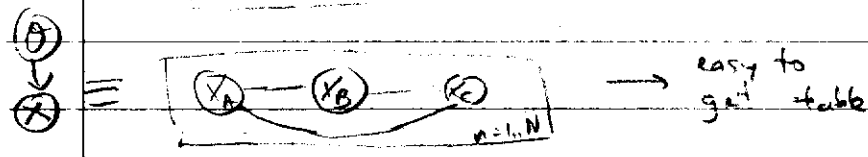
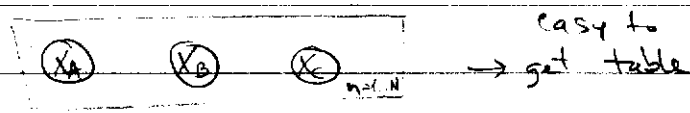


- Assignment 2 due tomorrow
- Assignment 3 on web tomorrow
- Project Info (2 people, April 1, May 6, on web tomorrow)
- Elimination Algo's, compute marg, conds, but only query node & lat time
- Learn tables, Junction tree: compute marg/conds but over several nodes & global.

ML for Graphical Models (fully observed)

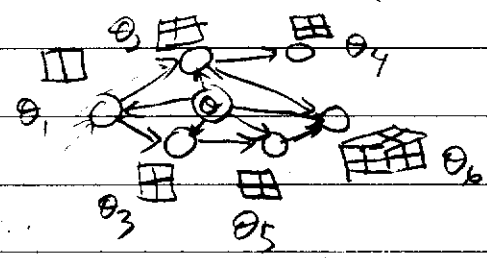
$G = (U, E)$   
 nodes vertices    edges



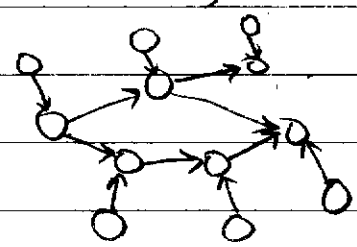
for directed graphical models → ML decouples

for undirected → decomposable (normalizer Z is hard to compute)

Directed  $p(x) = \prod_{u=1}^M p(x_u | x_{\pi_u})$  graph, no tables yet  
 $p(x|\theta) = \prod_{u=1}^M p(x_u | x_{\pi_u}, \theta_u)$   $\theta$  specifies tables



c.p.t.'s  
 can be varied indep.  
 So  $\theta$  breaks down



$p(x|\theta) = \prod_{u=1}^M p(x_u | x_{\pi_u}, \theta_u)$

IID have N observations of x (which is a set of M vars)

i.e. have N patients come in & observe

fly, headache, sinus fever, etc... For each

get N graphs:  $G^{(N)} = (U^{(N)}, E^{(N)})$

Replicates graph



to index nodes in super graph:  $X_{u,n}$

$u \in \mathcal{U}$ , the set of nodes

$n \in \{1, \dots, N\}$

$C$  is a set of vars, i.e.  $X_C$  for  $C \subseteq \mathcal{U}$

$X_{C,n}$  is the  $n$ th copy of the set of vars  $X_C$

$$X_{\mathcal{U}} = \{X_{u_1}, \dots, X_{u_N}\}$$

our dataset  $D = \{X_{u,1}, \dots, X_{u,N}\}$

or curly  $X$

Recall ML  $\theta^* = \underset{\theta}{\operatorname{argmax}} \log p(D|\theta)$

$$= \log \prod_n p(X_{\mathcal{U},n}|\theta)$$

$$= \log \prod_n \prod_u p(X_{u,n} | X_{\pi_u,n}, \theta_u) \quad \text{log likelihood}$$

$$= \sum_n \sum_u \log p(X_{u,n} | X_{\pi_u,n}, \theta_u) \equiv \ell(\theta; D)$$

$\ell(\theta)$  decouples, i.e.  $\frac{\partial}{\partial \theta_k} \ell(\theta; D) = \sum_n \frac{\partial}{\partial \theta_k} \log p(X_{u,n} | X_{\pi_u,n}, \theta_u)$

$\uparrow$  only consider single  $\theta_u$  at a time!

i.e. estimate each cpdf at a time.

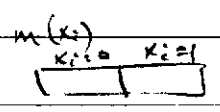
i.e. only look at a few vars at a time  $p(x_i | \pi_i)$

graphical models  $\Rightarrow$  memory efficiency

$\Rightarrow$  estimation efficiency

Estimating CPDFs ... first some notation...

$$\delta(x_i, x_j) = \begin{cases} 1 & \text{if } x_i = x_j \\ 0 & \text{otherwise} \end{cases}$$

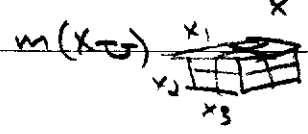


$$m(x_i) = \sum_n \delta(x_i, x_{i,n})$$

$\leftarrow$  # of data points where  $x_i$  assumes given val

$$m(X_{\mathcal{U}}) = \sum_n \delta(X_{\mathcal{U}}, X_{\mathcal{U},n})$$

$\leftarrow$  # of data points where  $X$  assumes value.



Sum over all  $X$  in  $\mathcal{U}$  except  $X$  in  $C$

$$m(X_C) = \sum_{x_U \in \mathcal{U}} m(x_U)$$

← # of times on a subset of vars.

i.e.  $m(x_1, x_2) = \sum_{x_3} m(x_1, x_2, x_3)$

$$m(x_i) = \sum_{x_j} \sum_{x_k} m(x_i, x_j, x_k)$$

$$\sum_{x_i} m(x_i) = N \quad (\text{total \# of data})$$

define:  $X_{\phi_i} := \{X_i, X_{\pi_i}\}$  node  $i$  & its parents (pdf)

$$\therefore m(X_{\phi_u}) = \sum_{x_U \in \mathcal{U}} m(x_U)$$

forming a table of counts over  $X_u$  & its parents.

our cond. prob. table

$\theta_u$  is a parameter vector. If  $\theta_u$  is over  $X_u$  by itself, it is a 1d table 

$x_u=0$	$x_u=1$
0.3	0.7

 if  $X_u | X_{\pi_u}$ 

$x_u=0$	$x_u=1$
$x_{\pi_u}=0$	
$x_{\pi_u}=1$	

$\therefore \theta_u$  is a table over  $X_{\phi_u}$   $\phi_u = \{u, \pi_u\}$

i.e.  $\theta_u(X_{\phi_u})$  like we had for coin flips

$$\theta(x=0)=0.3 \quad \theta(x=1)=0.7 \quad (\text{we called it } \theta_1 \text{ \& } \theta_2 \dots)$$

naturally  $\sum_{x_v} \theta_v(X_{\phi_v}) = \sum_{x_v} \theta_v(x_v, X_{\pi_v}) = 1$

i.e.  $p(x_v | X_{\pi_v}, \theta_v) = \theta_v(X_{\phi_v})$

$$\therefore p(x_{\mathcal{U}} | \theta) = \prod p(x_v | X_{\pi_v}, \theta_v) = \prod \theta_v(X_{\phi_v})$$

$$p(x_{\mathcal{U}}, n | \theta) = \prod_{x_{\mathcal{U}}} p(x_{\mathcal{U}} | \theta)^{S(x_{\mathcal{U}}, x_{\mathcal{U}}, n)}$$

over all possible  $\mathcal{U}$  all settings of the nodes  $X_1, \dots, X_M$

log-likelihood  $\log p(D | \theta) = \log \prod p(x_{\mathcal{U}}, n | \theta)$

$$= \sum_n \log p(x_{\mathcal{U}}, n | \theta)$$

$$= \sum_n \log \prod_{x_{\mathcal{U}}} p(x_{\mathcal{U}} | \theta)^{S(x_{\mathcal{U}}, x_{\mathcal{U}}, n)}$$

$$\begin{aligned}
 &= \sum_{\tau} \sum_{\phi} \delta(x_{\tau}, x_{\phi}, n) \log p(x_{\tau} | \theta) \\
 &= \sum_{\tau} m(x_{\tau}) \log p(x_{\tau} | \theta) \\
 &= \sum_{\tau} \sum_{\phi} m(x_{\tau}) \log \prod_{\phi} p(x_{\phi} | x_{\tau}, \theta) \\
 &= \sum_{\tau} m(x_{\tau}) \log \prod_{\phi} \Theta_{\phi}(x_{\phi}) \\
 &= \sum_{\tau} m(x_{\tau}) \sum_{\phi} \log \Theta_{\phi}(x_{\phi}) \\
 &= \sum_{\tau} \sum_{\phi} m(x_{\tau}) \log \Theta_{\phi}(x_{\phi}) \\
 &= \sum_{\tau} \sum_{\phi} \sum_{x_{\phi}} m(x_{\tau}) \log \Theta_{\phi}(x_{\phi}) \\
 &= \sum_{\tau} \sum_{\phi} \sum_{x_{\phi}} m(x_{\tau}, x_{\phi}) \log \Theta_{\phi}(x_{\phi})
 \end{aligned}$$

sufficient stats e-family property

(for each  $\Theta_{\nu}$ , it decouples, once again)

add Lagrangians  $l(\theta | D) = \sum_{\nu} \lambda_{\nu} (\sum_{x_{\nu}} \Theta_{\nu}(x_{\nu}) - 1)$

since  $\sum_{\nu} \Theta_{\nu}(x_{\nu}, x_{\tau_{\nu}}) = 1$

$2\Theta_{\nu}(x_{\nu}, x_{\tau_{\nu}})$

$\frac{m(x_{\nu}, x_{\tau_{\nu}})}{\Theta_{\nu}(x_{\nu}, x_{\tau_{\nu}})} - \lambda_{\tau_{\nu}} = 0$

$\Theta_{\nu}(x_{\nu}, x_{\tau_{\nu}}) = \frac{m(x_{\nu}, x_{\tau_{\nu}})}{\lambda_{\tau_{\nu}}}$

$\sum_{\nu} \Theta_{\nu}(x_{\nu}, x_{\tau_{\nu}}) = 1$

implicitly  $\Rightarrow \lambda_{\tau_{\nu}} = \sum_{\nu} m(x_{\nu}, x_{\tau_{\nu}}) = m(x_{\tau_{\nu}})$

$\hat{\Theta}_{\nu}(x_{\nu}, x_{\tau_{\nu}}) = \frac{m(x_{\nu}, x_{\tau_{\nu}})}{m(x_{\tau_{\nu}})}$

Simple algorithmically

example

$p(\text{fever}   \tau_{\nu})$	$x_2   x_1$																
$m(x_{\nu}, x_{\tau_{\nu}})$	<table border="1"> <tr><td>0</td><td>1</td><td>1</td><td>1</td></tr> <tr><td>1</td><td>0</td><td>1</td><td>1</td></tr> <tr><td>1</td><td>1</td><td>0</td><td>1</td></tr> <tr><td>1</td><td>1</td><td>1</td><td>0</td></tr> </table>	0	1	1	1	1	0	1	1	1	1	0	1	1	1	1	0
0	1	1	1														
1	0	1	1														
1	1	0	1														
1	1	1	0														
$m(x_{\tau_{\nu}})$	<table border="1"> <tr><td>1</td><td>1</td><td>1</td><td>1</td></tr> <tr><td>1</td><td>1</td><td>1</td><td>1</td></tr> <tr><td>1</td><td>1</td><td>1</td><td>1</td></tr> <tr><td>1</td><td>1</td><td>1</td><td>1</td></tr> </table>	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
1	1	1	1														
1	1	1	1														
1	1	1	1														
1	1	1	1														

Undirected Models factorize without being conditionals, but  $\frac{1}{2}$  <sup>undecoupled</sup> couples them.

$D = \{x_{\tau_1}, \dots, x_{\tau_n}\}$

$p(x_{\tau} | \theta) = \frac{1}{2} \prod_{\phi} \psi_{\phi}(x_{\phi})$        $z = \sum_{x_{\tau}} \prod_{\phi} \psi_{\phi}(x_{\phi})$

(recall previous  $m()$  notation...)

$p(x_{\tau}, n | \theta) = \prod_{\tau} p(x_{\tau} | \theta) \delta(x_{\tau}, x_{\tau}, n)$

$$\begin{aligned}
 l(\theta|D) &= \log p(D|\theta) = \log \prod_n p(x_{T_n}|\theta) \\
 &= \log \prod_n \prod_{k \in T_n} p(x_{T_n}|\theta) \delta(x_{T_n}, x_{T_n}) \\
 &= \sum_n \sum_{k \in T_n} \delta(x_{T_n}, x_{T_n}) \log p(x_{T_n}|\theta) \\
 &= \sum_{k \in T_n} m(x_{T_n}) \log p(x_{T_n}|\theta) \\
 &= \sum_{k \in T_n} m(x_{T_n}) \sum_c \log \psi_c(x_c) - \sum_{k \in T_n} m(x_{T_n}) \log 2 \\
 &= \sum_{k \in T_n} \sum_c m(x_c) \log \psi_c(x_c) - N \log 2
 \end{aligned}$$

use max likelihood to set  $\theta$  or the tables  $\psi_c(x_c)$ ...

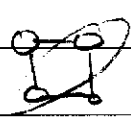
two cases: decomposable & non-decomposable.

Decomposable

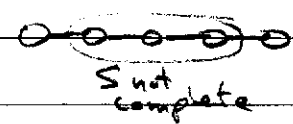
- divide graph into disjoint subsets A, B, S
- recursively, must have below for graph & sub-graphs

decomp. must recursive on subgraphs

if S separates A & B, S must be complete



S not complete

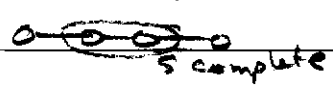


S not complete

all nodes fully intercon.

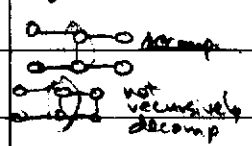


S complete



S complete

clique



not recursively decomp

if decomposable:

$$1) \psi_c(x_c) \leftarrow \tilde{p}(x_c) = \frac{1}{N} m(x_c)$$

$$2) \psi_c(x_c) \leftarrow \frac{\psi_c(x_c)}{P(X_c \cap X_d)}$$

$$\psi_c(x_c) \leftarrow \psi_d(x_d)$$

for all non-empty intersections

between pairs divide 1 of the intersecting cliques  
 (combo by the empirical marginal of intersection)

road to divide by middle for 3 nodes  
 $p(x_1, x_2)$   
 $p(x_2, x_3)$   
 $p(x_1, x_2, x_3)$

do example	count
$x_1, x_2, x_3$	110
	001
	111
	101
	000
	010
	001

more general

Non-decomposable (or otherwise)  $\Rightarrow$  I.P.F. algorithm

if decomposable, converges in finite steps

Iterative Prop. Fitting solution, properties, proof

IPF: { guess values for  $\psi_c(x_c)$   
 lock all  $\psi_c(x_c)$   
 for  $c=1, \dots, C$   $\psi_c(x_c)^{(t+1)} \leftarrow \psi_c(x_c)^{(t)} \frac{\hat{p}(x_c)}{p^{(t)}(x_c)}$   
 where  $\hat{p}(x_c) = \frac{1}{z} m(x_c)$   
 $p^{(t)}(x_c)$  is computed from current  
 $p(x) = \frac{1}{z} \prod_c \psi_c(x_c)^{(t)}$   
 $p(x_c) = \sum_{x \sim x_c} p(x)$

for decomposable graphs, IPF converges in finite iters

Note:  $z = \sum_x \prod_c \psi_c(x_c)^{(t)}$  so  $z$  can change with  $t$

No PE

$$\begin{aligned}
 p^{(t+1)}(x_c) &= \sum_{x \sim x_c} p^{(t+1)}(x) \\
 &= \sum_{x \sim x_c} \frac{1}{z^{(t+1)}} \prod_D \psi_D^{(t+1)}(x_D) \\
 &= \frac{1}{z^{(t+1)}} \sum_{x \sim x_c} \psi_c^{(t+1)}(x_c) \prod_{D \neq c} \psi_D^{(t+1)}(x_D) \\
 &= \frac{1}{z^{(t+1)}} \sum_{x \sim x_c} \psi_c^{(t)}(x_c) \frac{\hat{p}(x_c)}{p^{(t)}(x_c)} \prod_{D \neq c} \psi_D^{(t)}(x_D) \\
 &= \frac{z^{(t)}}{z^{(t+1)}} \frac{\hat{p}(x_c)}{p^{(t)}(x_c)} \sum_{x \sim x_c} \frac{1}{z^{(t)}} \prod_D \psi_D^{(t)}(x_D) \\
 &= \frac{z^{(t)}}{z^{(t+1)}} \frac{\hat{p}(x_c)}{p^{(t)}(x_c)} \sum_{x \sim x_c} p^{(t)}(x) \\
 &= \frac{z^{(t)}}{z^{(t+1)}} \frac{\hat{p}(x_c)}{p^{(t)}(x_c)} p^{(t)}(x_c)
 \end{aligned}$$

only  $\psi_c$  changed from  $t$  to  $t+1$

$$\begin{aligned}
 p^{(t+1)}(x_c) &= \frac{z^{(t)}}{z^{(t+1)}} \hat{p}(x_c) \\
 \sum_{x_c} p^{(t+1)}(x_c) &= \sum_{x_c} \frac{z^{(t)}}{z^{(t+1)}} \hat{p}(x_c) \\
 1 &= \frac{z^{(t)}}{z^{(t+1)}} \\
 \therefore z^{(t+1)} &= z^{(t)} = \text{constant with IPF iters.}
 \end{aligned}$$

sum both sides over  $x_c$