



Introduction to Machine Learning

# Project Phase 1

Mohammad Pouya Toroghi

Ilia Hashemi Rad

# Theory Questions

## Question 1

### Solution:

The **MM (Majorization-Minimization) algorithm** is designed to handle complex non-convex optimization problems by iteratively simplifying the problem using convex surrogate functions. The key idea is to replace the original, often non-convex, objective function  $l(\theta)$  with a simpler surrogate function  $Q(\theta, \theta^{(t)})$  at each iteration. This surrogate function satisfies the following conditions:

- **Tangency Condition:**  $Q(\theta^{(t)}, \theta^{(t)}) = l(\theta^{(t)})$ , meaning the surrogate and the original function coincide at the current parameter estimate.
- **Lower Bound Condition:**  $Q(\theta, \theta^{(t)}) \leq l(\theta)$  for all  $\theta$ , ensuring that the surrogate underestimates the original function.

These properties ensure that the MM algorithm produces a sequence of iterates  $\theta^{(t)}$  where each update improves the original objective function by solving a simpler problem. The update rule is:

$$\theta^{(t+1)} = \arg \max_{\theta} Q(\theta, \theta^{(t)}),$$

and this guarantees that the objective function  $l(\theta)$  increases monotonically at each step.

**Handling Non-Convexity:** For non-convex objective functions, directly maximizing  $l(\theta)$  can be difficult because of multiple local maxima. The MM algorithm tackles this by optimizing a surrogate function  $Q(\theta, \theta^{(t)})$ , which is often chosen to be convex. Convex functions are much easier to optimize because they have a unique global maximum, allowing the algorithm to avoid the difficulties associated with non-convex optimization.

Although the MM algorithm guarantees convergence to a local maximum, it may not reach the global maximum in the case of non-convex objectives. However, by carefully choosing surrogate functions that are convex, the algorithm simplifies the optimization problem, making it easier to find a good solution iteratively.

**Flexibility in Surrogate Function Selection:** The MM algorithm allows flexibility in constructing the surrogate function  $Q(\theta, \theta^{(t)})$ . Various techniques, such as Jensen's inequality, Cauchy-Schwarz inequality, or a second-order Taylor expansion, can be used to derive convex surrogates. This flexibility enables the algorithm to be adapted to different types of problems, improving the chances of efficient and effective optimization.

**Conclusion:** The MM algorithm handles non-convex optimization by iteratively constructing convex or simpler surrogate functions, optimizing them, and ensuring monotonic improvement in the original objective. While convergence is to a local maximum in non-convex settings, the flexibility in choosing surrogates and the ease of solving convex problems make the MM algorithm a powerful approach for complex optimization tasks.

## Question 2

### Solution:

The mixture model formula can be written as:

$$p(y; \theta) = \sum_{k=1}^K p_Z(z_k; \theta) p_{Y|Z}(y|Z = z_k; \theta)$$

This formula expresses the marginal likelihood of  $y$  by summing over all possible values of the latent variable  $Z = z_k$ . In a dataset with  $N$  observations, the likelihood for all data points  $\mathbf{y} = \{y^{(1)}, y^{(2)}, \dots, y^{(N)}\}$  becomes:

$$p_Y(\mathbf{y}; \theta) = \prod_{i=1}^N p_Y(y^{(i)}; \theta) = \prod_{i=1}^N \left( \sum_{k=1}^K p_{Y,Z}(y^{(i)}, Z^{(i)} = k; \theta) \right)$$

We can decompose the joint probability  $p_{Y,Z}(y^{(i)}, Z^{(i)} = k; \theta)$  using:

$$p_{Y,Z}(y^{(i)}, Z^{(i)} = k; \theta) = p_Z(z_k; \theta) p_{Y|Z}(y^{(i)} | Z^{(i)} = k; \theta)$$

Thus, the overall likelihood becomes:

$$p_Y(\mathbf{y}; \theta) = \prod_{i=1}^N \left( \sum_{k=1}^K p_Z(z_k; \theta) p_{Y|Z}(y^{(i)} | Z^{(i)} = k; \theta) \right)$$

This is equivalent to summing over all possible values of the latent variable  $Z^{(i)}$ . Now, the reason why it is easier to optimize  $p_{Y,Z}(y_n, z_n; \theta)$  rather than  $p_Y(y_n; \theta)$  is as follows:

Directly optimizing  $p_Y(y_n; \theta)$  is difficult because it involves summing over all possible latent variables  $Z$ , which makes the marginal likelihood complex and harder to optimize. On the other hand, optimizing the joint likelihood  $p_{Y,Z}(y_n, z_n; \theta)$  is simpler because:

1. The latent variables  $Z$  are treated as observed when optimizing the joint likelihood. This simplifies the problem by avoiding the need to marginalize over  $Z$ .
2. The Expectation-Maximization (EM) algorithm is commonly used to handle this. In the E-step, we compute the expected value of the latent variables, and in the M-step, we maximize the joint likelihood  $p_{Y,Z}(y_n, z_n; \theta)$ , which is easier to handle computationally.

Thus, optimizing  $p_{Y,Z}(y_n, z_n; \theta)$  is easier than optimizing  $p_Y(y_n; \theta)$  because we avoid the complexities of marginalizing over the latent variables.

### Question 3

#### Solution:

Variational Inference (VI) and the Expectation-Maximization (EM) algorithm are both used to approximate a true posterior distribution, but they employ different approaches and assumptions.

The EM algorithm iteratively maximizes the likelihood of observed data through two steps: the Expectation (E-step) and the Maximization (M-step). In the E-step, the algorithm computes the expected values of the latent variables, given the current estimates of the model parameters. In the M-step, the parameters are updated to maximize the expected log-likelihood of the observed data. EM assumes that latent variables are unobserved, and the goal is to estimate both the latent variables and the parameters. EM is widely used for problems where the direct optimization of the likelihood is difficult due to missing or incomplete data.

On the other hand, Variational Inference (VI) is a deterministic optimization method that minimizes the Kullback-Leibler (KL) divergence between the true posterior distribution and a simpler, tractable distribution (such as a Gaussian). Unlike EM, VI does not estimate specific latent variables but instead estimates an entire distribution over them. This is achieved by approximating the posterior distribution over the latent variables and parameters, allowing for fully Bayesian estimation. VI can thus be seen as a generalization of EM, moving from Maximum A Posteriori (MAP) estimation, which finds the most probable value of each parameter, to approximating the entire posterior distribution.

A key difference is that while EM provides point estimates for parameters, VI computes a full distribution, making it more flexible and applicable to complex probabilistic models where a full Bayesian treatment is desired. However, the trade-off is that VI requires specifying an approximating family of distributions, which introduces additional assumptions and might limit its flexibility depending on the choice of the distribution.

## Question 4

### Solution: 1:

Assume that we have a dataset of  $N$  observed  $d$ -dimensional points  $\{\mathbf{x}_i\}_{i=1}^N$ . The dataset follows a mixture of  $K$  Multivariate Normal distributions such that:

$$p(\mathbf{x}|\Theta) = \sum_k \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

where  $\pi_k$ ,  $\boldsymbol{\mu}_k$ , and  $\boldsymbol{\Sigma}_k$  are the prior probability, mean vector, and covariance matrix of the  $k$ -th Multivariate Gaussian model, respectively.

Our task is to estimate the set of parameters  $(\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  for each Multivariate Gaussian model using the dataset.

Clearly, in the first step, we should initialize the mentioned parameters. Here, we choose an algorithm to do so. In this algorithm, we initialize the parameters  $\pi_k$  uniformly:

$$\pi_k^0 = \frac{1}{K} \quad \forall k; \quad 1 \leq k \leq K$$

which means each point in the dataset is equally likely to belong to the  $k$ -th Multivariate Gaussian model. To initialize the Gaussian models' parameters  $\boldsymbol{\mu}_k$  and  $\boldsymbol{\Sigma}_k$ , we first divide the dataset randomly into  $K$  clusters with  $M = \frac{N}{K}$  data points. (Note that if the value of  $M$  is not an integer, we would divide the dataset into clusters such that all the clusters have approximately the same sizes.) Then, using the MLE method, we have:

$$\begin{cases} \boldsymbol{\mu}_k^0 = \frac{1}{M} \sum_{j=1}^M \mathbf{x}_j^{(k)} \\ \boldsymbol{\Sigma}_k^0 = \frac{1}{M} \sum_{j=1}^M (\mathbf{x}_j^{(k)} - \boldsymbol{\mu}_k^0)(\mathbf{x}_j^{(k)} - \boldsymbol{\mu}_k^0)^T \end{cases} \quad \forall k; \quad 1 \leq k \leq K$$

where  $\mathbf{x}_j^{(k)}$  denotes the  $j$ -th point of the  $k$ -th cluster.

### Solution: 2:

In this section, we are supposed to determine the likelihood of the complete dataset. According to the descriptions in the previous part, we have:

$$\begin{aligned} p(\mathcal{D}|\Theta) &= \prod_i p(\mathbf{x}_i, z_i|\Theta) = \prod_i \sum_k p(\mathbf{x}_i, z_i = k|\Theta) = \prod_i \sum_k p(z_i = k) p(\mathbf{x}_i|z_i = k, \Theta) = \prod_i \sum_k \pi_k \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \\ &= \prod_i \sum_k p(\mathbf{x}_i, z_i = k|\Theta) \frac{\pi_k \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{p(\mathbf{x}_i, z_i = k|\Theta)} = \prod_i \sum_k \omega_{i,k} \frac{\pi_k \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\omega_{i,k}} \end{aligned}$$

In which, using Bayes' Rule, we have:

$$\omega_{i,k} \triangleq p^*(z_i = k|\mathbf{x}_i, \Theta) = \frac{p(z_i = k) p(\mathbf{x}_i|z_i = k, \Theta)}{\sum_{k'} p(z_i = k') p(\mathbf{x}_i|z_i = k', \Theta)} = \frac{\pi_k \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k'} \pi_{k'} \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_{k'}, \boldsymbol{\Sigma}_{k'})}$$

Therefore, the complete log-likelihood function is determined as follows (this part is for the next section):

$$L(\mathcal{D}|\Theta) = \sum_i \log \sum_k \omega_{i,k} \frac{\pi_k \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\omega_{i,k}} \geq \sum_i \sum_k \omega_{i,k} \log \frac{\pi_k \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\omega_{i,k}} \quad *$$

**Solution:** 3:

Assume that the latent parameters at iteration  $t$  are defined as calculated previously:

$$\omega_{i,k}^t \triangleq p^*(z_i = k | \mathbf{x}_i, \Theta^t) = \frac{\pi_k^t \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k^t, \boldsymbol{\Sigma}_k^t)}{\sum_{k'} \pi_{k'}^t \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_{k'}^t, \boldsymbol{\Sigma}_{k'}^t)}.$$

Using the results from the previous parts and the reference book, we can define an instance iterative target function (the lower bound function) to work with:

$$\begin{aligned} Q(\Theta, \Theta^t) &= \sum_i \sum_k \omega_{i,k}^t \log \frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\omega_{i,k}^t} \\ &= \sum_i \sum_k \omega_{i,k}^t \log \left( \frac{\pi_k}{\omega_{i,k}^t \sqrt{(2\pi)^d |\boldsymbol{\Sigma}_k|}} \exp \left[ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right] \right). \end{aligned}$$

This function obviously satisfies the two constraints:  $Q(\Theta, \Theta^t) \leq L(\mathcal{D} | \Theta)$  (based on  $*$ ) and  $Q(\Theta^t, \Theta^t) = L(\mathcal{D} | \Theta^t)$  (based on the proof in the source).

We expand the result further to use it in subsequent steps:

$$Q(\Theta, \Theta^t) = \sum_i \sum_k \omega_{i,k}^t \left[ \log \pi_k - \log \omega_{i,k}^t - \frac{d}{2} \log 2\pi - \frac{1}{2} \log |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right]. \quad (1)$$

Now we proceed with the E-step and M-step:

1. **E-step:** The latent parameters to be estimated in each iteration of the E-step are determined in closed form:

$$\omega_{i,k}^t = \frac{\pi_k^t \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k^t, \boldsymbol{\Sigma}_k^t)}{\sum_{k'} \pi_{k'}^t \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_{k'}^t, \boldsymbol{\Sigma}_{k'}^t)},$$

where  $\omega$  is an  $N \times K$  matrix.

2. **M-step:** In this step, we update the  $\Theta$  parameters by maximizing  $Q(\Theta, \Theta^t)$  with respect to the parameters:

$$\Theta^{t+1} = \arg \max_{\Theta} Q(\Theta, \Theta^t).$$

Using the expanded equation in (1), we conclude:

$$\Theta^{t+1} = \arg \max_{\Theta} \sum_i \sum_k \omega_{i,k}^t \left[ \log \pi_k - \log \omega_{i,k}^t - \frac{d}{2} \log 2\pi - \frac{1}{2} \log |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right].$$

- Updating  $\boldsymbol{\mu}_k$ :

$$\boldsymbol{\mu}_k^{t+1} = \arg \max_{\boldsymbol{\mu}_k} Q(\Theta, \Theta^t).$$

To do so, by taking the partial derivative of  $Q(\Theta, \Theta^t)$  with respect to  $\boldsymbol{\mu}_k$ , we have:

$$\left. \frac{\partial Q(\Theta, \Theta^t)}{\partial \boldsymbol{\mu}_k} \right|_{\boldsymbol{\mu}_k = \boldsymbol{\mu}_k^{t+1}} = \sum_i \omega_{i,k}^t \left. \frac{\partial \left[ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right]}{\partial \boldsymbol{\mu}_k} \right|_{\boldsymbol{\mu}_k = \boldsymbol{\mu}_k^{t+1}} = 0. \quad (2)$$

On the other hand, it can be shown that if  $\mathbf{Y}$  is a symmetric matrix, for any two vectors  $\mathbf{x}$  and  $\mathbf{a}$  we have:

$$\frac{\partial (\mathbf{x} - \mathbf{a})^T \mathbf{Y} (\mathbf{x} - \mathbf{a})}{\partial \mathbf{x}} = -2\mathbf{Y} (\mathbf{x} - \mathbf{a}).$$

So, using the fact that the covariance matrix and consequently its inverse are symmetric, from (2) we conclude:

$$\sum_i \omega_{i,k}^t \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k^{t+1}) = 0 \Rightarrow \sum_i \omega_{i,k}^t \mathbf{x}_i = \boldsymbol{\mu}_k^{t+1} \sum_i \omega_{i,k}^t$$

$$\Rightarrow \mu_k^{t+1} = \frac{\sum_i \omega_{i,k}^t x_i}{\sum_i \omega_{i,k}^t}.$$

- Updating  $\Sigma_k$ :

$$\Sigma_k^{t+1} = \arg \max_{\Sigma_k} Q(\Theta, \Theta^t).$$

To do so, by taking the partial derivative of  $Q(\Theta, \Theta^t)$  with respect to  $\Sigma_k^{-1}$ , we have:

$$\left. \frac{\partial Q(\Theta, \Theta^t)}{\partial \Sigma_k^{-1}} \right|_{\Sigma_k = \Sigma_k^{t+1}} = \sum_i \omega_{i,k}^t \left. \frac{\partial \left[ -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (\mathbf{x}_i - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \mu_k) \right]}{\partial \Sigma_k^{-1}} \right|_{\Sigma_k = \Sigma_k^{t+1}} = 0. \quad (3)$$

On the other hand, for the symmetric matrices  $\Sigma$  and its inverse and for any vector  $a$ , we have:

$$\frac{\partial \log |\Sigma|}{\partial \Sigma^{-1}} = \frac{\partial \log |\Sigma|}{\partial |\Sigma|} \frac{\partial |\Sigma|}{\partial \Sigma^{-1}} = \frac{1}{|\Sigma|} (-|\Sigma| \Sigma^T) = -\Sigma^T,$$

$$\frac{\partial \mathbf{a}^T \Sigma^{-1} \mathbf{a}}{\partial \Sigma^{-1}} = (\mathbf{x} - \mu_k)(\mathbf{x} - \mu_k)^T.$$

So, using these results, from (3) we conclude:

$$\begin{aligned} \sum_i \omega_{i,k}^t \left[ \Sigma_k^{t+1} - (\mathbf{x}_i - \mu_k)(\mathbf{x}_i - \mu_k)^T \right] &= 0 \Rightarrow \sum_i \omega_{i,k}^t \Sigma_k^{t+1} = \sum_i \omega_{i,k}^t (\mathbf{x}_i - \mu_k)(\mathbf{x}_i - \mu_k)^T \\ \Rightarrow \Sigma_k^{t+1} &= \frac{\sum_i \omega_{i,k}^t (\mathbf{x}_i - \mu_k)(\mathbf{x}_i - \mu_k)^T}{\sum_i \omega_{i,k}^t}. \end{aligned}$$

- Updating  $\pi_k$ : Similarly, using the Lagrange multiplier method and the constraint  $\sum_k \pi_k = 1$ , it can be shown that:

$$\pi_k^{t+1} = \frac{1}{N} \sum_i \omega_{i,k}^t.$$

## Question 5

**Solution:** 1:

Assume we have a dataset of  $N$  observed  $d$ -dimensional points  $\{\mathbf{x}_i\}_{i=1}^N$ . The dataset follows a mixture of  $K$  Categorical distributions such that:

$$p(\mathbf{x}|\Theta) = \sum_k \pi_k \text{Cat}(\mathbf{x}|\theta_k)$$

where  $\pi_k$  and  $\theta_k$  are the prior probability and the parameter vector of the  $k$ th Categorical model. Each  $\theta = (\theta_{k,1}, \dots, \theta_{k,d})$  is a probability distribution over the labels  $\mathcal{A} = \{1, \dots, d\}$ .

For example, if  $x_i = c$  with  $1 \leq c \leq d$ , we can say  $p(x_i|\theta_k) = \theta_{k,c}$ .

Since each data point is  $d$ -dimensional, we can represent  $\mathbf{x}_i$  as a  $d$ -dimensional vector. For instance, if  $x_i = 3$ , we can express it by turning on the third component of the data point vector ( $\mathbf{x}_i = (0, 0, 1, 0, \dots, 0)$  and  $x_{i,3} = 1$ ). Our task is to estimate the set of parameters  $(\pi_k, \theta_k)$  for each Categorical model using the dataset. Clearly, the first step is to initialize these parameters. Here, we choose an algorithm to do so.

In this algorithm, we initialize the parameters  $\pi_k$  uniformly:

$$\pi_k^0 = \frac{1}{K} \quad \forall k; \quad 1 \leq k \leq K$$

This means each point in the dataset is equally likely to belong to the  $k$ th Categorical model.

To initialize the Categorical models' parameter vector  $\theta_k$ , we first divide the dataset randomly into  $K$  clusters with  $M = \frac{N}{K}$  data points. (If  $M$  is not an integer, we divide the dataset into clusters such that all clusters have approximately the same sizes.) Then, using the MLE method, we have:

$$\theta_{k,j}^0 = \frac{N_j^{(k)}}{M} \quad \forall k; \quad 1 \leq k \leq K$$

where  $N_j^{(k)}$  denotes the number of times the  $j$ th component of the data point vectors in the  $k$ th cluster is turned on. Equivalently,

$$\theta_k^0 = \frac{1}{M} \sum_{j=1}^M \mathbf{x}_i^{(k)}$$

### Solution: 2:

In this section, we aim to determine the complete dataset likelihood. According to the previous part, we have:

$$\begin{aligned} p(\mathcal{D}|\Theta) &= \prod_i p(\mathbf{x}_i, z_i|\Theta) = \prod_i \sum_k p(\mathbf{x}_i, z_i = k|\Theta) = \prod_i \sum_k p(z_i = k) p(\mathbf{x}_i|z_i = k, \Theta) = \prod_i \sum_k \pi_k \text{Cat}(\mathbf{x}_i|\theta_k) \\ &= \prod_i \sum_k p(\mathbf{x}_i, z_i = k|\Theta) \frac{\pi_k \text{Cat}(\mathbf{x}_i|\theta_k)}{p(\mathbf{x}_i, z_i = k|\Theta)} = \prod_i \sum_k \omega_{i,k} \frac{\pi_k \text{Cat}(\mathbf{x}_i|\theta_k)}{\omega_{i,k}} \end{aligned}$$

where, using Bayes' Rule, we have:

$$\omega_{i,k} \triangleq p^*(z_i = k|\mathbf{x}_i, \Theta) = \frac{p(z_i = k) p(\mathbf{x}_i|z_i = k, \Theta)}{\sum_{k'} p(z_i = k') p(\mathbf{x}_i|z_i = k', \Theta)} = \frac{\pi_k \text{Cat}(\mathbf{x}_i|\theta_k)}{\sum_{k'} \pi_{k'} \text{Cat}(\mathbf{x}_i|\theta_{k'})}$$

Therefore, the complete log likelihood function is determined as follows (this part is for the next section):

$$L(\mathcal{D}|\Theta) = \sum_i \log \sum_k \omega_{i,k} \frac{\pi_k \text{Cat}(\mathbf{x}_i|\theta_k)}{\omega_{i,k}} \geq \sum_i \sum_k \omega_{i,k} \log \frac{\pi_k \text{Cat}(\mathbf{x}_i|\theta_k)}{\omega_{i,k}} \quad *$$

### Solution: 3:

Assume that the latent parameters at iteration  $t$  are defined as calculated before:

$$\omega_{i,k}^t \triangleq p^*(z_i = k|\mathbf{x}_i, \Theta^t) = \frac{\pi_k^t \text{Cat}(\mathbf{x}_i|\theta_k^t)}{\sum_{k'} \pi_{k'}^t \text{Cat}(\mathbf{x}_i|\theta_{k'}^t)}$$

Using the results from previous parts and the source book, we define an iterative target function (the lower bound function):

$$\begin{aligned} Q(\Theta, \Theta^t) &= \sum_i \sum_k \omega_{i,k}^t \log \frac{\pi_k \text{Cat}(\mathbf{x}_i|\theta_k)}{\omega_{i,k}^t} \\ &= \sum_i \sum_k \omega_{i,k}^t \log \left( \frac{\pi_k \prod_{j=1}^d (\theta_{k,j})^{x_{i,j}}}{\omega_{i,k}^t} \right) \end{aligned}$$

This function satisfies two constraints:  $Q(\Theta, \Theta^t) \leq L(\mathcal{D}|\Theta)$  (based on  $*$ ) and  $Q(\Theta^t, \Theta^t) = L(\mathcal{D}|\Theta^t)$  (as proven in the source).

We expand the result further for use in the next steps:

$$Q(\Theta, \Theta^t) = \sum_i \sum_k \omega_{i,k}^t \left[ \log \pi_k - \log \omega_{i,k}^t + \underbrace{\sum_{j=1}^d x_{i,j} \log \theta_{k,j}}_{\mathbf{x}_i \cdot \log \theta_k} \right] \quad (1)$$

Now, we proceed to the E-step and M-step:

1. **E-step:** The latent parameters to be estimated in each iteration of the E-step are determined in a closed form:

$$\omega_{i,k}^t = \frac{\pi_k^t \text{Cat}(\mathbf{x}_i | \boldsymbol{\theta}_k^t)}{\sum_{k'} \pi_{k'}^t \text{Cat}(\mathbf{x}_i | \boldsymbol{\theta}_{k'}^t)}$$

where  $\omega$  is an  $N \times K$  matrix.

2. **M-step:** In the M-step, we update the  $\Theta$  parameters by maximizing  $Q(\Theta, \Theta^t)$ :

$$\Theta^{t+1} = \arg \max_{\Theta} Q(\Theta, \Theta^t)$$

Using the expanded equation in (1), we get:

$$\Theta^{t+1} = \arg \max_{\Theta} \sum_i \sum_k \omega_{i,k}^t \left[ \log \pi_k - \log \omega_{i,k}^t + \mathbf{x}_i \cdot \log \boldsymbol{\theta}_k \right]$$

$\pi_k$  and  $\boldsymbol{\theta}_k$  can be estimated separately:

- (a) Update prior probabilities  $\pi_k$ : We assume the Lagrange multiplier technique to preserve the constraint  $\sum_k \pi_k = 1$ . Hence,

$$\pi_k^{t+1} = \frac{1}{N} \sum_i \omega_{i,k}^t$$

- (b) Update Categorical parameters  $\boldsymbol{\theta}_k$ : Similarly, we use the Lagrange multiplier to preserve the constraint  $\sum_j \theta_{k,j} = 1$ . Hence,

$$\theta_{k,j}^{t+1} = \frac{\sum_i \omega_{i,k}^t x_{i,j}}{\sum_i \omega_{i,k}^t}$$

Finally, the algorithm iterates between these steps until convergence. Thus, the mixture of Categorical distributions model is successfully estimated using the EM algorithm.