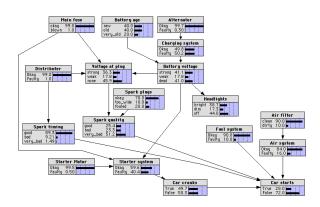
Cmput 466 / 551





# Learning Belief Net Parameters

Readings: ≈HTF 17

(Bayesian Networks without the Tears (Charniak))

# R Greiner University of Alberta

Some material taken from C Guesterin (CMU)



#### Introduce:

- Density estimation
- KL-divergence ... ≈ MLE
- Expectation Maximization
- Gibbs sampling

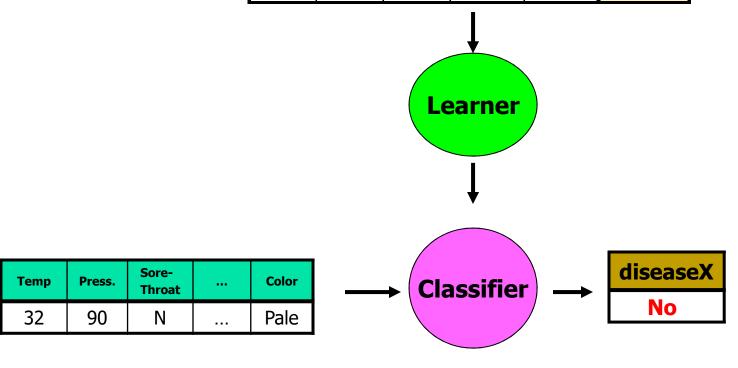
- Motivation
- What is a Belief Net?

<u>Jump</u>

- ...
- Learning a Belief Net
  - Goal?
  - Learning Parameters Complete Data
  - Learning Parameters Incomplete Data
  - Learning Structure

# Learning is ... Training a Classifier

Temp.	Press.	Sore Throat	 Colour	diseaseX
35	95	Y	 Pale	No
22	110	N	 Clear	Yes
:	:		:	:
10	87	N	 Pale	No



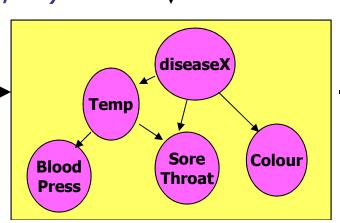
# Learning is ... Training a Model

Temp.	Press.	Sore Throat	 Colour	diseaseX
35	95	Y	 Pale	No
22	110	N	 Clear	Yes
:	:		:	•
10	87	N	 Pale	No

Learner

Then conditionalize, marginalize to answer *any question*:

Temp	Blood Press.	Sore- Throat	 Colour	diseaseX
32	90	N	 Pale	No



J	Н	В	P( j,b,h )
0	0	0	0.03395
0	0	1	0.0095
0	1	0	0.0003
0	1	1	0.1805
1	0	0	0.01455
1	0	1	0.038
1	1	0	0.00045
1	1	1	0.722

# Why Learn Belief Nets?

- Goal#1: Build a classifier
  - What is P(Cancer = + | HA = +, Fev = -, ... ) ?
  - Is P(Cancer = + | ... ) > P(Cancer = | ... )?
- Goal#2: Build a SET of classifiers
  - What is P(Cancer = + | HA = +, Fev = -, ... ) ?
  - What is P(Meningitis = | HA = +, Cold = 3, ...)?
  - What is  $P(HospStay = 3 \mid Smoke = 0.1, BNose = -1, ...)$ ?
- Goal#3: Build a model of the world!
  - ... all interrelations between all subsets of variables
  - Reveal (in)dependencies, connections, ...
  - Note: A completely accurate model will produce correct answers to EVERY P(X | Y ) query

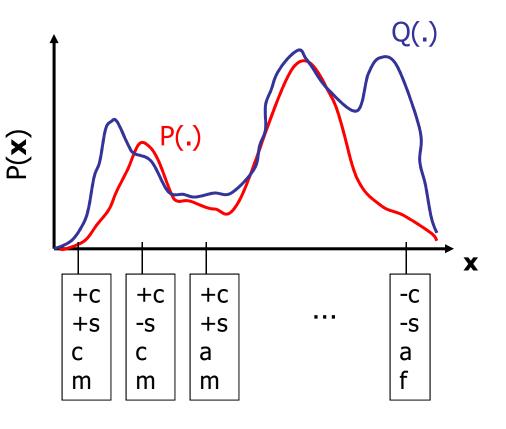




# **Generative Learning**

- Generative Learning:
  - Given (sample of) distribution, P(x)
  - Seek model Q(x)
     that matches P(x)

С	S	Α	•••	G	Р	Q
+	У	У	•••	m	0.3	0.2
+	У	У	•••	f	0.1	0.15
:	:	•		:	i	:
_	У	0	•••	f	0.01	0.02
÷	:	:		:	:	:



Note: no "y" vs "x" ...

# KL-Divergence ... ≈ MaxLikelihood

Seek the BN that minimizes KL-divergence

$$KL(D; BN) = \sum_{x} P_D(x) \ln \frac{P_D(x)}{P_{BN}(x)}$$

- KL-divergence ...
  - always ≥ 0
  - =0 iff distr's "identical"
  - not symmetric
- but... distrib'n *1* not known; Only have instances

$$S = \{d_r\}$$
  
drawn iid from  $\mathcal{D}$ 

 $BN^* = \operatorname{argmin}_{BN} K(\mathfrak{D}; BN)$ 

= argmax<sub>BN</sub>  $\sum_{x} P_{\mathcal{D}}(x) \ln P_{\mathcal{B}}(x)$ 

as  $\sum_{x} P_{D}(x) \ln P_{D}(x)$ 

 $\approx \operatorname{argmax}_{BN} \sum_{d \in S} \frac{1}{|S|} \ln P_{BN}^{\text{is independent of BN}}$ 

as S drawn from D

= argmax<sub>BN</sub>  $\prod_{d \in S} P_{BN}(d)$ 

=  $\operatorname{argmax}_{BN} P_{BN}(S)$ 



#### **Best Distribution**

If goal is BN that approximates 2:

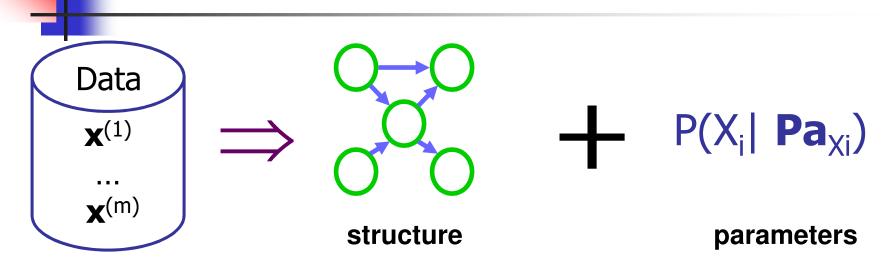
Find BN\* that maximizes likelihood of data S

$$\underset{BN}{\operatorname{arg\,min}} KL(\mathcal{D}; BN) \approx \underset{BN}{\operatorname{arg\,max}} P_{BN}(S)$$

- Approaches:
  - Frequentist: Maximize Likelihood
    - + tweaks to address overfitting: BDe, BIC, MDL, ...
  - Bayesian: Maximize a Posteriori

...

# Learning Belief Nets



# Structure Known Unknown Complete Easy NP-hard Missing Hard ... EM Very hard!!



# Typical (Benign) Assumptions

- 1. -Variables are discrete-
- 2. -Each-case -ε<sub>i</sub>-∈-S--is-complete
- Rows of CPtables are independent

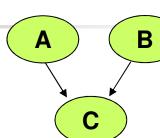
- 4. -Prior-p( $\Theta_{\overline{\chi}}$ -|-G)-is-uniform-- $\theta_{B|+a} \sim \text{Beta}(1,1)$ 

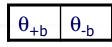
  - Later: relax Assumptions 1, 2, 4



# Learning the CPTs (Frequentist)







 $\theta_{+c|+a,+b}$ 

 $\theta_{+c|-a,+b}$ 

 $\theta_{+c|+\underline{a},\underline{-b}}$ 

 $\theta_{+c|-a,-b|}$ 

 $\theta_{-c|\underline{+a,+b}}$ 

 $\theta_{-c|-a,+b}$ 

 $\theta_{-c|+a,-b}$ 

 $\theta_{-c|-a,-b}$ 

#### Given

- Fixed structure
- over discrete variables { X<sub>i</sub> }
- Complete instances
- $\widehat{\Theta}$  = "empirical frequencies"
- Eg:

$$\theta_{+a} = 2 / (2+2) = 0.5$$

$$\theta_{-b} = 3 / (3+1) = 0.75$$

$$\theta_{+c|+a,-b} = 2 / (2+0) = 1.0$$

	A	В	С
$d_1$	1	0	1
d <sub>2</sub>	0	1	0
$d_3$	0	0	1
d₄	1	0	1

**WHY????** 





# One-Node Bayesian Net



• P(Heads) =  $\theta$ , P(Tails) =  $1-\theta$ 

$$\begin{array}{c|cccc} \hline & P(C=h) & P(C=t) \\ \hline & \theta & 1-\theta \end{array}$$

- Flips are i.i.d.:
  - Independent events
  - Identically distributed according to Binomial distribution
- Sequence S of  $\alpha_H$  Heads and  $\alpha_T$  Tails  $P(S \mid \theta) = \theta^{\alpha_H} (1 \theta)^{\alpha_T}$



#### **Maximum Likelihood Estimation**

- **Data:** Observed set S of  $\alpha_H$  Heads and  $\alpha_T$  Tails
- Hypothesis Space: Binomial distributions
- Learning θ is an optimization problem
  - What's the objective function?
- **MLE**: Choose  $\widehat{\boldsymbol{\theta}}$  that maximizes the probability of observed data:

$$\hat{\theta}$$
 =  $\underset{\theta}{\operatorname{arg max}} P(S | \theta)$   
 =  $\underset{\theta}{\operatorname{arg max}} ln[P(S | \theta)]$ 



# Simple "Learning" Algorithm

$$\hat{\theta}$$
 = arg max  $ln [P(S | \theta)]$   
= arg max  $ln [\theta^h (1 - \theta)^t]$ 

• Set derivative to zero: 
$$\frac{d}{d\theta} \ln P(|\mathcal{S}||\theta) = 0$$

$$\frac{\partial}{\partial \theta} \ln[\theta^h (1 - \theta)^t] = \frac{\partial}{\partial \theta} [h \ln \theta + t \ln (1 - \theta)^t] = \frac{h}{\theta} + \frac{-t}{(1 - \theta)}$$

$$\frac{h}{\theta} + \frac{-t}{(1-\theta)} = 0 \Rightarrow \theta = \frac{h}{t+h}$$
 so just average!!!

If 7 heads, 3 tails, set  $\hat{\theta} = 0.7$ 



### Factoids...

Recall that, for a Bayesian Network...
 For a COMPLETE instance, x = (x<sub>1</sub>, ..., x<sub>n</sub>)
 P(x) = product of CPtable values

 (one from each variable)

- In  $a^b = b \ln a$
- In  $(a \times b) = \ln a + \ln b$

$$\frac{\partial}{\partial \theta} \ln \theta = \frac{1}{\theta}$$

$$\frac{\partial}{\partial \theta} \ln \theta = \frac{1}{\theta}$$

$$\frac{\partial}{\partial \theta} \ln (1 - \theta) = \frac{-1}{(1 - \theta)}$$



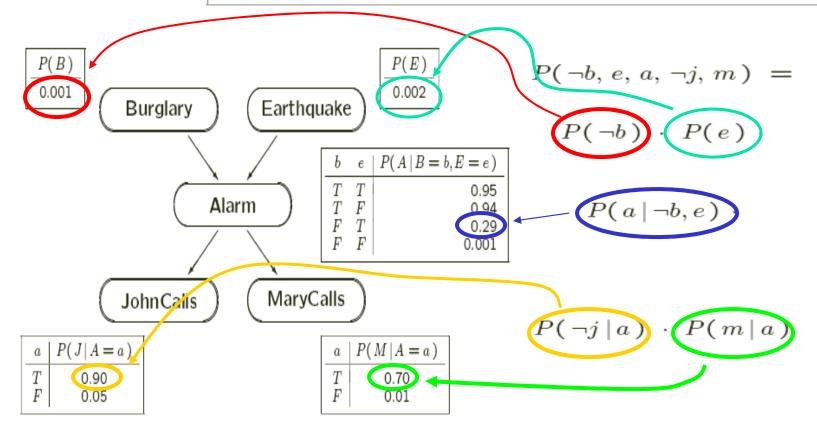
# Probability of Complete Instance

$$P(\neg b, e, a, \neg j, m) = P(\neg b) P(e|\neg b) P(a|e, \neg b) P(\neg j|a, e, \neg b) P(m|\neg j, a, e, \neg b)$$

$$P(\neg b) P(e) P(a|e, \neg b) P(\neg j|a) P(m|a)$$

$$0.99 \times 0.02 \times 0.29 \times 0.1 \times 0.70$$

Node independent of predecessors, given parents





# Likelihood of the Data (Frequentist)

- $\theta_{+a}$   $\theta_{-a}$
- A B
- $\theta_{+b}$   $\theta_{-b}$

- P(S |  $\Theta$ ) =  $\prod_r$  P(d<sub>r</sub> |  $\Theta$ )
- $P(d_1) = P_{\Theta}(+a, -b, +c)$ =  $P_{\Theta}(+a) P_{\Theta}(-b) P_{\Theta}(+c \mid +a, -b)$ =  $\Theta_{+a} \Theta_{-b} \Theta_{+c \mid +a, -b}$
- $P(d_2) = P_{\Theta}(-a, +b, -c)$ =  $P_{\Theta}(-a) P_{\Theta}(+b) P_{\Theta}(-c \mid -a, +b)$ =  $\Theta_{-a} \Theta_{+b} \Theta_{-c \mid -a, +b}$

$\theta_{+c +a,+b}$	$\theta_{-c +a,+b}$
$\theta_{+c -a,+b}$	$\theta_{-c -a,+b}$
$\theta_{+c +a,-b}$	$\theta_{-c +a,-b}$
$\theta_{+c -a,-b}$	$\theta_{-c -a,-b}$

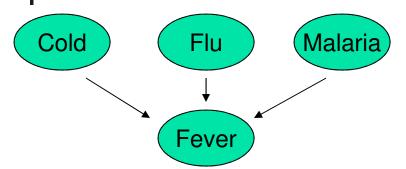
	A	В	С
$d_1$	1	0	1
d <sub>2</sub>	0	1	0
$d_3$	0	0	1
d <sub>4</sub>	1	0	1

$$\begin{array}{l} \bullet \ \mathsf{P}(\ \mathsf{S} \ | \ \Theta \ ) = \Theta_{+a}^{\phantom{+}2} \ \Theta_{-a}^{\phantom{-}2} \ \Theta_{+b}^{\phantom{-}1} \ \Theta_{-b}^{\phantom{-}3} \ \Theta_{+c|+a,+b}^{\phantom{-}0} \ \Theta_{+c|+a,+b}^{\phantom{-}0} \ \Theta_{+c|+a,-b}^{\phantom{-}2} \dots \\ = \Theta_{+a}^{\phantom{-}N_{+a}} \ \Theta_{-a}^{\phantom{-}N_{-a}} \ \Theta_{+b}^{\phantom{-}N_{+b}} \ \Theta_{-b}^{\phantom{-}N_{-b}} \ \Theta_{+c|+a,+b}^{\phantom{-}N_{+c|+a,+b}} \ \Theta_{+c|+a,-b}^{\phantom{-}0} \dots \\ = \prod_{ijk} \theta_{iik}^{N_{ijk}} \end{array}$$



# Example of Parameter θ<sub>ijk</sub>





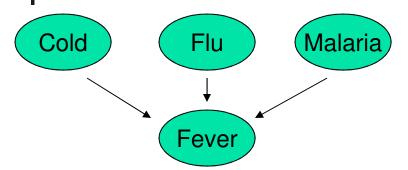
4th	$\Rightarrow$
•	

				P(Fever=?	Ca, Fiu, Maii)
	Cold	Flu	Malaria	True	False
	F	F	F	$ heta_{111}$	$ heta_{112}$
	F	F	Т	$ heta_{121}$	$ heta_{122}$
	F	Т	F	$\theta_{131}$	<i>6</i> <sub>132</sub>
-	F	T	I	$ heta_{141}$	▶ θ <sub>142</sub>
	Т	F	F	$ heta_{151}$	$\theta_{152}$
	Т	F	Т	$ heta_{161}$	$\theta_{162}$
	Т	Т	F	$ heta_{171}$	$\theta_{172}$
	Т	T	Т	$ heta_{181}$	$\theta_{182}$

- $\bullet \Theta_{ijk} = P(X_i = V_{ik} \mid Pa_i = pa_{ij})$ 
  - variable#1 -- here, "Fever"
  - 4th value of parents [Cold=F, Flu=T, Malaria=T]
  - 2nd value of Fever-node here, "Fever = FALSE"
- Note:  $\sum_{k} \Theta_{ijk} = 1$

# Example of Count N<sub>ijk</sub>





<b>4</b> th	$\Rightarrow$
-------------	---------------

			P(Fever=?	Cu, Fiu, Maii)
Cold	Flu	Malaria	True	False
F	F	F	N <sub>111</sub>	N <sub>112</sub>
F	F	Т	N <sub>121</sub>	N <sub>122</sub>
F	Т	F	N <sub>131</sub>	N <sub>132</sub>
F	T	T	N <sub>141</sub>	▶ N <sub>142</sub>
Т	F	F	N <sub>151</sub>	N <sub>152</sub>
Т	F	Т	N <sub>161</sub>	N <sub>162</sub>
Т	Т	F	N <sub>171</sub>	N <sub>172</sub>
Т	Т	Т	N <sub>181</sub>	N <sub>182</sub>

- N<sub>iik</sub> refers to ...
  - variable#1 -- here, "Fever"
  - 4<sup>th</sup> value of parents [Cold=F, Flu=T, Malaria=T]
  - 2nd value of Fever-node -- here, "Fever = FALSE"
- N<sub>ijk</sub> is number of data-tuples
   where variable#i = its k<sup>th</sup> value
   & parents(variable#i) = j<sup>th</sup> value

# Task#1:

# Fixed Structure, Complete Tuples

■ What are the ML values for Θ, given iid data  $S = \{ d_r \}, ...$ 

$$P(S \mid \Theta) = \prod_{d \in S} P(d \mid \Theta) = \prod_{d \in S} \prod_{[X_i = x_{ik}, Pa_i = pa_{ij}] \in d} \Theta_{ijk} =$$

$$\prod_{ijk} \Theta_{ijk}^{N_{ijk}} = \prod_{ij} \prod_{k} \Theta_{ijk}^{N_{ijk}}$$

- $\Theta^{(ML)}$  = argmax<sub> $\Theta$ </sub> { P(S |  $\Theta$  ) }

  - =  $\operatorname{argmax}_{\Theta} \{ \log P(S \mid \Theta) \}$ =  $\operatorname{argmax}_{\Theta} \{ \sum_{ij} \sum_{k} N_{ijk} \log \Theta_{ijk} ) \}$

 $\forall ij \sum_{k} \Theta_{iik} = 1$ 

#### **MLE Values**

- $\Theta^{(ML)} = \underset{\forall ij \ \sum_{k} \Theta_{ijk} = 1}{\operatorname{argmax}_{\Theta}} \left\{ \sum_{ij} \sum_{k} N_{ijk} \log \Theta_{ijk} \right\}$
- Notice  $\theta_{ij}$  is independent of  $\theta_{rs}$  when  $i \neq r$  or  $j \neq s$  ...  $\Rightarrow$  can solve each  $\sum_k N_{ijk} \log \theta_{ijk}$  individually!
- For each  $\sum_{k} N_{ijk} \log \theta_{ijk}$  ... as  $\sum_{k} \theta_{ijk} = 1$ , optimum is

$$\theta_{ijk} = \frac{N_{ijk}}{\sum_{k'} N_{ijk'}} = \frac{\#(X_i = v_{i,k} \& \mathbf{Pa}_i = \mathbf{pa}_{i,j})}{\#(\mathbf{Pa}_i = \mathbf{pa}_{i,j})}$$

- Observed Frequency Estimates!
- Undefined if  $\sum_k N_{ijk} = 0 \dots \#(\mathbf{Pa}_i = \mathbf{pa}_{i,j}) = 0$

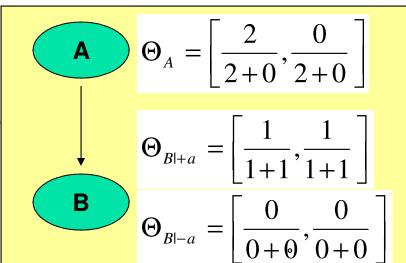
# •

# **Algorithm**

- ComputeMLE( graph G, data S): return MLE parameters  $[\theta_{ijk}]$
- Initialize N<sub>ijk</sub> ← 0
- Walk thru data \$
  - Whenever see [ X<sub>i</sub>=v<sub>ik</sub>, Pa<sub>i</sub>=pa<sub>ij</sub>], increment N<sub>ijk</sub> += 1
- Return parameters:  $\left|\theta_{ijk}\right| = \overline{\Sigma}$

$$\theta_{ijk} = \frac{N_{ijk}}{\sum_{r} N_{ijr}}$$





#### Buckets

$$N_{+a} = 0$$

$$N_{-a} = 0$$

$$N_{-a} = 0$$
 $N_{+b|+a} = 0$ 

$$N_{-b|+a} = 0$$

$$N_{+b|-a} = 0$$

$$N_{-b|-a} = 0$$

A	В	
+	+	
+	_	



0





### Problems with MLE

- 0/0 issues
- Do you really believe 0% if 0 / 0+2 ?
- Which is better?
  - 3 heads, 2 tails
  - 30 heads, 20 tails
  - 3E23 heads, 2E23 tails

- $\theta = 3/(3+2) = 0.6$
- $\theta = 30/(30+20) = 0.6$
- $\theta = 3E23/(3E3+2E23) = 0.6$
- What if you already know SOMETHING about the variable...

≈ 50/50 ...





# Bayesian Learning

$$P(\theta | S) \propto P(S | \theta) P(\theta)$$
 $\uparrow \qquad \uparrow \qquad \uparrow$ 
posterior

 $P(S | \theta) P(\theta) \rightarrow \uparrow \qquad \uparrow \qquad \uparrow$ 
likelihood prior

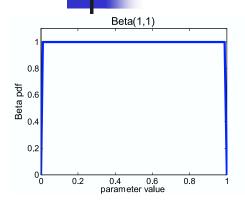
Likelihood function is simply Binomial:

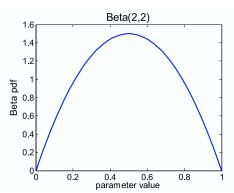
$$P(S \mid \theta) = \theta^{m_H} (1 - \theta)^{m_T}$$

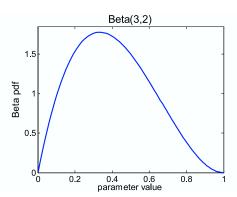
- What about prior?
  - Represent expert knowledge
  - Simple posterior form
- Conjugate priors:
  - Closed-form representation of posterior (more details soon)
  - For Binomial, conjugate prior is Beta distribution<sup>8</sup>

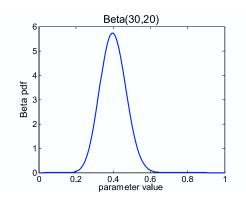


# Beta Prior Distribution – $P(\theta)$









• Prior: 
$$P(\theta) = \frac{\theta^{\alpha_H - 1} (1 - \theta)^{\alpha_T - 1}}{B(\alpha_H, \alpha_T)} \sim Beta(\alpha_H, \alpha_T)$$

• Likelihood function: 
$$P(\mathcal{D} \mid \theta) = \theta^{m_H} (1 - \theta)^{m_T}$$

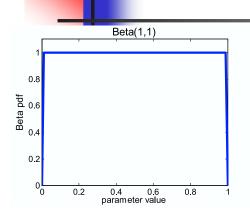
- Given X ~ Beta(a, b) :
  - Mean: a/(a + b)
  - Unimodal if a,b>1... here mode:  $\frac{a-1}{a-1}$

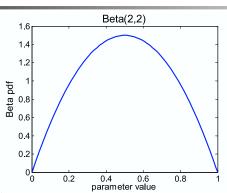
$$\frac{a-1}{a+b-2}$$

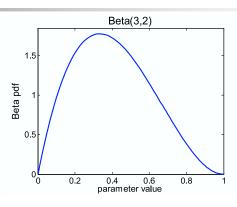
• Variance: 
$$\frac{a b}{(a+b)^2(a+b-1)}$$

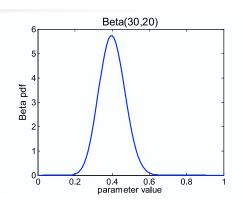
# Posterior distribution... from Beresti











$$P(\theta \mid \mathcal{D}) \propto P(\theta) P(\mathcal{D} \mid \theta)$$

Prior 
$$P(\theta)$$

Likelihood 
$$P(D|\theta)$$

$$= \Theta^{\alpha_H - 1} (1 - \Theta)^{\alpha_T - 1} \times \Theta^{m_H} (1 - \Theta)^{m_T}$$

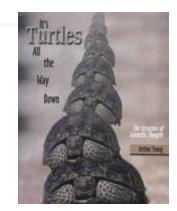
$$imes oxedot \Theta^{m_H} (1-\Theta)^{m_T}$$

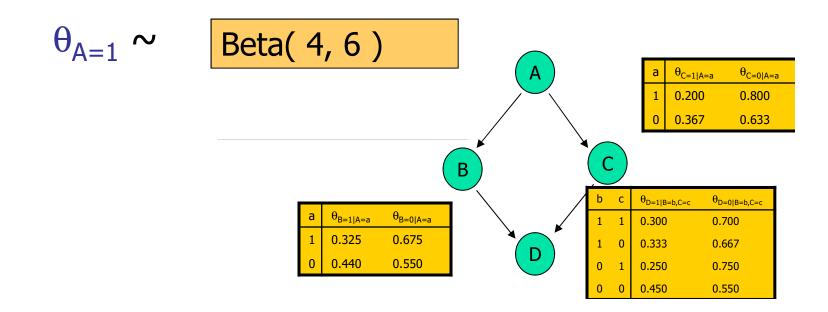
$$= \Theta^{\alpha_H + m_H - 1} (1 - \Theta)^{\alpha_T + m_T - 1}$$

$$\sim$$
 Beta $(\alpha_H + m_H, \alpha_T + m_T)$ 

## Distribution over Parameter

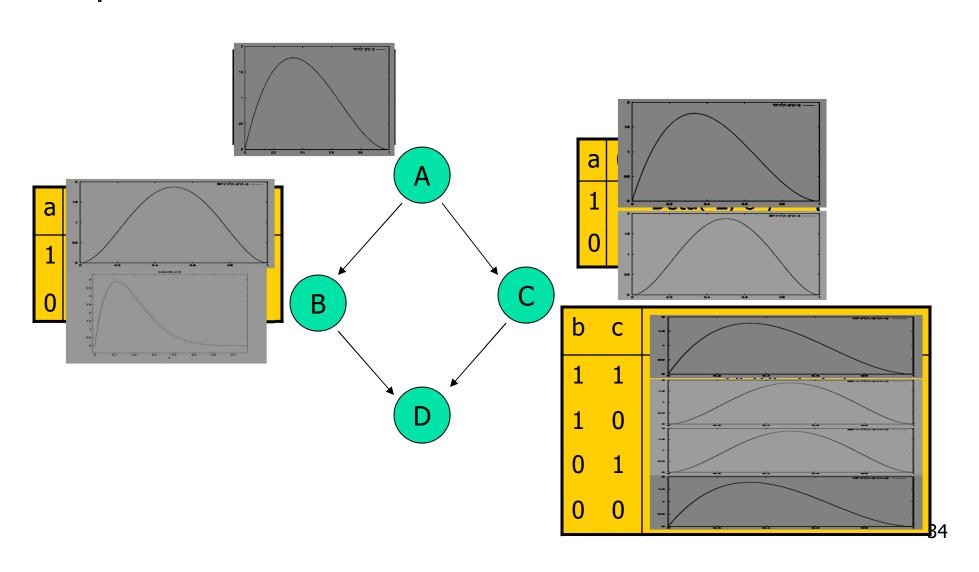
- What is "real" value of  $\theta_{A=1}$ ?
- If ...
  - uncertainty in expert opinion
  - limited training data only a distribution!



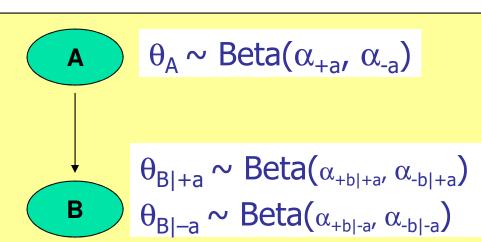




# Distribution over Parameters







#### Buckets

$$u_{+a} := \alpha_{+a}$$

• 
$$u_{-a}$$
 :=  $\alpha_{-a}$ 

$$u_{+b|+a} := \alpha_{+b|+a}$$

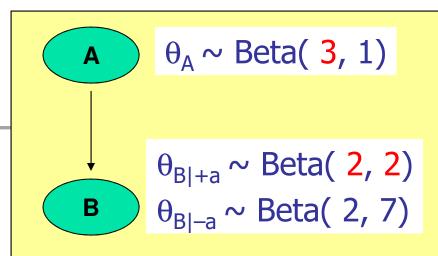
$$u_{-b|+a} := \alpha_{-b|+a}$$

$$u_{+b|-a} := \alpha_{+b|-a}$$

$$u_{-b|-a} := \alpha_{-b|-a}$$

A	В	
+	+	
+		





#### Buckets

$$u_{+a} := 1$$

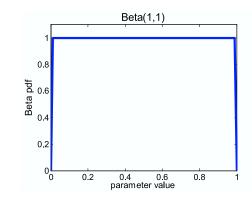
$$u_{+b|+a} := 1$$
 $u_{-b|+a} := 1$ 

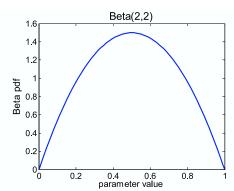
$$u_{-b|+a} := 1$$

$$u_{+bl-a} := 2$$

• 
$$u_{+b|-a} := 2$$
  
•  $u_{-b|-a} := 7$ 

A	В
+	+
+	







If you want POINT estimates...



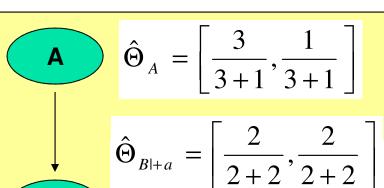
$$\mathbf{u}_{+a}$$
 :=  $\mathbf{1}'$ 

• 
$$u_{+b|+a} := 1$$
/
•  $u_{-b|+a} := 1$ /

$$u_{-bl+a} := 1$$

$$u_{+b|-a} := 2$$

$$u_{-b|-a} := 7$$



 $\hat{\Theta}_{B|-a} = \left| \frac{2}{2+7}, \frac{7}{2 + 7} \right|$ 

Α	В	
+	+	
+	_	

Note: no 0/0 issues!

# 4

#### **Beta Distribution**

Model row-parameter

$$\theta_{B|a=1} = \langle \theta_{b=0|a=1}, \theta_{b=1|a=1} \rangle$$

as Beta distribution

$$\bullet_{\mathsf{B}|\mathsf{A}=1} = \langle \theta_{\mathsf{B}=0|\mathsf{A}=1}, \ \theta_{\mathsf{B}=1|\mathsf{A}=1} \rangle \sim \mathsf{Beta}(\ 1,\ 1\ )$$

kinda like seeing 2 instances with  $\langle A=1 \rangle$ :

1	with	$\langle A=1,$	B=0	
1	with	$\langle A=1,$	$B=1\rangle$	

A	В	С	D
1	0	0	1
1	1	1	1
0	0	1	1
:	:	:	:



# Beta Distribution, II

 $\bullet_{B|A=1} = \langle \theta_{B=0|A=1}, \theta_{B=1|A=1} \rangle \sim \text{Beta}(1, 1)$ 

$$\Rightarrow \left| \mathbf{E}[\theta_{B=0|A=1}] = \widehat{\boldsymbol{\theta}}_{b|+a} \right| = \frac{1}{1+1} = 0.5$$

Now... observe data 5:

$$\begin{cases}
A & B & C & E \\
1 & 1 & 0 & 1 \\
1 & 1 & 1 & 1 \\
1 & 0 & 1 & 0 \\
1 & 0 & 1 & 0 \\
1 & 0 & 0 & 0 \\
1 & 0 & 0 & 1 \\
0 & 0 & 0 & 1 \\
\vdots & \vdots & \vdots & \vdots
\end{cases}$$

$$\begin{cases}
A & B & C & E \\
1 & 1 & 0 & 1 \\
1 & 1 & 1 & 1 \\
1 & 0 & 1 & 0 \\
1 & 0 & 0 & 0 \\
1 & 0 & 0 & 1 \\
\vdots & \vdots & \vdots & \vdots
\end{cases}$$

$$\begin{cases}
A & B & C & E \\
1 & 1 & 0 & 1 \\
1 & 0 & 1 & 0 \\
1 & 0 & 1 & 0 \\
0 & 0 & 0 & 1 \\
\vdots & \vdots & \vdots & \vdots
\end{cases}$$

$$\begin{cases}
A & \text{"}(A=1, B=0) \text{ "S} \\
A & \text{"}(A=1, B=0) \text{ "S} \\
A & \text{"}(A=1, B=0) \text{ "S}
\end{cases}$$



# Beta Distribution, III

$$\bullet_{\mathsf{B}|\mathsf{A}=1} = \langle \theta_{\mathsf{B}=0|\mathsf{A}=1}, \; \theta_{\mathsf{B}=1|\mathsf{A}=1} \rangle \sim \mathsf{Beta}(1,1)$$

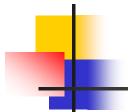
$$\Rightarrow \left| \mathbf{E}[\theta_{B=1|A=1}] = \widehat{\theta}_{+b|+a} \right| = \frac{1}{1+1} = 0.5$$

Then observe data S

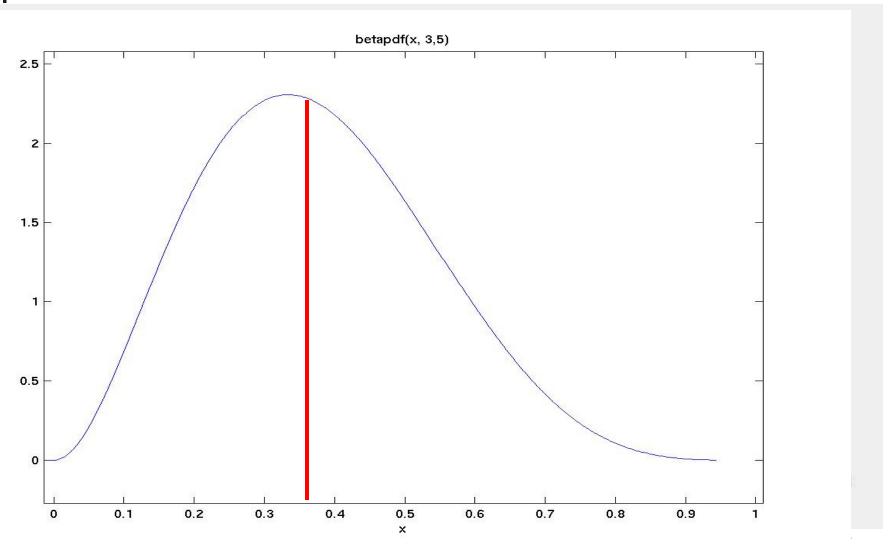
New distribution is

$$\theta'_{B|A=1} \sim Beta(1+2, 1+4) = Beta(3, 5)$$

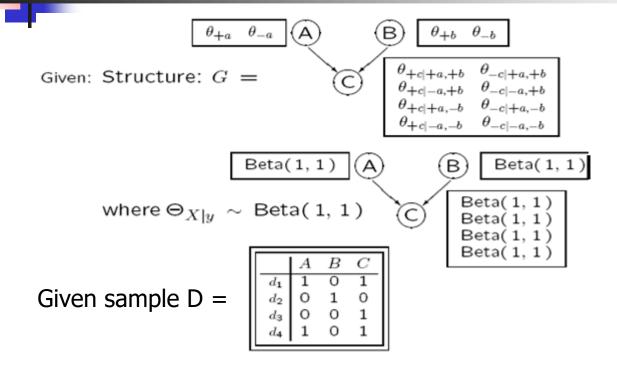
$$\Rightarrow E[\theta_{B=1|A=1} \mid S] = \hat{\theta}_{+b|+a} \mid S = \frac{3}{3+5} = 0.375$$



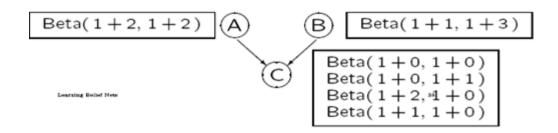
# $\theta_{B|A=1} \sim Beta(3,5)$ Distribution



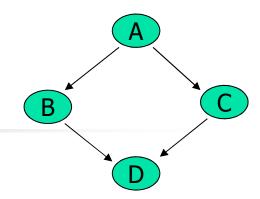
## Posterior Distribution of ⊕



Posterior distribution is...



## **Posterior Distribution**



- Initially:  $P(X_i | p_{ij}) \dots$  $\theta_{ij} \sim Dir(\alpha_{ij1}, \dots, \alpha_{ijr})$
- Data S includes
   N<sub>ijk</sub> examples including [ X<sub>i</sub>=v<sub>ik</sub>, Pa<sub>i</sub>=pa<sub>ij</sub>]
- Posterior  $\theta_{ij} \mid S \sim Dir(\alpha_{ij1} + N_{ij1}, ..., \alpha_{ijr} + N_{ijr})$
- Expected value

$$E[\theta_{ijk}] = \frac{N_{ijk} + \alpha_{ijk}}{\sum_{r} N_{ijr} + \alpha_{ijr}}$$

Compare to Frequentist:

$$\hat{\theta}_{ijk} = \frac{N_{ijk}}{\sum_{r} N_{ijr}}$$

# Algorithm

### ComputePosterior (graph G, data S, priors $[\alpha_{iik}]$ ): return posterior parameters [uiik]

- Initialize  $u_{ijk} \leftarrow \alpha_{ijk}$
- Walk thru data S
  - Whenever see [ X<sub>i</sub>=v<sub>ik</sub>, Pa<sub>i</sub>=pa<sub>ii</sub>], increment u<sub>iik</sub> += 1
- Set parameters:

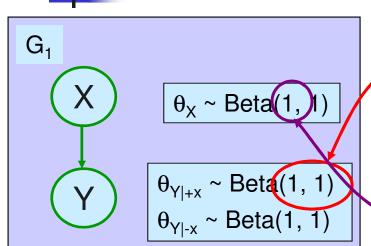
$$\theta_{ij}$$
 |S ~ Dir(  $u_{ij1}$ , ...,  $u_{ijr}$ )

If want expected value:  $E[\theta_{ijk}] = \frac{u_{ijk}}{\sum u_{...}}$ 

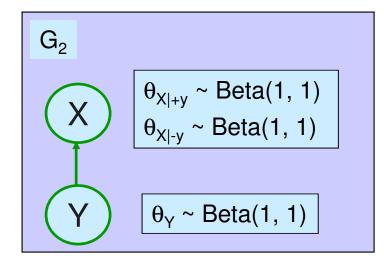
$$E[\theta_{ijk}] = \frac{u_{ijk}}{\sum_{r} u_{ijr}}$$



### **Priors for Parameters**



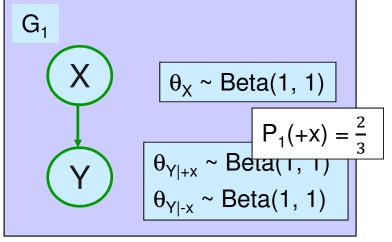
- Does this make sense?
  - EffectiveSampleSize( $\theta_{Y|+x}$ ) = 2
  - But only 1 example ~ "+x" ??



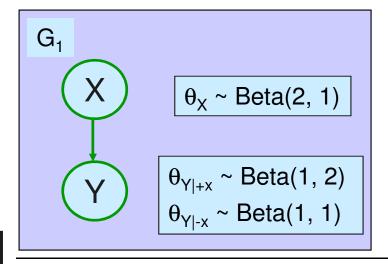
- J-Equivalent structure
- What happens after [+x, -y]?
  - Should be the same!!

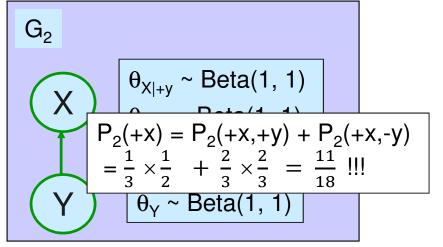


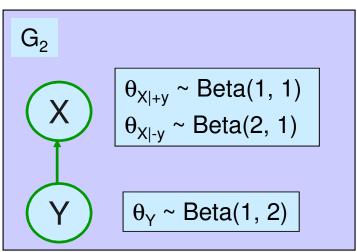
### **Priors for Parameters**



[+x, -y]

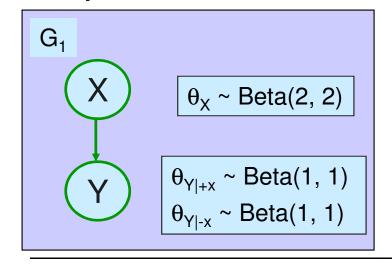




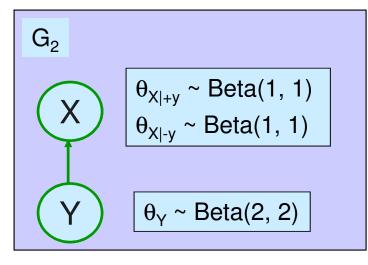




### **BDe Priors**



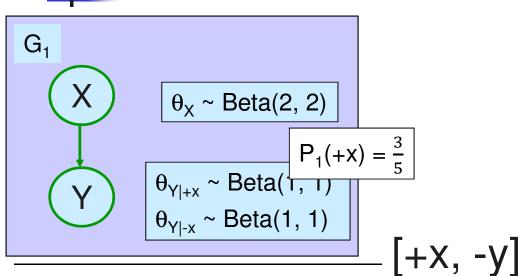
- This makes more sense:
  - EffectiveSampleSize( $\theta_{Y|+x}$ ) = 2
  - Now  $\approx$  2 examples  $\sim$  "+x"??

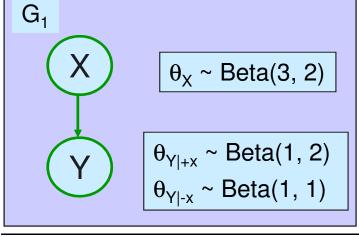


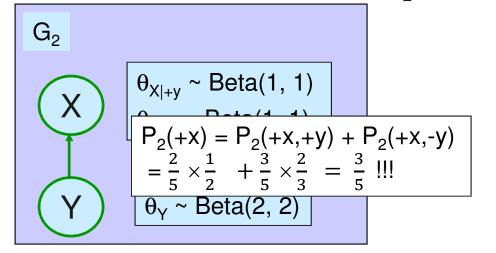
- J-Equivalent structure
- Now what happens after [+x, -y]?

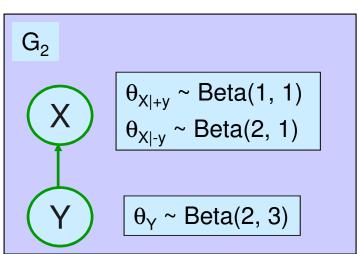


### **BDe Priors**









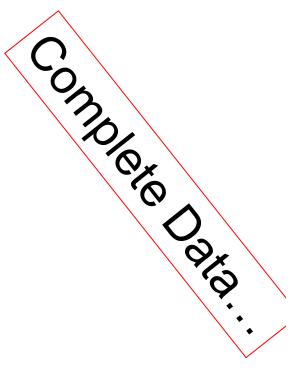
# BDe Prior

- View Dirichlet parameters as "fictitious samples"
  - equivalent sample size
- Pick a fictitious sample size m'
- For each possible family, define a prior distribution  $P(X_i, Pa_{X_i})$ 
  - Represent with a BN
  - Usually independent (product of marginals)
    - $P(X_i, Pa_{X_i}) = P'(X_i) \prod_{x_{j \in Pa[X_i]}} P'(x_j)$
    - P(  $\theta[x_i | Pa_{X_i} = u) = Dir( m' P'(x_i = 1, Pa_{X_i} = u), ..., m' P'(x_i = k, Pa_{X_i} = u))$
    - Typically, P'(X<sub>i</sub>) = uniform



- MLE:
  - score decomposes according to CPTs
  - optimize each CPT separately
- Bayesian parameter learning:
  - motivation for Bayesian approach
  - Bayesian prediction
- Bayesian learning for BN parameters
  - Global parameter independence
  - BDe if and only if score equivalence

  - Predictive distribution model averaging, for free!



# Outline

- Motivation
- What is a Belief Net?
- Learning a Belief Net
  - Goal?
  - Learning Parameters Complete Data
  - Learning Parameters Incomplete Data
    - Gradient Descent
    - EM
    - Gibbs
  - Learning Structure

# 4

## #2: Known structure, Missing data

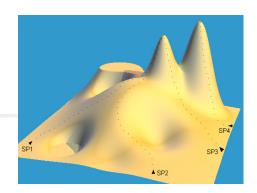
- To find good  $\Theta$ , need to compute  $P(\Theta \mid S, G)$
- Easy if ...Hard

$$S = \left\{ \begin{array}{cccc} c_1 \colon & \left\langle \begin{smallmatrix} * & & \dots & c_{1N} \\ c_2 \colon & \left\langle c_{21} & \dots & \begin{smallmatrix} * \\ * \\ \vdots & \left\langle \vdots & c_{ij} & \vdots \\ c_m \colon & \left\langle c_{m1} & \dots & c_{mN} \right\rangle \end{array} \right\} \text{ incomplete}$$

- What if S is incomplete ?
  - Some c<sub>ij</sub> = \*
  - $c_{iK} = * \forall i$  (ie,  $X_K$  never seen... "Hidden variables")
- Here:
  - Given fixed structure
  - Missing (Completely) At Random:
     Omission not correlated with value, etc.
- Approaches:
  - Gradient Ascent, EM, Gibbs sampling, ...



## **Gradient Ascent**



- Want to maximize likelihood
  - $\theta^{(MLE)} = \operatorname{argmax}_{\theta} L(\theta : S)$
- Unfortunately...
  - $L(\theta : S)$  is nasty, non-linear, multimodal fn
  - So...
- Gradient-Ascent
- ... 1<sup>st</sup>-order Taylor series expansion...

$$f_{\rm obj}(\theta) \approx f_{\rm obj}(\theta^0) + (\theta - \theta^0) \nabla f_{\rm obj}(\theta^0)$$

Need derivative!

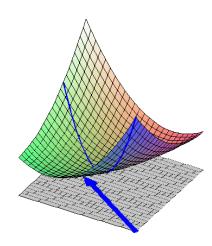
```
Procedure Gradient-Ascent ( \theta^1, // Initial starting point f_{\text{obj}}, // Function to be optimized \delta // Convergence threshold ) 1 \quad t \leftarrow 1 2 do 3 \quad \theta^{t+1} \leftarrow \theta^t + \sqrt{\nabla f_{\text{obj}}(\theta^t)} 4 t \leftarrow t+1 5 while \|\theta^t - \theta^{t-1}\| > \delta 6 return (\theta^t)
```



### **Issues with Gradient Ascent**

- Constraints
  - $\Theta_{iik} \in [0,1]$
  - $\sum_{\mathsf{r}} \Theta_{\mathsf{iir}} = 1$
  - But ...  $\Theta_{iik}$  +=  $\alpha \Delta \Theta_{iik}$  could violate constraints
  - Use  $\lambda_{ijk} = \log(\theta_{ijk})$ 
    - Find best  $\lambda_{ijk}$  ... ignore constraints ...
- Lots of Tricks for efficient ascent
  - Line Search
  - Conjugate Gradient
  - **...**

[See earlier notes on optimization]

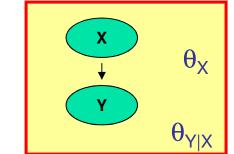


# **Expectation Maximization (EM)**

- EM is designed to find most likely θ, given incomplete data!
- Recall simple Maximization needs counts:

$$\#(+x, +y)$$
, ... for  $N_{+y|+x}$ , ...

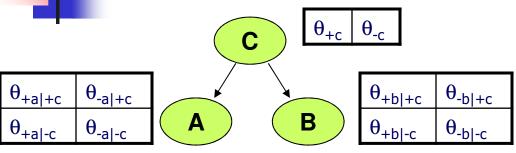
But is instance [?, +y] in ... #(+x, +y)? ... #(-x, +y)?



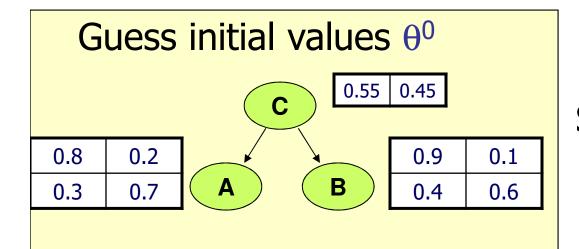
- Why not put it in BOTH... fractionally ?
  - What is weight of #(+x, +y)?
  - $P_{\theta}(+x + y)$ , based on current value of  $\theta$
- Compute "expected sufficient statistics":  $E_{\theta}[N_{ijk}]$

# 4

# EM Approach – E Step



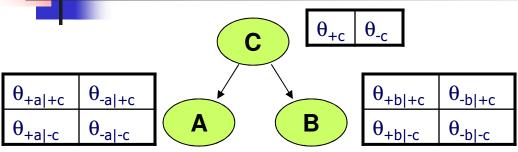
Camarala	Α	В	C
Sample S =	0(	0	1
	*	1	0
	9	*	1
	*	*	1



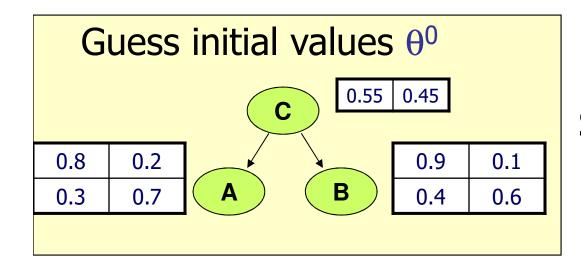
Set S(0) = 
$$\begin{bmatrix} A & B & C \\ 0 & 0 & 1 & 1.0 \\ 0 & 1 & 0 & 0.7 \\ 1 & 0 & 0.3 \\ 0 & 0 & 1 & 0.1 \\ 0 & 1 & 1 & 0.9 \\ \hline 0 & 0 & 1 & 0.2 \times 0.1 \\ 0 & 1 & 1 & 0.8 \times 0.1 \\ 1 & 1 & 1 & 0.8 \times 0.9 \end{bmatrix}$$

# 4

## EM Approach – E Step



Carranala	Α	В	С
Sample S =	0	0	1
	*	1	0
	0	*	1
	*	*	1



Set S(0) = 
$$\begin{vmatrix} A & B & C \\ 0 & 0 & 1 & 1.0 \\ 0 & 1 & 0 & 0.7 \\ 1 & 1 & 0 & 0.3 \\ 0 & 0 & 1 & 0.1 \\ 0 & 1 & 1 & 0.9 \\ 0 & 0 & 1 & 0.2 \times 0.1 \\ 0 & 1 & 1 & 0.2 \times 0.9 \\ + (0.8 \times 0.1) & 1 & 1 & 0.8 \times 0.9$$

$$E_{\theta^0}[N_{+b|+c}] = 0.9 + (0.2x0.9) + (0.8x0.9)$$
  
 $E_{\theta^0}[N_{-b|+c}] = 1 + 0.1 + (0.2x0.1) + (0.8x0.1)$ 



## EM Approach – M Step

Use fractional data:

$$S^{(0)} =$$

Α	В	С	
0	0	1	1.0
0	1	0	0.7
1	1	0	0.3
0	0	1	0.1
0	1	1	0.9
0	0	1	0.7 × 0.1
0	1	1	0.7 × 0.9
1	0	1	8.3 × 0.1
1	1	1	0.3 × 0.9
•			

	1	1.0	$\theta_{+a +c}$	$\theta_{-a +c}$	
	0	0.7			A
•	0	0.3	$\theta_{+a -c}$	$\theta_{-a -c}$	
	1	0.1			
• •	1	0.9			
	1	0.7 × 0.1			
	1	0.7 × 0.9			
	1	8.3 × 0.1			
		<u></u>			

New estimates:

$$E_{\Theta}[N_{+a|+c}]$$

$$(0.3 \times 0.1) + (0.3 \times 0.9)$$

$$= \frac{1 + E_{\Theta}[N_{-a|+c}]}{[(0.3 \times 0.1) + (0.3 \times 0.9)] + [1 + (0.1 + 0.9) + (0.7 \times 0.1) + (0.7 \times 0.9)]} = \frac{1 + E_{\Theta}[N_{-a|+c}]}{[(0.3 \times 0.1) + (0.3 \times 0.9)] + [1 + (0.1 + 0.9) + (0.7 \times 0.1) + (0.7 \times 0.9)]}$$

$$\hat{\theta}_{+c}^{(1)} = \frac{E_{\theta}[N_{-c}]}{E_{\theta}[N_{+c}] + E_{\theta}[N_{-c}]} = \frac{1.0 + (1.0) + (1.0)}{4} = 0.75$$

$$\hat{\theta}_{+b|+c}^{(1)} = \frac{E_{\theta}[N_{+b|+c}]}{E_{\theta}[N_{+b|+c}] + E_{\theta}[N_{-b|+c}]} = \frac{0.9 + (0.2 \times 0.9) + (0.8 \times 0.9)}{3} = 0.6$$

 $\theta_{+b|+c}$ 

 $\theta_{\text{+b|-c}}$ 

 $\theta_{-b|+c}$ 

 $\theta_{-b|-c}$ 

# EM Approach – M Step

Use fractional data:

$$S^{(0)} =$$

Α	В	С	
0	0	1	1.0
0	1	0	0.7
1	1	0	0.3
0	0	1	0.1
0	1	1	0.9
0	0	1	0.7 × 0.1
0	1	1	0.7 × 0.9
1	0	1	8.3 × 0.1
1	1	1	0.3 × 0.9

Α	В	C		
0	0	1	1.0	$\theta_{+a +c}$ $\theta_{-a +c}$
0	1	0	0.7	$\mathbf{A}$
1	1	0	0.3	$\theta_{+a -c}$ $\theta_{-a -c}$
0	0	1	0.1	
0	1	1	0.9	
0	0	1	0.7 × 0.1	
0	1	1	0.7 × 0.9	
1	0	1	0.3 × 0.1	
1	1	1	0.3 × 0.9	

•New estimates:

$$\hat{\theta}_{+a|+c}^{(1)}$$
 =

$$E_{\Theta}[N_{+a|+c}]$$

$$[E_{-a|+c}] + E_{\Theta}[N_{-a|+c}]$$

$$(0.3 \times 0.1) + (0.3 \times 0.9)$$

$$\frac{}{[(0.3\times0.1)+(0.3\times0.9)]+[1+(0.1+0.9)+(0.7\times0.1)+(0.7\times0.9)]}=0$$

 $\theta_{+b|+c}$ 

 $\theta_{+b|\underline{-c}}$ 

 $\theta_{-b|+c}$ 

 $\theta_{\text{-b|-c}}$ 

$$\hat{\theta}_{+c}^{(1)} =$$

$$\frac{L_{\theta}[I \setminus c]}{[N] + E_{\alpha}[N]}$$

$$=\frac{1.0+(1.0)+(1.0)}{1.0}=0.75$$

$$E_{\theta}[N_{+c}] + E_{\theta}[N_{-c}]$$

$$E_{\theta}[N_{+c}] + E_{\theta}[N_{-c}]$$

$$\frac{E_{\theta}[N_{+b|+c}]}{E_{\theta}[N_{+b|+c}] + E_{\theta}[N_{-b|+c}]} =$$

- **E-step**: estimate expected sufficient statistics (wrt missing values) using current  $\theta^{(t)}$  values
- **M-step**: compute new  $\theta^{(t+1)}$  values, using these expected sufficient statistics



## **EM Steps**

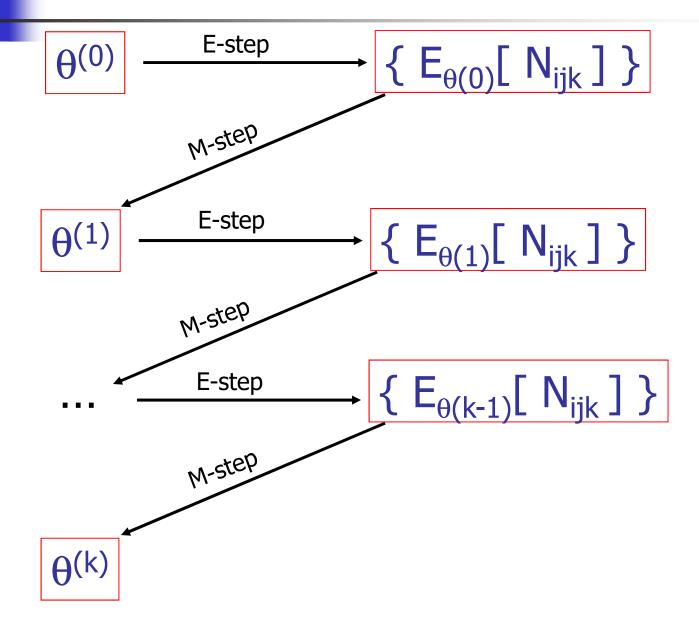
### E step:

- Given parameters  $\theta^{(t)}$
- find probability of each missing value
  - ... so get  $E_{\theta(t)}[N_{ijk}]$

### M step:

- Given completed (fractional) data
  - based on  $E_{\theta(t)}[N_{ijk}]$
- find max-likely parameters  $\theta^{(t+1)}$

## **EM Process**



# **EM Approach**

- Assign  $\Theta^{(0)} = \{\theta_{ijk}^{(0)}\}$  randomly.
- Iteratively, m = 0 ...

**E step:** Compute EXPECTED value of  $N_{ijk}$ , given  $\langle G, \theta^m \rangle$ 

$$\widehat{N}_{ijk} = E_{P(x|S,\theta^m,G)}(N_{ijk}) = \sum_{c_{\ell} \in S} P(x_i^k, pa_i^j \mid c_{\ell}, \theta^m, S)$$

**M step:** Update values of  $\theta^{m+1}$  based on  $\hat{N}_{ijk}$ 

$$\theta_{ijk}^{m+1} = \frac{\hat{N}_{ijk} + 0}{\sum_{k=1}^{r_i} (\hat{N}_{ijk} + 0)}$$

... until  $|\theta^{m+1} - \theta^m| \gtrsim 0$ .

• Return  $\theta^m$ 

1. This is ML computation; MAP is similar

"O" 
$$\rightarrow \alpha_{ijk}$$

- 2. Finds local optimum
- 4. Views each tupe with r "\*"s as O(2") partial tuples 3. Used for HMM

# 4

## Facts about EM ...

- Converges eventually
- When not converged: Always improves likelihood
  - L(  $\theta^{(t+1)} : S$  ) > L(  $\theta^{(t)} : S$  )
  - ... except at stationary points...
- For CPtable for Belief net:
  - Need to perform general BN inference
  - Use Clique-tree or ClusterGraph
     ... just needs one pass
     (as N<sub>iik</sub> depends on node+parents)

# Outline

- Motivation
- What is a Belief Net?
- Learning a Belief Net
  - Goal?
  - Learning Parameters Complete Data
  - Learning Parameters Incomplete Data
    - Gradient Descent
    - EM
    - Gibbs
  - Learning Structure

# Gibbs Sampling

ullet Let  $S^{(0)}$  be COMPLETED version of S, randomly filling-in each missing  $c_{ij}$ 

Let 
$$d_{ij}^{(0)}=c_{ij}$$
 If  $c_{ij}=*$ , then  $d_{ij}^{(0)}=\mathrm{Random}[\mathrm{\ Domain}(X_i)\ ]$ 

- For k = 0..
  - Compute  $\Theta^{(k)}$  from  $S^{(k)}$  [frequencies]
  - Form  $S^{(k+1)}$  by...
    - \* If  $c_{ij} \neq *$ ,  $d_{ij}^{k+1} \coloneqq c_{ij}$
    - \* If  $c_{ij}=*$  then

Let  $d_{ij}^{k+1}$  be random value for  $X_i$ , based on current distr  $\Theta^k$  over  $Z-X_i$ 

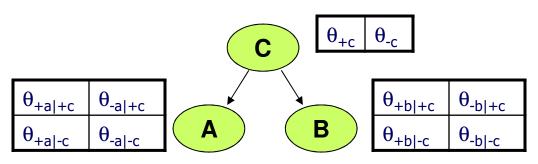
• Return average of these  $\Theta^{(k)}$ 's

Note: As  $\Theta^{(k)}$  based on COMPLETE DATA  $S^{(k)}$   $\Rightarrow \Theta^{(k)}$  can be computed efficiently!

"Multiple Imputation"



# Gibbs Sampling – Example



### New

$$S^{(1)} =$$

Flip 0.3-coin:

Flip 0.9-coin:

Flip 0.8-coin:

Flip 0.9-coin:

Α	В	С
0	0	1
0	1	0
0	1	1
1	1	1

Guess initial values $\theta^0$					
			0.55	0.45	
0.8	0.2			0.9	0.1
0.3	0.7	A	B	0.4	0.6

### Then

- Use  $S^{(1)}$  to get new  $\theta^{(2)}$  parameters
- Form new  $S^{(2)}$  by drawing new values from  $\theta^{(2)}$

# Gibbs Sampling (con't)

- Algorithm: Repeat
  - Given COMPLETE data  $S^{(i)}$ , compute new ML values for  $\{\theta_{ijk}^{(i+1)}\}$
  - Using NEW parameters, impute (new) missing values S(i+1)
- Q: What to return?

AVERAGE over **separated ⊕**(i)'s

- eg,  $\Theta^{(500)}$ ,  $\Theta^{(600)}$ ,  $\Theta^{(700)}$ , ...
- Q: When to stop?

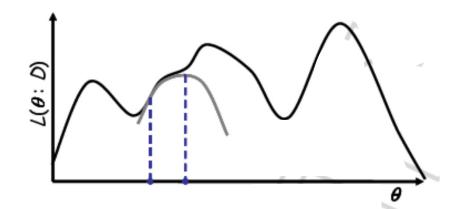
When distribution over ⊖(i)s has converged

- Comparison: Gibbs vs EM
  - + EM "splits" each instance
     ...into 2<sup>r</sup> parts if r \*'s
  - EM knows when it is done, and what to return



## General Issues

- All alg's are heuristic...
  - Starting values θ<sup>(0)</sup>
  - Stopping criteria
  - Escaping local maxima



So far, trying to optimize likelihood.
 Could try to optimize APPROXIMATION to likelihood...



# Summary of Approaches

- Gradient Ascent
- EM-based (many variants)
- Gibbs sampling
  - Multiple imputation