

## Machine Learning

### Structuring PDFs for Inference

- Inference: goal is to predict some variables given others  
 $x_1$ : flu  
 $x_2$ : fever  
 $x_3$ : sinus infection  
 $x_4$ : temperature  
 $x_5$ : sinus swelling  
 $x_6$ : headache
- Patient claims headache and high temperature. Does he have a flu?

Given findings variables  $X_f$  and unknown variables  $X_u$   
predict queried variables  $X_q$

- Classical approach: truth tables (slow) or logic networks
- Modern approach: probability tables (slow) or Bayesian networks (fast belief propagation, junction tree algorithm)

### Graphical Models & Bayes Nets

- Node: a random variable (discrete or continuous)
- Independent: no link  $p(x, y) = p(x)p(y)$
- Dependent: link  $p(x, y) = p(y | x)p(x)$
- Arrow: from parent to child (like causality, not exactly)
- Child: destination of arrow, response
- Parent: root of arrow, trigger  $\text{parents of child } i = pa_i = \pi_i$
- Graph: dependence/independence
- Graph: shows factorization of joint  
joint = products of conditionals  
 $p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i | pa_i) = \prod_{i=1}^n p(x_i | \pi_i)$
- DAG: directed acyclic graph

## Topic 14

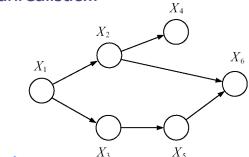
- Structuring Probability Functions for Storage
- Structuring Probability Functions for Inference
- Basic Graphical Models
- Graphical Models
- Parameters as Nodes

## Structuring PDFs for Storage

- Probability tables quickly grow if  $p$  has many variables  
 $p(x) = p(\text{flu?}, \text{headache?}, \dots, \text{temperature?})$
- For D true/false "medical" variables  $\text{table size} = 2^D$
- Exponential blow-up of storage size for the probability
- Example: 8x8 binary images of digits
- If multinomial with M choices, probabilities are how big?  
 $p(x) = p(\text{flu?})p(\text{headache?})\dots p(\text{temperature?})$
- As in Naïve Bayes or Multivariate Bernoulli, if words were independent things are much more efficient  
 $p(x) = p(\text{flu?})p(\text{headache?})\dots p(\text{temperature?})$
- For D true/false "medical" variables  $\text{table size} = 2 \times D$  (really even less than that...)

### From Logic Nets to Bayes Nets

- 1980's expert systems & logic networks became popular
- Problem: inconsistency, 2 paths can give different answers
- Problem: rules are hard, instead use soft probability tables
- $x_3 = x_1 \wedge x_2$   $p(x_3 | x_1, x_2)$
- $x_3 = 0$   $x_3 = 1$
- $x_3 = 0$   $x_3 = 1$
- These directed graphs are called Bayesian Networks



## Graphical Models & Bayes Nets

- Independence assumptions make probability tables smaller
- But real events in the world not completely independent!
- Complete independence is unrealistic...
- Graphical models use a graph to describe more subtle dependencies and independencies:  
...namely: conditional independencies (like causality but not exactly...)
- Directed Graphical Model, also called Bayesian Network use a directed acyclic graph (DAG).
- Neural Network = Graphical Function Representation
- Bayesian Network = Graphical Probability Representation

### Basic Graphical Models

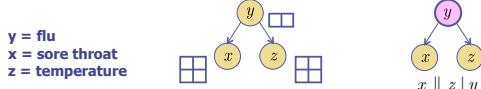
- Independence: all nodes are unlinked
- Shading: variable is 'observed', condition on it moves to the right of the bar in the pdf
- Examples of simplest conditional independence situations...  
 $p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i | pa_i) = \prod_{i=1}^n p(x_i | \pi_i)$
- Markov chain:  $p(x, y, z) = p(x)p(y | x)p(z | y)$
- Example binary events:  
 $x = \text{president says war}$   
 $y = \text{general orders attack}$   
 $z = \text{soldier shoots gun}$
- $x \perp\!\!\!\perp z | y$
- $p(x, y, z) = p(x)p(z)p(y | x, z)$
- Explaining away...
- $p(x, y, z) = p(x)p(z)p(y | x, z)$
- $p(x | y, z) = \frac{p(x, y, z)}{p(y, z)} = p(x | y)$

### Basic Graphical Models

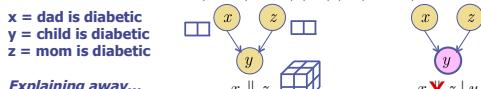
- 1 Cause, 2 effects:  $p(x, y, z) = p(y)p(x | y)p(z | y)$
- $y = \text{flu}$   
 $x = \text{sore throat}$   
 $z = \text{temperature}$
- $x \perp\!\!\!\perp z | y$
- 2 Causes, 1 effect:  $p(x, y, z) = p(x)p(z)p(y | x, z)$
- $x = \text{rain}$   
 $y = \text{wet driveway}$   
 $z = \text{car oil leak}$
- $x \perp\!\!\!\perp z | y$
- Explaining away...
- Each conditional is a mini-table (Multinomial or Bernoulli conditioned on parents)

## Basic Graphical Models

2) 1 Cause, 2 effects:  $p(x, y, z) = p(y)p(x|y)p(z|y)$



3) 2 Causes, 1 effect:  $p(x, y, z) = p(x)p(z)p(y|x, z)$



Explaining away...

- Each conditional is a mini-table  
(Multinomial or Bernoulli conditioned on parents)

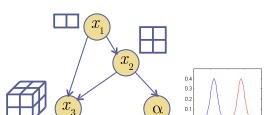
## Graphical Models

- Normalizing probability tables. Joint distributions sum to 1.
- BUT, conditionals sum to 1 for each setting of parents.

$p(x)$	<b>2-1</b> $\sum_{x=0}^1 p(x) = 1$	
$p(x, y)$	<b>4-1</b> $\sum_{x,y} p(x, y) = 1$	
$p(x, y, z)$	<b>8-1</b> $\sum_{x,y,z} p(x, y, z) = 1$ $y=0 \quad x=0 \quad z=0$ $x=0 \quad x=1$	
	<b>8-2</b> $\sum_x p(x y) = 1$ $\sum_x p(x y=0) = 1$ $\sum_x p(x y=1) = 1$	
	<b>8-3</b> $\sum_y p(x y, z) = 1$ $\sum_y p(x y=0, z=0) = 1$ $\sum_y p(x y=0, z=1) = 1$	
	<b>8-4</b> $\sum_z p(x y, z) = 1$ $\sum_z p(x y=0, z=0) = 1$ $\sum_z p(x y=0, z=1) = 1$	

## Continuous Conditional Models

- In previous slide,  $\theta$  and  $\alpha$  were a random variable in graph
- But,  $\theta$  and  $\alpha$  are continuous
- Network can have both discrete & continuous nodes
- Joint factorizes into conditionals that are either:
  - 1) discrete conditional probability tables
  - 2) continuous conditional probability distributions



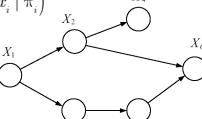
- Most popular continuous distribution = Gaussian

## Graphical Models

- Example: factorization of the following system of variables

$$\begin{aligned} p(x_1, \dots, x_n) &= \prod_{i=1}^n p(x_i | pa_i) = \prod_{i=1}^n p(x_i | \pi_i) \\ p(x_1, \dots, x_6) &= p(x_1) \dots \\ &= p(x_1)p(x_2 | x_1) \dots \\ &= p(x_1)p(x_2 | x_1)p(x_3 | x_1)p(x_4 | x_2) \dots \\ &= p(x_1)p(x_2 | x_1)p(x_3 | x_1)p(x_4 | x_2)p(x_5 | x_3) \dots \\ &= p(x_1)p(x_2 | x_1)p(x_3 | x_1)p(x_4 | x_2)p(x_5 | x_3)p(x_6 | x_2, x_5) \end{aligned}$$

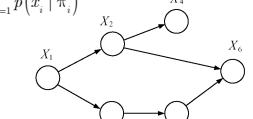
- How big are these tables (if binary variables)?



## Graphical Models

- Example: factorization of the following system of variables

$$\begin{aligned} p(x_1, \dots, x_n) &= \prod_{i=1}^n p(x_i | pa_i) = \prod_{i=1}^n p(x_i | \pi_i) \\ p(x_1, \dots, x_6) &= p(x_1)p(x_2 | x_1)p(x_3 | x_1)p(x_4 | x_2)p(x_5 | x_3)p(x_6 | x_2, x_5) \\ 2^6 &\quad 2^1 \quad 2^2 \quad 2^2 \quad 2^2 \quad 2^3 \end{aligned}$$

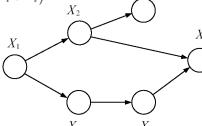


## Graphical Models

- Example: factorization of the following system of variables

$$\begin{aligned} p(x_1, \dots, x_n) &= \prod_{i=1}^n p(x_i | pa_i) = \prod_{i=1}^n p(x_i | \pi_i) \\ \bullet \text{Interpretation} \\ 1: \text{flu} & \\ 2: \text{fever} & \\ 3: \text{sinus infection} & \\ 4: \text{temperature} & \\ 5: \text{sinus swelling} & \\ 6: \text{headache} & \\ p(x_1, \dots, x_6) &= p(x_1)p(x_2 | x_1)p(x_3 | x_1)p(x_4 | x_2)p(x_5 | x_3)p(x_6 | x_2, x_5) \\ 2^6 - 1 &\quad 2^1 - 1 \quad 2^2 - 2 \quad 2^2 - 2 \quad 2^2 - 2 \quad 2^3 - 4 \end{aligned}$$

- 63 vs. 13 degrees of freedom



## Parameters as Nodes

- Consider the model variable  $\theta$  ALSO as a random variable



- But would need a prior distribution  $P(\theta)$ ... ignore for now

- Recall: Naïve Bayes, word probabilities are independent

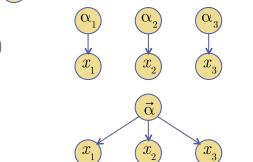


- Text: Multivariate Bernoulli

$$p(x | \vec{\alpha}) = \prod_{d=1}^{50000} \alpha_d^{x_d} (1 - \alpha_d)^{1-x_d}$$

- Text: Multinomial

$$p(X | \vec{\alpha}) = \frac{\sum_{m=1}^M X_m!}{\prod_{m=1}^M X_m!} \prod_{m=1}^M \alpha_m^{X_m}$$



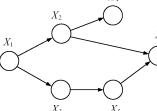
## Graphical Models

- In EM, we saw how to handle nodes that are: observed (shaded), hidden variables (E), parameters (M)

- But, only considered simple iid, single parent, structures
- More generally, have arbitrary DAG without loops

• Notation:

$$\begin{aligned} G &= \{X, E\} = \{\text{nodes / randomvars, edges}\} \\ X &= \{x_1, \dots, x_M\} \\ E &= \{(x_i, x_j) : i \neq j\} \\ X_c &= \{x_1, x_3, x_4\} = \text{subset} \end{aligned}$$



- Want to do 4 things with these graphical models:

- 1) Learn Parameters (to fit to data)
- 2) Query independence/dependence
- 3) Perform Inference (get marginals/max a posteriori)
- 4) Compute Likelihood (e.g. for classification)

## Graphical Models

- Graph factorizes probability:  $p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i | \pi_i)$

- Topological graph:

nodes are in order so that parents  $\pi$  come before children

$$\begin{aligned} p(x_1, \dots, x_6) &= p(x_1)p(x_2 | x_1) \\ &\times p(x_3 | x_1)p(x_4 | x_2) \\ &\times p(x_5 | x_3)p(x_6 | x_2, x_5) \end{aligned}$$

- Question? Which is the more general graph?



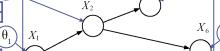
- Conditional probability tables can be chosen to make 'busier' graph look like simpler graph

# Machine Learning

- ## Topic 15

- Graphical Models
  - Maximum Likelihood for Graphical Models
  - Testing for Conditional Independence & D-Separation
  - Bayes Ball

## Maximum Likelihood CPTs

- Each conditional probability table  $\theta_i$  part of our parameters
  - Given table, have pdf  
 $p(X_U | \theta) = \prod_{i=1}^M p(x_i | \pi_i, \theta_i)$
  - Have M variables:  
 $X = \{x_1, \dots, x_M\}$
  - Have N x M dataset:  
 $\mathcal{D} = \{X_{U,1}, \dots, X_{U,N}\}$
  - Maximum likelihood:  
 $\theta^* = \arg \max_{\theta} \log p(\mathcal{D} | \theta)$   
 $= \arg \max_{\theta} \sum_{n=1}^N \log p(X_{U,n} | \theta)$   
 $= \arg \max_{\theta} \sum_{n=1}^N \log \prod_{i=1}^M p(x_{i,n} | \pi_{i,n}, \theta_i)$   
 $= \arg \max_{\theta} \sum_{n=1}^N \sum_{i=1}^M \log p(x_{i,n} | \pi_{i,n}, \theta_i)$

each  $\theta_i$  appears independently, can do ML for each CPT alone! efficient storage & efficient learning

## Maximum Likelihood CPTs

- Continuing:  $\ell(\theta) = \sum_{X_U} \sum_{i=1}^M m(X_U) \log p(x_i | \pi_i, \theta_i)$   
 $= \sum_{i=1}^M \sum_{x_i, \pi_i} m(X_U) \log p(x_i | \pi_i, \theta_i)$   
 $= \sum_{i=1}^M \sum_{x_i, \pi_i} m(x_i, \pi_i) \log p(x_i | \pi_i, \theta_i)$
  - Define:  $\theta(x_i, \pi_i) = p(x_i | \pi_i, \theta_i)$  Constraint:  $\sum_{x_i} \theta(x_i, \pi_i) = 1$
  - Now have above with Lagrange multipliers:  
 $\ell(\theta) = \sum_{i=1}^M \sum_{x_i} \sum_{\pi_i} m(x_i, \pi_i) \log \theta(x_i, \pi_i) - \sum_{i=1}^M \sum_{\pi_i} \lambda_{\pi_i} (\sum_{x_i} \theta(x_i, \pi_i) - 1)$   
 $\frac{\partial \ell(\theta)}{\partial \theta(x_i, \pi_i)} = \frac{m(x_i, \pi_i)}{\theta(x_i, \pi_i)} - \lambda_{\pi_i} = 0 \rightarrow \theta(x_i, \pi_i) = \frac{m(x_i, \pi_i)}{\lambda_{\pi_i}}$
  - Plug constraint:  $\sum_{x_i} \frac{m(x_i, \pi_i)}{\lambda_{\pi_i}} = 1 \rightarrow \lambda_{\pi_i} = \sum_{x_i} m(x_i, \pi_i) = m(\pi_i)$
  - Final solution (trivial!):  $\theta(x_i, \pi_i) = \frac{m(x_i, \pi_i) + \varepsilon}{m(\pi_i) + \varepsilon |I_i|}$

## Maximum Likelihood CPTs

$$\delta\left(X_{U,n}, X_{U,m}\right) = \begin{cases} 1 & \text{if } X_{U,n} = X_{U,m} \\ 0 & \text{otherwise} \end{cases}$$

Counts: # of times what's in the bracket appeared in data, for example:

$$N = \sum_{x_1} m(x_1) = \sum_{x_1} \left( \sum_{x_2} m(x_1, x_2) \right) = \sum_{x_1} \left( \sum_{x_2} \left( \sum_{x_3} m(x_1, x_2, x_3) \right) \right)$$

• So...  $i(\theta) = \sum_{i=1}^N \log p(X_{U,i} | \theta)$

$$\begin{aligned} &= \sum_{i=1}^N \log \prod_{X_U} p(X_U | \theta)^{\delta(X_U, X_{i,n})} \\ &= \sum_{i=1}^N \sum_{X_U} \delta(X_U, X_{i,n}) \log p(X_U | \theta) \\ &= \sum_{X_U} m(X_U) \log p(X_U | \theta) = \sum_{X_U} m(X_U) \log \prod_{i=1}^M p(x_i | \pi_i, \theta_i) \\ &= \sum_{X_U} \prod_{i=1}^M m(X_U) \log p(x_i | \pi_i, \theta_i) \end{aligned}$$

## Maximum Likelihood CPTs

- Let's try an example:
  - Compute the cpt
  - $p(x_3 | x_1)$
  - Using the formula:  $\theta(x_i, \pi_i) = \frac{m(x_i, \pi_i)}{m(\pi_i)}$
  - Note, here 0/0 = prior constant

PATIENT	FLU	FEVER	SINUS	TEMP	SWELL	HEAD
1	Y	Y	N	L	Y	Y
2	N	N	N	M	N	Y
3	Y	N	Y	H	Y	N
4	Y	N	Y	M	N	N

**Efficient**, only count over subset of variables in  $p(X_8 | X_4)$   
Not all  $p(x_1, \dots, x_6)$

## Learning Fully Observed Models

- Easiest scenario: we have observed all the nodes
  - Want to learn the probability tables from data...
  - Have N iid patients:

PATIENT	FLU	FEVER	SINUS	TEMP	SWELL	HEAD
1	Y	Y	N	L	Y	Y
2	N	N	N	M	N	Y
3	Y	N	Y	H	Y	N
4	Y	N	Y	M	N	N

  - Simplest case: least general graph  handle each dim individually as Bernoulli/Multinomial
  - 2<sup>nd</sup> Simplest case: most general, count each entry in pdf 
  - What about learning graphs in between?

Divide by total count  
Since  $\sum_{x_1} \dots \sum_{x_6} p(x) = 1$

## Maximum Likelihood CPTs

- Continuing:  $l(\theta) = \sum_{X_i} \sum_{m=1}^M m(X_i) \log p(x_i | \pi_i, \theta_i)$   
 $= \sum_{i=1}^M \sum_{x_i, \pi_i} \sum_{X_i=x_i, \pi_i} m(X_i) \log p(x_i | \pi_i, \theta_i)$   
 $= \sum_{i=1}^M \sum_{x_i, \pi_i} m(x_i, \pi_i) \log p(x_i | \pi_i, \theta_i)$
- Define:  $\theta(x_i, \pi_i) = p(x_i | \pi_i, \theta)$  Constraint:  $\sum_{x_i} \theta(x_i, \pi_i) = 1$
- Now have above with Lagrange multipliers:  
 $l(\theta) = \sum_{i=1}^M \sum_{x_i} \sum_{\pi_i} m(x_i, \pi_i) \log \theta(x_i, \pi_i) - \sum_{i=1}^M \sum_{x_i} \lambda_{\pi_i} (\sum_{x_i} \theta(x_i, \pi_i) - 1)$   
 $\frac{\partial l(\theta)}{\partial \theta(x_i, \pi_i)} = \frac{m(x_i, \pi_i)}{\theta(x_i, \pi_i)} - \lambda_{\pi_i} = 0 \rightarrow \theta(x_i, \pi_i) = \frac{m(x_i, \pi_i)}{\lambda_{\pi_i}}$
- Plug constraint:  $\sum_{x_i} \frac{m(x_i, \pi_i)}{\lambda_{\pi_i}} = 1 \rightarrow \lambda_{\pi_i} = \sum_{x_i} m(x_i, \pi_i) = m(\pi_i)$
- Final solution (trivial!!):  $\theta(x_i, \pi_i) = \frac{m(x_i, \pi_i)}{m(\pi_i)}$

## Conditional Dependence Tests

- Another thing we would like to do with a graphical model: Check conditional independencies...
    - *"Is Temperature Independ. of Flu Given Fever?"*
    - *"Is Temperature Independ. of Sinus Infection Given Fever?"*
  - Try computing & simplify marginals of  $p(\mathbf{x})$ 

$$p(\bar{X}) = p(x_1) p(x_2 | x_1) p(x_3 | x_1) p(x_4 | x_2) p(x_5 | x_3) p(x_6 | x_2, x_5)$$

$$p(x_4 | x_1, x_2, x_3) = \frac{p(x_1, x_2, x_3, x_4)}{p(x_1, x_2, x_3)} = \sum_{x_1} \sum_{x_2} \sum_{x_3} p(\bar{X})$$

$$= \frac{p(x_1) p(x_2 | x_1) p(x_3 | x_1) p(x_4 | x_2)}{p(x_1) p(x_2 | x_1) p(x_3 | x_1)}$$

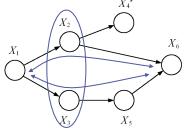
$$= p(x_4 | x_2) \quad \longleftrightarrow \quad x_4 \perp\!\!\!\perp x_1, x_3 | x_2$$
  - In this case it was easy, what if checking:  $x_1 \parallel x_6 | x_2, x_3$
  - Hard to compute  $p(x_6 | x_2, x_3, x_4)$  want efficient algorithm...

## D-Separation & Bayes Ball

- There is a graph algorithm for checking independence
- Intuition: separation or blocking of some nodes by others

- Example:

if nodes  $x_2, x_3$  "block" path from  $x_1$  to  $x_6$   
we might say that  
 $x_1 \perp\!\!\!\perp x_6 | x_2, x_3$



This is not exact for directed graphs (true for Undirected)

We need more than just simple Separation

Need D-Separation (directed separation)

D-Separation is computed via the Bayes Ball algorithm

Use to prove general statements over subsets of vars:

$$X_A \perp\!\!\!\perp X_B | X_C$$

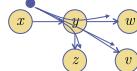
## Bayes Ball Algorithm

- The algorithm:
  - 1) Shade nodes  $X_A$
  - 2) Place a ball at each node in  $X_A$
  - 3) Bounce balls around graph according to some rules
  - 4) If no ball reaches  $X_B$ , then  $X_A \perp\!\!\!\perp X_B | X_C$  is true (else false)

Balls can travel along/against arrows

Pick any incoming & outgoing path

Test each to see if ball goes through or bounces back



Look at canonical sub-graphs & leaf cases for rules...

## Bayes Ball Algorithm

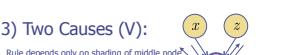
- 1) Markov Chain:
 

Rule depends only on shading of middle node

Ball stops  $x \perp\!\!\!\perp z | y$
- 2) Two Effects:
 

Rule depends only on shading of middle node

Ball stops  $x \perp\!\!\!\perp z | y$

Go Through  $x \not\perp\!\!\!\perp z$
- 3) Two Causes (V):
 

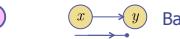
Rule depends only on shading of middle node

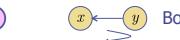
Go Through  $x \not\perp\!\!\!\perp z | y$

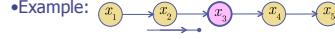
Ball stops  $x \perp\!\!\!\perp z$

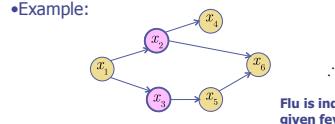
## Bayes Ball Algorithm

- Also need to look at special 'leaf' cases:

Bounces back   Ball is stopped

Ball is stopped   Bounces back

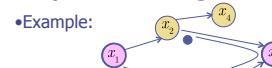
Example:   $x_1 \perp\!\!\!\perp x_6 | x_3$

Example: 

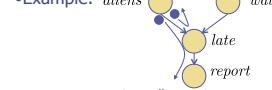
$x_1, x_2, x_4 \text{ Stopped}$ $x_1, x_2, x_5 \text{ Stopped}$ $x_1, x_3, x_5 \text{ Stopped}$ $\therefore x_1 \perp\!\!\!\perp x_6   \{x_2, x_3\}$	<b>Flu is independent of headache given fever &amp; sinus infection!</b>
---	--

## Bayes Ball Algorithm

- Example:


 **$x_2, x_6, x_5$  Goes Through Because of V-structure**  
 $\therefore x_2 \not\perp\!\!\!\perp x_3 | \{x_1, x_6\}$

- Example: aliens -> watch -> late -> report


**aliens  $\perp\!\!\!\perp$  watch | report**  
**aliens  $\not\perp\!\!\!\perp$  watch | report**  

Ball bounces back from report leaf and goes to right if report is shaded.  
 Bob is waiting for Alice but can't know if she is late. Instead a security guard says if she is. She can be late if aliens abduct her or Bob's watch is ahead (daylight savings time). Guard reports she is late.  
 If watch is ahead,  $p(\text{alien}=\text{true})$  goes down, they are dependent.

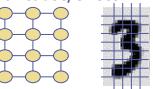
# Machine Learning

## Topic 16

- Undirected Graphs
- Undirected Separation
- Inferring Marginals & Conditionals
- Moralization
- Junction Trees
- Triangulation

## Undirected Graphs

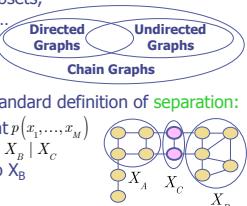
- Separation is *much easier* for **undirected graphs**
- But, what are undirected graphs and why use them?
- Might be hard to call vars parent/child or cause/effect
- Example: Image pixels
- Each pixel is Bernoulli = {0,1}
- Where 0=dark, 1=bright
- Have probability over all pixels  $p(x_{11}, \dots, x_{1M}, \dots, x_{M1}, \dots, x_{MM})$
- Bright pixels have Bright neighbors
- Nearby pixels dependent, so connect with links
- Get a graphical model that looks like a grid
- But who is parent? No parents really, just probability
- Grid models are called **Markov Random Fields**
- Used in vision, physics (lattice, spin, or Ising models), etc.



Tony Jebara, Columbia University

## Undirected Graphs

- Undirected & directed not subsets,
- **Chain Graphs** are a superset...
- Some distributions behave as undirected graphs, some as directed, some as both
- Undirected graphs use the standard definition of **separation**:  
an undirected graph says that  $p(x_1, \dots, x_M)$  satisfies any statement  $X_A \perp\!\!\!\perp X_B | X_C$  if no paths can go from  $X_A$  to  $X_B$  unless they go through  $X_C$
- Thus, undirected graphs obey the **general Markov property**
- Recall the simple Markov property  
 $x_1 - x_2 - x_3 \quad x_1 \perp\!\!\!\perp x_3 | x_2 \Rightarrow p(x_1 | x_2, x_3) = p(x_1 | x_2)$



## Hammersley Clifford Theorem

Theorem [HC]: any distribution that obeys the **Markov property**

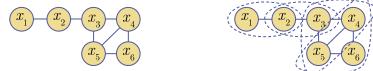
$$p(x_i | X_{U(i)}) = p(x_i | X_{Ne(i)}) \quad \forall i \in U$$

can be written as a product of terms over all maximal cliques

$$p(X_U) = p(x_1, \dots, x_M) = \frac{1}{Z} \prod_{c \in C} \psi_c(X_c)$$

**Clique:** a subset of nodes that are all pair-wise adjacent

**Maximal clique:** cannot add more variables and still be a clique  
Each c is a maximal clique of variables  $X_c$  in the graph  
C is the set of all maximal cliques



## Undirected Graph Functions

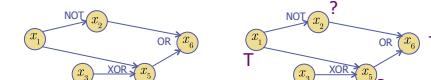
- Probability for undirected factorizes as a product of small non-negative **Potential Functions** over cliques in the graph  
 $p(X) = p(x_1, \dots, x_M) = \frac{1}{Z} \prod_{c \in C} \psi_c(X_c)$
- Normalizing term  $Z = \sum_X \prod_{c \in C} \psi_c(X_c)$  makes  $p(X)$  sum to 1
- Potentials  $\psi$  are non-negative un-normalized functions over cliques (subgroups of fully inter-connected variables)
- Use only **maximal cliques** since small  $\psi$  absorb into larger  $\psi$   
 $\psi(x_2, x_3) \psi(x_2) \rightarrow \psi(x_2, x_3) = \begin{bmatrix} 1 & 2 \\ 5 & 0 \end{bmatrix}$
- $p(X) = \frac{1}{Z} \psi(x_1, x_2) \psi(x_2, x_3) \psi(x_3, x_4, x_5) \psi(x_4, x_5, x_6)$

## Undirected Separation Examples

- Example:  
  
 $x \perp\!\!\!\perp y | \{w, z\}$   
 $w \perp\!\!\!\perp z | \{x, y\}$
- Example:  
  
 $x \perp\!\!\!\perp y | \{w\}$   
 $x \text{ } \cancel{\perp\!\!\!\perp} y | \{w, z\}$   
**Directed can't do it!**  
Must be acyclic  
Will have at least one V structure and ball goes through
- Example:  
  
 $x \perp\!\!\!\perp z | y$   
 $x \text{ } \cancel{\perp\!\!\!\perp} z | y$   
**Undirected can't do it!**

## Logical Inference

- Classic logic network: nodes are binary
- Arrows represent AND, OR, XOR, NAND, NOR, NOT etc.
- Inference: given observed binary variables, predict others



- Problems: uncertainty, conflicts and inconsistency
- Could get  $x_3=T$  and  $x_3=F$  following two different paths
- We need a way to enforce consistency and combine conflicting statements via probabilities and Bayes rule!

## Probabilistic Inference

- Replace logic network with Bayesian network
- Tables represent AND, OR, XOR, NAND, NOR, NOT etc.
- Probabilistic Inference: given observed binary variables, predict marginals over others

NOT		$x_3=f \quad x_3=t$ $x_1=f \quad [0, 1]$ $x_1=t \quad [1, 0]$
XOR		$x_1=f \quad x_1=t$ $x_2=f \quad [0, 1]$ $x_2=t \quad [1, 0]$ $x_3=f \quad x_3=t$ $x_3=f \quad [1, 0]$ $x_3=t \quad [0, 1]$
soft NOT		$x_3=f \quad x_3=t$ $x_1=f \quad [1, 0]$ $x_1=t \quad [0, 1]$

## Probabilistic Inference

- Two types of inference with a probability distribution:

$$p(X) = p(x_1, \dots, x_M) \text{ with queries } X_F \subseteq X \text{ given evidence } X_E \subseteq X$$

### Marginal Inference:

$$p(X_F | X_E) = \frac{p(X_F, X_E)}{p(X_E)} = \sum_{X \setminus X_F \cup X_E} p(X) / \sum_{X \setminus X_E} p(X)$$

or...  $p(x_i | X_E) \forall x_i \in X_F$

### Maximum a posteriori (MAP) inference:

$$\arg \max_{X_F} p(X_F | X_E)$$

...for now we focus on marginal inference

## Traditional Marginal Inference

- Marginal inference problem: given graph and probability function  $p(X) = p(x_1, \dots, x_M)$  for any subsets of variables

find

$$p(X_F | X_E) = \frac{p(X_F, X_E)}{p(X_E)}$$

- So, we basically compute both marginals and divide

- But finding marginals can take exponential work!

- A problem for both directed & undirected graphs:

$$\begin{aligned} p(x_j, x_k) &= \sum_{x_1} \sum_{x_2} \dots \sum_{x_M} \prod_{i=1}^M p(x_i | x_i) \\ p(x_j, x_k) &= \sum_{x_1} \sum_{x_2} \dots \sum_{x_M} \frac{1}{Z} \prod_{c \in C} \psi_c(X_c) \end{aligned}$$

- Graphs have efficient storage, learning, Bayes Ball...

- Graphs can also be used to perform efficient inference!

- Junction Tree Algorithm: method to efficiently find marginals

## Traditional Marginal Inference

- Example: brute force inference on a directed graph...

- Given a directed graph structure & filled-in CPTs

- We would like to efficiently compute arbitrary marginals

- Or we would like to compute arbitrary conditionals

$$\begin{aligned} p(X) &= p(x_1) p(x_2 | x_1) p(x_3 | x_1) p(x_4 | x_2) p(x_5 | x_3) p(x_6 | x_2, x_5) \\ p(x_1, x_3) &= p(x_1) p(x_3 | x_1) \\ p(x_1, x_6) &= \sum_{x_2, x_5} p(x_1) p(x_2 | x_1) p(x_3 | x_1) p(x_4 | x_2) p(x_5 | x_1) p(x_6 | x_2, x_5) \\ p(x_1 | x_6) &= \sum_{x_2, x_5} \frac{p(X)}{\sum_{x_2, x_5} p(x_1) p(x_2 | x_1) p(x_3 | x_1) p(x_4 | x_2) p(x_5 | x_1) p(x_6 | x_2, x_5)} \end{aligned}$$

- For example, we may have some evidence, i.e.  $x_6 = \text{TRUE}$

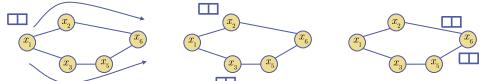
$$p(x_1 | x_6 = \text{TRUE}) = \frac{\sum_{x_2, x_5} p(x_1) p(x_2 | x_1) p(x_3 | x_1) p(x_4 | x_2) p(x_5 | x_1) p(x_6 | x_2, x_5)}{\sum_{x_2, x_5} p(x_1) p(x_2 | x_1) p(x_3 | x_1) p(x_4 | x_2) p(x_5 | x_1) p(x_6 | x_2, x_5)}$$

- This is tedious & does not exploit the graph's efficiency

Tony Jebara, Columbia University

## Efficient Marginals & Inference

- Another idea is to use some efficient graph algorithm
- Try sending messages (small tables) around the graph



- Hopefully these somehow settle down and equal marginals

$$\hat{p}(x_1, x_6) = \sum_{x_2, x_3, x_4, x_5} p(X)$$

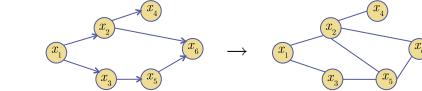
$$\text{AND marginals are self-consistent } \sum_{x_1} \hat{p}(x_1, x_6) = \sum_{x_2} \hat{p}(x_2, x_6)$$

- Note: can't just return conditionals since they can be inconsistent  $\sum_{x_1} \hat{p}(x_6 | x_1) \neq \sum_{x_2} \hat{p}(x_6 | x_2)$
- Junction Tree Algorithm must find consistent marginals

## Junction Tree Algorithm

- An algorithm that achieves fast inference, by doing message passing on undirected graphs.

- We first convert a directed graph to an undirected one



- Then apply the efficient Junction Tree Algorithm:

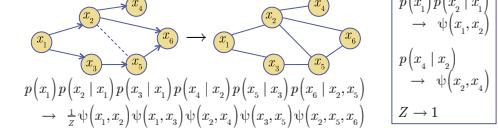
- 1) Moralization
- 2) Introduce Evidence
- 3) Triangulate
- 4) Construct Junction Tree
- 5) Propagate Probabilities (Junction Tree Algorithm)

## Moralization

- Converts directed graph into undirected graph
- By moralization, marrying the parents:

- 1) Connect nodes that have common children

- 2) Drop the arrow heads to get undirected



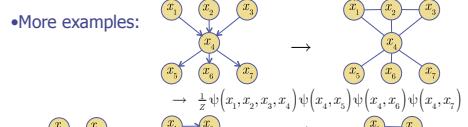
- Note: moralization resolves coupling due to marginalizing

- moral graph is more general (loses some independencies)



## Moralization

- More examples:



- More general graph less efficient but same inference:

$$\begin{aligned} p(x_1) &= \sum_{x_2, x_3} p(x_1, x_2, x_3) \\ &= \sum_{x_2} p(x_1 | x_2) p(x_3) \\ &= \sum_{x_2} p(x_1 | x_2) p(x_2) \end{aligned}$$

same

## Introducing Evidence

- Given moral graph, note what is observed  $X_E \rightarrow \bar{X}_E$

$$p(X_F | X_E = \bar{X}_E) \equiv p(X_F | \bar{X}_E)$$

- If we know this is always observed at  $X_E \rightarrow \bar{X}_E$ , simplify...

- Reduce the probability function since those  $X_E$  fixed

- Only keep probability function over remaining nodes  $X_F$

- Only get marginals and conditionals with subsets of  $X_F$

$$\begin{aligned} p(X_F | \bar{X}_E) &\propto \frac{1}{Z} \psi(x_1, x_2, x_3, x_4) \psi(x_4, x_5) \psi(x_4, x_6) \psi(x_4, x_7) \\ &\propto \frac{1}{Z} \tilde{\psi}(x_1, x_2, x_3, x_4) \tilde{\psi}(x_4, x_5) \tilde{\psi}(x_4, x_6) \tilde{\psi}(x_4, x_7) \\ &\text{Replace potential functions with slices} \quad \begin{array}{|c|c|} \hline 0.3 & 0.13 \\ \hline 0.12 & 0.1 \\ \hline \end{array} \\ &\text{But, need to recompute different normalization Z...} \end{aligned}$$

## Introducing Evidence

- Recall undirected separation, observing  $X_E$  separates others

$$p(X_F, X_E) = \psi(x_1, x_2, x_3, x_4, x_5, x_6, x_7)$$

$$p(X_F | \bar{X}_E) = \psi(x_1, x_2, x_3, x_4, x_5, x_6, x_7)$$

- But, need to recompute new normalization ...

$$\begin{aligned} p(X_F | \bar{X}_E) &\propto \frac{1}{Z} \tilde{\psi}(x_1, x_2, x_3, x_4) \tilde{\psi}(x_5, x_6, x_7) \\ &\propto \tilde{p}(X_F) = \frac{1}{Z} \tilde{\psi}(x_1, x_2, x_3, x_4) \tilde{\psi}(x_5, x_6, x_7) \end{aligned}$$

- Just avoid Z & normalize at the end when we are querying individual marginals and conditionals as subsets of  $X_F$

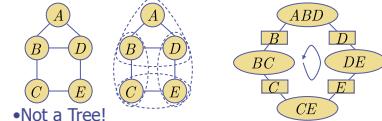
$$\begin{aligned} \tilde{p}(x_2) &= \frac{\sum_{x_1, x_3, x_4, x_5, x_6, x_7} \tilde{\psi}(x_1, x_2, x_3, x_4) \tilde{\psi}(x_5, x_6, x_7) \tilde{\psi}(x_2)}{\sum_{x_1, x_2, x_3, x_4, x_5, x_6, x_7} \tilde{\psi}(x_1, x_2, x_3, x_4) \tilde{\psi}(x_5, x_6, x_7) \tilde{\psi}(x_2)} \end{aligned}$$

## Junction Trees

- Given moral graph want to build **Junction Tree**:
    - each node is a **clique** ( $\psi$ ) of variables in moral graph
    - edges connect cliques of the potential functions
    - unique path between nodes & root node (tree)
    - between adjacent clique nodes, create **separators** ( $\phi$ )
    - separator nodes contain intersection of variables
- undirected**    **cliques**    **clique tree**    **junction tree**
- $$p(X) = \frac{1}{z} \psi(A, B, D) \psi(B, C, D) \psi(C, D, E)$$
- 

## Triangulation

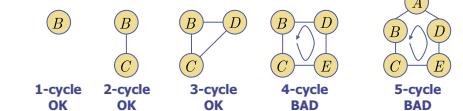
- Problem: imagine the following undirected graph



- Not a Tree!
- To ensure Junction Tree is a tree (no loops, etc.) before forming it must first **Triangulate** moral graph before finding the cliques...
- Triangulating gives more general graph (like moralization)
- Adds links to get rid of cycles or loops
- Triangulation: Connect nodes in moral graph until no chordless cycle of 4 or more nodes remains in the graph

## Triangulation

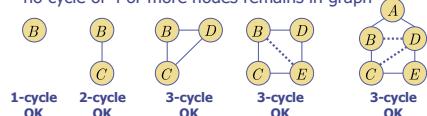
- Triangulation: Connect nodes in moral graph such that no cycle of 4 or more nodes remains in graph



- So, *add links*, but many possible choices...
- HINT: Try to keep largest clique size small (makes junction tree algorithm more efficient)
- Sub-optimal triangulations of moral graph are Polynomial
- Triangulation that minimizes largest clique size is NP
- But, OK to use a suboptimal triangulation (slower JTA...)

## Triangulation

- Triangulation: Connect nodes in moral graph such that no cycle of 4 or more nodes remains in graph



- So, *add links*, but many possible choices...
- HINT: Try to keep largest clique size small (makes junction tree algorithm more efficient)
- Sub-optimal triangulations of moral graph are Polynomial
- Triangulation that minimizes largest clique size is NP
- But, OK to use a suboptimal triangulation (slower JTA...)

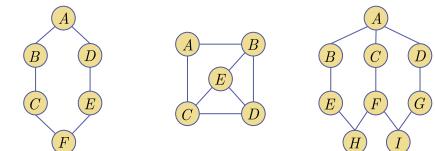
# Machine Learning

## 4771

### Topic 17

- Triangulation Examples
- Running Intersection Property
- Building a Junction Tree
- The Junction Tree Algorithm

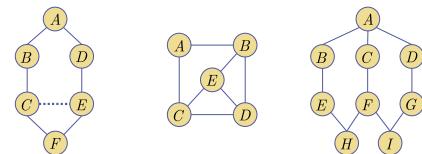
### Triangulation Examples



- Cycle:** A closed (simple) path, with no repeated vertices other than the starting and ending vertices
- Chordless Cycle:** a cycle where no two non-adjacent vertices on the cycle are joined by an edge.
- Triangulated Graph:** a graph that contains no chordless cycle of four or more vertices (aka a Chordal Graph).

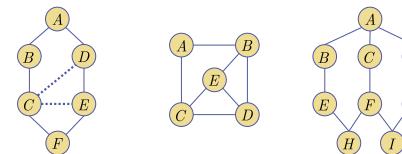
Tony Jebara, Columbia University

### Triangulation Examples



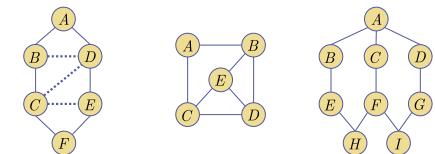
Tony Jebara, Columbia University

### Triangulation Examples



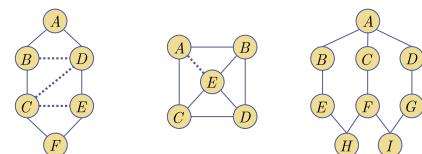
Tony Jebara, Columbia University

### Triangulation Examples



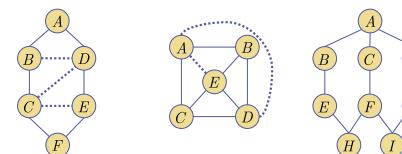
Tony Jebara, Columbia University

### Triangulation Examples



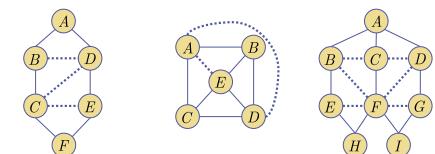
Tony Jebara, Columbia University

### Triangulation Examples

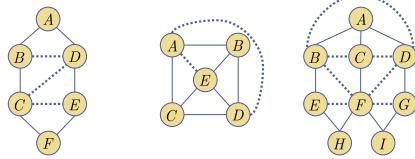


Tony Jebara, Columbia University

### Triangulation Examples

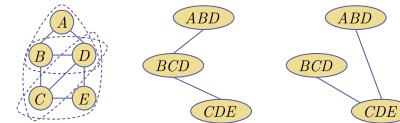


## Triangulation Examples



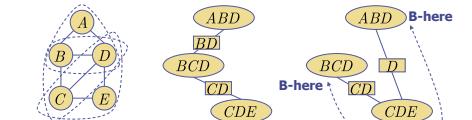
## Running Intersection Property

- Junction Tree must satisfy **Running Intersection Property**
- RIP: On unique path connecting clique  $V$  to clique  $W$ , all other cliques share nodes in  $V \cap W$



## Running Intersection Property

- Junction Tree must satisfy **Running Intersection Property**
- RIP: On unique path connecting clique  $V$  to clique  $W$ , all other cliques share nodes in  $V \cap W$



HINT: Junction Tree has largest total separator cardinality

$$|\Phi| = |\phi(B, D)| + |\phi(C, D)| = 2 + 2$$

$$|\Phi| = |\phi(C, D)| + |\phi(D)| = 2 + 1$$

## Forming the Junction Tree

- Goal: connect  $k$  cliques into a tree...  $k^{k-2}$  possibilities!
- For each, check Running Intersection Property, too slow...
- Theorem: a valid (RIP) Junction Tree connection is one that maximizes the cardinality of the separators

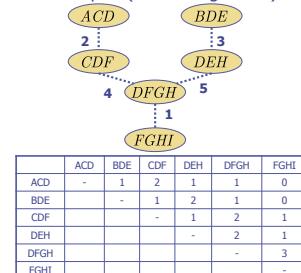
$$\begin{aligned} JT^* &= \arg \max_{\text{TREE STRUCTURES}} |\Phi| \\ &= \arg \max_{\text{TREE STRUCTURES}} \sum_S |\phi(X_S)| \end{aligned}$$

- Use very fast **Kruskal algorithm**:

- Init Tree with all cliques unconnected (no edges)
- Compute size of separators between all pairs
- Connect the two cliques with the biggest separator cardinality which doesn't create a loop in current Tree (maintains Tree structure)
- Stop when all nodes are connected, else goto 3

## Kruskal Example

- Start with unconnected cliques (after triangulation)



## Junction Tree Probabilities

- We now have a valid Junction Tree!
- What does that mean?
- Recall probability for undirected graphs:  
 $p(X) = p(x_1, \dots, x_m) = \frac{1}{Z} \prod_c \psi(X_c)$
- Can write junction tree as potentials of its cliques:  
 $p(X) = \frac{1}{Z} \prod_c \psi(X_c)$
- Alternatively: clique potentials over separator potentials:  
 $p(X) = \frac{1}{Z} \prod_s \phi(X_s)$
- This doesn't change/do anything! Just less compact...
- Like *de-absorbing* smaller cliques from maximal cliques:

$$\tilde{\psi}(A, B, D) = \frac{\psi(A, B, D)}{\phi(B, D)} \quad \text{...gives back original formula if } \phi(B, D) \triangleq 1$$

## Junction Tree Probabilities

- Can quickly convert directed graph into this form:

$$p(X) = \frac{1}{Z} \prod_c \psi(X_c)$$

- Example:

$$\begin{aligned} p(X) &= p(x_1) p(x_2 | x_1) p(x_3 | x_2) p(x_4 | x_3) \\ &= p(x_1) p(x_2) p(x_3 | x_1) p(x_4 | x_3) \\ &= p(x_1, x_2) p(x_3, x_4 | x_1, x_3) \\ &= \frac{1}{Z} \frac{p(x_1, x_2) p(x_3, x_4 | x_1, x_3)}{1 \times 1} \\ &= \frac{1}{Z} \frac{\psi(x_1, x_2) \psi(x_3, x_4) \phi(x_3)}{\phi(x_2) \phi(x_3)} \end{aligned}$$

By inspection, can just cut & paste CPTs as clique and separator potential functions

## Junction Tree Algorithm

- Running the JTA converts clique potentials & separator potentials into marginals over their variables ... and does not change  $p(X)$
- Don't want just normalization!  
 $\frac{\psi(A, B, D)}{\sum_{A, B, D} \psi(A, B, D)} \neq p(A, B, D)$
- These marginals should all agree & be consistent  
 $\psi(A, B, D) \rightarrow p(A, B, D) \rightarrow \sum_A p(A, B, D) = \tilde{p}(B, D) \rightarrow \phi(B, D) \rightarrow p(B, D) \rightarrow \psi(B, C, D) \rightarrow p(B, C, D) \rightarrow \sum_B p(B, C, D) = \tilde{p}(B, D)$  ALL EQUAL
- Consistency: all distributions agree on submarginals
- JTA sends messages between cliques & separators dividing each by the others marginals until consistency...

## Junction Tree Algorithm

- Send message from each clique to its separators of what it thinks the submarginal on the separator is.
- Normalize each clique by incoming message from its separators so it agrees with them

$$\begin{array}{ccc} AB & \xrightarrow{\psi_{AB}} & BC \\ V = \{A, B\} & S = \{B\} & W = \{B, C\} \end{array}$$

If agree:  $\sum_{V \setminus S} \psi_V = \phi_S = p(S) = \phi_S = \sum_{W \setminus S} \psi_W \quad \text{...Done!}$

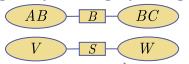
Else: Send message From V to W... Now they Agree...Done!

$$\begin{array}{lcl} \phi_S^* = \sum_{V \setminus S} \psi_V^* & \sum_{V \setminus S} \psi_V^* = \sum_{W \setminus S} \frac{\phi_S}{\phi_S^*} \psi_W^* \\ \psi_W^* = \frac{\phi_S^*}{\phi_S} \psi_W^* & = \frac{\phi_S^*}{\phi_S} \sum_{V \setminus S} \psi_V^* \\ \psi_V^* = \psi_V & = \phi_S^* = \sum_{W \setminus S} \psi_W^* \end{array}$$

## Junction Tree Algorithm

- When "Done", all clique potentials are marginals and all separator potentials are submarginals!
- Note that  $p(X)$  is unchanged by message passing step:

$$\begin{aligned}\phi_s^* &= \sum_{V,S} \psi_V \\ \psi_w^* &= \frac{\phi_s^*}{\phi_s} \psi_w \\ \psi_v^* &= \psi_v\end{aligned}$$



$$p(X) = \frac{1}{Z} \frac{\psi_v^* \psi_w^*}{\phi_s^*} = \frac{1}{Z} \frac{\psi_v \frac{\phi_s^*}{\phi_s} \psi_w}{\phi_s^*} = \frac{1}{Z} \frac{\psi_v \psi_w}{\phi_s}$$

- Potentials set to conditionals (or slices) become marginals!

$$\begin{aligned}\psi_{AB} &= p(B | A) p(A) \\ &= p(A, B) \\ \psi_{BC} &= p(C | B) \\ \phi_B &= 1\end{aligned} \longrightarrow \begin{aligned}\phi_B^* &= \sum_A \psi_{AB} = \sum_A p(A, B) = p(B) \\ \psi_{BC}^* &= \frac{\phi_s^*}{\phi_s} \psi_{BC} = \frac{p(B)}{1} p(C | B) = p(B, C)\end{aligned}$$

## Machine Learning

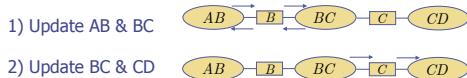
### Topic 18

- The Junction Tree Algorithm
- Collect & Distribute
- Algorithmic Complexity
- ArgMax Junction Tree Algorithm



### JTA with many cliques

- Problem: what if we have more than two cliques?



- Problem: AB has not heard about CD!  
After BC updates, it will be inconsistent for AB
- Need to iterate the pairwise updates many times
- This will eventually converge to consistent marginals
- But, inefficient... can we do better?

### JTA: Collect & Distribute

• Use tree recursion rather than iterate messages mindlessly!

```

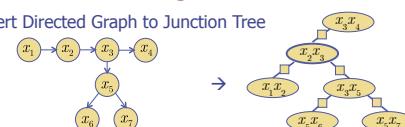
initialize(DAG){ Pick root
    Set all variables as:  $\psi_{C_i} = p(x_i | \pi_i)$ ,  $\phi_S = 1$ 
    collectEvidence(node) {
        for each child of node {
            update1(node,collectEvidence(child));
        }
        return(node);
    }
    distributeEvidence(node) {
        for each child of node {
            update2(child,node);
            distributeEvidence(child);
        }
    }
    update1(node w,node v) {  $\phi_{V\cap W}^* = \sum_{V\setminus(V\cap W)} \psi_V$ ,  $\psi_W = \frac{\phi_{V\cap W}}{\phi_{V\cap W}} \psi_W$  }
    update2(node w,node v) {  $\phi_{V\cap W}^{**} = \sum_{V\setminus(V\cap W)} \psi_V$ ,  $\psi_W = \frac{\phi_{V\cap W}^{**}}{\phi_{V\cap W}} \psi_W$  }
    normalize() {  $p(X_C) = \frac{1}{\sum_C \psi_C^*}$   $\forall C$ ,  $p(X_S) = \frac{1}{\sum_S \phi_S^{**}}$   $\forall S$  }
}

```

## Algorithmic Complexity

- The 5 steps of JTA are all efficient:

- 1) Moralization  
Polynomial in # of nodes
- 2) Introduce Evidence (fixed or constant)  
Polynomial in # of nodes (convert pdf to slices)
- 3) Triangulate (Tarjan & Yannakakis 1984)  
Suboptimal=Polynomial, Optimal=NP
- 4) Construct Junction Tree (Kruskal)  
Polynomial in # of cliques
- 5) Junction Tree Algorithm (Init,Collect,Distribute,Normalize)  
Polynomial (linear) in # of cliques, Exponential in Clique Cardinality



- Convert Directed Graph to Junction Tree
- Initialize separators to 1 (and Z=1) and set clique tables to the CPTs in the Directed Graph

$$p(X) = p(x_1)p(x_2 | x_1)p(x_3 | x_2)p(x_4 | x_3)p(x_5 | x_4)p(x_6 | x_5)p(x_7 | x_6)$$

$$p(X) = \frac{1}{Z} \prod_{S \in \mathcal{C}} \psi(S)$$

• Run Collect, Distribute, Normalize  
• Get valid marginals from all  $\psi, \phi$  tables

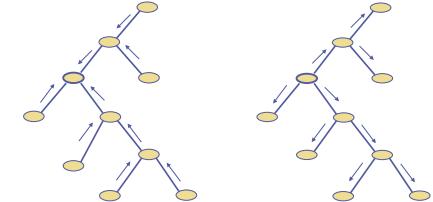
### Review: Junction Tree Algorithm

- Send message from each clique to its separators of what it thinks the submarginal on the separator is.
- Normalize each clique by incoming message from its separators so it agrees with them

	$V = \{A, B, C, D\}$ $S = \{B\}$ $W = \{B, C\}$	
<b>If agree:</b> $\sum_{V \setminus S} \psi_V = \phi_S = p(S) = \phi_S = \sum_{W \setminus S} \psi_W$ ...Done!		
<b>Else:</b> Send message From V to W...      Send message From W to V...      Now they Agree...Done!		
$\phi_S^* = \sum_{V \setminus S} \psi_V$ $\psi_W^* = \frac{\phi_S}{\phi_S^*} \psi_W$ $\psi_V^* = \psi_V$	$\phi_S^{**} = \sum_{W \setminus S} \psi_W^*$ $\psi_V^* = \frac{\phi_S^{**}}{\phi_S} \psi_V$ $\psi_W^* = \psi_W$	$\sum_{V \setminus S} \psi_V^* = \sum_{V \setminus S} \frac{\phi_S^*}{\phi_S} \psi_V^*$ $= \frac{\phi_S^*}{\phi_S} \sum_{V \setminus S} \psi_V^*$ $= \phi_S^{**} = \sum_{W \setminus S} \psi_W^*$

### Junction Tree Algorithm

- JTA: 1) Initialize 2) Collect 3) Distribute 4) Normalize



- Note: leaves do not change their  $\psi$  during collect
- Note: the first cliques collect changes are parents of leaves
- Note: root does not change its  $\psi$  during distribute

### JTA with Extra Evidence

- If extra evidence is observed, must slice tables accordingly
- Example:  $p(A, B, C, D) = \frac{1}{Z} \psi_{AB} \psi_{BC} \psi_{CD}$

$$Z = 1$$

$\psi_{AB} = \begin{bmatrix} 8 & 4 \\ 3 & 1 \end{bmatrix}$	$\psi_{BC} = \begin{bmatrix} 2 & 3 \\ 1 & 1 \end{bmatrix}$	$\psi_{CD} = \begin{bmatrix} 1 & 4 \\ 1 & 1 \end{bmatrix}$
$A=0$	$B=0$	$C=0$
$A=1$	$B=1$	$C=1$
$B=0$	$C=0$	$D=0$
$B=1$	$C=1$	$D=1$

- You are given evidence: A=0. Replace table with slices...

$$\psi_{AB} \rightarrow \begin{bmatrix} 8 & 4 \\ 0 & 0 \end{bmatrix} \quad \psi_{BC} \rightarrow \begin{bmatrix} 2 & 3 \\ 0 & 0 \end{bmatrix} \quad \psi_{CD} \rightarrow \begin{bmatrix} 1 & 4 \\ 0 & 0 \end{bmatrix}$$

- JTA now gives  $\psi, \phi$  as marginals conditioned on evidence

$$p(B | A=0) = \frac{\psi_{B|0}^*}{\sum_B \psi_{AB}^*}, \quad p(B, C | A=0) = \frac{\psi_{BC|0}^*}{\sum_{B,C} \psi_{BC}^*}, \quad p(C, D | A=0) = \frac{\psi_{CD|0}^*}{\sum_{C,D} \psi_{CD}^*}$$

- All denominators equal the new normalizer Z'

$$Z' = p(EVIDENCE) = \sum_B \psi_{AB}^* = \sum_{B,C} \psi_{BC}^* = \sum_{C,D} \psi_{CD}^*$$

## ArgMax Junction Tree Algorithm

- We can also use JTA for finding the max not the sum over the joint to get argmax of marginals & conditionals
  - Say have some evidence:  $p(X_p, \bar{X}_E) = p(x_1, \dots, x_n, \bar{x}_{n+1}, \dots, \bar{x}_N)$
  - Most likely (highest p)  $X_p^* = \arg \max_{x_p} p(X_p, \bar{X}_E)$
  - What is most likely state of patient with fever & headache?
- $$\begin{aligned} p_p^* &= \max_{x_1, x_2, x_3, x_4, x_5, x_6} p(x_1 = 1, x_2, x_3, x_4, x_5, x_6 = 1) \\ &= \max_{x_4} p(x_2 | x_1 = 1) p(x_1 = 1) \max_{x_5} p(x_3 | x_1 = 1) \\ &\quad \max_{x_6} p(x_4 | x_2) \max_{x_3} p(x_3 | x_2) p(x_6 = 1 | x_2, x_5) \end{aligned}$$
- Solution:** replace sum with max inside JTA update code
- $$\phi_{V \cap W}^* = \max_{V \setminus (V \cap W)} \psi_V, \quad \psi_W = \frac{\phi_{V \cap W}^*}{\phi_{V \cap W}} \psi_W \quad \phi_{V \cap W}^{**} = \max_{V \setminus (V \cap W)} \psi_V, \quad \psi_W = \frac{\phi_{V \cap W}^{**}}{\phi_{V \cap W}} \psi_W$$
- Final potentials are **max marginals**:  $\psi^{**}(X_c) = \max_{V \setminus C} p(X)$
- Highest value in potential is most likely:  $X_c^* = \arg \max_C \psi^{**}(X_c)$

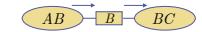
## ArgMax Junction Tree Algorithm

- Why do I need the ArgMax junction tree algorithm?
  - Can't I just compute marginals using the Sum algorithm and then find the highest value in each marginal???
  - No!! Here's a counter-example:
- |                 |   |
|-----------------|---|
| $p(x_1, x_2) =$ | $x_1 \quad x_1 \quad x_1$                   |
|                 | $A \quad B \quad C$                         |
|                 | $x_2 = 0 \quad [ .14 \quad .05 \quad .27 ]$ |
|                 | $x_2 = 1 \quad [ .24 \quad .20 \quad .10 ]$ |
- Most likely is  $x_1^* = C$  and  $x_2^* = 0$
- But the sub-marginals  $p(x_1)$  and  $p(x_2)$  do not reveal this...

$$p(x_1) = \begin{bmatrix} A & B & C \end{bmatrix} \quad p(x_2) = \begin{bmatrix} x_2 = 0 & [.46] \\ x_2 = 1 & [.54] \end{bmatrix}$$

- The marginals would *falsely* imply that  $x_1^* = A$  and  $x_2^* = 1$

## Example



- Note that products are element-wise
  - Let us send a regular JTA message from AB to BC
- $$\begin{aligned} \psi_{AB} &= \begin{bmatrix} 8 & 4 \\ 3 & 1 \end{bmatrix} \begin{matrix} A=0 \\ B=1 \end{matrix} & \phi_B &= \begin{bmatrix} 1 & \\ 1 & \end{bmatrix} \begin{matrix} B=0 \\ B=1 \end{matrix} & \psi_{BC} &= \begin{bmatrix} 2 & 3 \\ 1 & 1 \end{bmatrix} \begin{matrix} B=0 \\ C=1 \end{matrix} \\ \phi_B^* &= \sum_A \psi_{AB} = \sum_A \begin{bmatrix} 8 & 4 \\ 3 & 1 \end{bmatrix} = \begin{bmatrix} 11 & 5 \\ 11 & 1 \end{bmatrix} = \begin{bmatrix} 11 \\ 5 \end{bmatrix} \begin{matrix} B=0 \\ B=1 \end{matrix} & & & & \\ \psi_{BC}^* &= \frac{\phi_B^*}{\phi_B} \psi_{BC} = \begin{bmatrix} 11 \\ 5 \end{bmatrix} \begin{bmatrix} 2 & 3 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 22 & 33 \\ 5 & 5 \end{bmatrix} \begin{matrix} B=0 \\ C=1 \end{matrix} & & & & \end{aligned}$$
- If argmax JTA, just change the separator update to:
- $$\phi_B^* = \max_A \psi_{AB} = \max_A \begin{bmatrix} 8 & 4 \\ 3 & 1 \end{bmatrix} = \begin{bmatrix} 8 & 4 \\ 8 & 4 \end{bmatrix} = \begin{bmatrix} 8 & \\ 8 & \end{bmatrix} \begin{matrix} B=0 \\ B=1 \end{matrix}$$

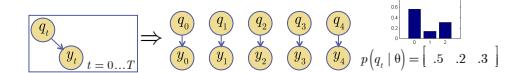
## Machine Learning

### Topic 19

- Hidden Markov Models
- HMMs as State Machines & Applications
- HMMs Basic Operations
- HMMs via the Junction Tree Algorithm

### Hidden Markov Models

- A great application of Junction Tree Algorithm and EM
- Recall mixture of Gaussians model on IID data

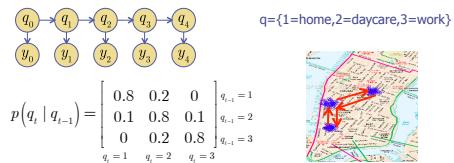


- Example: location data of a single parent as a mixture of Gaussians
- Parent has 3 internal states:  $q = \{\text{home}, \text{daycare}, \text{work}\}$
- Based on  $q$ , sample from appropriate Gaussian mean and covariance to get  $y = (\text{latitude}, \text{longitude})$



### Hidden Markov Models

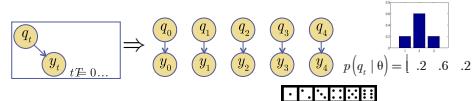
- Parent drops child at daycare before & after work. Not IID!



- Have dependence on previous state
- Can't go straight from home to work!
- Now, order of  $y_0, \dots, y_T$  matters (in IID order doesn't matter)

### Hidden Markov Models

- Consider mixture of multinomials (dice)  $y = \{1, 2, 3, 4, 5, 6\}$



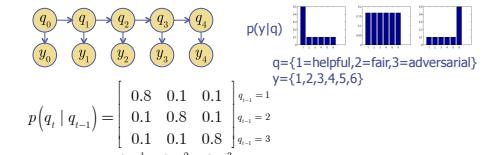
- Example: a crooked casino croupier using mixture of dice.
- You win if he rolls 1,2,3. You lose if he rolls 4,5,6.
- Croupier has 3 internal states  $q = \{\text{helpful}, \text{fair}, \text{adversarial}\}$
- Based on  $q$ , sample different 'dice' multinomial



### Hidden Markov Models

- But if the dealer has a memory or mood? Not IID!

5646166166 4321534161414341634 1113114121



- If you tip, dealer starts to like you and rolls the helpful die
- Dealer has a memory of his mood and last type of die  $q_{t-1}$
- Will often use same die for  $q_t$  as was rolled before...
- Now, order of  $y_0, \dots, y_T$  matters (if IID order doesn't matter)

### Hidden Markov Models

- Since next choice of the dice depends on previous one...



- Add left-right arrows. This is a **hidden Markov model**

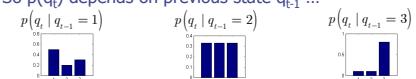
•Markov: *future*  $\perp$  *past*  $\mid$  *present*

$$p(q_t | q_{t-1}, q_{t-2}, \dots, q_1, q_0) = p(q_t | q_{t-1})$$

- From graph, have the following general formula:

$$p(X_T) = p(q_0) \prod_{t=1}^T p(q_t | q_{t-1}) \prod_{t=1}^T p(y_t | q_t)$$

- So  $p(q_t)$  depends on previous state  $q_{t-1}$  ...



### HMMs as State Machines

- HMMs have two variables: **state**  $q$  and **emission**  $y$

- Typically, we don't know  $q$  (hidden variable 1,2,3,4)

- HMMs are like **stochastic automata** or finite state machines...

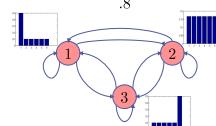
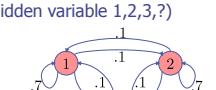
next state depends

on previous one...

(helpful, fair, adversarial)

- Can't observe state  $q$  directly, just a random related emission  $y$  outcome (dice roll) so...

doubly-stochastic automaton



### HMM Applications

- Speech Recognition

phonemes from  
audio cepstral vectors



Ba-ra-  
kk-oo-  
oo-d  
ah

- Language Parsing

parts of speech  
from words



Noun  
Verb  
Noun  
John  
Ate  
Pizza

- Genomics

splice site from  
gene sequence



-Intron- -Exon- -Promoter-

GATTACATTACCAACCATACG

## HMMs: Parameters

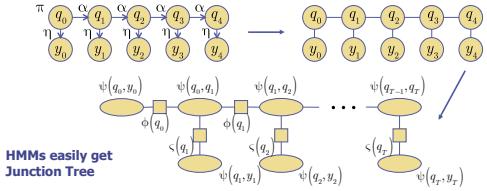
- We focus on HMMs with: discrete state  $q$  (of size M)  
discrete emission  $y$  (of size N)
  - Input will be arbitrary length string:  $y_0, \dots, y_T$
  - The pdf or (complete) likelihood is:
$$p(q, y) = p(q_0) \prod_{i=1}^T p(q_i | q_{i-1}) \prod_{i=0}^T p(y_i | q_i)$$
  - We don't know hidden states, the incomplete likelihood is:
$$p(y) = \sum_{q_0} \dots \sum_{q_T} p(q, y)$$
  - Assume HMM is stationary, tables are repeated:  $\theta = \{\pi, \eta, \alpha\}$

$p(q_i   q_{i-1}) = \prod_{i=1}^M \prod_{j=1}^M \alpha_j^{q_{i-1}^j q_i^j}$	$\sum_{j=1}^M \alpha_j = 1$		$M \times M$
$p(y_i   q_i) = \prod_{i=1}^M \prod_{j=1}^N [\eta_{ij}]^{y_i^j}$	$\sum_{j=1}^N \eta_{ij} = 1$		$M \times N$
$p(q_0) = \prod_{i=1}^M [\pi_i]^{q_0^i}$	$\sum_{j=1}^M \pi_j = 1$		$M$

# HMMs: Basic Operations

- Would like to do 3 basic things with our HMMs:
    - 1) Evaluate: given  $y_0, \dots, y_T$  &  $\theta$  compute  $p(y_0, \dots, y_T | \theta)$
    - 2) Decode: given  $y_0, \dots, y_T$  &  $\theta$  find  $q_0, \dots, q_T$  or  $p(q_0), \dots, p(q_T)$
    - 3) Max Likelihood: given  $y_0, \dots, y_T$  learn parameters  $\theta$

- Typically use Baum-Welch ( $\alpha$ - $\beta$  algo)... JTA is more general:



HMMs easily get  
Junction Tree

## HMMs: JTA Init & Verify

- **Init:**  $\psi(q_0, y_0) = p(q_0) p(y_0 | q_0)$     $\psi(q_1, q_{-1}) = p(q_{1+1} | q_1) = \alpha_{q_{-1}, q_1} \psi(q_1, y_1) = p(y_1 | q_1)$
  - **Collect up** (this time it actually doesn't change the zetas)  
 $\zeta^*(q_t) = \sum_{y_t} \psi(q_t, y_t) = \sum_y p(y_t | q_t) = 1$     $\psi^*(q_{t-1}, q_t) = \sum_{y_t} \psi(q_{t-1}, q_t, y_t) = \psi(q_{t-1}, q_t)$
  - **Collect left-right via phi's: change backbone to marginals**  
 $\hat{\phi}^*(q_t) = \sum_{y_t} \psi(q_t, y_t) = p(q_t)$     $\hat{\psi}^*(q_a, q_b) = \frac{1}{\zeta^*(q_b)} \psi(q_a, q_b) = p(q_a | q_b)$   
 $\hat{\phi}^*(q_t) = \sum_{y_t} \psi^*(q_{t-1}, q_t) = p(q_t)$     $\hat{\psi}^*(q_{t-1}, q_t) = \frac{1}{\zeta^*(q_t)} \psi^*(q_{t-1}, q_t) = p(q_t | q_{t-1})$
  - **Distribute:**  
 $\zeta^*(q_t) = \sum_{y_t} \psi^*(q_{t-1}, q_t) = \sum_{y_t} p(q_{t-1}, q_t) = p(q_t)$   
 $\psi^*(q_t, y_t) = \frac{1}{\zeta^*(q_t)} \psi^*(q_t, y_t) = \frac{p(q_t)}{\zeta^*(q_t)} p(y_t | q_t) = p(y_t, q_t)$    ...done!

## Machine Learning

## Topic 20

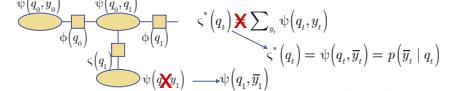
- HMMs with Evidence
- HMM Collect
- HMM Evaluate
- HMM Distribute
- HMM Decode
- HMM Parameter Learning via JTA & EM

## HMMs: JTA with Evidence

- If  $y$  sequence is observed (in problems 1,2,3) get evidence:

$$p(q, \bar{y}) = p(q_0) \prod_{t=1}^T p(q_t | q_{t-1}) \prod_{t=0}^T p(\bar{y}_t | q_t)$$

- The potentials turn into slices:

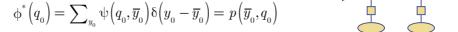


- Next, pick a root, for example rightmost one:  $\psi(q_{T-1}, q_T)$

- Collect all zeta separators bottom up:



- Collect leftmost phi separator to the right:



## HMMs: Collect with Evidence

- Now, we will collect (\*) along the backbone left to right
  - Update each clique with its left and bottom separators:
- 
- $\psi^*(q_t, q_{t+1}) = \frac{\phi^*(q_t) \zeta^*(q_{t+1})}{\sum_y \psi^*(q_t, y)} \psi(q_t, q_{t+1}) = \phi^*(q_t) p(\bar{y}_{t+1} | q_{t+1}) \alpha_{q_t, q_{t+1}}$
- $\phi^*(q_{t+1}) = \sum_y \psi^*(q_t, q_{t+1}) = \sum_y \phi^*(q_t) p(\bar{y}_{t+1} | q_{t+1}) \alpha_{q_t, q_{t+1}}$
- Keep going along chain until right most node
- Note: above formula for phi is recursive, could use as is.
- Property: recall we had  $\phi^*(q_0) = p(\bar{y}_0, q_0)$
- $\phi^*(q_t) = \sum_{q_{t+1}} p(\bar{y}_0, q_0) p(\bar{y}_1 | q_1) p(q_1 | q_0) = p(\bar{y}_0, \bar{y}_1, q_t)$
- $\phi^*(q_2) = \sum_{q_3} p(\bar{y}_0, \bar{y}_1, q_1) p(\bar{y}_2 | q_2) p(q_2 | q_1) = p(\bar{y}_0, \bar{y}_1, \bar{y}_2, q_2)$
- $\phi^*(q_{t+1}) = \sum_{q_{t+2}} p(\bar{y}_0, \dots, \bar{y}_t, q_t) p(\bar{y}_{t+1} | q_{t+1}) p(q_{t+1} | q_t) = p(\bar{y}_0, \dots, \bar{y}_{t+1}, q_{t+1})$

## HMMs: Evaluate with Evidence

- Say we are solving the first HMM problem:
  - 1) **Evaluate:** given  $y_0, \dots, y_T$  &  $\theta$  compute  $p(y_0, \dots, y_T | \theta)$
  - If we want to compute the likelihood, we are already done!
  - We really just need to do collect (not even distribute).
  - From previous slide we had:
  - $\phi^*(q_{t+1}) = \sum_{q_t} p(\bar{y}_0, \dots, \bar{y}_t, q_t) p(\bar{y}_{t+1} | q_{t+1}) p(q_{t+1} | q_t) = p(\bar{y}_0, \dots, \bar{y}_{t+1}, q_{t+1})$
  - Collect 'til root (rightmost node):  $\psi^*(q_{T-1}, q_T) = p(\bar{y}_0, \dots, \bar{y}_T, q_{T-1}, q_T)$  its normalizer is  $p(\text{EVIDENCE})!$
- 
- Hypothetical  $\phi^*(q_T)$
- Or use hypothetical  $\phi^*(q_T) = p(\bar{y}_0, \dots, \bar{y}_T, q_T)$
- Can compute the likelihood just by marginalizing this phi
- $p(\bar{y}_0, \dots, \bar{y}_T) = \sum_{q_T} p(\bar{y}_0, \dots, \bar{y}_T, q_T) = \sum_{q_T} \phi^*(q_T)$
- So, adding up the entries in last  $\phi^*$  gives us the likelihood

## HMMs: Distribute with Evidence

- Back to collecting... say just finished collecting to the root with our last update formula:

$$\psi^*(q_{T-1}, q_T) = \frac{\phi^*(q_{T-1})}{1} \zeta^*(q_T) = \phi^*(q_{T-1}) p(\bar{y}_T | q_T) \alpha_{q_{T-1}, q_T}$$

- Now, we distribute (\*\*\*) along the backbone right to left
- Have first \*\*\* for root (stays the same):  $\psi^{**}(q_{T-1}, q_T) = \psi^*(q_{T-1}, q_T)$
- Start going to the left from there:



for t=T-1 to 0

- a)  $\phi^{**}(q_t) = \sum_{q_{t+1}} \psi^{**}(q_t, q_{t+1})$
- b)  $\zeta^{**}(q_{t+1}) = \sum_{q_t} \psi^{**}(q_t, q_{t+1})$
- c)  $\psi^{**}(q_t, q_{t+1}) = \frac{\phi^{**}(q_{t+1})}{\zeta^{**}(q_{t+1})} \psi^*(q_t, q_{t+1})$
- d)  $\psi^{**}(q_t, q_{t+1}) = \frac{\zeta^{**}(q_t)}{\zeta^{**}(q_{t+1})} \psi^*(q_t, q_{t+1})$

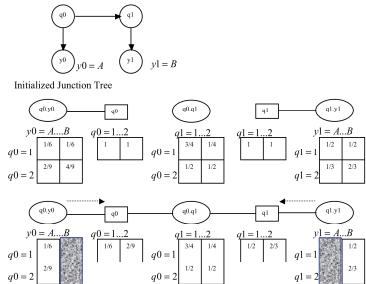
## HMM Example

You are given the parameters of a 2-state HMM. You observed the input sequence AB (from a 2-symbol alphabet A or B). In other words, you observe two symbols from your finite state machine, A and then B. Using the junction tree algorithm, evaluate the likelihood of this data  $p(y)$  given your HMM and its parameters. Also compute (or decoding) the individual marginals of the states after the evidence from this sequence is observed:  $p(q_0|y)$  and  $p(q_1|y)$ . The parameters for the HMM are provided below. They are the initial state prior  $p(q_0)$ , the state transition matrix given by  $p(q_t|q_{t-1})$ , and the emission matrix  $p(y_t|q_t)$ , respectively.

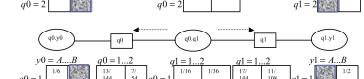
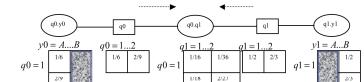
$$\pi = p(q_0) = \begin{bmatrix} 1 & 2 \\ 1/3 & 2/3 \end{bmatrix}$$

$$A^T = p(q_t | q_{t-1}) = \begin{bmatrix} 1 & 2 \\ 1 & 2 \end{bmatrix} \quad \eta^T = p(y_t | q_t) = \begin{bmatrix} 1 & 2 \\ A & B \end{bmatrix} \quad A = \begin{bmatrix} 1/2 & 1/3 \\ 1/4 & 1/2 \end{bmatrix}$$

## HMM Example



## HMM Example



$$\text{So the likelihood } p(y) = \frac{13}{144} + \frac{7}{54} + \frac{1}{16} + \frac{1}{18} + \frac{1}{36} + \frac{2}{27} = \frac{17}{144} + \frac{11}{108} = \frac{43}{144} = 0.2199$$

$$p(q_0 = 1 | y) = \frac{13/144 + 7/54}{13/144 + 7/54} = \frac{95}{144} = 0.65$$

$$p(q_1 = 1 | y) = \frac{17/144}{17/144 + 11/108} = \frac{51}{95} = 0.5368$$

$$p(q_2 = 1 | y) = \frac{11/108}{17/144 + 11/108} = \frac{44}{95} = 0.4632$$

## HMMs: Marginals & MaxDecoding

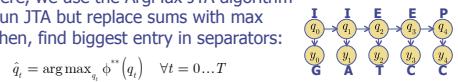
- Now that JTA is finished, we have the following:

$$\phi^{**}(q_t) \propto p(q_t | \bar{y}_1, \dots, \bar{y}_T) \quad \varsigma^{**}(q_{t+1}) \propto p(q_{t+1} | \bar{y}_1, \dots, \bar{y}_T)$$

$$\psi^{**}(q_t, q_{t+1}) \propto p(q_t, q_{t+1} | \bar{y}_1, \dots, \bar{y}_T)$$

- The separators define a distribution over the hidden states
- This gives the probability the DNA symbol  $y_t$  was  $q_t = \{I, E, P\}$
- We've done 2) **Decode**: given  $y_0, \dots, y_T$  &  $\theta$  find  $p(q_0), \dots, p(q_T)$

- Can also do 2) **Decode**: given  $y_0, \dots, y_T$  &  $\theta$  find  $q_0, \dots, q_T$
- We can also decode to find the most likely path  $q_0 \dots q_T$
- Here, we use the ArgMax JTA algorithm
- Run JTA but replace sums with max
- Then, find biggest entry in separators:



## HMMs: EM Learning

- Finally 3) **Max Likelihood**: given  $y_0, \dots, y_T$  learn parameters  $\theta$

•Recall max likelihood:  $\hat{\theta} = \arg \max_{\theta} \log p(\bar{y} | \theta)$

- If observe  $q_t$ , it's easy to maximize the *complete* likelihood:

$$l(\theta) = \log(p(q, y))$$

$$= \log(p(q_0) \prod_{t=1}^T p(q_t | q_{t-1}) \prod_{t=0}^T p(\bar{y}_t | q_t))$$

$$= \log p(q_0) + \sum_{t=1}^T \log p(q_t | q_{t-1}) + \sum_{t=0}^T \log p(\bar{y}_t | q_t)$$

$$= \log \prod_{i=1}^M [\pi_i]^{\eta_i^k} + \sum_{t=1}^T \log \prod_{i=1}^M \prod_{j=1}^N [q_{t-1}^i q_t^j \alpha_{ij}]^{\eta_{t-1}^i \eta_t^j} + \sum_{t=0}^T \log \prod_{i=1}^M \prod_{j=1}^N [\eta_{t-1}^i]^{\eta_t^j}$$

$$= \sum_{i=1}^M q_0^i \log \pi_i + \sum_{t=1}^T \sum_{i=1}^M q_{t-1}^i q_t^j \log \alpha_{ij} + \sum_{t=0}^T \sum_{i=1}^M \sum_{j=1}^N q_t^j \eta_t^j \log \eta_{t-1}^i$$

Introduce Lagrange & take derivatives

$$\sum_{i=1}^M \pi_i = 1 \quad \sum_{j=1}^N \alpha_{ij} = 1 \quad \sum_{i=1}^M \eta_{t-1}^i = 1$$

$$\hat{\pi}_i = q_0^i \quad \hat{\alpha}_{ij} = \frac{\sum_{t=0}^{T-1} q_t^i q_{t+1}^j}{\sum_{i=1}^M \sum_{t=0}^{T-1} q_t^i q_{t+1}^j} \quad \hat{\eta}_{t-1}^i = \frac{\sum_{t=0}^T q_t^i y_t^j}{\sum_{i=1}^M \sum_{t=0}^T q_t^i y_t^j}$$

## HMMs: EM Learning

- But, we don't observe the  $q$ 's, incomplete...

$p(\bar{y} | \theta) = \sum_y p(q, \bar{y} | \theta) = \sum_{q_0} \dots \sum_{q_T} p(q_0) \prod_{i=1}^T p(q_i | q_{i-1}) \prod_{t=0}^T p(\bar{y}_t | q_t)$

- EM**: Max expected complete likelihood given current  $p(q)$

$$E\{l(\theta)\} = E_{p(q_0, \dots, q_T | y)} \{ \log p(q, y) \} = \sum_{q_0} \dots \sum_{q_T} p(q | y) \log p(q, y)$$

$$= E \left\{ \sum_{i=1}^M q_0^i \log \pi_i + \sum_{t=1}^T \sum_{i=1}^M \sum_{j=1}^N q_{t-1}^i q_t^j \log \alpha_{ij} + \sum_{t=0}^T \sum_{i=1}^M \sum_{j=1}^N q_t^j \eta_t^j \log \eta_{t-1}^i \right\}$$

$$= \sum_{i=1}^M E\{q_0^i\} \log \pi_i + \sum_{t=1}^T \sum_{i=1}^M E\{q_{t-1}^i q_t^j\} \log \alpha_{ij} + \sum_{t=0}^T \sum_{i=1}^M \sum_{j=1}^N E\{q_t^j\} \eta_t^j \log \eta_{t-1}^i$$

- M-step is maximizing as before:

$$\hat{\pi}_i = E\{q_0^i\} \quad \hat{\alpha}_{ij} = \frac{\sum_{t=0}^{T-1} E\{q_t^i q_{t+1}^j\}}{\sum_{k=1}^M \sum_{t=0}^{T-1} E\{q_t^i q_{t+1}^k\}} \quad \hat{\eta}_{t-1}^i = \frac{\sum_{t=0}^T E\{q_t^i\} y_t^j}{\sum_{i=1}^N \sum_{t=0}^T E\{q_t^i\} y_t^j}$$

- What are  $E\{x\}$ 's?  $E_{p(x)}\{x\} = \sum_x p(x)x = \sum_x p(x)\delta(x=x) = p(x)$

- Our JTA  $\psi$  &  $\phi$  marginals! (JTA is the E-Step for given  $\theta$ )

$$E\{q_t^i q_{t+1}^j\} = p(q_t = i, q_{t+1} = j | \bar{y}) = \frac{\psi^{**}(q_{t-1}, q_{t+1})}{\sum_y \psi^{**}(q_{t-1}, q_{t+1})} \quad E\{q_t^i\} = p(q_t = i | \bar{y}) = \frac{\phi^{**}(q_{t-1})}{\sum_y \phi^{**}(q_{t-1})}$$