

## Structure Learning 2

1

## Structure Scores

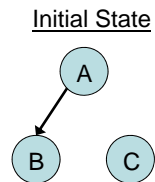
- Searching for highest-scoring network structure is intractable
- Need to resort to heuristic search (eg. hillclimbing)
- Need:
  1. Search space
  2. Scoring function
  3. Search procedure

2

# Structure Scores

## 1. Search space

- Start with initial state (eg. disconnected graph or randomly generated one)

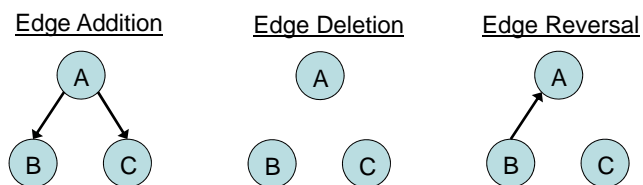


3

# Structure Scores

## 1. Search space

- Move to a neighboring state by applying an operator:



Can only perform an operator if it doesn't lead to a cycle!

4

## Structure Scores

### 2. Scoring function:

- Two general classes of scoring functions:
  1. Likelihood scoring functions
  2. Bayesian scoring functions
- More about this in a bit...assume we have a scoring function for now

5

## Structure Scores

### 3. Search procedure

- Greedy search: pick the best scoring neighboring state to move to
- Repeat until convergence
- Converges to a local optimum

Tricks for dealing with this:  
random restart, simulated  
annealing, tabu search and  
data perturbation

6

## Structure Scores

7

## Likelihood Scores

$$\begin{aligned} & \max_{\mathcal{G}, \boldsymbol{\theta}_{\mathcal{G}}} L(\langle \mathcal{G}, \boldsymbol{\theta}_{\mathcal{G}} \rangle; \mathcal{D}) \\ &= \max_{\mathcal{G}} \left[ \max_{\boldsymbol{\theta}_{\mathcal{G}}} L(\langle \mathcal{G}, \boldsymbol{\theta}_{\mathcal{G}} \rangle; \mathcal{D}) \right] \\ &= \max_{\mathcal{G}} [L(\langle \mathcal{G}, \hat{\boldsymbol{\theta}}_{\mathcal{G}} \rangle; \mathcal{D})] \end{aligned}$$

Graph structure that  
maximizes the likelihood

Maximum likelihood estimates  
of parameters

$$\text{score}_L(\mathcal{G}; \mathcal{D}) = l(\hat{\boldsymbol{\theta}}_{\mathcal{G}}; \mathcal{D})$$

Log likelihood

8

## Likelihood Scores

Let  $M$  be the number of samples. We use the notation  $M[x]$  to be the count of  $x$  in the data.

Let  $\hat{P}$  be the empirical distribution observed in the data. Eg.

- $M[x, y] = M \cdot \hat{P}(x, y)$
- $M[y] = M \cdot \hat{P}(y)$

Note that:

- $\hat{\theta}_{y|x} = \hat{P}(y|x)$
- $\hat{\theta}_y = \hat{P}(y)$

9

## Likelihood Scores

Mutual Information

$$\begin{aligned} I_{\hat{P}}(X; Y) &= \sum_{x,y} \hat{P}(x,y) \log \frac{\hat{P}(x,y)}{\hat{P}(x)\hat{P}(y)} \\ &= \frac{1}{M} \sum_{x,y} M[x,y] \log \left( \frac{M[x,y]}{M[x]M[y]} \right) \end{aligned}$$

10

## Likelihood Scores

Claim:

$$\begin{aligned} \text{score}_L(\mathcal{G}: \mathcal{D}) &= M \sum_{i=1}^n \mathbf{I}_{\hat{p}}(X_i; \text{Parents}(X_i, \mathcal{G})) - M \sum_{i=1}^n H_{\hat{p}}(X_i) \\ &= M \sum_{i=1}^n [\mathbf{I}_{\hat{p}}(X_i; \text{Parents}(X_i, \mathcal{G})) - H_{\hat{p}}(X_i)] \end{aligned}$$

11

## Likelihood Scores

Proof:

$$\begin{aligned} l(\hat{\theta}_{\mathcal{G}}: \mathcal{D}) &= \sum_{i=1}^n \left[ \sum_{\mathbf{u}_i \in \text{Val}(\text{Parents}(X_i, \mathcal{G}))} \sum_{x_i} M[x_i, \mathbf{u}_i] \log \hat{\theta}_{x_i | \mathbf{u}_i} \right] \\ &= M \sum_{i=1}^n \left[ \frac{1}{M} \sum_{\mathbf{u}_i} \sum_{x_i} M[x_i, \mathbf{u}_i] \log \hat{\theta}_{x_i | \mathbf{u}_i} \right] \end{aligned}$$

12

## Likelihood Scores

$$\begin{aligned}
 & \frac{1}{M} \sum_{\mathbf{u}_i} \sum_{x_i} M[x_i, \mathbf{u}_i] \log \hat{\theta}_{x_i | \mathbf{u}_i} \\
 &= \sum_{\mathbf{u}_i} \sum_{x_i} \hat{P}(x_i, \mathbf{u}_i) \log \hat{P}(x_i | \mathbf{u}_i) \\
 &= \sum_{\mathbf{u}_i} \sum_{x_i} \hat{P}(x_i, \mathbf{u}_i) \log \left( \frac{\hat{P}(x_i, \mathbf{u}_i)}{\hat{P}(\mathbf{u}_i)} \cdot \frac{\hat{P}(x_i)}{\hat{P}(x_i)} \right) \\
 &= \sum_{\mathbf{u}_i} \sum_{x_i} \hat{P}(x_i, \mathbf{u}_i) \log \left( \frac{\hat{P}(x_i, \mathbf{u}_i)}{\hat{P}(\mathbf{u}_i) \hat{P}(x_i)} \right) + \sum_{x_i} \left( \sum_{\mathbf{u}_i} \hat{P}(x_i, \mathbf{u}_i) \right) \log \hat{P}(x_i) \\
 &= I_{\hat{P}}(X_i; \mathbf{U}_i) - \sum_{x_i} \hat{P}(x_i) \log \frac{1}{\hat{P}(x_i)} \\
 &= I_{\hat{P}}(X_i; \mathbf{U}_i) - H_{\hat{P}}(X_i)
 \end{aligned}$$

Note that if  $\text{Parents}(X_i, \mathcal{G}) = \emptyset$ , then  $I_{\hat{P}}(X_i; \text{Parents}(X_i, \mathcal{G})) = 0$

13

## Likelihood Scores

What are the implications of

$$I_{\hat{P}}(X_i; \mathbf{U}_i) - H_{\hat{P}}(X_i)$$

Depends on network structure (because  $\mathbf{U}_i = \text{Parents}(X_i, \mathcal{G})$ ). Only need to maximize this.

Does not depend on network structure

The likelihood of a network measures how informative  $\text{Parents}(X_i)$  are about  $X_i$

14

# Likelihood Scores

An alternate representation:

$$\frac{1}{M} \text{score}_L(G; \mathcal{D}) = \underbrace{H_{\hat{p}}(X_1, \dots, X_n)}_{\text{Does not depend on network structure}} - \sum_{i=1}^n \underbrace{I_{\hat{p}}(X_i; \{X_1, \dots, X_{i-1}\} - \text{Parents}(X_i, G) \mid \text{Parents}(X_i, G))}_{\text{Depends on network structure}}$$

Measures to what extent the Markov properties of the graph are violated in the data (fewer violations  $\Rightarrow$  larger score)

15

# Problems with Likelihood Score

Never prefers a simpler network over a more complex one eg.



$$\text{score}_L(G_1; \mathcal{D}) \geq \text{score}_L(G_0; \mathcal{D})$$

16



## Problems with Likelihood Score

- Exhibits a conditional independence only if it holds exactly in the **empirical distribution**
  - Due to noise, this almost never happens
- Learns a fully connected graph
  - Overfits the training data and does not generalize well to unseen cases
- Needs a penalty for learning overly complex structures

17

## Bayesian Scoring

18

## Bayesian Score

- Bayesian philosophy: if you are uncertainty about something, put a distribution over it
- In structure learning, we have uncertainty over the **structure** and the **parameters**
- We will have two prior distributions:
  - Structure prior  $P(\mathcal{G})$
  - Parameter prior  $P(\boldsymbol{\theta}_{\mathcal{G}}|\mathcal{G})$

19

## Bayesian Score

Recall: 
$$P(\mathcal{G}|D) = \frac{P(D|\mathcal{G})P(\mathcal{G})}{P(D)} = \alpha P(D|\mathcal{G})P(\mathcal{G})$$

$$score_B(\mathcal{G}: D) = \log P(D|\mathcal{G}) + \log P(\mathcal{G})$$

Marginal Likelihood  
(dominates the score)

Structure prior

$$P(D|\mathcal{G}) = \int_{\boldsymbol{\theta}_{\mathcal{G}}} P(D|\boldsymbol{\theta}_{\mathcal{G}}, \mathcal{G}) P(\boldsymbol{\theta}_{\mathcal{G}}|\mathcal{G}) d\boldsymbol{\theta}_{\mathcal{G}}$$

“Averages” out  $P(D|\boldsymbol{\theta}_{\mathcal{G}}, \mathcal{G})$  over the distribution of  $\boldsymbol{\theta}_{\mathcal{G}}$ . Contrast this with maximum likelihood which finds the  $\boldsymbol{\theta}_{\mathcal{G}}$  that maximizes the likelihood of the data

20

## Bayesian score

- How does the Bayesian score improve over the likelihood score?
  - By avoiding overfitting
- Likelihood score commits to a single  $\hat{\theta}$  value
- Bayesian score works with a distribution of  $\theta_G$  and averages  $P(\mathcal{D}|\theta_G, \mathcal{G})$  over this distribution
  - Results in an **expected likelihood**

21

## Marginal Likelihood (Single Variable case)

- Suppose we have a single binary random variable  $X$
- Let the prior distribution over the parameters of  $X$  be *Dirichlet*( $\alpha_1, \alpha_0$ )
- Let the data  $\mathcal{D} = \{x[1], \dots, x[M]\}$  have  $M[1]$  heads and  $M[0]$  tails
- Maximum likelihood value given  $\mathcal{D}$

$$\text{is: } P(\mathcal{D}|\hat{\theta}) = \left(\frac{M[1]}{M}\right)^{M[1]} \cdot \left(\frac{M[0]}{M}\right)^{M[0]}$$

22

## Marginal Likelihood (Single Variable case)

What about the marginal likelihood?

$$P(\mathcal{D}|\mathcal{G}) = \int_{\theta_{\mathcal{G}}} \underbrace{P(\mathcal{D}|\theta_{\mathcal{G}}, \mathcal{G})}_{\left(\frac{M[1]}{M}\right)^{M[1]} \cdot \left(\frac{M[0]}{M}\right)^{M[0]}} \underbrace{P(\theta_{\mathcal{G}}|\mathcal{G})}_{\text{Dirichlet}(\alpha_1, \alpha_0)} d\theta_{\mathcal{G}}$$

Shorthand: let  $p_i = \frac{M[i]}{M}$  and  $\alpha = \alpha_0 + \alpha_1$

23

## Marginal Likelihood (Single Variable case)

$$\begin{aligned} P(\mathcal{D}|\mathcal{G}) &= \int_{\theta_{\mathcal{G}}} p_1^{M[1]} p_0^{M[0]} \frac{\Gamma(\alpha_0 + \alpha_1)}{\Gamma(\alpha_0)\Gamma(\alpha_1)} p_1^{(\alpha_1-1)} p_0^{(\alpha_0-1)} d\theta_{\mathcal{G}} \\ &= \frac{\Gamma(\alpha)}{\Gamma(\alpha_0)\Gamma(\alpha_1)} \int_{\theta_{\mathcal{G}}} p_1^{(M[1]+\alpha_1-1)} p_0^{(M[0]+\alpha_0-1)} d\theta_{\mathcal{G}} \end{aligned}$$

Note:

$$\begin{aligned} &\int_{\theta_{\mathcal{G}}} \text{Beta}(M[1] + \alpha_1, M[0] + \alpha_0) d\theta_{\mathcal{G}} = 1 \\ &\Rightarrow \int_{\theta_{\mathcal{G}}} \frac{\Gamma(\alpha + M)}{\Gamma(M[1] + \alpha_1)\Gamma(M[0] + \alpha_0)} p_1^{(M[1]+\alpha_1-1)} p_0^{(M[0]+\alpha_0-1)} d\theta_{\mathcal{G}} = 1 \\ &\Rightarrow \int_{\theta_{\mathcal{G}}} p_1^{(M[1]+\alpha_1-1)} p_0^{(M[0]+\alpha_0-1)} d\theta_{\mathcal{G}} = \frac{\Gamma(\alpha + M)}{\Gamma(M[1] + \alpha_1)\Gamma(M[0] + \alpha_0)} \end{aligned}$$

24

## Marginal Likelihood (Single Variable case)

$$P(\mathcal{D}|\mathcal{G}) = \int_{\theta_{\mathcal{G}}} p_1^{M[1]} p_0^{M[0]} \frac{\Gamma(\alpha_0 + \alpha_1)}{\Gamma(\alpha_0)\Gamma(\alpha_1)} p_1^{(\alpha_1-1)} p_0^{(\alpha_0-1)} d\theta_{\mathcal{G}}$$

$$= \frac{\Gamma(\alpha)}{\Gamma(\alpha_0)\Gamma(\alpha_1)} \int_{\theta_{\mathcal{G}}} p_1^{(M[1]+\alpha_1-1)} p_0^{(M[0]+\alpha_0-1)} d\theta_{\mathcal{G}}$$

Note that the Gamma function is as follows:

$$\Gamma(1) = 1$$

$$\Gamma(x + 1) = x\Gamma(x)$$

ie. it is a continuous generalization of the factorial:  
 $\Gamma(n + 1) = n!$

Note:  $\int_{\theta_{\mathcal{G}}} \text{Beta}(M[1] + \alpha_1, M[0] + \alpha_0) d\theta_{\mathcal{G}} = 1$

$$\Rightarrow \int_{\theta_{\mathcal{G}}} \frac{\Gamma(\alpha + M)}{\Gamma(M[1] + \alpha_1)\Gamma(M[0] + \alpha_0)} p_1^{(M[1]+\alpha_1-1)} p_0^{(M[0]+\alpha_0-1)} d\theta_{\mathcal{G}} = 1$$

$$\Rightarrow \int_{\theta_{\mathcal{G}}} p_1^{(M[1]+\alpha_1-1)} p_0^{(M[0]+\alpha_0-1)} d\theta_{\mathcal{G}} = \frac{\Gamma(\alpha + M)}{\Gamma(M[1] + \alpha_1)\Gamma(M[0] + \alpha_0)}$$

25

## Marginal Likelihood (Single Variable case)

$$= \frac{\Gamma(\alpha)}{\Gamma(\alpha_0)\Gamma(\alpha_1)} \frac{\Gamma(\alpha + M)}{\Gamma(M[1] + \alpha_1)\Gamma(M[0] + \alpha_0)}$$

$$P(\mathcal{D}|\mathcal{G}) = \frac{\Gamma(\alpha)}{\Gamma(\alpha + M)} \cdot \frac{\Gamma(\alpha_1 + M)}{\Gamma(\alpha_1)} \cdot \frac{\Gamma(\alpha_0 + M)}{\Gamma(\alpha_0)}$$

We can easily generalize to a **multinomial distribution** over the space of values  $x^1, \dots, x^k$  with a prior

**Dirichlet**( $\alpha_1, \dots, \alpha_k$ ):

$$P(\mathcal{D}|\mathcal{G}) = \frac{\Gamma(\alpha)}{\Gamma(\alpha + M)} \cdot \prod_i^k \frac{\Gamma(\alpha_i + M[x^i])}{\Gamma(\alpha_i)}$$

26

## Bayesian Scoring

### Global parameter independence:

Let  $\mathcal{G}$  be a Bayesian network structure with parameters  $\theta = (\theta_{X_1|Pa(X_1)}, \dots, \theta_{X_n|Pa(X_n)})$ .

The distribution  $P(\theta)$  satisfies global parameter independence if it has the form:

$$P(\theta) = \prod_{i=1}^n P(\theta_{X_i|Pa(X_i)})$$

27

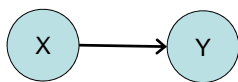
## Bayesian Scoring

### Local parameter independence:

Let  $X$  be a variable with parents  $U$ . We say that distribution  $P(\theta_{X|U})$  satisfies local parameter independence if:

$$P(\theta_{X|U}) = \prod_{u} P(\theta_{X|u})$$

Example:



X	Y	P(Y X)
0	0	$\theta_{Y x^0}$
0	1	
1	0	$\theta_{Y x^1}$
1	1	

Only one of the  $\theta$ s applies, depending on the value of  $x$ . In other words, the  $\theta$ s don't affect each other

28

# Bayesian Scoring

Now suppose there are two binary random variables X and Y. Let  $\mathcal{G}_0$  be a graph with X and Y and no edges



$$P(\mathcal{D}|\mathcal{G}_0) = \int_{\Theta_X \times \Theta_Y} P(\mathcal{D}|\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \mathcal{G}_0) P(\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y|\mathcal{G}_0) d[\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y]$$

1. Decompose likelihood in terms of each variable

$$P(\mathcal{D}|\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \mathcal{G}_0) = \prod_{i=1}^M P(x[m]|\boldsymbol{\theta}_X, \mathcal{G}_0) P(y[m]|\boldsymbol{\theta}_Y, \mathcal{G}_0)$$

2. Global Parameter Independence:  $P(\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y|\mathcal{G}_0) = P(\boldsymbol{\theta}_X|\mathcal{G}_0)P(\boldsymbol{\theta}_Y|\mathcal{G}_0)$

29

# Bayesian Scoring

$$\begin{aligned} P(\mathcal{D}|\mathcal{G}_0) &= \int_{\Theta_X \times \Theta_Y} P(\mathcal{D}|\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y, \mathcal{G}_0) P(\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y|\mathcal{G}_0) d[\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y] \\ &= \int_{\Theta_X} \int_{\Theta_Y} \prod_{m=1}^M P(x[m]|\boldsymbol{\theta}_X, \mathcal{G}_0) P(y[m]|\boldsymbol{\theta}_Y, \mathcal{G}_0) P(\boldsymbol{\theta}_X|\mathcal{G}_0) P(\boldsymbol{\theta}_Y|\mathcal{G}_0) d[\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y] \\ &= \int_{\Theta_X} \int_{\Theta_Y} \left( \prod_{m=1}^M P(x[m]|\boldsymbol{\theta}_X, \mathcal{G}_0) P(\boldsymbol{\theta}_X|\mathcal{G}_0) \right) \left( \prod_{m=1}^M P(y[m]|\boldsymbol{\theta}_Y, \mathcal{G}_0) P(\boldsymbol{\theta}_Y|\mathcal{G}_0) \right) d[\boldsymbol{\theta}_X, \boldsymbol{\theta}_Y] \end{aligned}$$

Integral of a product of independent functions is the product of integrals:

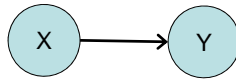
$$= \underbrace{\left( \int_{\Theta_X} \prod_{m=1}^M P(x[m]|\boldsymbol{\theta}_X, \mathcal{G}_0) P(\boldsymbol{\theta}_X|\mathcal{G}_0) d\boldsymbol{\theta}_X \right)}_{\text{Term for X}} \underbrace{\left( \int_{\Theta_Y} \prod_{m=1}^M P(y[m]|\boldsymbol{\theta}_Y, \mathcal{G}_0) P(\boldsymbol{\theta}_Y|\mathcal{G}_0) d\boldsymbol{\theta}_Y \right)}_{\text{Term for Y}}$$

Note: decomposes into one term for each random variable

30

## Bayesian Scoring

Now suppose there are two binary random variables X and Y and let  $\mathcal{G}_{X \rightarrow Y}$  be the graph below:



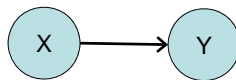
X	P(X)
0	$\theta_x$
1	

X	Y	P(Y X)
0	0	$\theta_{Y x^0}$
0	1	
1	0	$\theta_{Y x^1}$
1	1	

31

## Bayesian Scoring

Now suppose there are two binary random variables X and Y and let  $\mathcal{G}_{X \rightarrow Y}$  be the graph below:



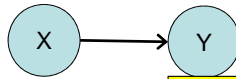
$$\begin{aligned}
 P(\mathcal{D}|\mathcal{G}_{X \rightarrow Y}) &= \left( \int_{\Theta_X} \prod_{m=1}^M P(x[m]|\theta_x, \mathcal{G}_{X \rightarrow Y}) P(\theta_x|\mathcal{G}_{X \rightarrow Y}) d\theta_x \right) \cdot \\
 &\left( \int_{\Theta_{Y|x^0}} \prod_{m:x[m]=x^0}^M P(y[m]|\theta_{Y|x^0}, \mathcal{G}_{X \rightarrow Y}) P(\theta_{Y|x^0}|\mathcal{G}_{X \rightarrow Y}) d\theta_{Y|x^0} \right) \cdot \\
 &\left( \int_{\Theta_{Y|x^1}} \prod_{m:x[m]=x^1}^M P(y[m]|\theta_{Y|x^1}, \mathcal{G}_{X \rightarrow Y}) P(\theta_{Y|x^1}|\mathcal{G}_{X \rightarrow Y}) d\theta_{Y|x^1} \right)
 \end{aligned}$$

32



# Bayesian Scoring

Now suppose there are two binary random variables X and Y and let  $\mathcal{G}_{X \rightarrow Y}$  be the graph below:



$$P(\mathcal{D}|\mathcal{G}_{X \rightarrow Y}) = \left( \int_{\theta_X} \prod_{m=1}^M P(x[m]|\theta_X, \mathcal{G}_{X \rightarrow Y}) \right) \left( \int_{\theta_{Y|x^0}} \prod_{m:x[m]=x^0}^M P(y[m]|\theta_{Y|x^0}, \mathcal{G}_{X \rightarrow Y}) P(\theta_{Y|x^0}|\mathcal{G}_{X \rightarrow Y}) \right) \left( \int_{\theta_{Y|x^1}} \prod_{m:x[m]=x^1}^M P(y[m]|\theta_{Y|x^1}, \mathcal{G}_{X \rightarrow Y}) P(\theta_{Y|x^1}|\mathcal{G}_{X \rightarrow Y}) d\theta_{Y|x^1} \right)$$

One term for each parameter family. Each term has a closed form solution

$$P(\mathcal{D}|\mathcal{G}) = \frac{\Gamma(\alpha)}{\Gamma(\alpha + M)} \cdot \prod_i^k \frac{\Gamma(\alpha_i + M[x^i])}{\Gamma(\alpha_i)}$$

33

# Bayesian Scoring

The general case: let  $\mathcal{G}$  be a network structure, and let  $P(\theta_{\mathcal{G}}|\mathcal{G})$  be a parameter prior satisfying **global parameter independence**. Then:

$$P(\mathcal{D}|\mathcal{G}) = \prod_{i=1}^n \int_{\theta_{X_i|Pa(X_i)}} \prod_{m=1}^M P(x_i[m]|pa(X_i)[m], \theta_{X_i|Pa(X_i)}, \mathcal{G}) P(\theta_{X_i|Pa(X_i)}|\mathcal{G}) d\theta_{X_i|Pa(X_i)}$$

If  $P(\theta_{\mathcal{G}})$  also satisfies local parameter independence, then

$$P(\mathcal{D}|\mathcal{G}) = \prod_{i=1}^n \prod_{\mathbf{u}_i \in Val(Pa^{\mathcal{G}}(X_i))} \int_{\theta_{X_i|\mathbf{u}_i}} \prod_{m:\mathbf{u}_i[m]=\mathbf{u}_i}^M P(x_i[m]|\mathbf{u}_i)[m], \theta_{X_i|\mathbf{u}_i}, \mathcal{G}) P(\theta_{X_i|\mathbf{u}_i}|\mathcal{G}) d\theta_{X_i|\mathbf{u}_i}$$

34

## Bayesian Scoring

If we have a Bayesian network with Dirichlet priors where  $P(\theta_{X_i|pa(X_i)}|\mathcal{G})$  has hyperparameters  $\{\alpha^{\mathcal{G}}_{x^j_i|u_i} : j = 1, \dots, |X_i|\}$  then

$$P(\mathcal{D}|\mathcal{G}) = \prod_{i=1}^n \prod_{u_i \in \text{Val}(Pa^{\mathcal{G}}(X_i))} \frac{\Gamma(\alpha^{\mathcal{G}}_{x_i|u_i})}{\Gamma(\alpha^{\mathcal{G}}_{x_i|u_i} + M[u_i])} \prod_{x_i^j \in \text{Val}(X_i)} \left[ \frac{\Gamma(\alpha^{\mathcal{G}}_{x_i^j|u_i}) + M[x_i^j, u_i]}{\Gamma(\alpha^{\mathcal{G}}_{x_i^j|u_i})} \right]$$

Where:

$$\alpha^{\mathcal{G}}_{x_i|u_i} = \sum_j \alpha^{\mathcal{G}}_{x_i^j|u_i}$$

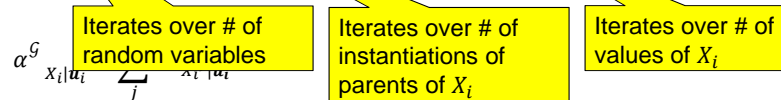
35

## Bayesian Scoring

If we have a Bayesian network with Dirichlet priors where  $P(\theta_{X_i|pa(X_i)}|\mathcal{G})$  has hyperparameters  $\{\alpha^{\mathcal{G}}_{x^j_i|u_i} : j = 1, \dots, |X_i|\}$  then

$$P(\mathcal{D}|\mathcal{G}) = \prod_{i=1}^n \prod_{u_i \in \text{Val}(Pa^{\mathcal{G}}(X_i))} \frac{\Gamma(\alpha^{\mathcal{G}}_{x_i|u_i})}{\Gamma(\alpha^{\mathcal{G}}_{x_i|u_i} + M[u_i])} \prod_{x_i^j \in \text{Val}(X_i)} \left[ \frac{\Gamma(\alpha^{\mathcal{G}}_{x_i^j|u_i}) + M[x_i^j, u_i]}{\Gamma(\alpha^{\mathcal{G}}_{x_i^j|u_i})} \right]$$

Where:



36

## Bayesian Scoring

If we use a Dirichlet parameter prior for all parameters in our network, then, because  $M \rightarrow \infty$  (proof omitted), we have:

$$\log P(\mathcal{D}|\mathcal{G}) = l(\hat{\theta}_{\mathcal{G}}; \mathcal{D}) - \frac{\log M}{2} \underbrace{\text{Dim}[\mathcal{G}]}_{\substack{\text{\# of independent} \\ \text{parameters in } \mathcal{G}}} + O(1)$$

This is the Bayesian Information Criterion (BIC) score

37

## Bayesian Scoring

This is the Bayesian Information Criterion (BIC) score:

$$\text{score}_{BIC}(\mathcal{G}; \mathcal{D}) = \underbrace{l(\hat{\theta}_{\mathcal{G}}; \mathcal{D})}_{\text{Fit to data}} - \underbrace{\frac{\log M}{2} \text{Dim}[\mathcal{G}]}_{\text{Model complexity}} + O(1)$$

Can also interpret this as the # of bits to encode the model and the data given the model (minimum description length)

38

## Bayesian Scoring

$$score_{BIC}(\mathcal{G}; \mathcal{D}) = M \sum_{i=1}^n I_{\hat{p}}(X_i; Pa(X_i)) - M \sum_{i=1}^n H_{\hat{p}}(X_i) - \frac{\log M}{2} Dim[\mathcal{G}]$$

Things to note:

- Entropy term  $M \sum_{i=1}^n H_{\hat{p}}(X_i)$  can be ignored (doesn't depend on graph structure)
- Trades off fit to data and model complexity
  - The stronger the dependence of a variable on its parents, the higher the score (grows linearly)
  - The more complex the network, the lower the score (grows logarithmically)
- As  $M$  grows, the score pays more attention to the data fit

39

## Bayesian Scoring

Assume that our data are generated by some distribution  $P^*$  for which the network  $\mathcal{G}^*$  is a perfect map.

We say that a scoring function is **consistent** if the following properties hold as the amount of data  $M \rightarrow \infty$ , with probability that approaches 1 (over all possible choices of data set  $\mathcal{D}$ ):

- The structure  $\mathcal{G}^*$  will maximize the score
- All structures  $\mathcal{G}$  that are not I-equivalent to  $\mathcal{G}^*$  will have strictly lower score

40

## Bayesian Scoring

- The BIC score (and the Bayesian score) is **consistent** [proof omitted]
- In practice though, the BIC score tends to have a very strong preference for simpler structures

41

## Structure Priors

Recall that

$$\text{score}_B(\mathcal{G}: \mathcal{D}) = \underbrace{\log P(\mathcal{D} | \mathcal{G})}_{\substack{\text{Grows linearly with the number of} \\ \text{examples (dominates the score)}}} + \underbrace{\log P(\mathcal{G})}_{\substack{\text{Structure prior (stays} \\ \text{constant). Only} \\ \text{matters for small} \\ \text{sample sizes}}}$$

Grows linearly with the number of examples (dominates the score)

Structure prior (stays constant). Only matters for small sample sizes

42

## Structure Priors

- Typically assign uniform priors over structures
- If you can provide an informed structure prior, you could penalize edges in the graph:
  - $P(\mathcal{G}) \propto c^{|\mathcal{G}|}$  (where  $c < 1$  and  $|\mathcal{G}|$  is the number of edges)
- Mathematically convenient to have structure prior with structure modularity:
  - $P(\mathcal{G}) \propto \prod_i P(\text{Pa}(X_i) = \text{Pa}^{\mathcal{G}}(X_i))$

Uses local properties not global properties  
of the graph

43