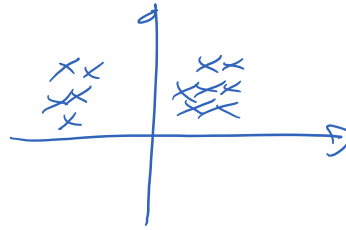


- today:
- unsupervised learning
  - K-means & EM

Unsupervised Learning

here  $X$  without any labels ✓



Consider the Gaussian mixture model (GMM)  
(can be obtained from FLD)

$$Y \sim \text{Mult}(\pi) \quad \pi \in \Delta_K$$

[extension of FLD]  
to multiple classes

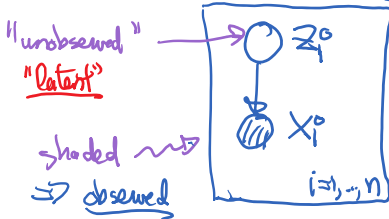
$$X | Y=j \sim N(\mu_j, \Sigma)$$

$$p(x) = \sum_y p(x,y) = \sum_y p(x|y)p(y) = \sum_{j=1}^K \pi_j N(x | \mu_j, \Sigma)$$

"GMM model"

more generally  
can have  $\Sigma_j$  per class

graphical model for this "latent variable model"



"plate" = repetition



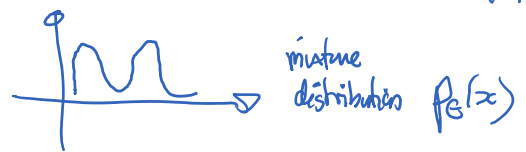
GMM

$$z_i \sim \text{Mult}(\pi)$$

$$x_i | z_i \sim N(x_i | \mu_{z_i}, \Sigma_{z_i})$$

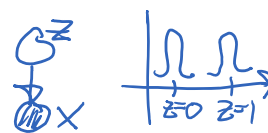
two views on  $p(x)$

unstructured



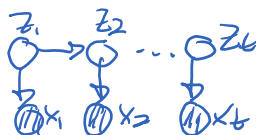
structured

latent variable model



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

(later in class, we will add  
time structure: HMM)



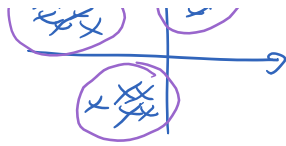
K-means

→ to do clustering i.e. group data



we want to get a cluster assignment for every data point  $x_i$

represent  $z_i, j=1$  to mean + put  $x_i$  belongs to cluster  $i$



represent  $z_{i,j} = 1$  to mean that  $x_i$  belongs to cluster  $j$   
 $j = 1, \dots, k$   
 # of clusters (specified in advance for k-means)

- applications:
- vector quantization (compression)
  - in computer vision: use k-means to get "bag of visual words" representation of image patches
  - many many others!

k-mean alg.: → can be derived as a block-coordinate minimization alg. of objective fct.:

(distortion measure) →  $J(z, \mu) \triangleq \sum_{i=1}^n \|x_i - \mu_{z_i}\|_2^2 = \sum_{j=1}^k \left( \sum_{i=1}^n z_{i,j} \|x_i - \mu_j\|_2^2 \right)$

cluster assignment  $z_1, \dots, z_n \in \text{corners of } \Delta_k$  ("one-hot encoding")

cluster index represented by  $z_i$

cluster centroids  $\mu_1, \dots, \mu_k \in \mathbb{R}^d$

$x_1, x_2, x_3, x_4$  (block 1)     $x_3, x_4$  (block 2)

alg.: 1) initialize  $\mu^{(1)}$

2) iterate until convergence:

"E" step:  $z^{(t+1)} = \underset{z \in \text{valid ass.}}{\text{argmin}} J(z, \mu^{(t)})$

⇒  $z_{i,j^*}^{(t+1)} = 1$  for  $j^* = \underset{j}{\text{argmin}} \|x_i - \mu_j^{(t)}\|$

"M" step:  $\mu^{(t+1)} = \underset{\mu \in \mathbb{R}^{d \times k}}{\text{argmin}} J(z^{(t+1)}, \mu)$

⇒  $\mu_j^{(t+1)} = \frac{\sum_i z_{i,j} x_i}{\sum_i z_{i,j}}$  empirical mean of cluster

Visualization: <https://www.naftaliharris.com/blog/visualizing-k-means-clustering/>

17h32

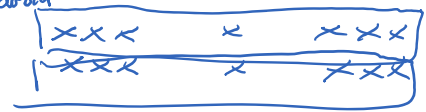
properties of k-means:

- converge in finite # of iterations to a local min
- NP hard in general to compute the global minimum in  $z$

k-means++: clever initialization scheme which guarantees that alg. is within  $\log k$  of global opt. (w.h.p.)  
 to avoid

guarantees that alg. is within  $\log k$  of global opt. (w.h.p.)

→ idea: spread out as much as possible the initial means



### 3) choice of $k$ ?

• one heuristic is:  $J(\mu, z, k) = \sum_{i=1}^n \sum_{j=1}^k z_{ij} \|x_i - \mu_j\|^2 + \lambda k$

→ we'll see later in class

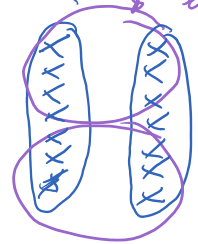
"non-parametric" models where " $k$ " is basically infinite and can get  $p(k|data)$

e.g. Dirichlet process mixture model

↑ hyperparameter

k-means solution

4) k-mean is very sensitive on distance measure: it assumes spherical clusters



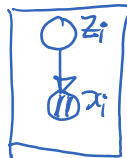
Side reference I mentioned: see <https://icml.cc/2012/papers/291.pdf> for interpreting regularized K-means as approximate inference in a Dirichlet process mixture model... [by Kulis & Jordan]

↳ GMM fixes that problem

Mahalanobis distance  $d_{\Sigma}(x, x') \stackrel{\text{def}}{=} \sqrt{(x-x')^T \Sigma^{-1} (x-x')}$

### EM - maximum likelihood in latent variable model

setup



$z$  latent variable

$x$  observed variable

log-likelihood  $\log p(x_{1:n}; \theta) = \log \left( \prod_{i=1}^n p(x_i; \theta) \right)$

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

$= \sum_{i=1}^n \log \left[ \sum_{z_i} p(x_i, z_i; \theta) \right]$

problem? → gives multi-modal opt. problem (non-convex)

options for MLE in latent variable model

$f$  is convex  $\Leftrightarrow -f$  is concave

1) do gradient ascent on a non-concave obj.

2) EM alg. → block-coordinate ascent on auxiliary fct. which lower bounds  $\log p(x_{1:n}; \theta)$   
nice interpretation in terms of filling "missing data"

i.e. E step → fill  $z$  with "soft-values"

M step → max w.r. to  $\theta$  for fully observed model



trick overview:

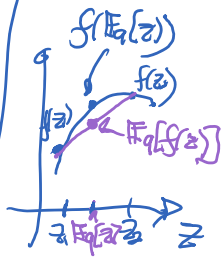
Jensen's inequality

trick overview:

$$\log \sum_z p(x,z) = \log \sum_z q(z) \frac{p(x,z)}{q(z)}$$

$$= \log \mathbb{E}_q \left[ \frac{p(x,z)}{q(z)} \right]$$

Jensen's inequality  
 $\mathbb{E}_q[f(g(z))] \leq f(\mathbb{E}_q(g(z)))$   
 when  $f$  is concave



Jensen's inequality trick:

$$\mathbb{E}_q \left[ \log \frac{p(x,z)}{q(z)} \right] = \sum_z q(z) \log p(x,z) - \sum_z q(z) \log q(z)$$

$$\triangleq \mathcal{J}(q, \theta) \triangleq \underbrace{\mathbb{E}_q \left[ \log p(x,z; \theta) \right]}_{\text{"expected complete log-likelihood"}} + \underbrace{H(q)}_{\text{"entropy of q"}}$$

we have  $\log p(x; \theta) \geq \mathcal{J}(q, \theta) \quad \forall q \in \mathcal{G}$

EM algorithm: E step:  $q_{t+1} \triangleq \arg \max_{q \in \mathcal{G}} \mathcal{J}(q, \theta_t) \Rightarrow q_{t+1}(z) = p(z|x; \theta_t)$

M step:  $\theta_{t+1} \triangleq \arg \max_{\theta} \mathcal{J}(q_{t+1}, \theta)$

$= \arg \max_{\theta} \mathbb{E}_{q_{t+1}(z)} \left[ \log p(x,z; \theta) \right]$

this is another MLE problem, but for complete information

(often, replace  $z$  with  $\mathbb{E}_q[z]$  in this expression)

from Jensen's inequality

\* we had  $\log p(x; \theta) \geq \mathcal{J}(q, \theta)$

$\log(\mathbb{E}_q[g(z)]) \geq \mathbb{E}_q[\log(g(z))]$

in Jensen's ineq, you get a strict ineq unless the dist. for  $g(z)$  is degenerate (i.e. takes only one value)

(when  $f$  is strictly concave)

i.e. when  $g(z) = \text{constant}$

[log is " " " ]

$g(z) = \text{constant}$  i.e.  $\frac{p(x,z)}{q(z)} = \text{const.} \quad \forall z \Rightarrow q(z) \propto p(x,z)$

$\mathcal{J}(q_{t+1}, \theta_t) = \log p(x; \theta_t) \geq \mathcal{J}(q, \theta_t) \quad \forall q$

i.e.  $q(z) = p(z|x; \theta)$

$\Rightarrow q_{t+1}$  maximizes  $\mathcal{J}(q, \theta_t)$  w.r. to  $q$

i.e.  $\arg \max_q \mathcal{J}(q, \theta_t) = p(z|x; \theta_t) = q_{t+1}$

and  $\mathcal{J}(q_{t+1}, \theta_t) = \log p(x; \theta_t)$

makes equality  $\mathcal{J}(q, \theta) = \log p(x; \theta)$