

today: finish EM algorithm + GMM MLE  
 graph theory & DGM

properties of EM algorithm

recall: block coordinate ascent on  $J(q, \theta) \equiv \log p(x; \theta) \quad \forall \theta \neq q$

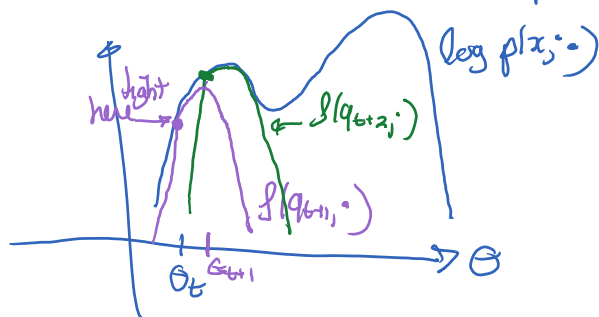
note:  $\log p(x; \theta_t) = J(q_{t+1}, \theta_t)$

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

properties:

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

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



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

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

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

for GMM

could use k-mean++ to initialize the  $\mu$ 's

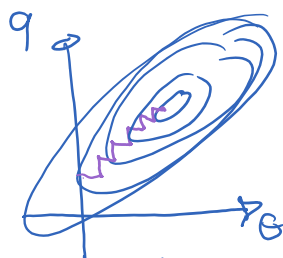
c)  $J(q, \theta) = \mathbb{E}_q \left[ \log \frac{p(x, z; \theta)}{q(z)} \right]$

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

$\log p(x; \theta) - J(q, \theta) \} \text{KL}(q || p(\cdot|x; \theta))$

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

we will revisit this for variational inference  $q \in \mathcal{Q}$



block-coordinate method can sometimes be slow

for GMM model:



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

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

short-hand to say  $z_{i,j} = 1$

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

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

$z = z_{1:n}$

exercise:

$$p(z|x) = \prod_i p(z_i|x) = \prod_i p(z_i|x_i)$$

complete log-likelihood:

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

Gaussian

multinomial

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

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

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

$$\mathbb{E}_q [z_{i,j}] = q(z_{i,j}=1) \text{ [marginal dist.]}$$

during EM,  $q_{t+1}(z) = p(z|x; \Theta_t)$

weight  $\hat{\pi}_{i,j}^t \triangleq p(z_{i,j}=1 | x_i; \Theta_t) = q_{t+1}(z_{i,j}=1)$

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

$$= \prod_i p(z_i | x_i; \Theta_t)$$

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

$$\propto p(x_i | z_i; \Theta_t) p(z_i; \Theta_t)$$

Gaussian

$\prod_{z_i}$

$$\hat{\pi}_{i,j}^t = q_{t+1}(z_{i,j}=1) = \frac{\pi_j^{(t)} N(x_i | \mu_j^{(t)}, \Sigma_j^{(t)})}{\sum_{l=1}^k \pi_l^{(t)} N(x_i | \mu_l^{(t)}, \Sigma_l^{(t)})} \left\{ \begin{array}{l} p(x_i, z_{i,j}=1 | \Theta^{(t)}) \\ p(z_i | \Theta^{(t)}) \end{array} \right.$$

E step for GMM: computing  $\hat{\pi}_{i,j}^{(t)}$  for  $i=1, \dots, n$  using  $\Theta^{(t)}$

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

exercise:  $\hat{\pi}_i^{(t+1)} = \sum_j \hat{\pi}_{i,j}^{(t)}$  "soft-counts"

exercise :

M step  
for EM  
for GMM

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

"soft-counts"

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

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

- initialize: e.g.  $\mu_j^{(0)}$  from k-means++  
 $\Sigma_j^{(0)}$  big spherical covariance  $\Sigma_j^{(0)} = \sigma^2 I$   
 $\pi_j^{(0)}$ : proportions from k-means++ big

- EM step in GMM with fixed  $\Sigma_j = \sigma^2 I$  with  $\sigma^2 \rightarrow 0$   
 $\rightsquigarrow$  get k-means alg?

16403

Graphical model

graph. model  $\rightsquigarrow$  proba theory + C.S.  
 $\downarrow$   $\downarrow$   
 R.V. graph

graph  $\rightarrow$  efficient data structure

eg.  $X_1, \dots, X_n$  R.V.  
 $X_i \in \{0, 1\}$

$n=100s$   
 $\Rightarrow 2^{100}$  #'s table  $\rightarrow$  intractable

QMR



Graph theory review:

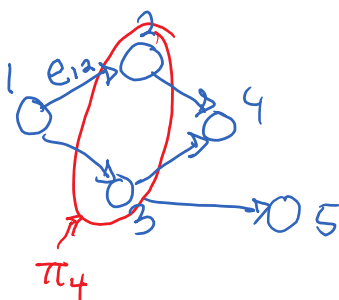
directed graph  
 "digraph"

$G = (V, E)$

$V = \{1, \dots, n\}$  "nodes/vertices"

$E \subseteq V \times V$  "directed edges"

$e_{12} = (1, 2)$



directed path

$1 \rightsquigarrow 4$  compatible seq. of edges  
 $(1, 2), (2, 4)$

$\pi_4$   
 $\pi_i \triangleq \{j \in V : \exists (j, i) \in E\}$   
 set of parents of  $i$

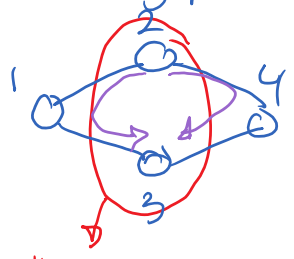
seq. of edges  
 $(1,2), (2,4)$   
 or  
 $(1,3), (3,4)$

children( $i$ )  $\triangleq \{j \in V : \exists (i, j) \in E\}$

(note: no self loop in this class i.e.  $|e|=2$ )

undirected graph:  $G = (V, E)$

where elements of  $E$  are 2-sets  
 (set of 2 elements)  
 thus we have  $\{i, j\} = \{j, i\}$



vs.  $(i, j) \neq (j, i)$  [order matters]

undirected path 2 and 3

"neighbors" of node 4 or 1

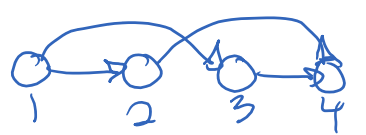
$N(i)$  for neighbors of  $i$   
 $= \{j \in V : \exists (i, j) \in E\}$

$\oplus$  neighbors replace the parents/children terminology from a digraph

def: DAG = directed acyclic graph = digraph with no cycles

def: an ordering  $I: V \rightarrow \{1, \dots, |V|\}$  is said to be topological for  $G$   
 iff nodes in  $\pi_i$  appear before  $i$  in  $I \forall i$

$$(j \in \pi_i \Rightarrow I(j) < I(i))$$

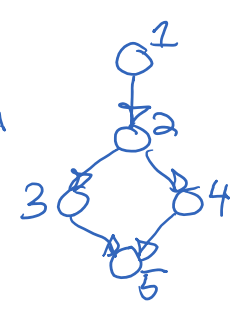


$\rightarrow$  if top ordering  $\Rightarrow$  all edges go from left to right  
 [no "back edge"]

prop: digraph  $G$  is a DAG  $\Leftrightarrow \exists$  a topological ordering of  $G$

proof:  $\Leftarrow$ ) trivial: no back edge  $\Rightarrow$  no cycle

$\Rightarrow$  use DFS algorithm to label nodes in decreasing order when have no children



top. sort  
 $\rightarrow$  finding a topological ordering

to construct a top. sort.  
 in  $O(|E| + |V|)$

Notation for graph models

$n$  discrete R.V.  $X_1, \dots, X_n$   
 $\hookrightarrow$  discrete R.V. for simplicity

conditioned dist. concept is

↳ discrete R.V. for simplicity  
 $V \leftarrow$  set of vertices  
 one R.V. variable per node

conditional dist. concept is  
tricky to formalize for cfs. R.V.  
 (see Borel-Kolmogorov paradox)

joint  $p(X_1=x_1, \dots, X_n=x_n) \stackrel{\text{short-hand}}{=} p(x_1, \dots, x_n) \quad x = x_{1:n}$   
 $= p(x_V) \stackrel{(P)}{=} p(x)$

for any  $A \subseteq V$   
 $p(x_A) = P\{X_i=x_i : i \in A\} = \sum_{x_{A^c} \in \mathcal{R}} p(x_A, x_{A^c})$   
subset of "subscripts"  $x_{A^c} \in \mathcal{R}$  summing over all possible values of  $x_{A^c}$  can take i.e.  $(x_i)_{i \in A^c} = V \setminus A$

Question: is  $p(x_1, x_2, x_4) \stackrel{?}{=} p(x_2, x_1, x_4)$ ?  
 yes! usually in this class, we use "typed convention"

revisit cond. indep.:

let  $A, B, C \subseteq V$

(\*)  $X_A \perp\!\!\!\perp X_B \mid X_C$

(F)  $\Leftrightarrow p(x_A, x_B \mid x_C) = p(x_A \mid x_C) p(x_B \mid x_C)$   $\forall x_A, x_B, x_C$  s.t.  $p(x_C) > 0$

(C)  $\Leftrightarrow p(x_A \mid x_B, x_C) = p(x_A \mid x_C)$   $\forall x_A, x_B, x_C$  s.t.  $p(x_C) > 0$

(\*) "marginal independence":  $X_A \perp\!\!\!\perp X_B \mid \emptyset$   
or  $X_A \perp\!\!\!\perp X_B$

$p(x_A, x_B) = p(x_A) p(x_B)$

2 facts about cond. indep.:

1) can repeat variables in statement (for convenience)

$X \perp\!\!\!\perp Y, Z \mid Z, W$  is fine to say  
 is equivalent  $X \perp\!\!\!\perp Y \mid Z, W$   
does not do anything

2) decomposition:  $X \perp\!\!\!\perp (Y, Z) \mid W \Rightarrow X \perp\!\!\!\perp Y \mid W$   
 $X \perp\!\!\!\perp Z \mid W$

(\*) pairwise indep.  $\not\Rightarrow$  mutual indep.  $\rightarrow$  see lecture 3

$Z = X \oplus Y$   
XOR

(\*) chain rule  $p(x_V) = \prod_{i=1}^n p(x_i \mid x_{1:i-1})$  last cond.  $p(x_n \mid x_{1:n-1})$  table with  $2^n$  entries

(always true)  $\bar{i}=1$

assumption in DBM:

$$p(x_v) = \prod_{i=1}^n p(x_i | x_{\pi_i}) \rightarrow \text{table of } 2^{\max\{k_i, l_i\}}$$

$X_0 \perp\!\!\!\perp X_{1:\bar{i}-1} \mid X_{\pi_i}$