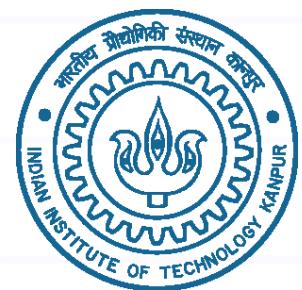


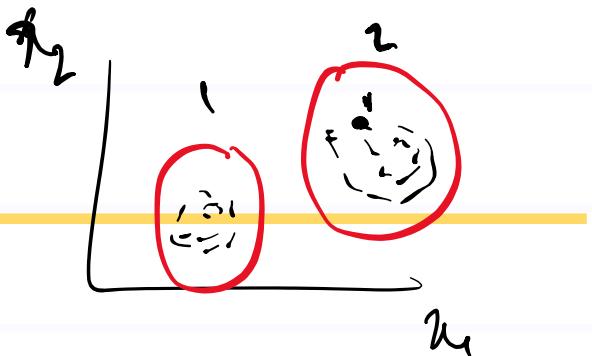
Gaussian Mixture Model

EE698V - Machine Learning for Signal Processing

Vipul Arora



K-means Clustering



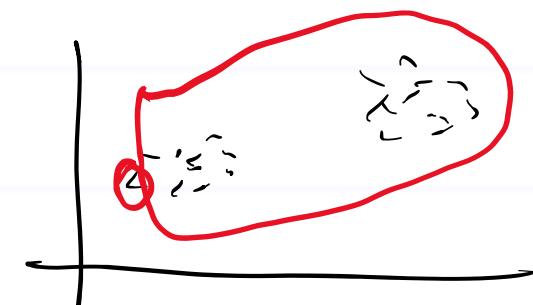
- Distortion measure: $\mathcal{L} = \sum_{n,k} r_{nk} \|x_k - \mu_k\|^2$

- k is the cluster index ✓

- n is the sample index ✓

- $r_{nk} = 1$ if sample n belongs to cluster k (hard clustering) r_{nk}
membership

- \mathcal{L} avoids very large (+ singular) clusters



Optimal value of r_{nk} and μ_k

$$\mathcal{L} = \sum_{n,k} r_{nk} \|x_n - \mu_k\|^2$$

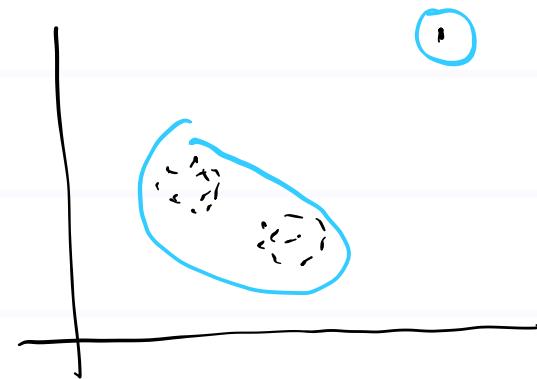
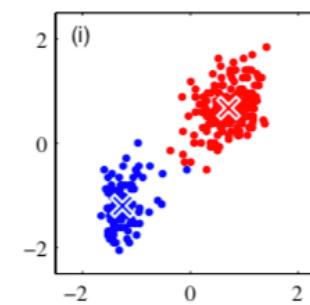
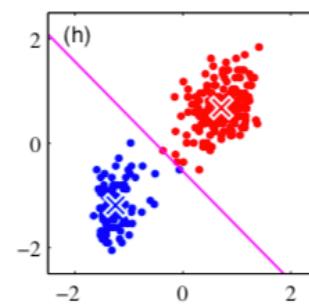
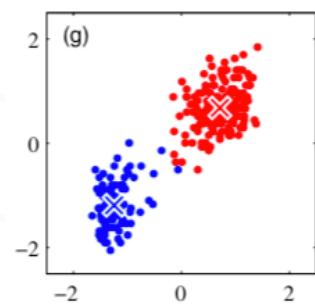
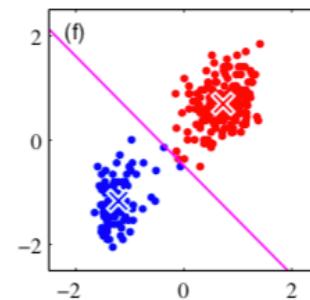
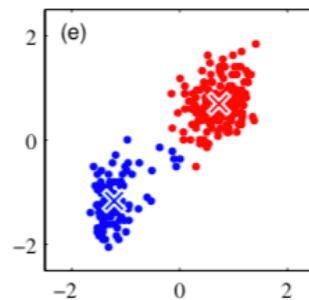
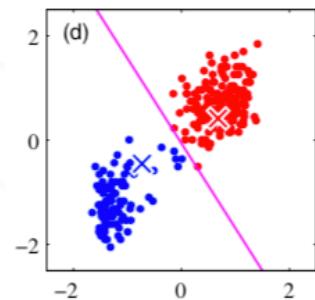
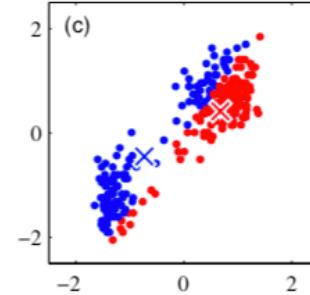
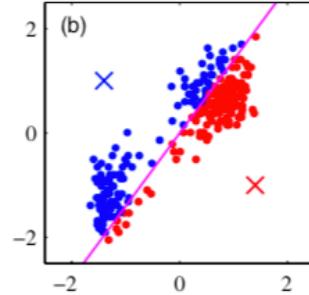
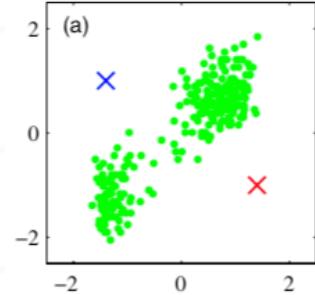
Given r_{nk}

$$\frac{\partial \mathcal{L}}{\partial \mu_{k'}} = 2 \sum_n r_{nk'} (x_{k'} - \mu_{k'}) = 0 \Rightarrow \mu_{k'} = \frac{\sum_n r_{nk'} x_{k'}}{\sum_n r_{nk'}}$$

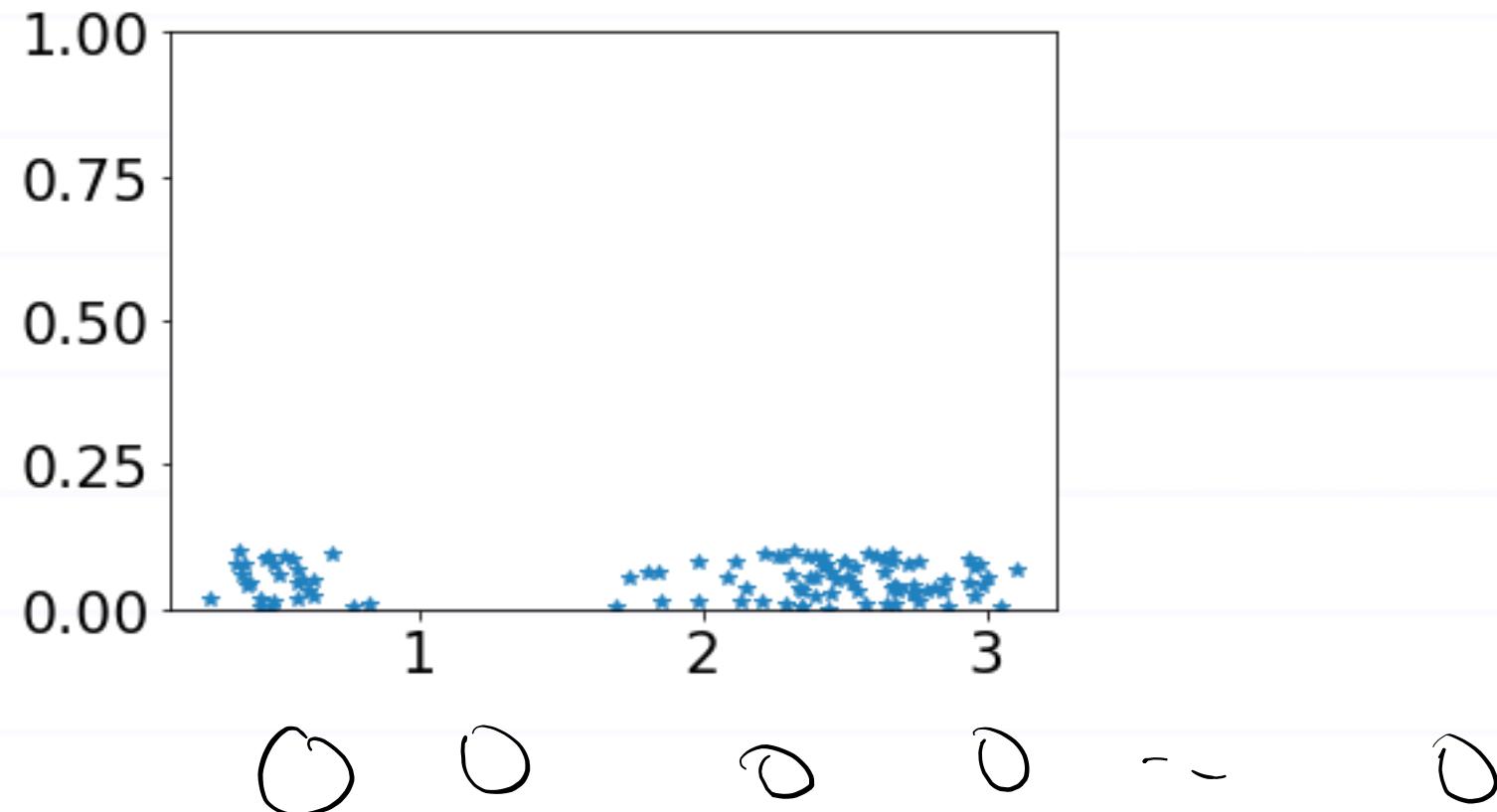
Given $\mu_{k'}$

$$r_{nk} = \begin{cases} 1 & , \text{ if } k = \operatorname{argmin}_{k'} \|x_n - \mu_{k'}\| \\ 0 & , \text{ o.w.} \end{cases}$$

K-means iteration



How would you model this data?



Latent Variable Models

Reference: PRML Chapter 9

Latent Variables

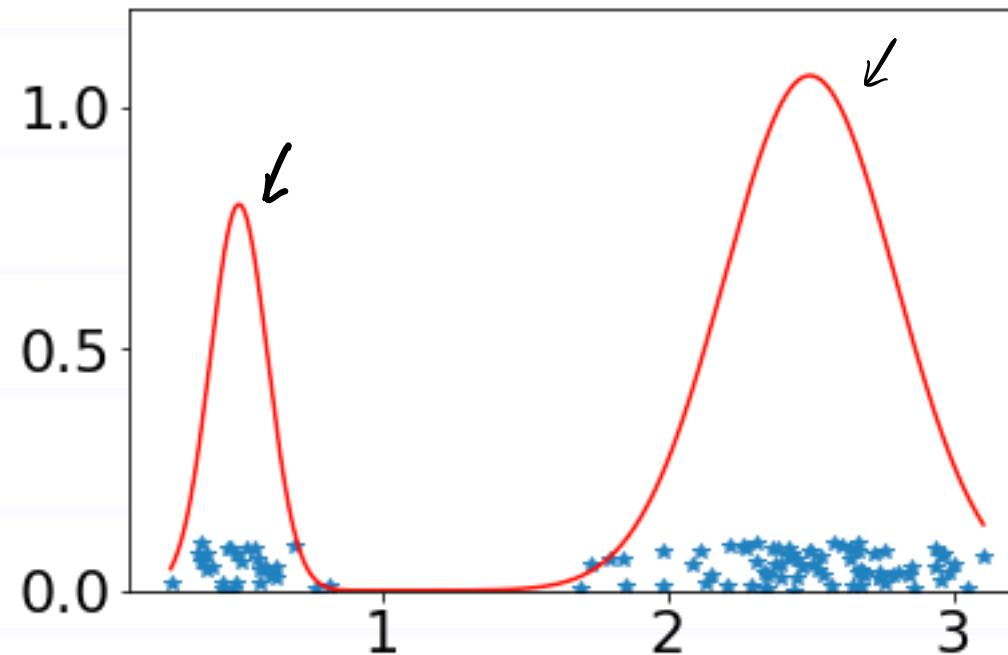
- Latent variables help complicated distributions to be formed from simpler components
- $p(x) = p(x|z=1)p(z=1) + p(x|z=2)p(z=2)$

$$p(x) = \sum_z p(x|z)p(z)$$

$\underbrace{\sum_z p(x,z)}_{\text{discrete r.v. } z \in \{1, 2\}}$ $\xrightarrow{\text{gauss}} \quad \xrightarrow{\text{p.m.f.}}$

Gaussian Mixture Model

$$p(x) = \sum_k p(z_k) \mathcal{N}(x; \mu_k, \sigma_k); z_k \in \{0,1\}, \sum_k z_k = 1$$



$$z = \begin{pmatrix} z_0 \\ z_1 \end{pmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}^T \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

one hot vector

Gaussian Mixture Model

- $p(\mathbf{x}) = \sum_z p(z)p(\mathbf{x}|z)$
 - $z \in \{1, 2, \dots, K\}$ is a discrete r.v.
 - $\mathbf{x} \in \mathbb{R}^D$

Equivalently,

- $p(\mathbf{x}) = \sum_k p(z_k = 1)p(\mathbf{x}|z_k)$
- z_k is k th element of a 1-hot vector. It is 1 if \mathbf{x} belongs to k th cluster
- $p(\mathbf{x}) = \sum_k \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$

for a given sample

$$z_1 = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \quad \text{3rd cluster}$$

$$z_3 = 1 ; z_i = 0 \text{ for all } i \neq 3$$

$$z_k = \delta_{k3}$$

Estimating the parameters

- Given the samples $\{s_n; n = 1, \dots, N\}$
 - the learnable parameters are $\theta = \{\pi_k, \mu_k, \Sigma_k; k = 1, \dots, K\}$
 - K is a hyperparameter

MLE for μ_k

$$\ln \mathcal{L} = \sum_n \ln p(s_n | \theta) = \sum_n \ln \left(\sum_k \pi_k N(s_n | \mu_k, \Sigma_k) \right)$$

$\frac{1}{\sqrt{2\pi^D |\Sigma|}}$

$e^{-\frac{1}{2}(s_n - \mu_k)^T \Sigma^{-1} (s_n - \mu_k)}$

Differentiate w.r.t. μ_k :

$$\frac{\partial \ln \mathcal{L}}{\partial \mu_k} = \sum_n \underbrace{\frac{\pi_{k'} N(s_n | \mu_{k'}, \Sigma_{k'})}{\sum_k \pi_k N(s_n | \mu_k, \Sigma_k)} \cdot \left(-\Sigma_{k'}^{-1} (s_n - \mu_{k'}) \right)}_{= \gamma_{nk'}} = 0$$

$\rightarrow p(z_{k'}=1) \text{ for } s_n$

- prob. that s_n belongs to the k^{th} cluster (soft assignment)
- k-means does hard assignment

MLE for μ_k

$$\Rightarrow \sum_n r_{nk} \mu_k = \sum_n r_{nk} s_n$$

$$\Rightarrow \hat{\mu}_k = \frac{\sum_n r_{nk} s_n}{N_k}$$

$N_k = \sum_n r_{nk}$: effective # samples assigned to cluster k

MLE for Σ_k

Differentiating $\ln L$ w.r.t $\Sigma_{k'}$

$$\frac{\partial \ln L}{\partial \Sigma_{k'}} = \sum_n \gamma_{n k'} \left(-\frac{1}{2|\Sigma|_k'^2} \frac{\partial |\Sigma|}{\partial \Sigma_{k'}} + \frac{1}{|\Sigma|} \frac{\partial}{\partial \Sigma} \left((\mathbf{s}_n - \boldsymbol{\mu}_{k'})^T \Sigma^{-1} (\mathbf{s}_n - \boldsymbol{\mu}_{k'}) \right) \right) = 0$$

\Rightarrow

$$|\Sigma| (\Sigma^{-1})^T$$

$$= \text{tr} \{ (\mathbf{s}_n - \boldsymbol{\mu}_{k'}) (\mathbf{s}_n - \boldsymbol{\mu}_{k'})^T \Sigma^{-1} \}$$

$$- (\Sigma^{-1})^T (\mathbf{s}_n - \boldsymbol{\mu}_{k'}) (\mathbf{s}_n - \boldsymbol{\mu}_{k'})^T (\Sigma^{-1})^T$$

$$\Rightarrow \sum_n \gamma_{n k'} \left[(\Sigma^{-1})^T - (\Sigma^{-1})^T (\mathbf{s}_n - \boldsymbol{\mu}_{k'}) (\mathbf{s}_n - \boldsymbol{\mu}_{k'})^T (\Sigma^{-1})^T \right] = 0$$

$$\sigma^2 = \frac{1}{n} \sum_n (x_n - \mu)^2$$

\Rightarrow

$$\Sigma_{k'} = \frac{1}{N_{k'}} \sum_n \gamma_{n k'} (\mathbf{s}_n - \boldsymbol{\mu}_{k'}) (\mathbf{s}_n - \boldsymbol{\mu}_{k'})^T$$

MLE for π_k

Differentiate $[\ln L + \lambda(\sum_k \pi_k - 1)]$ w.r.t. $\pi_{k'}$,

$$\frac{\partial [\cdot]}{\partial \pi_k} = \sum_n \frac{N(s_n | \mu_{k'}, \Sigma_{k'})}{\sum_k \pi_k N(s_n | \mu_k, \Sigma_k)} + \lambda = 0 \quad \dots (1)$$

Multiply with $\pi_{k'}$ and sum over k'

$$\Rightarrow \sum_n 1 + \lambda = 0 \Rightarrow \lambda = -N$$

Substituting in (1),

$$\sum_n \frac{\gamma_{nk'}}{\pi_{k'}} + (-N) = 0$$

$$\Rightarrow \pi_{k'} = \frac{N_k}{N}$$

MLE for GMM

- Now none of the estimates is independent of each other
- So follow iterative approach, as in K-means
 1. Initialize all parameters
 2. ASSIGN: calculate γ_k
 3. UPDATE θ
 4. Stop if converged, else go to step 2

Given L, find optimal θ :-

- $\frac{\partial L}{\partial \theta} = 0$, solve for θ possible if $\frac{\partial L}{\partial \theta} = 0$ is solvable for θ
- if θ 's are interdependent, use iterative update (in tandem).
 $\frac{\partial L}{\partial \theta}$
- If $\frac{\partial L}{\partial \theta}$ not solvable, use Grad. desc.

Expectation Maximization Algorithm

Reference: PRML Section 9.3

EM Algo

$$\ln p(S|\Theta) \xrightarrow{\substack{\text{all the} \\ \text{data samples}}} = \ln \left\{ \sum_z p(S, z | \Theta) \right\}$$
$$= \ln \left\{ \sum_z p(z | \Theta) p(S | z, \Theta) \right\}$$

$$\ln p(S | \Theta) = \sum_n \ln \left\{ \sum_k \pi_k \underbrace{N(s_n | \mu_k, \Sigma_k)}_{\text{exp, but ln can't act on it } \because \text{ of } \sum_k} \right\}$$

How to get rid of sum inside log?

EM Algo

- $p(S, Z | \theta)$, without marginalizing over Z , is exp.
- How to avoid marginalization?
- Assume a value of Z , and later average over all Z .

$$Q = \sum_Z \underbrace{p(z|s, \theta^{\text{old}})}_{\text{arg.}} \ln \underbrace{p(s, z|\theta)}_{\text{exp.}}$$

$$\underset{\theta}{\operatorname{argmax}} Q(\theta; \theta^{\text{old}})$$

But again $\frac{\partial Q}{\partial \theta}$ will include difficult terms \therefore take θ^{old} outside \ln .

$$Q = E_{z|s, \theta^{\text{old}}} [\ln p(s, z|\theta)]$$

EM Algo

$$p(s, z | \theta) = \prod_n \prod_k \left\{ \pi_k N(s_n | \mu_k, \Sigma_k) \right\}^{z_{nk}}$$

Latent variable for
cluster assignment

$$\ln p(s, z | \theta) = \sum_n \sum_k z_{nk} \ln \left\{ \pi_k N(s_n | \mu_k, \Sigma_k) \right\}$$

$$\therefore Q = \sum_z p(z | s, \theta^{old}) \ln p(s, z | \theta)$$

$$= \sum_n \sum_k \gamma_{nk} \ln \left\{ \pi_k N(s_n | \mu_k, \Sigma_k) \right\}$$

EM Algo

Iterative Algorithm :

1. Initialize θ^{old}
2. ASSIGN : Evaluate $p(z|x, \theta^{\text{old}})$, i.e., γ_{nk} E-step
3. UPDATE PARAMETERS : μ_k, Σ_k, π_k
by maximizing Q M-step
4. Stop if converged, else go to Step 2

Further Reading (Not mandatory)

- PRML Section 9.4