

today: • finish k-means  
 • EM (★) & GMM (☺)

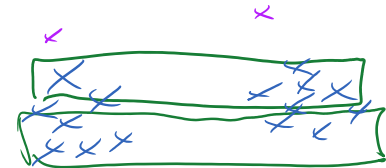
properties of k-means:

K-means demo: <http://web.stanford.edu/class/ee103/visualizations/kmeans/kmeans.html>

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

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

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



3) choice of k?

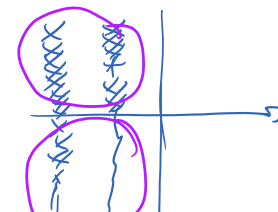
• one heuristic is:  $J(\mu, z, k) = \sum_{i=1}^n \sum_{j=1}^k z_{ij} \|x_i - \mu_j\|^2 + \lambda k$   
 hyperparameter

→ 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

4) k-mean is very sensitive on distance measure: if assumes spherical cluster



14/04

↳ GMM fixes that



## EM - maximum likelihood in latent variable model



$$\begin{aligned}\log\text{-likelihood } \log p(x_{1:n}; \theta) &= \log \prod_{i=1}^n p(x_i; \theta) \\ &= \sum_{i=1}^n \log p(x_i; \theta) \\ &= \sum_{i=1}^n \log \left[ \sum_{z_i} p(x_i, z_i; \theta) \right]\end{aligned}$$

problem?

→ gives multi-modal difficult optimization problem (non-convex)

options for ML in latent variable model

- 1) do gradient ascent on non-convex objective
- 2) EM alg. → block-coordinate ascent on auxiliary function which lower bounds  $\log p(x_{1:n}; \theta)$

nice interpretation in terms of filling "missing data"

i.e. E step  $\rightarrow$  finds  $z$  with "soft-values"

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

trick overview:

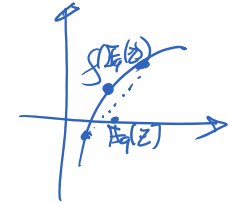
$$\log \sum_z p(x, z) = \log \sum_z q(z) \frac{p(x, z)}{q(z)}$$

$f(\cdot)$  here is  $\log(\cdot)$

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

$$\begin{aligned} &\geq \mathbb{E}_q \left[ \log \frac{p(x, z)}{q(z)} \right] = \sum_z q(z) \log p(x, z) - \sum_z q(z) \log q(z) \\ &\stackrel{\text{Jensen's inequality trick?}}{\cong} \mathcal{L}(q, \theta) \triangleq \mathbb{E}_{q(z)} [\log p(x, z; \theta)] + H(q) \end{aligned}$$

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



"entropy of  $q$ "

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

we get equality in Jensen if  $\frac{p(x, z)}{q(z)} = \text{constant w.r. to } z \Rightarrow \frac{p(x, z)}{q(z)} = c$   
 $\mathbb{E}_q [c] = c$

(explanation:  $f(z) = \log(z)$  is a \*strictly\* concave function; so this means Jensen's inequality becomes an equality only if the R.V. is constant -- this implies that  $p(x, z)/q(z)$  should not depend on  $z$  for this to happen)

i.e.  $q^*(z) \propto p(x, z)$

$$\Rightarrow q^*(z) = \frac{p(x, z)}{\sum_z p(x, z)} = p(z|x)$$

$$\begin{aligned} \log \mathbb{E}_q [c] &= \log c \\ \mathbb{E}_q [\log c] &= \log c \end{aligned}$$

this means that  $\underset{q \in \text{distributions over } z}{\text{argmax}} \mathcal{L}(q, \theta) = p(z|x; \theta)$

inference

$q(z)$  distributions over  $z$

inference

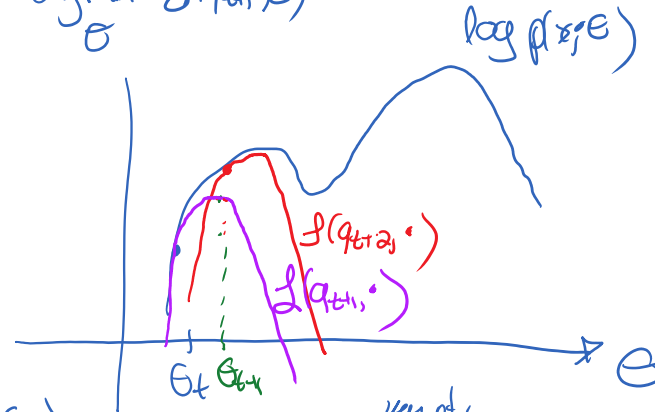
EM algorithm: E step:  $q_{t+1} \triangleq \arg \max_q \mathcal{J}(q, \theta_t) \Rightarrow q_{t+1}(z) = p(z|x; \theta_t)$

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

block-coordinate ascent on  $\mathcal{J}(q, \theta) \leq \log p(x; \theta)$

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

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



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

14/55

properties:

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

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

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

like K-means, initialization is crucial

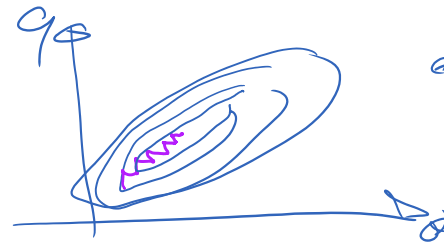
$\rightarrow$  usually random restarts

for GMM, could use K-mean++ to initialize the  $\mu$ 's

c)  $\mathcal{J}(q, \theta) = \mathbb{E}_q[\log p(x, z; \theta)]$   $n(z|x; \theta)$

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

$$\begin{aligned} \log p(x, \theta) - \mathcal{J}(q, \theta) &= -\mathbb{E}_q \left[ \log \frac{p(x, z; \theta)}{q(z) p(z; \theta)} \right] \\ &= \mathbb{E}_q \left[ \log \frac{q(z)}{p(z; \theta)} \right] \triangleq \text{KL}(q(\cdot) \| p(\cdot | x, \theta)) \\ &\quad (\text{KL divergence}) \end{aligned}$$



where conjugate gradient could be faster

$$M \text{ step: } \Theta_{t+1} = \underset{\Theta \in \Theta}{\text{argmax}} \underbrace{\mathbb{E}_{q_{t+1}} \left[ \log p(z, x; \theta) \right]}_{\text{expected complete log-likelihood}} + \text{cst.}$$

→ this is another ML problem, but for complete information

(usually replace  $z$  with  $\mathbb{E}_q[z]$ )

for GMM model:



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

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

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

notation here  $x = x_{1:n}$

$\downarrow$   
 $\mathbb{R}^n$

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

↑  
shorthand for  $z_{i,j} = 1$

notation here  $x = x_{1:n}$   
 $z = z_{1:n}$

complete log-likelihood

$$\begin{aligned} \log p(x, z; \theta) &= \sum_{i=1}^n \left[ \log p(x_i | z_i; \theta) + \log p(z_i; \theta) \right] \\ &= \sum_{i=1}^n \left[ \sum_{j=1}^k z_{i,j} \log N(x_i | \mu_j, \Sigma_j) + \sum_{j=1}^k z_{i,j} \log \pi_j \right] \end{aligned}$$

Gaussian log-likelihood

multinomial

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

↓  
 $\mathbb{E}_q[z_{i,j}] = q(z_{i,j} = 1)$  marginal distribution

E step is computing  $q_{t+1}(z) \triangleq p(z | x; \theta_t)$

$$\propto p(x | z; \theta_t) p(z; \theta_t)$$

$$\prod_{i=1}^n (p(x_i | z_i; \theta_t) p(z_i; \theta_t))$$

$$\Rightarrow q_{t+1}(z_i) \propto p(x_i | z_i; \theta_t) p(z_i; \theta_t)$$

↓

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

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

$$= q_{t+1}(z_{i,j} = 1) = \frac{\pi_j^{(t)} N(x_i | \mu_j^{(t)}, \Sigma_j^{(t)})}{\sum_{k=1}^k \pi_k^{(t)} N(x_i | \mu_k^{(t)}, \Sigma_k^{(t)})} \} p(x_i, z_{i,j} = 1 | \theta^{(t)})$$

$$\frac{\sum_{k=1}^k \pi_k^{(t)} N(x_i | \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_{k=1}^k \pi_k^{(t)} N(x_i | \mu_k^{(t)}, \Sigma_k^{(t)})} \} p(x_i | \theta^{(t)})$$

$$\prod_{l=1}^k \pi_l^{(t)} N(x | \mu_l^{(t)}, \Sigma_l^{(t)}) \quad \} \quad p(x | \theta^{(t)})$$

E step: compute  $\tau_{ij}^t$  for  $i=1, \dots, n$   
 $j=1, \dots, k$

M step:  $\max_{\{\mu_j, \Sigma_j, \pi_j\}} \sum_{i=1}^n \sum_{j=1}^k \tau_{ij}^t [\log p(x_i | \mu_j, \Sigma_j) + \log \pi_j]$

exercise:  $\hat{\pi}_j^{(t+1)} = \sum_i \tau_{ij}^t$  "soft count"

$\hat{\mu}_j^{(t+1)} = \frac{\sum_{i=1}^n \tau_{ij}^t x_i}{(\sum_{i=1}^n \tau_{ij}^t)}$  "soft-cluster assignment" [weighted empirical mean]

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

initializations: e.g.  $\mu_j^{(0)}$  from k-means++

$\Sigma_j^{(0)}$  big spherical variance

$\pi_j^{(0)}$  proportions k-means++

$\Sigma_j^{(0)} = \sigma^2 I$   
 $\uparrow$   
 big