

# Gaussian Mixture Model

Goal: Approximate the distribution using the data

## Gaussian mixture distribution

The distribution has  $k$  clusters, cluster  $j$  has probability  $\pi_j$ .

$\text{Prob}\{x \in \text{cluster } j\} = \pi_j$ ,  $j=1, 2, \dots, k$ , where  $\sum_{j=1}^k \pi_j = 1$

Define one-hot vectors,  $r_j \in \{0, 1\}^k$ , where  $j=1, 2, \dots, k$  to

identify the clusters:  $(r_j)_I = \delta_{ji}$ ,  $I=1, 2, \dots, k$

$x \in \text{cluster } j \Leftrightarrow r = r_j$ ,  $j=1, 2, \dots, k$

Within each cluster, we assume the distribution is Gaussian  $p(x|r=r_j) = p_g(x; z_j, \Sigma_j)$

where  $p_g(x; z_j, \Sigma_j) = (2\pi)^{-d/2} |\Sigma_j|^{-1/2} e^{-\frac{1}{2}(x-z_j)^T \Sigma_j^{-1}(x-z_j)}$

• Law of total probability:  $p(x) = \sum_{j=1}^k p(x|r=r_j) p(r_j) = \sum_{j=1}^k \pi_j p_g(x; z_j, \Sigma_j)$

parameter estimation (observed data  $\{x_i\}_{i=1}^N \rightarrow \{\pi_j, z_j, \Sigma_j\}_{j=1}^k$ )

Log-likelihood:  $l(\{\pi_j, z_j, \Sigma_j\}) = \log p(x_1, \dots, x_N)$

$$= \log \prod_{i=1}^N p(x_i) = \sum_{i=1}^N \log \left[ \sum_{j=1}^k \pi_j p_g(x_i; z_j, \Sigma_j) \right]$$

$$l = \sum_{i=1}^N \log \left[ \sum_{j=1}^k \pi_j p_g(x_i; z_j, \Sigma_j) \right]$$

$$0 = \nabla_{z_j} l = - \sum_{i=1}^N \left[ \frac{\pi_j p_g(x_i; z_j, \Sigma_j)}{\sum_{e=1}^k \pi_e p_g(x_i; z_e, \Sigma_e)} \right] \cdot \Sigma_j^{-1} (x_i - z_j)$$

prior:  $p(r=r_j) \cdot p(x_i|r=r_j)$  conditional

$$\text{where } r_{ij} = \frac{\pi_j p_g(x_i; z_j, \Sigma_j)}{\sum_{e=1}^k \pi_e p_g(x_i; z_e, \Sigma_e)} = \underset{\substack{\uparrow \\ \text{posterior}}}{p(r=r_j | x_i)}$$

$$\downarrow$$

$$\sum_{e=1}^k p(x, x_i) = p(x_i) \quad \text{marginal}$$

$r_{ij}$ : Probability that  $x_i \in \text{cluster } j$  after observing the data point.

$$\left( \sum_{j=1}^k r_{ij} = 1 \right)$$

Cont.

$$0 = \nabla_{z_j} \ell = - \sum_{i=1}^N r_{ij} \cdot \Sigma_j^{-1} (\underline{x}_i - \underline{z}_j)$$

$$0 = \Sigma_j^{-1} \sum_{i=1}^N r_{ij} (\underline{x}_i - \underline{z}_j)$$

$$0 = \sum_{i=1}^N r_{ij} \underline{x}_i - \left( \sum_{i=1}^N r_{ij} \right) \underline{z}_j$$

$$\Rightarrow \underline{z}_j = \frac{\sum_{i=1}^N r_{ij} \underline{x}_i}{\sum_{i=1}^N r_{ij}} \quad (\text{weighted average of all samples})$$

$\propto$  (effective # of sample in  $j$ -th cluster  $N_j$ )

Maximization conditions:

- Setting the derivative of  $\ell$  w.r.t.  $\underline{z}_j$  to zero gives

$$\underline{z}_j = \frac{1}{N_j} \sum_{i=1}^N \gamma_{ij} \underline{x}_i, \quad j = 1, \dots, K$$

where

$$\gamma_{ij} = \frac{\pi_j p_g(\underline{x}_i; \underline{z}_j, \Sigma_j)}{\sum_{l=1}^K \pi_l p_g(\underline{x}_i; \underline{z}_l, \Sigma_l)} = p(r_j | \underline{x}_i),$$

$$N_j = \sum_{i=1}^N \gamma_{ij}.$$

$$\Sigma_j = \frac{1}{N_j} \sum_{i=1}^N r_{ij} (\underline{x}_i - \underline{z}_j)(\underline{x}_i - \underline{z}_j)^T, \quad j = 1, \dots, K$$

- $$\begin{cases} \max \{\pi_j\} \ell = \sum_{i=1}^N \log \left[ \sum_{j=1}^K \pi_j p_g(\underline{x}_i; \underline{z}_j, \Sigma_j) \right] \\ \text{s.t. } \sum_{j=1}^K \pi_j = 1 \end{cases}$$

$$\mathcal{L} = \ell + \lambda \left( \sum_{j=1}^K \pi_j - 1 \right) \quad \lambda: \text{Lagrange multiplier}$$

$$0 = \frac{\partial \mathcal{L}}{\partial \pi_j} = \sum_{i=1}^N \left[ \frac{p_g(\underline{x}_i; \underline{z}_j, \Sigma_j)}{\sum_{l=1}^K \pi_l p_g(\underline{x}_i; \underline{z}_l, \Sigma_l)} \right] + \lambda \quad \text{ratio} = 1$$

Multiply both sides by  $\pi_j$ ;

$$\left[ \sum_{j=1}^K \pi_j \right] \lambda = - \sum_{i=1}^N \frac{\sum_{j=1}^K \pi_j p_g(\underline{x}_i; \underline{z}_j, \Sigma_j)}{\sum_{l=1}^K \pi_l p_g(\underline{x}_i; \underline{z}_l, \Sigma_l)} = - \sum_{i=1}^N 1 = -N$$

$$\Rightarrow \lambda = -N$$

$$0 = \frac{\partial \mathcal{L}}{\partial \pi_j} = \sum_{i=1}^N \left[ \frac{\pi_j p_g(\underline{x}_i; \underline{z}_j, \Sigma_j)}{\sum_{l=1}^K \pi_l p_g(\underline{x}_i; \underline{z}_l, \Sigma_l)} \right] - N \pi_j$$

$$0 = \sum_{i=1}^N r_{ij} - N \pi_j \Rightarrow \pi_j = \frac{\sum_{i=1}^N r_{ij}}{N} = \frac{N_j}{N}, \quad j = 1, 2, \dots, K$$

## Summary

$$\left\{ \begin{array}{l} Z_j = \frac{1}{N_j} \sum_{i=1}^N \gamma_{ij} X_i \quad (1) \\ \Sigma_j = \frac{1}{N_j} \sum_{i=1}^N \gamma_{ij} (X_i - Z_j)(X_i - Z_j)^T \quad (2) \\ \pi_j = \frac{N_j}{N} \quad (3) \end{array} \right. \quad \text{where} \quad \left\{ \begin{array}{l} \gamma_{ij} = \frac{\pi_j p_g(X_i; Z_j, \Sigma_j)}{\sum_{l=1}^K \pi_l p_g(X_i; Z_l, \Sigma_l)} = p(r_j | X_i) \quad (4) \\ N_j = \sum_{i=1}^N \gamma_{ij} \quad (5) \end{array} \right.$$

## Algorithm:

- Given the dataset  $\{\mathbf{x}_i\}_{i=1}^N$  in  $R^d$ .
- Initialize  $\{\pi_j, \mathbf{z}_j, \Sigma_j\}_{j=1}^K$ :

$$\pi_j = \frac{1}{K}, \quad \mathbf{z}_j \in R^d, \quad \Sigma_j \in R^{d \times d}, \quad j = 1, \dots, K$$

- Repeat the following two steps until convergence:
  - Update  $\{\gamma_{ij}, N_j\}$  according to Eqs. (4), (5).
  - Update  $\{\pi_j, \mathbf{z}_j, \Sigma_j\}$  according to Eqs. (1) – (3).

## Clustering:

- Given a data point  $\mathbf{x}$ , compute the posterior probability  $p(\mathbf{r} = \mathbf{r}_j | \mathbf{x})$  - the probability that the data belongs to cluster  $j$  after observing the sample.
- Bayes rule:

$$p(\mathbf{r} = \mathbf{r}_j | \mathbf{x}) = \frac{p(\mathbf{x} | \mathbf{r} = \mathbf{r}_j) p(\mathbf{r} = \mathbf{r}_j)}{p(\mathbf{x})}$$

where

$p(\mathbf{x} | \mathbf{r} = \mathbf{r}_j) = p_g(\mathbf{x}; \mathbf{z}_j, \Sigma_j)$  : class conditional

$p(\mathbf{r} = \mathbf{r}_j) = \pi_j$  : prior

$$p(\mathbf{x}) = \sum_{l=1}^K \pi_l p_g(\mathbf{x}; \mathbf{z}_l, \Sigma_l) : \text{margin}$$

$\{\pi_j, \mathbf{z}_j, \Sigma_j\} \rightarrow$  compute posterior probability  $p(r=r_j | x)$

Assign  $x$  to the cluster with the largest probability

$$K = \operatorname{argmax}_j p(r=r_j | x)$$