# MIXTURE MODELS & EM

K-means clustering

data set $\{x_1, \dots, x_N\}$      $x_n \in \mathbb{R}^D$

partition the data set into $K$ clusters

Goal: to find [1] an assignment of data points to clusters and [2] a set of vector $\{\mu_k\}$, such that the sum of the squares of the distances of each data point to its closest vector $\mu_k$, is a minimum.

$x_n \longrightarrow r_{nk} \in \{0, 1\}$ , $k = 1, \dots, K$    (1-of-$K$ coding scheme)

$$r_{nk} = \begin{cases} 1, & x_n \text{ in cluster } K, \\ 0, & \text{otherwise.} \end{cases}$$

$\underset{\{\mu_k\}\,\{r_{nk}\}}{\text{minimize}}$    $J = \sum_{n=1}^{N} \sum_{k=1}^{K} r_{nk} \| x_n - \mu_k \|^2$

See Fig 9.1

Two-stage optimization

E step    $r_{nk}$ :    $r_{nk} = \begin{cases} 1, & \text{if } k = \underset{j}{\arg\min} \| x_n - \mu_j \|^2 \\ 0, & \text{otherwise} \end{cases}$

M step    $\mu_k$ :    $\dfrac{\partial J}{\partial \mu_k} = 0 \;\Rightarrow\; -2 \sum_{n=1}^{N} r_{nk} (x_n - \mu_k) = 0$

$\mu_k = \sum_{n=1}^{N} r_{nk} x_n \Big/ \sum_{n=1}^{N} r_{nk}$    (cluster mean)

---

Sequential Update

$$\mu_k^{new} = \mu_k^{old} + \lambda_n \frac{\partial J(x_n)}{\partial \mu_k}$$

$$= \mu_k^{old} + \lambda_n (-2 r_{nk})(x_n - \mu_k^{old})$$

$$= \mu_k^{old} + \eta_{nk} (x_n - \mu_k^{old})$$

Robustness to Outliers?    $l_2$-norm is not robust to outliers

K-medoids    $\tilde{J} = \sum_{n=1}^{N} \sum_{k=1}^{K} r_{nk}\, \mathcal{V}(x_n, \mu_k)$

outliers

general dissimilarity measure

Examples of K-means clustering
     image segmentation
     vector quantization

# Mixtures of Gaussians

## Generative Models

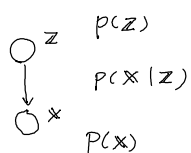latent variables    $p(z_k = 1) = \pi_k$

$p(z)$
$p(x|z)$
$P(x)$

$$0 \le \pi_k \le 1 \quad , \quad \sum_{k=1}^{K} \pi_k = 1$$

$z$ uses a 1-of-K coding scheme

$$\begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \Big\} k$$

$$p(z) = \prod_{k=1}^{K} \pi_k^{z_k}$$

$$p(z_k = 1) = 1 \cdots 1 \cdot \pi_k \cdot 1 \cdots 1$$

Since $p(x|z_k = 1)$ is assumed to be a Gaussian,

$$p(x \mid z_k = 1) = \mathcal{N}(x \mid \mu_k, \Sigma_k)$$

Considering $z_k$ as a selector:

$$p(x \mid z) = \prod_{k=1}^{K} \mathcal{N}(x \mid \mu_k, \Sigma_k)^{z_k}$$

the marginal distribution of $x$

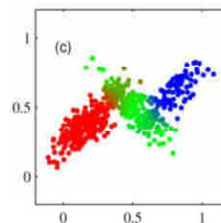$$p(x) = \sum_{z} p(x, z) = \sum_{z} p(z)\, p(x|z)$$

$$= \sum_{z} \prod_{k=1}^{K} \pi_k^{z_k} \prod_{k=1}^{K} \mathcal{N}(x \mid \mu_k, \Sigma_k)^{z_k} = \sum_{k=1}^{K} \pi_k \mathcal{N}(x \mid \mu_k, \Sigma_k)$$

⇑

sum over all possible states of $z \in \left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \\ \vdots \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \end{pmatrix}, \cdots, \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix} \right\}$

---

## responsibility

$$\gamma(z_k) \equiv p(z_k = 1 \mid x) = \frac{p(z_k=1)\, p(x \mid z_k = 1)}{\sum_{j=1}^{K} p(z_j=1)\, p(x \mid z_j = 1)}$$



$$= \frac{\pi_k\, \mathcal{N}(x \mid \mu_k, \Sigma_k)}{\sum_{j=1}^{K} \pi_j\, \mathcal{N}(x \mid \mu_j, \Sigma_j)}$$
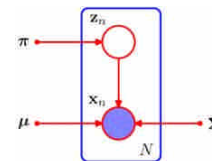
## Maximum Likelihood

Given $N$ observations $\{x_1, \ldots, x_N\}$, we want to find $\mu_k, \Sigma_k, \pi_k$ that maximize the likelihood function
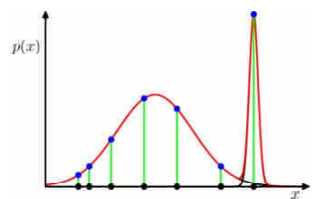
log-likelihood

$$L \equiv \ln p(X \mid \{\pi_k\}, \{\mu_k\}, \{\Sigma_k\})$$

$$= \ln \prod_{n=1}^{N} \sum_{k=1}^{K} \pi_k \mathcal{N}(x_n \mid \mu_k, \Sigma_k)$$

$$= \sum_{n=1}^{N} \ln \left\{ \sum_{k=1}^{K} \pi_k \mathcal{N}(x_n \mid \mu_k, \Sigma_k) \right\}$$

⇑

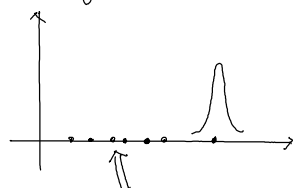maximize it w.r.t. $\mu_k, \Sigma_k, \pi_k$

Singularities in the Likelihood Function of Mixtures of Gaussians



$$\mathcal{N}(x_n \mid \mu_j = x_n, \, \sigma_j^2 I)$$
$$= \frac{1}{\sqrt{2\pi}} \frac{1}{\sigma_j} \longrightarrow \infty \quad \text{as } \sigma_j \to 0$$

↑ one Gaussian collapses onto a specific data point
the other takes care of the remaining data points

For a single Gaussian model, such a situation would not happen.



Suppose the Gaussian collapses onto a specific data point

↑ these data points have zero likelihoods
and thus the overall likelihood goes to zero.

Maximum Likelihood $\quad \frac{\partial L}{\partial \mu_k} = 0. \quad \frac{\partial L}{\partial \Sigma_k} = 0, \quad \frac{\partial L}{\partial \pi_k} = 0$

Recall $\quad \mathcal{N}(x \mid \mu, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp\left\{ -\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu) \right\}$

① $\quad 0 = \frac{\partial L}{\partial \mu_k} = \frac{\partial}{\partial \mu_k} \sum_{n=1}^{N} \ln \left\{ \sum_{k=1}^{K} \pi_k \mathcal{N}\left( x_n \mid \mu_k, \Sigma_k \right) \right\}$

$$= \sum_{n=1}^{N} \frac{\pi_k \frac{\partial \mathcal{N}(x_n \mid \mu_k, \Sigma_k)}{\partial \mu_k}}{\sum_{j=1}^{K} \pi_j \mathcal{N}(x_n \mid \mu_j, \Sigma_j)}$$

derivative w.r.t. $\mu$.
(a useful Gaussian identity)

$$\frac{\partial \mathcal{N}}{\partial \mu} = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp\left\{ -\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu) \right\} \cdot \left( \frac{-1}{2} \cdot 2 \cdot \Sigma^{-1}(x-\mu) \cdot -1 \right)$$

$$= \mathcal{N} \, \Sigma^{-1} (x-\mu)$$

Plugged into $\frac{\partial L}{\partial \mu_k}$, we have

$$0 = \sum_{n=1}^{N} \frac{\pi_k \mathcal{N}(x_n \mid \mu_k, \Sigma_k)}{\sum_{j=1}^{K} \pi_j \mathcal{N}(x_n \mid \mu_j, \Sigma_j)} \Sigma_k^{-1} (x_n - \mu_k)$$

multiplying both sides by $\Sigma_k$, we obtain

$$\boxed{\mu_k = \frac{1}{N_k} \sum_{n=1}^{N} \gamma(z_{nk}) \, x_n} \quad , \quad N_k = \sum_{n=1}^{N} \gamma(z_{nk})$$

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(x_n \mid \mu_k, \Sigma_k)}{\sum_{j=1}^{K} \pi_j \mathcal{N}(x_n \mid \mu_j, \Sigma_j)}$$

②

Next, we set $\frac{\partial L}{\partial \Sigma_k}$ to zero to get $\Sigma_k$

$$0 = \frac{\partial L}{\partial \Sigma_k} = \frac{\partial}{\partial \Sigma_k} \sum_{n=1}^{N} \ln \left\{ \sum_{k=1}^{K} \pi_k \, \mathcal{N}(x_n | \mu_k, \Sigma_k) \right\}$$

$$= \sum_{n=1}^{N} \frac{\pi_k \, \frac{\partial}{\partial \Sigma_k} \mathcal{N}(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^{K} \pi_j \, \mathcal{N}(x_n | \mu_j, \Sigma_j)}$$

derivative w.r.t. $\Sigma$

$$\frac{\partial \mathcal{N}}{\partial \Sigma} = \frac{1}{(2\pi)^{D/2}} \left( \frac{\partial}{\partial \Sigma} |\Sigma|^{-\frac{1}{2}} \right) \exp\left\{ -\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu) \right\}$$

$$+ \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \frac{\partial}{\partial \Sigma} \exp\left\{ -\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu) \right\}$$

$$= -\frac{1}{2} \Sigma^{-1} \mathcal{N} - \Sigma^{-2} \left( -\frac{1}{2}(x-\mu)(x-\mu)^T \right) \mathcal{N}$$

matrix calculus

$$\frac{\partial}{\partial \Sigma} |\Sigma|^{-\frac{1}{2}} = -\frac{1}{2} |\Sigma|^{-\frac{3}{2}} \frac{\partial}{\partial \Sigma} |\Sigma|$$

$$= -\frac{1}{2} |\Sigma|^{-\frac{3}{2}} |\Sigma| \Sigma^{-1}$$

$$= -\frac{1}{2} |\Sigma|^{-\frac{1}{2}} \Sigma^{-1}$$

$$\frac{\partial}{\partial \Sigma} \exp\left\{ -\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu) \right\} = \exp\left\{ \right\} \left( -\Sigma^{-2} \right)\left( -\frac{1}{2}(x-\mu)(x-\mu)^T \right)$$

$$\frac{\partial}{\partial A} x^T A x = \frac{\partial}{\partial A} Tr(x^T A x) = \left( \frac{\partial}{\partial A} x^T A \right) x^T = x x^T$$

$$0 = \frac{\partial}{\partial \Sigma_k} \sum_{n=1}^{N} \ln \sum_{k=1}^{K} \pi_k \, \mathcal{N}(x_n | \mu_k, \Sigma_k)$$

$$= \sum_{n=1}^{N} \frac{\pi_k \, \mathcal{N}(x_n | \mu_k, \Sigma_k)}{\sum_{s=1}^{K} \pi_j \, \mathcal{N}(x_n | \mu_j, \Sigma_j)} \left( -\frac{1}{2} \Sigma_k^{-1} + \frac{1}{2} \Sigma_k^{-2}(x_n - \mu_k)(x - \mu_k)^T \right)$$

multiplying both sides by $2\Sigma_k^2$

$$\boxed{ \Sigma_k = \frac{1}{N_k} \sum_{n=1}^{N} \gamma(z_{nk})(x_n - \mu_k)(x_n - \mu_k)^T }$$

③

$$\tilde{\mathcal{L}} = \sum_{n=1}^{N} \ln \sum_{k=1}^{K} \pi_k \, \mathcal{N}(x_n | \mu_k, \Sigma_k) + \lambda \left( \sum_{k=1}^{K} \pi_k - 1 \right)$$

⇑ given by the constraint

$$\frac{\partial \tilde{\mathcal{L}}}{\partial \pi_k} = 0$$

$$0 = \sum_{n=1}^{N} \frac{\mathcal{N}(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^{K} \pi_j \, \mathcal{N}(x_n | \mu_j, \Sigma_j)} + \lambda$$

⇓ multiplying by $\pi_k$, summing over $k$

$$0 = \sum_{n=1}^{N} 1 + \sum_{k=1}^{K} \pi_k \lambda$$

⇓

$$\lambda = -N$$

$$\pi_k \, N = \sum_{n=1}^{N} \frac{\pi_k \, \mathcal{N}(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^{K} \pi_j \, \mathcal{N}(x_n | \mu_j, \Sigma_j)}$$

$$\boxed{\pi_k = \frac{N_k}{N}} \qquad \Bigg| \quad N_k = \sum_{n=1}^{N} \gamma(z_{nk})$$

## SUMMARY

E step :

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

M step :

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

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

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

where

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