

# Lecture 10 - scribbles

Friday, October 6, 2017

13:21

today: finish EM & GMM  
 ✓ checked graphical model

EM continuation:

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

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

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

$$= KL(q(\cdot) \parallel p(\cdot|x; \theta)) \quad (KL\text{-divergence})$$

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

Operationally:

$$E \text{ step: } q_{t+1} = \underset{\substack{q \in \text{all dist.} \\ \text{over } z}}{\text{argmax}} J(q, \theta_t) = \underset{q}{\text{argmin}} KL(q \parallel p(\cdot|x; \theta_t))$$

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

[inference]

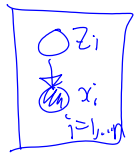
$$M \text{ step} \quad \theta_{t+1} = \underset{\theta \in \Theta}{\text{argmax}} J(q_{t+1}, \theta)$$

$$\mathbb{E}_{q_{t+1}} [\log p(x, z; \Theta)] + \text{cst.}$$

"expected complete log-likelihood"

need to solve another ML problem, but for complete information (i.e.  $z$  is known)

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

complete log-likelihood

$$\begin{aligned} \log p(x, z; \Theta) &= \sum_{i=1}^n [\log p(x_i | z_i; \Theta) + \log p(z_i; \Theta)] \\ &\stackrel{\text{Gaussian log-likelihood}}{\approx} \sum_{i=1}^n \left[ \sum_{j=1}^K z_{ij} \log N(x_i | \mu_j, \Sigma_j) \right] \stackrel{\text{multinomial}}{+} \sum_{j=1}^K z_{ij} \log \pi_j \end{aligned}$$

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

$\downarrow$   
 $q(z_{ij}=1)$  (marginal distribution)

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

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

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

marginal

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

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

$$= \frac{\pi_j^{(t)} N(x_i | \mu_j^{(t)}, \Sigma_j^{(t)})}{\sum_{l=1}^k \pi_l^{(t)} N(x_i | \mu_l^{(t)}, \Sigma_l^{(t)})} \} p(x_i, z_{i,j}=1 | \theta_t)$$

$$\} p(x_i | \theta_t)$$

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

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

exercise:  $N_j^{(t+1)} \approx \sum_i \gamma_{i,j}^t$  "soft-count"

$$\mu_j^{(t+1)} = \frac{\sum_{i=1}^n \gamma_{i,j}^t x_i}{\left( \sum_{i=1}^n \gamma_{i,j}^t \right)}$$

soft cluster assignment

[weighted empirical mean]

$$\Sigma_j^{(t+1)} = \frac{\sum_{i=1}^n \gamma_{i,j}^t (x_i - \mu_j^{(t+1)})(x_i - \mu_j^{(t+1)})^T}{\sum_{i=1}^n \gamma_{i,j}^t}$$

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

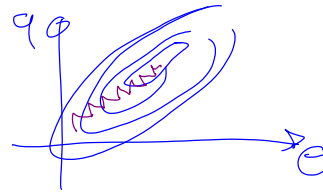
$\Sigma_j^{(0)}$  high spherical covariance

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

$\pi_i^{(6)}$  proportions from K-means++

→ one problem with EM:

Sometimes make very slow progress



→ where 'conjugate gradient' method can use

Aside 9

- with  $\mathbf{b}$ -data, true MLE gives the correct parameters  
but it's intractable to find them (non-convex)
- instead, use method of moments to recover parameters "cheaply"  
with guarantees  $\Rightarrow$

see. [Hsu & Kakade TCS 2013]

<https://arxiv.org/abs/1206.5766>



## Graphical mode

graph model  $\mapsto$  probability + C.S.

## R.V. and graphs

GMR

arrang  $\rightarrow$  efficient data structure

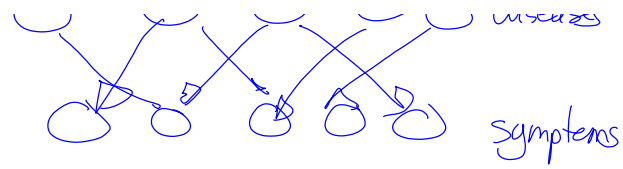
○ ○ ○ ○ ○ disease

graph: system with state -

eg.  $X_1, \dots, X_n$  R.V.'s

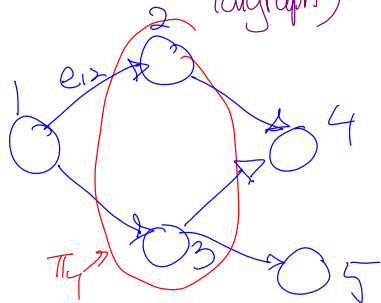
$X_i \in \{0, 1\}$   $n \rightarrow 100$

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



## Graph theory

directed graph  $G = (V, E)$   
(digraph)



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

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

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

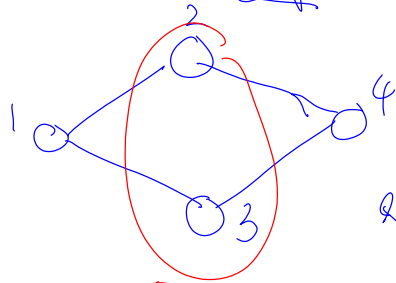
$e_{12} = (1, 2)$

directed path  $1 \rightarrow 4$

$(1, 2), (2, 4)$

undirected graph: elements of  $E$  are 2-sets

(note: no self loops)



$\nearrow$   
"neighbors" of node 1

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

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

$\nwarrow$  here we have (undirected) path from 2 to 3

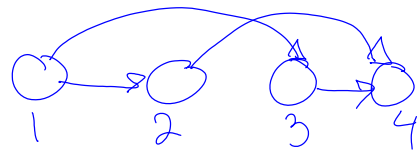
neighbors replace the parents/children terminology from digraphs

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

def: an ordering  $I: V \rightarrow \{1, \dots, n\}$  is said topological for  $G$

iff nodes in  $\pi_i$  appear before  $i$  in  $I \quad \forall i$

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



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

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

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

$\Rightarrow$ ) use DFS algorithm to construct a topological sorting in  $O(|E| + |V|)$

general issues in this class

A) representation  $\begin{cases} \rightarrow \text{DGM} \\ \rightarrow \text{UGM} \end{cases}$

parameterization  $\rightarrow$  exponential family

} probabilities

B) inference computing  $p(\mathcal{Z}_Q | \mathcal{Z}_F)$   
 "query" "evidence"

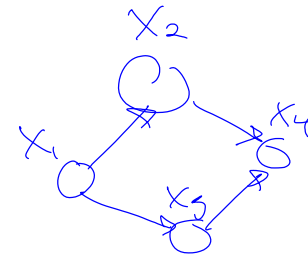
→ sum-product alg.  
junction

C) statistical estimation → MLE  
maximum entropy  
method of moments

Notation:

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

$V$  a set of vertices  
one R.V. per node



joint  $p(X_1=x_1, X_2=x_2, \dots, X_n=x_n) \stackrel{\text{shortcut}}{=} p(x_1, \dots, x_n)$   
 $= p(x_V) \stackrel{(\text{v})}{=} p(x)$

recall.  
 $p(x_{1:n}) = p(x)$

for any  $A \subseteq V$  marginal on  $x_A$

$p(x_A) = P\{X_i=x_i : i \in A\} = \sum_{x_{A^c}} p(x_A, x_{A^c})$   
 subset of "subscripts"  $x_{\{1,2,4\}}$  or  $\{x_1, x_2, x_4\}$   $A^c = V \setminus A$   $\{x_i : i \in V \setminus A\}$   
 summing over all possible values of

revisit conditional independence:

let  $A, B, C \subseteq V$

$X_A \perp\!\!\!\perp X_B \mid X_C$

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

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

"marginal independence" :  $X_A \perp\!\!\!\perp X_B \mid \emptyset$

3 facts about C.I.:

1) can repeat variables e.g.  $X \perp\!\!\!\perp Y, Z \mid Z, W$  is fine to say

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

3) trick: extra conditioning on both sides of equation maintain true statement

$$\text{e.g. } p(x, y) = p(x|y) p(y) \quad (\text{always true})$$

$$\Rightarrow p(x, y|z) = p(x|y, z) p(y|z) \quad \parallel \parallel$$

(\*) pairwise independence  $\not\Rightarrow$  mutual independence

e.g.  $X_3 = X_1 \text{ xor } X_2$  with  $X_1, X_2 \sim \text{Bernoulli}(\frac{1}{2})$

here  $X_1 \perp\!\!\!\perp X_2$ ,  $X_2 \perp\!\!\!\perp X_3$ ,  $X_1 \perp\!\!\!\perp X_3$

but  $X_1 \not\perp\!\!\!\perp (X_2, X_3)$

chain rule:  
(always true)

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

last conditional is  $p(x_n | x_{1:n-1})$

label with  $x_{1:n}$  entries



directed graph  
model

$$p(x_v) = \prod_{i=1} p(x_i | \pi_i)$$

parents of  $i$  in graph  $G$

→ tables of  $2^{\max_i |\pi_i| + 1}$