}_{p(\mathbf{x};\boldsymbol{\theta})}[\mathcal{T}(\mathbf{x})] = \mathbb{E}_{p(\mathbf{x};\boldsymbol{\theta})}[(\mathbf{x}, \mathbf{x}\mathbf{x}^T)] = (\boldsymbol{\mu}, \boldsymbol{\Sigma} + \boldsymbol{\mu}\boldsymbol{\mu}^T) \quad (37)$$

Now we need to compute the left side of Equation (36) to update $\boldsymbol{\theta} = (\boldsymbol{\mu}, \boldsymbol{\Sigma})$. We need to use the result of the conditional probability for a multivariate Gaussian. Assume

$$p(\mathbf{x}_n^{\text{mis}}, \mathbf{x}_n^{\text{obs}} \mid \boldsymbol{\theta}^t) = \mathcal{N}\left(\begin{pmatrix} \boldsymbol{\mu}_t^{\text{mis}} \\ \boldsymbol{\mu}_t^{\text{obs}} \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_t^{mm} & \boldsymbol{\Sigma}_t^{mo} \\ \boldsymbol{\Sigma}_t^{om} & \boldsymbol{\Sigma}_t^{oo} \end{pmatrix}\right). \quad (38)$$

Then, by standard results on the conditional distribution of a multivariate Gaussian, we have

$$p(\mathbf{x}_n^{\text{mis}} \mid \mathbf{x}_n^{\text{obs}}, \boldsymbol{\theta}^t) = \mathcal{N}(\mathbf{m}_n^t, \mathbf{V}_n^t) \quad (39)$$

$$\mathbf{m}_n^t := \boldsymbol{\mu}_t^{\text{mis}} + \boldsymbol{\Sigma}_t^{mo} (\boldsymbol{\Sigma}_t^{oo})^{-1} (\mathbf{x}_n^{\text{obs}} - \boldsymbol{\mu}_t^{\text{obs}}) \quad (40)$$

$$\mathbf{V}_n^t := \boldsymbol{\Sigma}_t^{mm} - \boldsymbol{\Sigma}_t^{mo} (\boldsymbol{\Sigma}_t^{oo})^{-1} \boldsymbol{\Sigma}_t^{om} \quad (41)$$

Now we can compute the left side of (36)

$$\sum_{n=1}^N \mathbb{E}_{p(\mathbf{x}_n^{\text{mis}} \mid \mathbf{x}_n^{\text{obs}}, \boldsymbol{\theta}^t)} [\mathcal{T}(\mathbf{x}_n^{\text{mis}}, \mathbf{x}_n^{\text{obs}})] \quad (42)$$

$$= \sum_{n=1}^N \mathbb{E}_{p(\mathbf{x}_n^{\text{mis}} \mid \mathbf{x}_n^{\text{obs}}, \boldsymbol{\theta}^t)} \left[\left((\mathbf{x}_n^{\text{mis}}, \mathbf{x}_n^{\text{obs}}), (\mathbf{x}_n^{\text{mis}}, \mathbf{x}_n^{\text{obs}})(\mathbf{x}_n^{\text{mis}}, \mathbf{x}_n^{\text{obs}})^T \right) \right] \quad (43)$$

$$= \sum_{n=1}^N \left((\mathbf{m}_n^t, \mathbf{x}_n^{\text{obs}}), \begin{pmatrix} \mathbf{V}_n^t + \mathbf{m}_n^t (\mathbf{m}_n^t)^T & \mathbf{m}_n^t (\mathbf{x}_n^{\text{obs}})^T \\ \mathbf{x}_n^{\text{obs}} (\mathbf{m}_n^t)^T & \mathbf{x}_n^{\text{obs}} (\mathbf{x}_n^{\text{obs}})^T \end{pmatrix} \right) \quad (44)$$

Thus, to update $\boldsymbol{\theta}^{t+1} = (\boldsymbol{\mu}^{t+1}, \boldsymbol{\Sigma}^{t+1})$, we have

$$\boldsymbol{\mu}^{t+1} = \frac{1}{N} \sum_{n=1}^N \text{Reordered}(\mathbf{m}_n^t, \mathbf{x}_n^{\text{obs}}), \quad (45)$$

$$\boldsymbol{\Sigma}^{t+1} = \frac{1}{N} \sum_{n=1}^N \text{Reordered} \left(\begin{pmatrix} \mathbf{V}_n^t + \mathbf{m}_n^t (\mathbf{m}_n^t)^T & \mathbf{m}_n^t (\mathbf{x}_n^{\text{obs}})^T \\ \mathbf{x}_n^{\text{obs}} (\mathbf{m}_n^t)^T & \mathbf{x}_n^{\text{obs}} (\mathbf{x}_n^{\text{obs}})^T \end{pmatrix} - \boldsymbol{\mu}^{t+1} (\boldsymbol{\mu}^{t+1})^T \right). \quad (46)$$

The Reordered operation here is to make sure that variables are in the same order because different n may have different set of missing variables.

For more details on the exponential family, see https://en.wikipedia.org/wiki/Exponential_family or Section I.3.4 of [1].