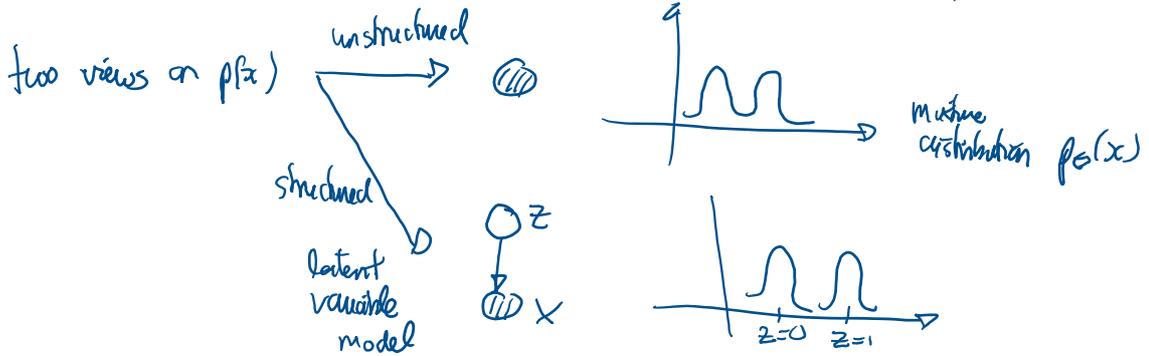
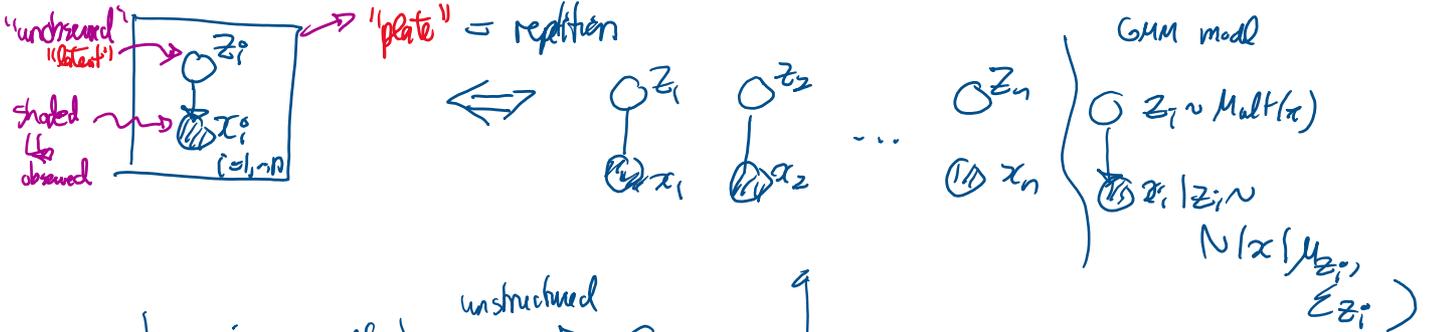


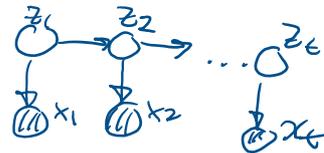
today: • k-means
• EM for GMM

graphical model for this "latent variable model"



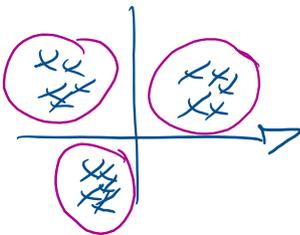
$$p(x) = \sum_z p_G(x|z) p_G(z)$$

(later in class, we will add fine structure:



HMM)

k-means → ^{alg.} to clustering i.e. group data



we want to get a cluster assignment for every data pt. x_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 derive as block-coordinate minimization dg.

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

cluster assignment → z_1, \dots, z_k
 \in corners of Δ_k ("one hot encoding")
 cluster centroids → $\mu_1, \dots, \mu_k \in \mathbb{R}^d$
 cluster index represented by z_i

- alg.:
- 1) initialize $\mu^{(1)}$
 - 2) iterate until convergence

"E step": $z^{(t+1)} = \underset{z \in \text{valid assign.}}{\text{argmin}} J(z, \mu^{(t)})$
 $\Rightarrow z_{ij}^{(t+1)} = 1$ for $j^* = \underset{j}{\text{argmin}} \|x_i - \mu_j^{(t)}\|$

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

$\Rightarrow \mu_j^{(t+1)} = \frac{\sum_i z_{ij} x_i}{\left(\sum_i z_{ij} \right)}$ empirical mean of cluster j

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

properties:

- 1) converge in finite # of iterations to a local min
- 2) NP hard in general to compute global min in z

note: run for multiple random initializations

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

→ idea: spread as much as possible the initial means

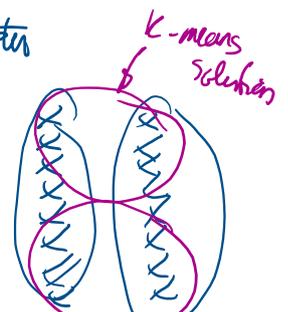


3) choice of k ?

Side reference I mentioned: • one heuristic is $J(z, \mu, k) = \sum_i \sum_j z_{ij} \|x_i - \mu_j\|^2 + \lambda k$
 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]
 hyperparameter

4) k-means is very sensitive in distance measure: it assumes spherical clusters

↳ GMM fixes this problem (partially)



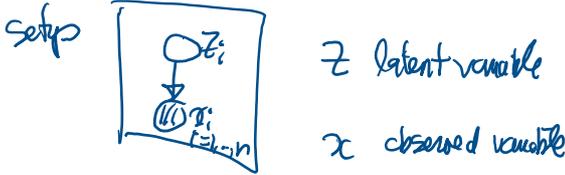
↳ GMM uses this problem (partially)



Ma halandris distance

$$d_{\Sigma}(x, x') \triangleq \sqrt{(x-x')^T \Sigma^{-1} (x-x')}$$

EM - maximum likelihood in latent variable model



$$\begin{aligned} \log\text{-likelihood } \log p(x_1, \dots, x_n; \Theta) &= \log \left(\prod_i p(x_i; \Theta) \right) \\ &= \sum_{i=1}^n \log p(x_i; \Theta) \\ &= \sum_{i=1}^n \log \left(\sum_{z_i} p(x_i, z_i; \Theta) \right) \end{aligned}$$

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

options MLE in latent variable model

- 1) do gradient ascent on a non-convex obj.
- 2) EM alg. → block coordinate ascent on a convex obj. that lower bounds $\log p(x_{1:n}; \Theta)$
 nice interpretation in terms of filling "missing data"
 ie. E step → fill z with "soft-values"
 M step → max u.r. to Θ for fully observed model $\begin{matrix} \odot z_i \\ \downarrow \\ \odot x_i \end{matrix}$

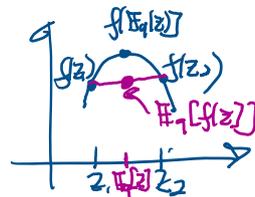
trick overview:

$$\begin{aligned} \log \sum_z p(x, z) &= \log \sum_z q(z) \frac{p(x, z)}{q(z)} \\ &= \log \mathbb{E}_q \left[\frac{p(x, z)}{q(z)} \right] \\ &\stackrel{\text{Jensen's ineq. trick?}}{\geq} \mathbb{E}_q \left[\log \frac{p(x, z)}{q(z)} \right] \\ &= \sum_z q(z) \log p(x, z) - \sum_z q(z) \log q(z) \end{aligned}$$

Jensen's inequality

$$\mathbb{E}_q [f(g(z))] \leq f(\mathbb{E}_q [g(z)])$$

when f is concave



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

13436

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

EM algorithm: E step: $q_{t+1} \triangleq \underset{q \in \text{distribution } \mathcal{Z}}{\text{argmax}} \mathcal{J}(q, \theta_t) \Rightarrow$ turns out $q_{t+1}(z) = p(z|x; \theta_t)$

M step: $\theta_{t+1} \triangleq \underset{\theta}{\text{argmax}} \mathcal{J}(q_{t+1}, \theta)$
 $= \underset{\theta}{\text{argmax}} \mathbb{E}_{q_{t+1}} [\log p(x, z; \theta)]$
 this is another MLE problem, but for complete information
 (often, replace z with $\mathbb{E}_{q_{t+1}}[z]$ in this expression)

E step derivation: by Jensen's inequality

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

in Jensen's inequality, you get a strict inequality (when f is strictly concave)

above $g(z) = \frac{p(x, z)}{q(z)} \stackrel{!}{=} \text{constant} \quad \forall z$
 $\Rightarrow q(z) \propto p(x, z)$

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

unless R.V. is deterministic

\hookrightarrow here $g(z)$ being deterministic \Rightarrow constant
 i.e. $g(z) = \text{constant}$

i.e. $\boxed{q^*(z) = p(z|x; \theta)}$

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

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

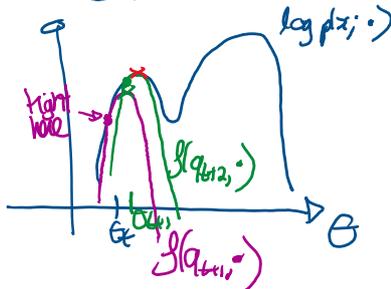
i.e. $\underset{q}{\text{argmax}} \mathcal{J}(q, \theta_t)$

$= p(z|x; \theta_t) \triangleq q_{t+1}$

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

properties of EM algorithm

a) $\log p(x; \theta_{t+1}) \geq \log p(x; \theta_t)$



proof: $\log p(x; \theta_{t+1}) \geq \mathcal{J}(q_{t+1}, \theta_{t+1})$

$\stackrel{\text{by M step}}{\geq} \mathcal{J}(q_{t+1}, \theta_t)$
 $\geq \log p(x; \theta_t)$

b) θ_t in EM converges to a stationary pt. of $\log p(x; \theta)$

i.e. $\nabla_{\theta} \log p(x; \theta) \Big|_{\theta} = 0$

$$\text{i.e. } \nabla_{\theta} \log p(x; \theta) \Big|_{\hat{\theta}} = 0$$

like k-means, initialization is crucial
 → usually do random centroids

• for GMM, could use k-means to initialize μ 's

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

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

$$= \mathbb{E}_q \left[\log \frac{q(z)}{p(z; \theta)} \right] \leq \text{KL}(q(\cdot) \parallel p(\cdot | x; \theta))$$

KL divergence

$$\left. \begin{matrix} \log p(x; \theta) \\ \mathcal{J}(q, \theta) \end{matrix} \right\} \text{KL}(q \parallel p(\cdot | x; \theta))$$

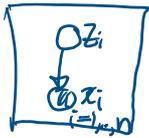


Block-coordinate method sometimes slow

we will revisit this for variational inference

$q \in \mathcal{Q}$
 (simple distributions)

for GMM model:



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

$$x_i | z_i = j \sim N(\mu_j, \Sigma_j)$$

↑ shorthand to say $z_{ij} = 1$

$$\theta = (\pi, \{\mu_j\}_{j=1}^k, \{\Sigma_j\}_{j=1}^k)$$

$$\text{notation: } x = x_{1:n}$$

$$z = z_{1:n}$$

$$\text{exercice } p(z|x) = \prod_{i=1}^n p(z_i|x) = \prod_{i=1}^n p(z_i|x_i)$$

complete log-likelihood

$$\log p(x, z; \theta) = \sum_{i=1}^n \log p(x_i | z_i; \theta) + \log p(z_i; \theta)$$

$$= \sum_{i=1}^n \left[\sum_{j=1}^k z_{ij} \log N(x_i | \mu_j, \Sigma_j) + \sum_{j=1}^k z_{ij} \log \pi_j \right]$$

↑ Gaussian ↑ multinomial

$$\mathbb{E}_q \left[\log p(x, z; \theta) \right] = \sum_{i=1}^n \left[\sum_{j=1}^k \mathbb{E}_q[z_{ij}] (\log N(x_i | \mu_j, \Sigma_j) + \log \pi_j) \right]$$

$$\mathbb{E}_q [1\{z_{ij}=1\}] = q(z_{ij}=1)$$

↑ marginal

$$\text{during EM: } q_{t+1}(z) = p(z|x; \theta_t)$$

$$\text{weight } \gamma_{i,j}^t \triangleq p(z_{i,j}=1|x; \theta_t) = q_{t+1}(z_{i,j}=1)$$

$\rightarrow q(z_i) = \prod_{j=1}^k \pi_j \mathcal{N}(x_i | \mu_j, \Sigma_j)$ (marginal of z_i) weight $\tau_{ij}^{(t)} \triangleq p(z_{ij}=1 | x_i; \Theta^t) = q_{t_{ij}}(z_{ij}=1)$

E-step: compute $q_{t_{ij}}(z) \triangleq p(z | x_i; \Theta^t)$
 $= \prod_{j=1}^k p(z_{ij} | x_i; \Theta^t)$
 $\Rightarrow q_{t_{ij}}(z_i) = p(z_i | x_i; \Theta^t)$
 $\propto p(x_i | z_i; \Theta^t) p(z_i; \Theta^t)$

$$\tau_{ij}^{(t)} = q_{t_{ij}}(z_{ij}=1) = \frac{\pi_j^{(t)} \mathcal{N}(x_i | \mu_j^{(t)}, \Sigma_j^{(t)})}{\sum_{l=1}^k \pi_l^{(t)} \mathcal{N}(x_i | \mu_l^{(t)}, \Sigma_l^{(t)})}$$

$\left\{ \begin{array}{l} \pi_j \\ p(x_i, z_i | \Theta^t) \\ p(x_i | \Theta^t) \end{array} \right.$

E step for GMM: compute $\tau_{ij}^{(t)}$ for $i=1, \dots, n$ using $\Theta^{(t)}$

M step: $\max_{\{\mu_j, \Sigma_j, \pi_j\}} \sum_{i,j} \tau_{ij}^{(t)} [\log p(x_i | \mu_j, \Sigma_j) + \log \pi_j]$

Exercise:

$$\hat{\pi}_j^{(t+1)} = \frac{\sum_i \tau_{ij}^{(t)}}{n}$$

$$\hat{\mu}_j^{(t+1)} = \frac{\sum_i \tau_{ij}^{(t)} x_i}{\sum_i \tau_{ij}^{(t)}}$$

$$\hat{\Sigma}_j^{(t+1)} = \frac{\sum_i \tau_{ij}^{(t)} (x_i - \hat{\mu}_j^{(t+1)}) (x_i - \hat{\mu}_j^{(t+1)})^T}{\sum_i \tau_{ij}^{(t)}}$$

$\tau_{ij}^{(t)}$ soft-counts
 $\tau_{ij}^{(t)}$ big spherical covariance $\Sigma_j^{(t)} = \sigma^2 I$
 $\pi_j^{(t)}$ from proportions from k-means

- initialize eg. $\mu_j^{(0)}$ from k-means++
- EM step in GMM with $\Sigma_j = \sigma^2 I$ with $\sigma^2 \rightarrow 0$
 \rightarrow get k-means by "hard EM"