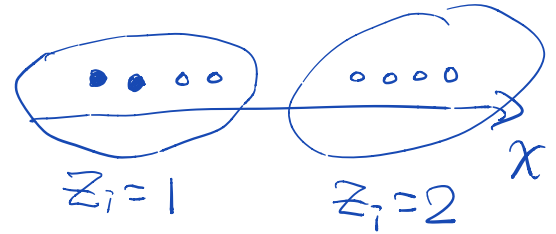


## Expectation Maximization (EM)

**Qiang Liu**  
UT Austin

# Clustering

- **Inputs:**  $n$  objects (data points)  $\{x_i\}_{i=1}^n$  and a number  $K$  of clusters.
- **Goal:** Group the points into several groups
  - Deciding a group ID,  $z_i \in \{1, \dots, K\}$ , for each data point  $x_i$ .



## Probabilistic Modeling of Clustering

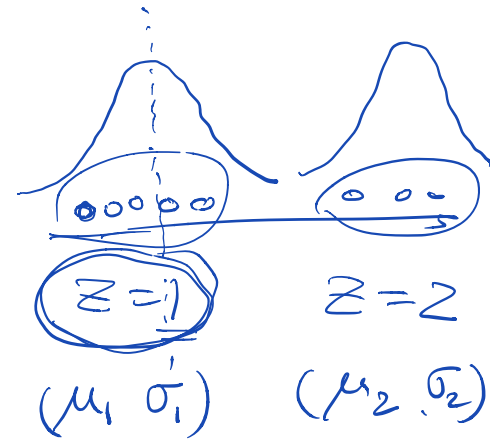
### Probabilistic Approach:

- Assume a joint distribution  $p(x, z | \theta)$  that generates data and labels  $(x, z)$ .
- Estimate parameter  $\theta$ .
- Infer the labels from the posterior distribution  $p(z | x; \theta)$ .

$$P(x, z) = P(z) P(x | z)$$

$$\textcircled{1} \quad \underline{P(z=k)} = \pi_k \quad \left( \begin{array}{l} \pi_k \geq 0 \\ \sum_k \pi_k = 1 \end{array} \right)$$

$$\textcircled{2} \quad P(x | z=k) = \mathcal{N}(x | \mu_k, \sigma_k^2)$$



$$\begin{aligned} P(x, \underline{z=k}) &= P(z=k) P(x | z=k) \\ &= \underline{\pi_k} \underline{\mathcal{N}(x | \mu_k, \sigma_k^2)} \end{aligned}$$

$$\theta = \underline{\underline{\{\pi_k, \mu_k, \sigma_k^2\}_{k=1}^K}}}$$

$$\underline{P(x, z | \theta)}$$

# MLE with Complete Information

- If we observe both data and labels  $\{x_i, z_i\}_{i=1}^n$  (complete information), we can estimate the parameter  $\theta$  by maximizing the **joint probability**:

$$\max_{\theta} \sum_{i=1}^n \log p(x_i, z_i | \theta)$$

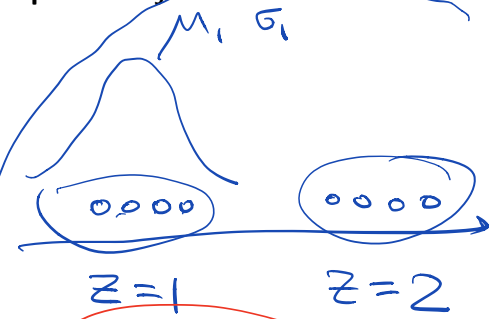
$$P(x_i, z_i = k | \theta) = \pi_k \mathcal{N}(x_i | \mu_k, \sigma_k^2)$$

$$\sum_{i=1}^n \sum_{k=1}^K \mathbb{I}(z_i = k) \log P(x_i, z_i = k | \theta)$$

$$= \sum_{i=1}^n \sum_{k=1}^K \mathbb{I}(z_i = k) \log(\pi_k \mathcal{N}(x_i | \mu_k, \sigma_k^2))$$

$$= \sum_{i=1}^n \sum_{k=1}^K \mathbb{I}(z_i = k) (\log \pi_k + \log \mathcal{N}(x_i | \mu_k, \sigma_k^2))$$

$$\pi_k = \frac{\sum_{i=1}^n \mathbb{I}(z_i = k)}{n}$$



$$\pi_1 = \frac{\#(z=1)}{n}$$

$$= \frac{1}{2}$$

$$\mu_1 = \frac{\sum_{i=1}^n \mathbb{I}(z_i = 1) x_i}{\sum_{i=1}^n \mathbb{I}(z_i = 1)}$$

$$\sigma_1^2 = \text{Var}(\{x_i | z_i = 1\})$$

# Clustering and Mixture Models

□ If we only observe data  $\{x_i\}_{i=1}^n$  (incomplete information), we shall estimate  $\theta$  by maximizing the marginal probability:

$$\sum_{i=1}^n \log p(x_i | \theta) = \sum_{i=1}^n \log \sum_{z_i} p(x_i, z_i | \theta)$$

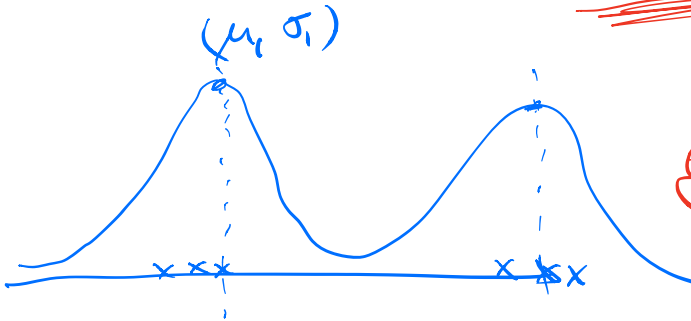
$$\sum_{i=1}^n \log P(x_i, z_i | \theta)$$

$$P(x | \theta) = \sum_k P(x, z=k | \theta)$$

$$= \sum_k \pi_k \mathcal{N}(x | \mu_k, \sigma_k^2)$$

mixture  
of Gaussian  
distributions

$$\sum_{i=1}^n \log P(x_i | \theta) = \sum_{i=1}^n \log \left( \sum_{k=1}^K \pi_k \mathcal{N}(x | \mu_k, \sigma_k^2) \right)$$



$$\max \{\pi_k, \mu_k, \sigma_k\}_{k=1}^K$$

$$P(x | \theta) = \left( \frac{1}{2} \right) \exp \left( -\frac{(x - \mu_1)^2}{2\sigma_1^2} \right) \left( \frac{1}{\sqrt{2\pi} \sigma_1} \right)$$

$$+ \frac{1}{2} \exp \left( -\frac{(x - \mu_2)^2}{2\sigma_2^2} \right) \left( \frac{1}{\sqrt{2\pi} \sigma_2} \right)$$



## Expectation Maximization: Algorithm Procedure

**Expectation Maximization (EM)** is an iterative method for maximizing the marginal likelihood.

□ Initialize parameter  $\theta_0$ .

□ For iteration  $t$ :

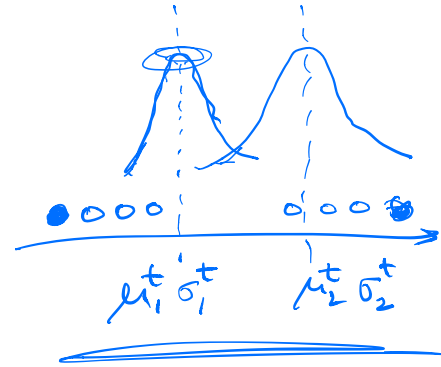
○ Given  $\theta_t$ , "impute" the missing labels by drawn samples from the posterior distribution

$$z_i \sim p(z | x_i; \theta_t).$$

$$\max_{\theta} \sum_{i=1}^n \log p(x_i | \theta)$$

○ Update  $\theta$  by maximizing the expected joint likelihood:

$$\theta^{t+1} = \arg \max_{\theta} \sum_{i=1}^n \mathbb{E}_{z_i \sim p(\cdot | x_i; \theta_t)} [\log p(x_i, z_i | \theta)].$$



$$\theta_t = [\pi_k^t, \mu_k^t, \sigma_k^t]_{k=1}^K$$

$$\begin{aligned} \underline{P(z_i = k | x_i, \theta_t)} &= \frac{P(x_i, z_i = k | \theta_t)}{\sum_l P(x_i, z_i = l | \theta_t)} \\ &= \frac{\pi_k \mathcal{N}(x_i | \mu_k^t, \sigma_k^t)}{\sum_l \pi_l \mathcal{N}(x_i | \mu_l^t, \sigma_l^t)} \end{aligned}$$

$$\begin{aligned} P(z_i = 1 | x_i, \theta_t) &= \frac{\pi_1 \exp\left(-\frac{(x_i - \mu_1)^2}{2\sigma_1^2}\right) / (\sqrt{2\pi} \sigma_1)}{\sum_l \pi_l \exp\left(-\frac{(x_i - \mu_l)^2}{2\sigma_l^2}\right) / (\sqrt{2\pi} \sigma_l)} \end{aligned}$$

EM algorithm:

$$\theta^{t+1} = \arg \max_{\theta} \sum_{i=1}^n \mathbb{E}_{z_i \sim p(\cdot | x_i; \theta_t)} [\log p(x_i, z_i | \theta)]. \triangleq L(\theta)$$

Define  $\gamma_{ik}^t = P(z_i = k | x_i, \theta_t)$

$$L(\theta) = \sum_{i=1}^n \sum_{k=1}^K \gamma_{ik}^t \log P(x_i, z_i = k | \theta)$$

$$\Rightarrow \left\{ \begin{aligned} \mu_k^{t+1} &= \frac{\sum_{i=1}^n \gamma_{ik}^t}{n} \\ \mu_k^{t+1} &= \left( \sum_{i=1}^n \gamma_{ik}^t x_i \right) / \left( \sum_{i=1}^n \gamma_{ik}^t \right) \\ \sigma_k^{t+1} &= \dots \end{aligned} \right.$$