

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

properties of EM algorithm

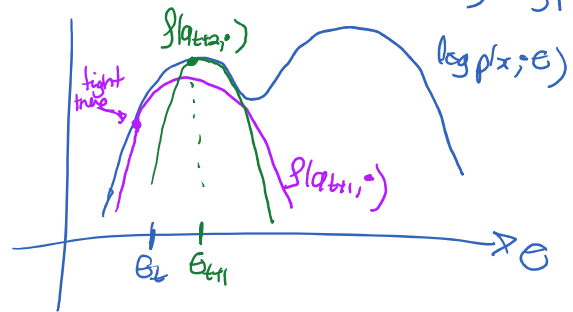
recall: block coordinate ascent on  $f(q, \theta) \leq \log p(x; \theta) \forall q$

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

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

properties:

a)  $\log p(x; \theta_{t+1}) \geq \log p(x; \theta_t)$        $\log p(x; \theta_{t+1}) \geq \mathcal{L}(q_{t+1}, \theta_{t+1})$   
 $\geq \mathcal{L}(q_{t+1}, \theta_t) = \log p(x; \theta_t)$



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

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

like K-means, initialization is crucial  
 → usually use random restarts

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

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

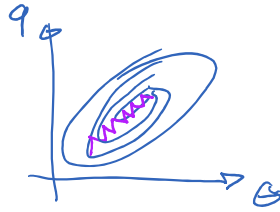
$$\begin{aligned} \log p(x; \theta) - \mathcal{L}(q, \theta) &= - \mathbb{E}_q \left[ \log \frac{p(x, z; \theta)}{q(z) p(x; \theta)} \right] \\ &= \mathbb{E}_q \left[ \log \frac{q(z)}{p(z|x; \theta)} \right] \triangleq \text{KL}(q(\cdot) \parallel p(\cdot|x; \theta)) \end{aligned}$$

KL divergence

(we will revisit this)

for variational inference  $q \in \mathcal{Q}$

block-coordinate method can be slow



for GMM model:



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

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

↑ shorthand to say  $z_{i,j}=1$

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

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

$$Z = Z_{1:n}$$

complete log-likelihood:

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

↓ Gaussian
↓ multinomial

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

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

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

weight  $\tau_{ij}^t \triangleq p(z_{i,j}=1 | x_i; \Theta^t) = q_{t+1}(z_{i,j}=1)$

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

$$\propto p(x | z; \Theta^t) p(z; \Theta^t)$$

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

$$\Rightarrow q_{t+1}(z_i) \propto \underbrace{p(x_i | z_i; \Theta^t)}_{\text{Gaussian}} \underbrace{p(z_i; \Theta^t)}_{\pi_{z_i}}$$

$$\tau_{ij}^t = q_{t+1}(z_{i,j}=1) = \frac{\pi_j^{(t)} N(x_i | \mu_j^{(t)}, \Sigma_j^{(t)})}{\sum_{l=1}^k \pi_l^{(t)} N(x_i | \mu_l^{(t)}, \Sigma_l^{(t)})} \left. \vphantom{\frac{\pi_j^{(t)} N(x_i | \mu_j^{(t)}, \Sigma_j^{(t)})}{\sum_{l=1}^k \pi_l^{(t)} N(x_i | \mu_l^{(t)}, \Sigma_l^{(t)})}} \right\} p(x_i, z_{i,j}=1 | \Theta^t)$$

$$\left. \vphantom{\frac{\pi_j^{(t)} N(x_i | \mu_j^{(t)}, \Sigma_j^{(t)})}{\sum_{l=1}^k \pi_l^{(t)} N(x_i | \mu_l^{(t)}, \Sigma_l^{(t)})}} \right\} p(x_i | \Theta^{(t)})$$

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

$$j = 1, \dots, k$$

$$M \text{ step} : \max_{\{\mu_j, \Sigma_j, \pi_j\}} \sum_i \sum_j \gamma_{ij}^t [\log p(x_i | \mu_j, \Sigma_j) + \log \pi_j]$$

M step for EM for GMM

Exercise :

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

"soft-count"

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

soft-cluster assignment

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

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

$\Sigma_j^{(0)}$  big spherical variance  $\Sigma_j^{(0)} = \sigma^2 I$   
 $\rho$  big

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

EM step in GMM with fixed  $\Sigma_j = \sigma^2 I$  with  $\sigma \rightarrow 0$

$\Leftrightarrow$  k-means

14h55

## Graphical model

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

graph  $\rightsquigarrow$  efficient data structure

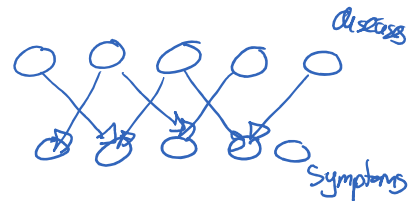
e.g.  $X_1, \dots, X_n$  R.V.'s

$$X_i \in \{0, 1\}$$

$$n \approx 100$$

$\Rightarrow 2^{100}$  #'s table  $\rightsquigarrow$  untractable

QMR



## Graph theory review:

directed graph

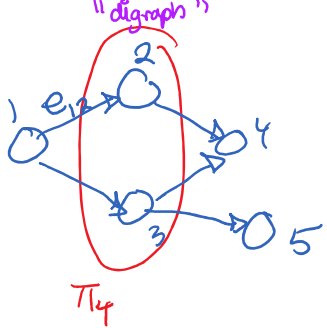
"digraph"

$$G = (V, E)$$

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

$$E \subseteq V \times V \text{ "directed edges"}$$

undirected graph  $G = (V, E)$



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

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

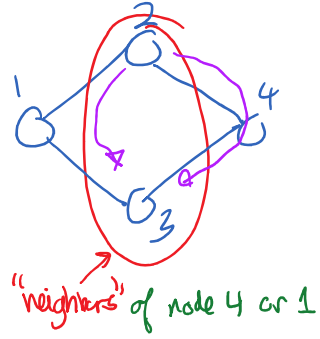
$e_{1,2} = (1, 2)$

directed path  $1 \rightsquigarrow 4$   $\rightarrow$  compatible sequence of edges  
 $(1, 2), (2, 4)$   
 or  
 $(1, 3), (3, 4)$

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

elements of  $E$  are 2-sets  
 (set of 2 elements)

note: no self loop  $\Rightarrow |e| = 2$



thus we have  $\{i, j\} = \{j, i\}$

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

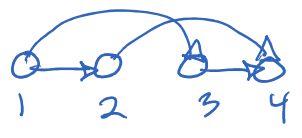
undirected path  $2 \rightsquigarrow 3$

neighbors replace the parents/children terminology from digraphs

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

def: an ordering  $I: V \rightarrow \{1, \dots, n\}$  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 edges"]

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

Notation for graph models:

$n$  discrete R.V.  $X_1, \dots, X_n$

discrete R.V. for simplicity

conditional dist. concepts tricky to formalize

see Borel-Kolmogorov paradox

$V$  set of vertices

one R.V. per node shorthand

$$\text{joint } p(X_1=x_1, \dots, X_n=x_n) \stackrel{!}{=} p(x_1, \dots, x_n) \quad \mathcal{X} = x_{1:n}$$

$$= p(x_V) \stackrel{(\ast)}{=} p(z)$$

for any  $A \subseteq V$

$$p(x_A) = P\{X_i=x_i : i \in A\} = \sum_{x_{A^c}} p(x_A, x_{A^c})$$

subsets of 'subscripts'

summing over all possible values of  $\{x_i : i \in V \setminus A\}$

$$x_{\{1,2,4\}} \leftrightarrow \{x_1, x_2, x_4\}$$

Question: is  $p(x_1, x_2, x_4) = p(x_2, x_1, x_4)$ ?

yes usually in this class ("typed convention")

revisit cond. indep.:

let  $A, B, C \subseteq V$

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

(F)  $\Leftrightarrow p(x_A, x_B \mid x_C) = p(x_A \mid x_C) p(x_B \mid x_C) \quad \forall x_A, x_B, x_C \text{ st. } p(x_C) > 0$

(C)  $\Leftrightarrow p(x_A \mid x_B, x_C) = p(x_A \mid x_C)$

(\*) "marginal indep.":  $X_A \perp\!\!\!\perp X_B \mid \emptyset$   
 $X_B$

$\forall x_B, x_C \text{ st. } p(x_B, x_C) > 0$

2 facts about C.I.:

1) can repeat variables in statement (for convenience)  $X \perp\!\!\!\perp Y, Z \mid Z, W$  is fine to say  
↑  
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

$$Z = X \oplus Y$$

(\*) chain rule (always true)

$$p(x_V) = \prod_{i=1}^n p(x_i \mid x_{1:i-1})$$

assumption in DCM

last cond.  $p(x_n \mid x_{1:n-1})$  table with  $2^n$  entries

assumption in DGM

$$= \prod_{i=1}^n p(x_i | x_{\pi_i}) \rightarrow \text{tables of } 2^{\max_{i: |x_{\pi_i}|+1} \text{ entries}}$$

## Directed graphical model

let  $G = (V, E)$  be a DAG

a directed graphical model (DGM) (associated with  $G$ ) (aka Bayesian network)

is a family of distributions over  $X_V$

here  $\mathcal{L}(G) = \{ p \text{ is a dist. over } X_V : \exists \text{ legal factors } f_i \text{ s.t. } \}$  (not nec. unique)

$$p(x_V) = \prod_{i=1}^n f_i(x_i | x_{\pi_i})$$

(potential sets)

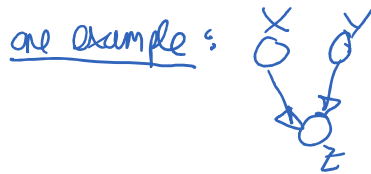
$$f_i: \Omega_{x_i} \times \Omega_{x_{\pi_i}} \rightarrow [0, 1]$$

$$\text{s.t. } \forall i \sum_{x_i} f_i(x_i | x_{\pi_i}) = 1 \quad \forall x_{\pi_i}$$

$f_i$  is like a CPT (conditional prob table)

terminology: if we can write  $p(x_V) = \prod_i f_i(x_i | x_{\pi_i})$  this determines  $G$

then we say that " $p$  factorizes according to  $G$ "



$$p \in \mathcal{L}(G) \Leftrightarrow p(x, y, z) = f_x(x) f_y(y) f_z(z | x, y)$$

very soon, we show that  $p(x_i | x_{\pi_i}) = f_i(x_i | x_{\pi_i})$