

Clustering and Mixture Models

Wenting Tu

SHUFE, SIME

Machine Learning and Deep Learning

Course No. 1638

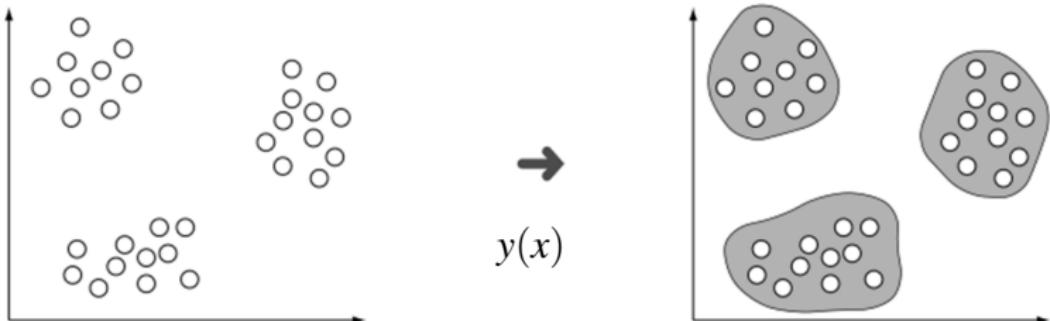
Outline

Clustering

Mixture Models

Extension

Clustering



x_n

$r_{nk} \in \{0, 1\}, k = 1, \dots, K$

K-means

$$\boldsymbol{\mu}^*, \mathbf{r}^* = \arg \min \sum_{k=1}^K \sum_{n=1}^N r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$

$$J = \sum_{k=1}^K r_{1k} \|\mathbf{x}_1 - \boldsymbol{\mu}_k\|^2 + \sum_{k=1}^K r_{2k} \|\mathbf{x}_2 - \boldsymbol{\mu}_k\|^2 + \cdots + \sum_{k=1}^K r_{Nk} \|\mathbf{x}_N - \boldsymbol{\mu}_k\|^2$$

$$\rightarrow r_{nk}^* = \begin{cases} 1 & \text{if } k = \arg \min_j \|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2 \\ 0 & \text{otherwise} \end{cases}$$

K-means

$$\boldsymbol{\mu}^*, \mathbf{r}^* = \arg \min \sum_{k=1}^K \sum_{n=1}^N r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$

$$J = \sum_{n=1}^N r_{n1} \|\mathbf{x}_n - \boldsymbol{\mu}_1\|^2 + \sum_{n=1}^N r_{n2} \|\mathbf{x}_n - \boldsymbol{\mu}_2\|^2 + \cdots + \sum_{n=1}^N r_{nK} \|\mathbf{x}_n - \boldsymbol{\mu}_K\|^2$$

$$\rightarrow \boldsymbol{\mu}_k^* = \frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}}$$

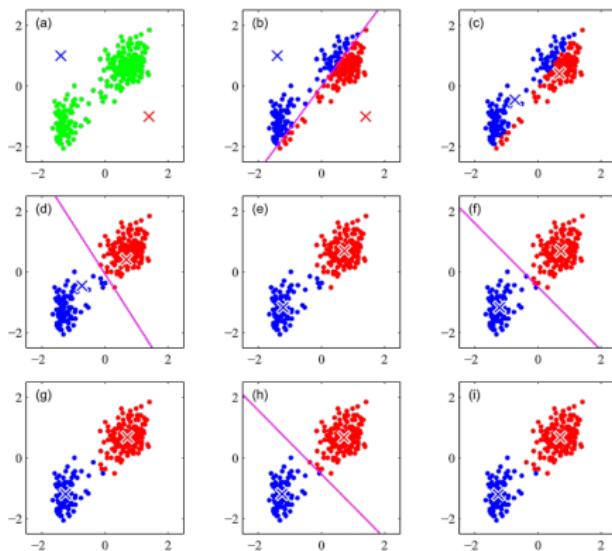
K-means

- choose some initial values for the μ_k
- repeated until convergence
- minimize J with respect to the r_{nk}

$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|x_n - \mu_j\|^2 \\ 0 & \text{otherwise} \end{cases}$$

- minimize J with respect to the μ_k

$$\mu_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}}$$



K-medoids

$$\tilde{J} = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \mathcal{V}(\mathbf{x}_n, \boldsymbol{\mu}_k)$$

Others

- Density-based Clustering
- Hierarchical Clustering

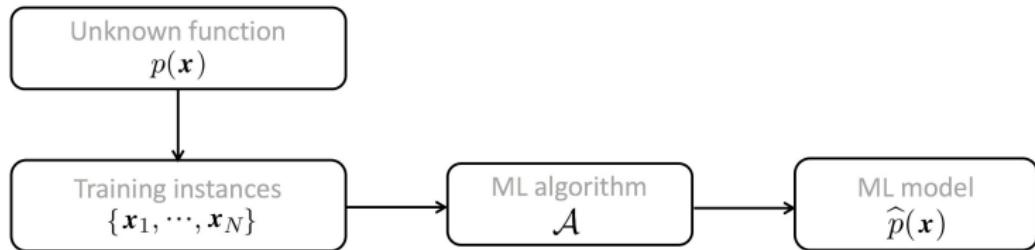
Outline

Clustering

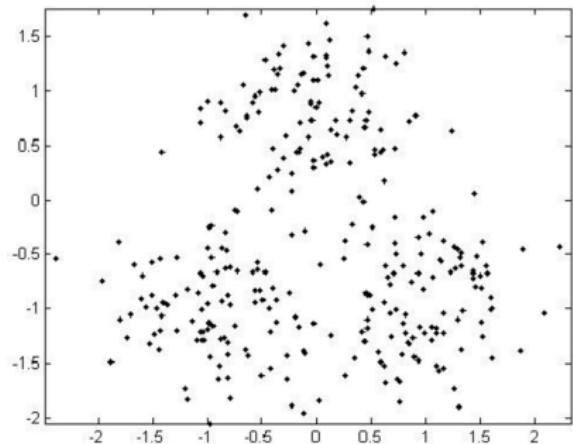
Mixture Models

Extension

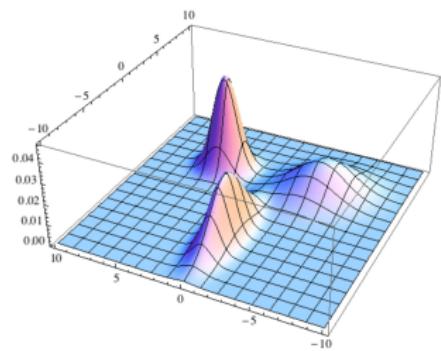
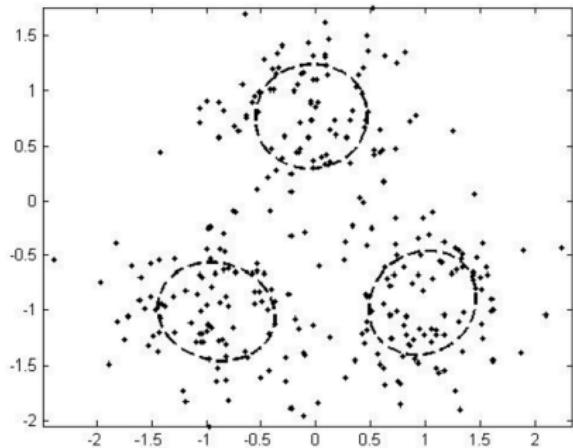
Unsupervised Learning



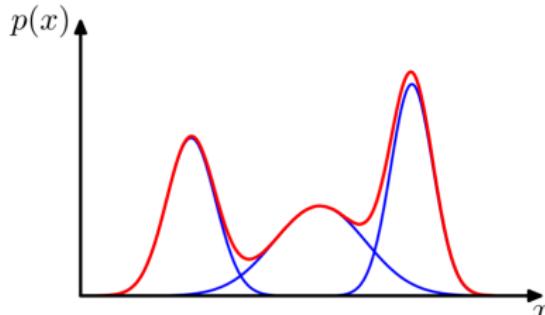
Mixtures of Gaussians



Mixtures of Gaussians



Mixtures of Gaussians



$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad \sum_{k=1}^K \pi_k = 1 \quad 0 \leq \pi_k \leq 1$$

mixing coefficient component of the mixture

$$\boldsymbol{\pi} \equiv \{\pi_1, \dots, \pi_K\} \quad \boldsymbol{\mu} \equiv \{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K\} \quad \boldsymbol{\Sigma} \equiv \{\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K\}$$

K-means

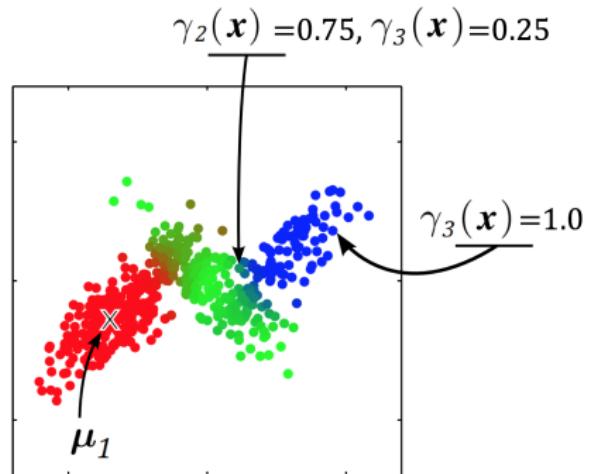
- GMM for soft clustering

$$\gamma_k(\mathbf{x}) \equiv p(k|\mathbf{x})$$

$$= \frac{p(k)p(\mathbf{x}|k)}{\sum_l p(l)p(\mathbf{x}|l)}$$

$$= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_l \pi_l \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l)}$$

(responsibilities)



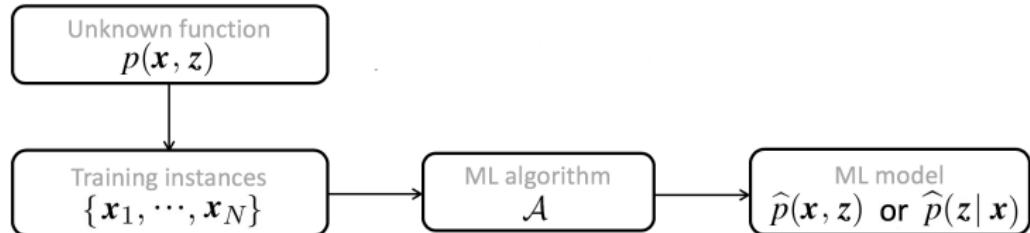
Maximum likelihood Solution

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

Note: the maximum likelihood solution for the parameters no longer has a closed-form analytical solution.

Latent variable

Introduce a K-dimensional binary random variable



$$z_k \in \{0, 1\}, \sum_k z_k = 1$$

$$p(\mathbf{x} | z_k = 1) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad p(z_k = 1) = \pi_k$$

$$p(\mathbf{x} | \mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k} \quad p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$$

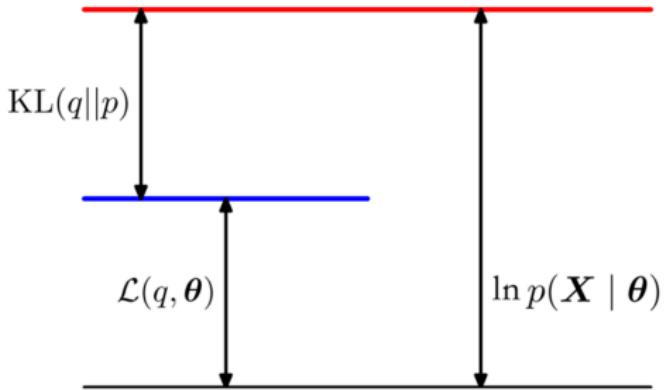
$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z}) p(\mathbf{x} | \mathbf{z}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Expectation Maximization

$$\ln p(\mathbf{X} \mid \boldsymbol{\theta}) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z} \mid \boldsymbol{\theta})}{q(\mathbf{Z})} \right\} - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z} \mid \mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \right\}$$

1st item: Evidence Lower Bound (ELBO), denoted as $\mathcal{L}(q, \boldsymbol{\theta})$

2nd item: Kullback-Leibler divergence, denoted as $\text{KL}(q \parallel p)$



$$\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}, \boldsymbol{\theta}^{(3)} \dots \longrightarrow \ln p(\mathbf{X} \mid \boldsymbol{\theta}^{(t+1)}) \geq \ln p(\mathbf{X} \mid \boldsymbol{\theta}^{(t)})$$

Expectation Maximization

$$\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) = \ln p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}) + \ln p(\mathbf{X} | \boldsymbol{\theta})$$

$$\ln p(\mathbf{X} | \boldsymbol{\theta}) = \ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) - \ln p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta})$$

$$\ln p(\mathbf{X} | \boldsymbol{\theta}) = \{\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) - \ln q(\mathbf{Z})\} - \{\ln p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}) - \ln q(\mathbf{Z})\}$$

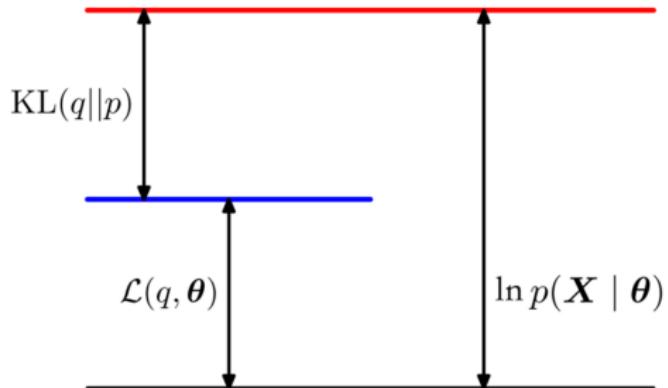
$$\ln p(\mathbf{X} | \boldsymbol{\theta}) = \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})}{q(\mathbf{Z})} \right\} - \ln \left\{ \frac{p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \right\}$$

$$\sum_{\mathbf{Z}} q(\mathbf{Z}) \ln p(\mathbf{X} | \boldsymbol{\theta}) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})}{q(\mathbf{Z})} \right\} - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \right\}$$

$$\ln p(\mathbf{X} | \boldsymbol{\theta}) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})}{q(\mathbf{Z})} \right\} - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \right\}$$

Expectation Maximization

$$\ln p(\mathbf{X} \mid \boldsymbol{\theta}) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z} \mid \boldsymbol{\theta})}{q(\mathbf{Z})} \right\} - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z} \mid \mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \right\}$$



$$\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}, \boldsymbol{\theta}^{(3)} \dots \longrightarrow \ln p(\mathbf{X} \mid \boldsymbol{\theta}^{(t+1)}) \geq \ln p(\mathbf{X} \mid \boldsymbol{\theta}^{(t)})$$

Expectation Maximization

$$\begin{aligned} & \ln p(X | \boldsymbol{\theta}^{(t)}) \\ &= \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(X, \mathbf{Z} | \boldsymbol{\theta}^{(t)})}{q(\mathbf{Z})} \right\} + \left(- \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z} | X, \boldsymbol{\theta}^{(t)})}{q(\mathbf{Z})} \right\} \right) \end{aligned}$$

$$= \mathcal{L}(q, \boldsymbol{\theta}^{(t)}) + KL(q(\mathbf{Z}) \| p(\mathbf{Z} | X, \boldsymbol{\theta}^{(t)}))$$

$$q^t(\mathbf{Z}) = p(\mathbf{Z} | X, \boldsymbol{\theta}^{(t)}) \rightarrow KL = 0$$

$$= \mathcal{L}(q^t, \boldsymbol{\theta}^{(t)})$$

$$= \sum_{\mathbf{Z}} q^t(\mathbf{Z}) \ln \left\{ \frac{p(X, \mathbf{Z} | \boldsymbol{\theta}^{(t)})}{q^t(\mathbf{Z})} \right\}$$

$$= \sum_{\mathbf{Z}} q^t(\mathbf{Z}) \ln p(X, \mathbf{Z} | \boldsymbol{\theta}^{(t)}) - \sum_{\mathbf{Z}} q^t(\mathbf{Z}) \ln q^t(\mathbf{Z})$$

$$\boldsymbol{\theta}^{(t+1)} = \arg \max_{\boldsymbol{\theta}} \mathcal{L}(q^t, \boldsymbol{\theta})$$

$$= \arg \max_{\boldsymbol{\theta}} \sum_{\mathbf{Z}} q^t(\mathbf{Z}) \ln p(X, \mathbf{Z} | \boldsymbol{\theta}) \rightarrow \mathcal{L}(q^t, \boldsymbol{\theta}^{(t)}) \leq \mathcal{L}(q^t, \boldsymbol{\theta}^{(t+1)})$$

$$\leq \mathcal{L}(q^t, \boldsymbol{\theta}^{(t+1)})$$

$$\leq \mathcal{L}(q^t, \boldsymbol{\theta}^{(t+1)}) + KL(q^t(\mathbf{Z}) \| p(\mathbf{Z} | X, \boldsymbol{\theta}^{(t+1)}))$$

$$= \ln p(X | \boldsymbol{\theta}^{(t+1)})$$

Expectation Maximization

- Choose an initial setting for the parameters θ^{old}
- **E step** Evaluate $p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}})$
- **M step** Evaluate θ^{new}

$$\theta^{\text{new}} = \arg \max \mathcal{Q}(\theta, \theta^{\text{old}})$$

where

$$\mathcal{Q}(\theta, \theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z} | \theta)$$

- Check for convergence of either the log likelihood or the parameter values. If the convergence criterion is not satisfied, then let

$$\theta^{\text{old}} \leftarrow \theta^{\text{new}}$$

and return to **E step**

EM for Mixtures of Gaussians

$$\begin{aligned}\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) &= \sum_{\mathbf{Z}} \ln p(\mathbf{X}, \mathbf{Z} \mid \boldsymbol{\theta}) p(\mathbf{Z} \mid \mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \\ &= \sum_{z_1} \cdots \sum_{z_N} \left\{ \sum_{n=1}^N \ln p(\mathbf{x}_n, z_n \mid \boldsymbol{\theta}) \prod_{n=1}^N p(z_n \mid \mathbf{x}_n, \boldsymbol{\theta}^{\text{old}}) \right\}\end{aligned}$$

EM for Mixtures of Gaussians

$$\begin{aligned} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) &= \sum_{\mathbf{Z}} \ln p(X, \mathbf{Z} \mid \boldsymbol{\theta}) p(\mathbf{Z} \mid X, \boldsymbol{\theta}^{\text{old}}) \\ &= \sum_{z_1} \cdots \sum_{z_N} \left\{ \sum_{n=1}^N \ln p(\mathbf{x}_n, z_n \mid \boldsymbol{\theta}) \prod_{n=1}^N p(z_n \mid \mathbf{x}_n, \boldsymbol{\theta}^{\text{old}}) \right\} \end{aligned}$$

$$\frac{\partial \mathcal{Q}}{\partial \pi_k} = 0 \quad \text{s.t. } \sum_{k=1}^K \pi_k = 1$$

$$\pi_k^{\text{new}} = \frac{N_k}{N}$$

$$\gamma(z_{nk}) = \frac{\pi_k^{\text{old}} \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k^{\text{old}}, \boldsymbol{\Sigma}_k^{\text{old}})}{\sum_{j=1}^K \pi_j^{\text{old}} \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_j^{\text{old}}, \boldsymbol{\Sigma}_j^{\text{old}})}$$

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

EM for Mixtures of Gaussians

$$\begin{aligned} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) &= \sum_{\mathbf{Z}} \ln p(X, \mathbf{Z} \mid \boldsymbol{\theta}) p(\mathbf{Z} \mid X, \boldsymbol{\theta}^{\text{old}}) \\ &= \sum_{z_1} \cdots \sum_{z_N} \left\{ \sum_{n=1}^N \ln p(\mathbf{x}_n, z_n \mid \boldsymbol{\theta}) \prod_{n=1}^N p(z_n \mid \mathbf{x}_n, \boldsymbol{\theta}^{\text{old}}) \right\} \end{aligned}$$

$$\frac{\partial \mathcal{Q}}{\partial \boldsymbol{\mu}_k} = 0 \quad \frac{\partial \mathcal{Q}}{\partial \boldsymbol{\Sigma}_k} = 0$$

$$\boldsymbol{\mu}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

$$\boldsymbol{\Sigma}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}})^T$$

EM for Mixtures of Gaussians

- Initialize μ_k, Σ_k, π_k
- **E step.** $\gamma(z_{nk}) = \frac{\pi_k^{\text{old}} \mathcal{N}(x_n | \mu_k^{\text{old}}, \Sigma_k^{\text{old}})}{\sum_{j=1}^K \pi_j^{\text{old}} \mathcal{N}(x_n | \mu_j^{\text{old}}, \Sigma_j^{\text{old}})}$
- **M step.** $N_k = \sum_{n=1}^N \gamma(z_{nk}), \pi_k^{\text{new}} = \frac{N_k}{N}$
 $\mu_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) x_n$
 $\Sigma_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (x_n - \mu_k^{\text{new}}) (x_n - \mu_k^{\text{new}})^T$
- Check for convergence
If the convergence criterion is not satisfied, then let

$$(\pi_k^{\text{old}}, \mu_k^{\text{old}}, \Sigma_k^{\text{old}}) \leftarrow (\pi_k^{\text{new}}, \mu_k^{\text{new}}, \Sigma_k^{\text{new}})$$

and return to **E step**

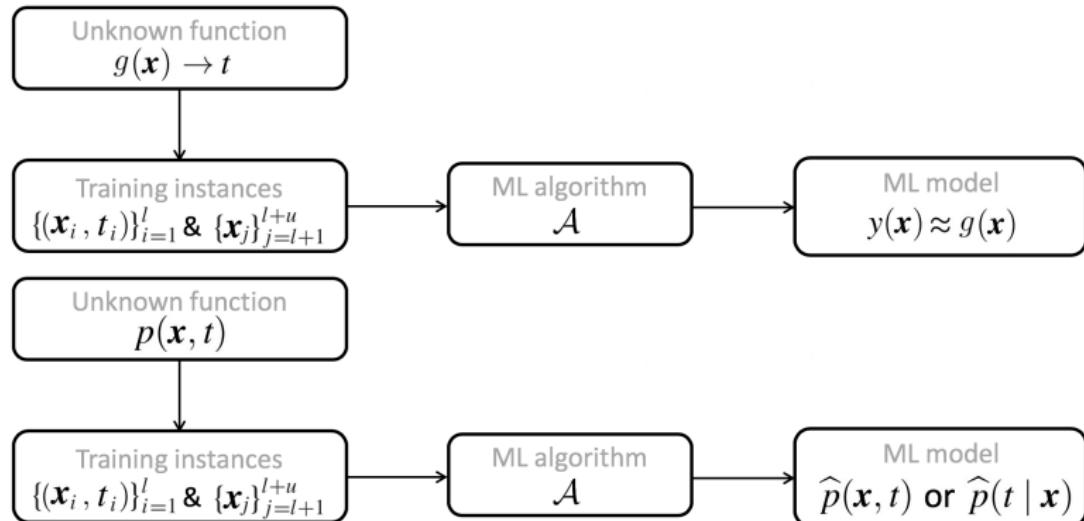
Outline

Clustering

Mixture Models

Extension

Semi-supervised Learning (SSL)



Thanks

Some images and slides are from the internet.
If related to copyright, please contact me.

[tu.wenting@mail.shufe.edu.cn](mailto:tudongtian@163.com)