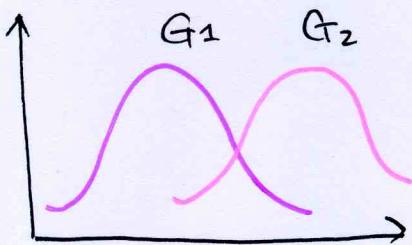


Mixture Models and EM - Expectation Maximization

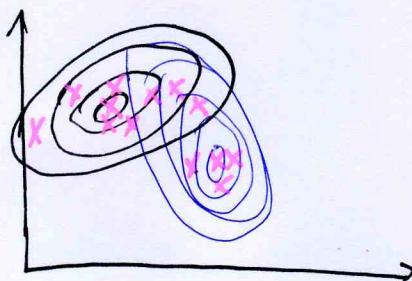
① mixture model Revisited



$$x \sim \pi_1 G_1 + \pi_2 G_2$$

$$G_i \sim N(\mu_i, \sigma_i)$$

$$\sum_i \pi_i = 1$$



With the probability of π_i to select the cluster i , with the probability density function (Pdf) f_1 and f_2 to sample the data.

HMM can be considered as a dynamic mixture model. (See previous lecture notes)

② K-Means Algorithm Revisited

Identifying groups (clusters) of data points in a multi-dimensional space. Given K as the cluster number. We hope to minimize the an objective function (also called distortion measure)

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

~ (belonging to the cluster or not)

where r_{nk} is an indicator that:

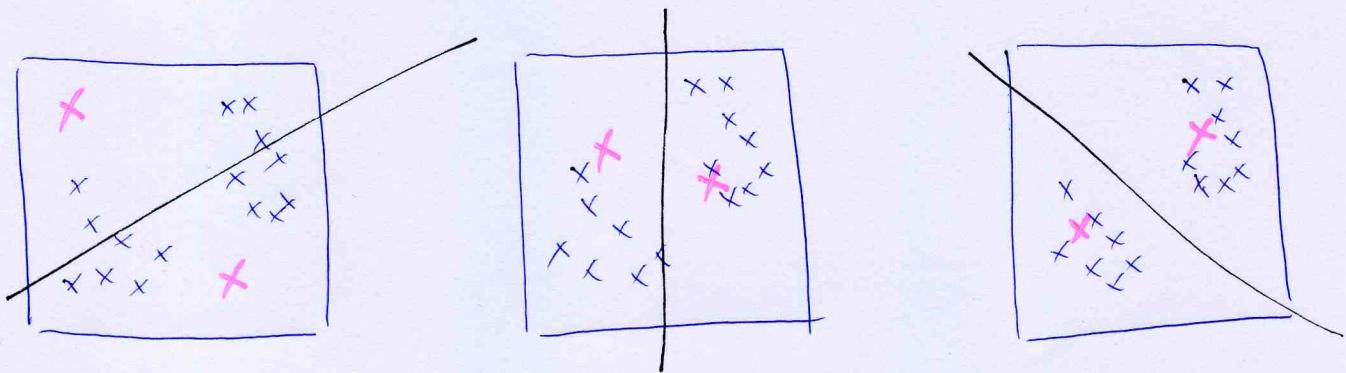
$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|x_n - \mu_j\|^2 \\ 0 & \text{otherwise} \end{cases}$$

If r_{nk} is fixed. (the clusters have existed!)

$$\frac{\partial J}{\partial \mu_k} = 2 \sum_{n=1}^N r_{nk} (x_n - \mu_k) = 0$$

$$\Rightarrow \mu_k = \frac{\sum_{n=1}^N r_{nk} x_n}{\sum_{n=1}^N r_{nk}}$$

↗ The number of elements in cluster k .
 ↗ all the elements in cluster k



Above figures are some schematic illustration of K-Means.

In calculating M , we call it E-step (Expectation)
given a set of new centers, we need to decide r_{nk}
(the way of assigning each data point x_n) in order to
maximize the objective function J . This is called
M-step (Maximization).

③ Given a general form of Gaussian Mixture Model

$$P(x) = \sum_z P(z) P(x|z) = \sum_{k=1}^K \pi_k N(x|\mu_k, \Sigma_k)$$

where $N(x|\mu_k, \Sigma_k) = \frac{1}{(2\pi)^{D/2} |\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)\right\}$

$$\begin{aligned} P(z_k=1|x) &= \frac{P(z_k=1) P(x|z_k=1)}{\sum_{k=1}^K P(z_k=1) P(x|z_k=1)} \\ &= \frac{\pi_k N(x|\mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x|\mu_j, \Sigma_j)} = r(z_k) \end{aligned}$$

$z_k : k=1 \dots K$ can be regarded as a random variable

$z_k = \begin{cases} 0 & \text{for the data point } x \text{ belonging to cluster } k \\ 1 & \end{cases}$

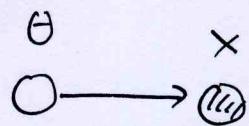
z is a multinomial random variable. e.g.

$z = [0 \ 1 \ 0 \ 0 \ 0 \ 0]$, $z = [0 \ 0 \ 0 \ 1 \ 0 \ 0]$ for $K=6$.

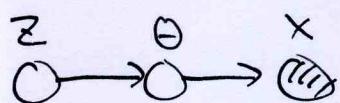


latent Variables

④ Bayesian and Frequentist



We can regard θ as the parameter.
It is just a parameter estimation problem



θ in Bayesian can also be treated as a random variable, which can be conditional on another random variable. Here, we refer to them as "latent variables", or hidden variables. Because they cannot be directly observed.

⑤ Maximum Likelihood of Gaussian Mixture Model (see previous lecture notes). In multi-dimensional case:

$$\mathcal{L} = \ln P(X | \pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k) \right\}$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mu_k} &= - \sum_{n=1}^N \frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_j \pi_j N(x_n | \mu_j, \Sigma_j)} \sum_k (x_n - \mu_k) \\ &= - \sum_{n=1}^N \gamma(z_{nk}) \underbrace{\sum_k (x_n - \mu_k)}_{\text{constant}} = 0 \end{aligned}$$

See Appendix

$$\begin{aligned} \Rightarrow \cancel{\sum_k} \sum_{n=1}^N \gamma(z_{nk}) &= \sum_{n=1}^N \sum_k (x_n - \mu_k) \\ &= \sum_k \left(- \sum_{n=1}^N \gamma(z_{nk}) x_n + \sum_{n=1}^N \gamma(z_{nk}) \cancel{\mu_k} \right) = 0 \\ \Rightarrow \mu_k &= \left(\sum_{n=1}^N \gamma(z_{nk}) x_n \right) \cdot \frac{1}{\sum_{n=1}^N \gamma(z_{nk})} \rightarrow \frac{1}{N_k} \\ \Rightarrow \boxed{\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) x_n} \end{aligned}$$

$$N_k = \sum_{n=1}^N \gamma(z_{nk}) = \sum_{n=1}^N P(z=k | x_n)$$

$$\frac{\partial \mathcal{L}}{\partial \Sigma_k} = 0 \Rightarrow \sum_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (x_n - \mu_k)(x_n - \mu_k)^T$$

Now, Let's consider the third parameter π_k with the constraint of $\sum_i \pi_i = 1$. Using lagrange multiplier

$$\mathcal{L} = \ln P(x|\mu, \Sigma, \pi) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right)$$

$$\frac{\partial \mathcal{L}}{\partial \pi_k} = 0 \Rightarrow \sum_{n=1}^N \frac{N(x_n | \mu_k, \Sigma_k)}{\sum_j N(x_n | \mu_j, \Sigma_j)} + \lambda = 0$$

$$\Rightarrow \pi_k = \frac{N_k}{N}$$

⑥ In Summary. The EM Algorithm for the Gaussian Mixture Model is as the following:

1. Initialize, μ_k , Σ_k and mixing coefficient π_k to evaluate the initial value of log likelihood.

2. E-step. evaluate the responsibilities using the current parameter values

$$r(z_{nk}) = \frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_n | \mu_j, \Sigma_j)} \text{ for each } x_n \quad n=1 \dots N$$

3. M-step: Re-estimate the parameter using current responsibilities

$$\mu_k^{new} = \frac{1}{N_k} \sum_{n=1}^N r(z_{nk}) x_n$$

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

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

$$\text{where, } N_k = \sum_{n=1}^N r(z_{nk})$$

4. Evaluate the log likelihood

$$\ln P(x|\mu, \Sigma, \pi) = \sum_{n=1}^N \ln \left(\sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k) \right)$$

where $N(x_n | \mu, \Sigma)$ is the pdf of multi-dim gaussian.

$$x \sim N(\mu_k, \Sigma_k)$$

$$P(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma_k|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) \right\}$$

~~Integrate over μ_k~~

$$\begin{aligned} J(x_1, \dots, x_N) &= \ln P(x_1, \dots, x_N | \mu_k, \Sigma_k) \\ &= \ln \prod_{n=1}^N P(x_n | \mu_k, \Sigma_k) = \sum_{n=1}^N \ln P(x_n | \mu_k, \Sigma_k) \\ &= -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\Sigma_k| - \\ &\quad - \frac{1}{2} \sum_{n=1}^N (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) \end{aligned}$$

$$\begin{aligned} \frac{\partial J}{\partial \mu_k} &= \left(-\frac{ND}{2} \ln(2\pi) \right)' - \left(\frac{N}{2} \ln |\Sigma_k| \right)' - \left(\frac{1}{2} \sum_{n=1}^N (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) \right)' \\ &= -\sum_{n=1}^N (x_n - \mu_k) = 0 \Rightarrow \mu_k = \frac{1}{N} \sum_{n=1}^N x_n \end{aligned}$$

For mixture model where $K = 1 \dots K$.

$$J = \underbrace{\sum_{k=1}^K \ln \{ \pi_k N(x_n | \mu_k, \Sigma_k) \}}_{P(x | \pi, \mu, \Sigma)} \quad P(x | \pi, \mu, \Sigma) =$$

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

$$\begin{aligned} F &= \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k f_k \right\} \quad \boxed{\frac{\partial F}{\partial f_k} = \frac{\pi_k}{\pi_k f_k}} \\ &= \sum_{n=1}^N \ln \{ \pi_1 f_1 + \pi_2 f_2 \dots \pi_K f_K \} \end{aligned}$$

$$\frac{\partial F}{\partial f_k} = \frac{\pi_k}{\pi_1 f_1 + \dots + \pi_K f_K}$$

$$\begin{aligned} \frac{f_k}{\partial \mu} \dots \frac{\partial F}{\partial \mu} &= \frac{\pi_K f_K}{\pi_1 f_1 \dots \pi_K f_K} \frac{\partial f_K}{\partial \mu} \\ &= r(\dots) \sum_k (x_n - \mu_k) \end{aligned}$$

Appendix 1