

## Multivariate Gaussian mixture model

We use multivariate Gaussian mixture model to describe the expression of genes. The joint probability of gene expression value is

$$p(g_1, \dots, g_D) = \sum_{k=1}^K \pi_k N(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

where  $\{g_1, \dots, g_D\}$  are genes,  $\pi_k$  is the weight of the  $k$ th component and  $\sum_{k=1}^K \pi_k = 1$ ,  $\boldsymbol{\mu}_k$  and  $\boldsymbol{\Sigma}_k$  are the mean vector and the covariance matrix of the  $k$ th component respectively .

## Parameter estimation

Suppose we have a set of genes  $\{g_1, \dots, g_D\}$  in one module and there are  $N$  expression data points  $\{\mathbf{e}_1, \dots, \mathbf{e}_N\}$ , this data set can be represented as an  $N \times D$  matrix  $\mathbf{E}$ . We assume the  $N$  data points are independent from the same multivariate Gaussian mixture distribution. The log-likelihood of the observed data is :

$$\ln p(\mathbf{E} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(\mathbf{e}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

The training process is to find the maximum likelihood estimate of  $(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ . An elegant and powerful method for handling this task is the *Expectation – Maximization* (EM) algorithms.

**E step:** Evaluate the responsibilities using the current parameter values

$$\gamma(z_{nk}) = \frac{\pi_k N(\mathbf{e}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j N(\mathbf{e}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

$\gamma(z_{nk})$  is the posterior probability that component  $k$  is responsible for generating  $\mathbf{e}_n$ .

**M step:** Re-estimate the parameters using the current responsibilities

$$\begin{aligned} \boldsymbol{\mu}_k^{\text{new}} &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{e}_n \\ \boldsymbol{\Sigma}_k^{\text{new}} &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{e}_n - \boldsymbol{\mu}_k^{\text{new}})(\mathbf{e}_n - \boldsymbol{\mu}_k^{\text{new}})^T \\ \pi_k^{\text{new}} &= \frac{N_k}{N} \end{aligned}$$

where

$$N_k = \sum_{n=1}^N \gamma(z_{nk})$$