

Readings:

K&F: 15.1, 15.2, 15.3, 15.4, 15.5

Structure Learning: the good, the bad, the ugly

in BNR

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

September 29th, 2006

1

Understanding the uniform prior over the parameters

■ Likelihood: $P(\mathcal{D} | \theta) = \theta^{m_H} (1 - \theta)^{m_T}$

■ Prior: Uniform, Beta(1,1)

■ $E[\theta] = \frac{m_H + 1}{m_H + m_T + 2}$

■ Why???

$$P(\theta | \mathcal{D}) \propto P(\theta) \cdot P(\mathcal{D} | \theta) = P(\theta) \cdot \theta^{m_H} (1 - \theta)^{m_T}$$

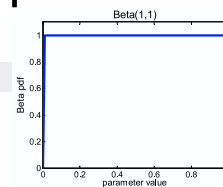
uniform $P(\theta) = 1$ $\Rightarrow \theta^{m_H} (1 - \theta)^{m_T} = \text{Beta}(m_H + 1, m_T + 1)$

$$E[\theta] = \int_0^1 \theta \cdot \theta^{m_H} (1 - \theta)^{m_T} d\theta = \frac{m_H + 1}{m_H + m_T + 2}$$

suppose: $m_H = 1, m_T = 0$
 $\theta_{\text{mean}} = \frac{2}{3}$ $\theta_{\text{mode}} = 1$

mode: $\frac{m_H}{m_H + m_T}$

■ Intuition: Prior knowledge is “don’t know anything”, so I don’t want to commit too fast



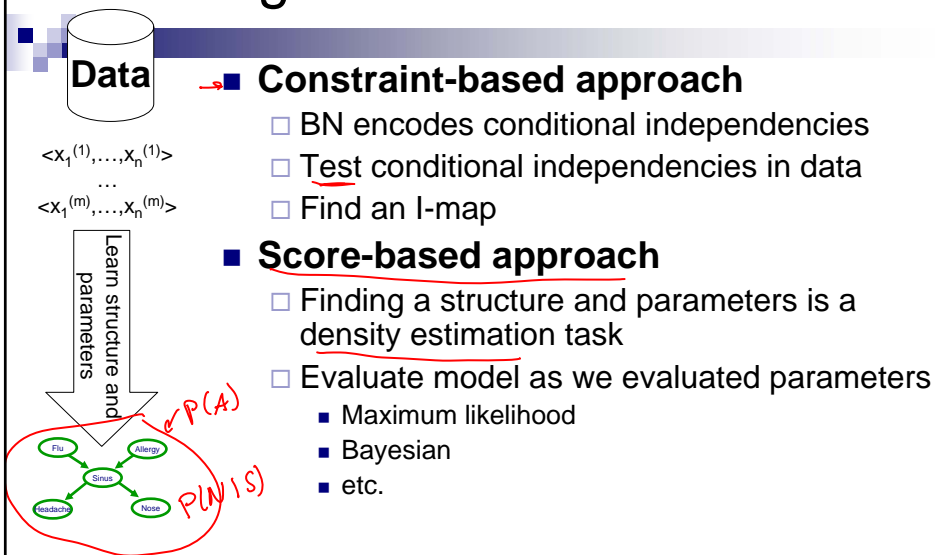
Where are we with learning BNs?

- Given structure, estimate parameters
 - Maximum likelihood estimation MLE
 - Bayesian learning
- What about learning structure?

10-708 - ©Carlos Guestrin 2006

3

Learning the structure of a BN



10-708 - ©Carlos Guestrin 2006

4

Remember: Obtaining a P-map?

- Given the independence assertions that are true for P
 - Obtain skeleton
 - Obtain immoralities
- From skeleton and immoralities, obtain every (and any) BN structure from the equivalence class

■ Constraint-based approach:

- Use Learn PDAG algorithm
- Key question: **Independence test**

10-708 - ©Carlos Guestrin 2006

5

Independence tests

- Statistically difficult task!

- Intuitive approach: **Mutual information**

$$I(X_i, X_j) = \sum_{x_i, x_j} P(x_i, x_j) \log \frac{P(x_i, x_j)}{P(x_i)P(x_j)}$$

- Mutual information and independence:

- X_i and X_j independent if and only if $I(X_i, X_j) = 0$

suppose $X_i \perp X_j$ $P(x_i, x_j) = P(x_i) \cdot P(x_j)$

$$\frac{P(x_i) \cdot P(x_j)}{P(x_i) \cdot P(x_j)} = 1; \log 1 = 0$$

- Conditional mutual information:

$$(X_i \perp X_j | Z)$$

$$I(X_i, X_j | Z) = \sum_{x_i, x_j, z} P(x_i, x_j, z) \cdot \log \frac{P(x_i, x_j | z)}{P(x_i | z) \cdot P(x_j | z)}$$

$$I(X_i, X_j | Z) = 0 \Leftrightarrow (X_i \perp X_j | Z); \text{ since } P(x_i, x_j | z) = P(x_i | z) \cdot P(x_j | z)$$

10-708 - ©Carlos Guestrin 2006

6

Independence tests and the constraint based approach

Using the data D

Empirical distribution: $\hat{P}(x_i, x_j) = \frac{\text{Count}(x_i, x_j)}{m}$

Mutual information: $\hat{I}(X_i, X_j) = \sum_{x_i, x_j} \hat{P}(x_i, x_j) \log \frac{\hat{P}(x_i, x_j)}{\hat{P}(x_i) \hat{P}(x_j)}$

Similarly for conditional MI

Use learning PDAG algorithm:

When algorithm asks: $(X \perp Y | U)$?

picking t
is about statistical significance

$\hat{I}(X, Y | U) > t$?
 $\leq t$?
 ≤ 0 ?

$t \rightarrow$ magic param.

no

yes!

Many other types of independence tests

See reading...

10-708 - ©Carlos Guestrin 2006

7

Score-based approach

Data

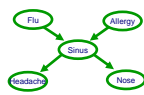
$\langle X_1^{(1)}, \dots, X_n^{(1)} \rangle$

\dots

$\langle X_1^{(m)}, \dots, X_n^{(m)} \rangle$

Possible structures

DAG



$F \rightarrow A$
 $S \rightarrow H \rightarrow N$
 $S \rightarrow F \rightarrow A$
 $H \rightarrow N$

$S \rightarrow F \rightarrow A$
 $H \rightarrow N$

Learn parameters

what score

Score structure

log space

-10 000

-12 000

-20 000

-10 500

many many

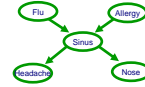
10-708 - ©Carlos Guestrin 2006

8

Information-theoretic interpretation of maximum likelihood

- Given structure, log likelihood of data:

$$\begin{aligned} \log P(\mathcal{D} | \theta_{\mathcal{G}}, \mathcal{G}) &= \sum_{j=1}^m \sum_{i=1}^n \log P(X_i = x_i^{(j)} | \text{Pa}_{X_i} = \mathbf{x}^{(j)} [\text{Pa}_{X_i}]) \\ &= \sum_{i=1}^n \sum_{j=1}^m \log P(X_i = x_i^{(j)} | \text{Pa}_{X_i} = \mathbf{x}^{(j)} [\text{Pa}_{X_i}]) \\ &= m \sum_i \sum_{x_i, u} \underbrace{\text{Count}(X_i = x_i, \text{Pa}_{X_i} = u)}_m \log P(X_i = x_i | \text{Pa}_{X_i} = u) \\ &\quad \underbrace{\hat{P}(X_i = x_i, \text{Pa}_{X_i} = u)} \end{aligned}$$



10-708 - ©Carlos Guestrin 2006

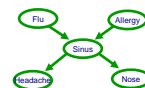
9

Information-theoretic interpretation of maximum likelihood 2

- Given structure, log likelihood of data:

$$\begin{aligned} \log \hat{P}(\mathcal{D} | \theta, \mathcal{G}) &= m \sum_i \sum_{x_i, \text{Pa}_{X_i, \mathcal{G}}} \hat{P}(x_i, \text{Pa}_{X_i, \mathcal{G}}) \log \hat{P}(x_i | \text{Pa}_{X_i, \mathcal{G}}) \\ &= m \sum_i (-\hat{H}(x_i | \text{Pa}_{X_i})) \\ &= -m \sum_i \hat{H}(x_i | \text{Pa}_{X_i}) \\ &\quad \downarrow \sum_i \hat{H}(x_i | \text{Pa}_{X_i}) \\ \text{learn structure} &\Rightarrow \text{for } x_i \text{ find very informative parents} \end{aligned}$$

Conditional Entropy:
 $H(X|U)$
 $= - \sum_{x_i, u} p(x_i, u) \log p(x_i | u)$
 information of U about X
 OR
 uncertainty of X after observing U



10-708 - ©Carlos Guestrin 2006

10

Decomposable score $\overset{I(X_i; \text{Pa}_{X_i, G})}{=} H(X_i) - H(X_i | \text{Pa}_{X_i, G})$

- Log data likelihood $\neg (X_i \perp \text{Pa}_{X_i})$ [very not indep.]

$$\uparrow \log \hat{P}(D | \theta, G) = m \sum_i \hat{I}(X_i, \text{Pa}_{X_i, G}) - \underbrace{m \sum_i \hat{H}(X_i)}_{\text{doesn't graph structure } G}$$

- Decomposable score:

- Decomposes over families in BN (node and its parents)
- Will lead to significant computational efficiency!!!
- $\text{Score}(G : D) = \sum_i \text{FamScore}(X_i | \text{Pa}_{X_i} : D)$

for MLE: $\text{FamScore}(X_i | \text{Pa}_{X_i} : D) = m [\hat{I}(X_i, \text{Pa}_{X_i}) - \hat{H}(X_i)]$

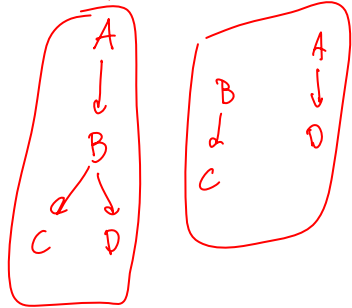
10-708 - ©Carlos Guestrin 2006

11

How many trees are there?

Nonetheless – Efficient optimal algorithm finds best tree

Trees: each node has ~~zero or one~~ parent



trees, not forests

- one root picked
- pick children
- for each child a nother tree



$$\sim 2^{\Theta(n \lg n)}$$

→ really large

10-708 - ©Carlos Guestrin 2006

12

Scoring a tree 1: I-equivalent trees

$\log \hat{P}(\mathcal{D} | \theta, \mathcal{G}) = M \sum_i \hat{I}(x_i, \text{Pa}_{x_i, \mathcal{G}}) - M \sum_i \hat{H}(X_i)$

I-equiv, Same score
 $\text{Score}(\text{Tree}) = \sum_{i,j \in \text{edges}} I(x_i, x_j) - C$

$X \rightarrow Z \rightarrow Y$
 $I(X, \emptyset) - H(X) + I(X, Z) - H(Z) + I(Z, Y) - H(Y)$

$I(X, Z) + I(Z, Y) + \underbrace{H(X) + H(Y) + H(Z)}_{\text{constant } C}$

$X \leftarrow Z \rightarrow Y$
 $I(X, Z) - H(X) + I(Z, \emptyset) - H(Z) + I(Y, Z) - H(Y)$

$X \leftarrow Z \leftarrow Y$
 $I(X, Z) - H(X) + I(Z, Y) - H(Z) + I(Y, \emptyset) - H(Y)$

not a tree

10-708 - ©Carlos Guestrin 2006

13

Scoring a tree 2: similar trees

$\log \hat{P}(\mathcal{D} | \theta, \mathcal{G}) = M \sum_i \hat{I}(x_i, \text{Pa}_{x_i, \mathcal{G}}) - M \sum_i \hat{H}(X_i)$

$X \rightarrow Z \rightarrow Y$
 $I(X, Z) + I(Z, Y) - C$

$Z \rightarrow Y \rightarrow X$
 $I(Z, Y) + I(Y, X) - C$

difference
 $I(X, Z) - I(Y, X)$

only change a few families
 → score only changed for these families

$\text{Score}(\text{Tree}) = \sum_{i,j \in \text{edges}} I(x_i, x_j)$
 best tree heuristic edges \Rightarrow maximum spanning tree

10-708 - ©Carlos Guestrin 2006

14

Chow-Liu tree learning algorithm 1

- For each pair of variables X_i, X_j

- Compute empirical distribution:

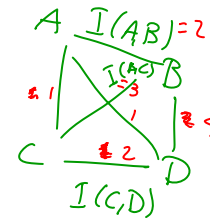
$$\hat{P}(x_i, x_j) \stackrel{MLE}{=} \frac{\text{Count}(x_i, x_j)}{m}$$

- Compute mutual information:

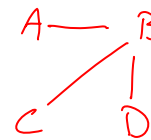
$$\hat{I}(X_i, X_j) = \sum_{x_i, x_j} \hat{P}(x_i, x_j) \log \frac{\hat{P}(x_i, x_j)}{\hat{P}(x_i) \hat{P}(x_j)}$$

- Define a graph

- Nodes X_1, \dots, X_n
 - Edge (i, j) gets weight $\hat{I}(X_i, X_j)$



MST:



10-708 - ©Carlos Guestrin 2006

15

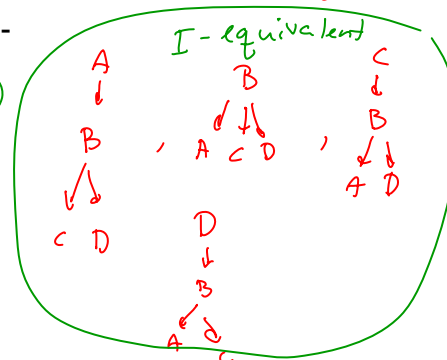
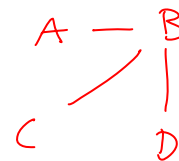
Chow-Liu tree learning algorithm 2

$$\log \hat{P}(\mathcal{D} \mid \theta, \mathcal{G}) = M \sum_i \hat{I}(x_i, \text{Pa}_{x_i, \mathcal{G}}) - M \sum_i \hat{H}(X_i)$$

- Optimal tree BN

- Compute maximum weight spanning tree
 - Directions in BN: pick any node as root, breadth-first-search defines directions

Even though $2^{\theta(n \log n)}$ trees,
found best
in poly time.



10-708 - ©Carlos Guestrin 2006

16

Can we extend Chow-Liu 1

■ OPTIMAL Tree augmented naïve Bayes (TAN) [Friedman et al. '97]

- Naïve Bayes model overcounts, because correlation between features not considered
- Same as Chow-Liu, but score edges with:

$$\hat{I}(X_i, X_j | C) = \sum_{c, x_i, x_j} \hat{P}(c, x_i, x_j) \log \frac{\hat{P}(x_i, x_j | c)}{\hat{P}(x_i | c) \hat{P}(x_j | c)}$$

TAN:



each x_i (except for root) has two parents C and some x_j

Naïve Bayes
 $x_1 \dots x_n (x_i | x_j | C)$
what if
 $x_1 = x_2$
double count
Explain x_2 by:

decreases or avoids double counting

10-708 - ©Carlos Guestrin 2006

17

Can we extend Chow-Liu 2

■ (Approximately learning) models with tree-width up to k

- [Narasimhan & Bilmes '04]
- But, $O(n^{k+1})$...
 - and more subtleties

10-708 - ©Carlos Guestrin 2006

18

What you need to know about learning BN structures so far

- Decomposable scores
 - Maximum likelihood *Score*
 - Information theoretic interpretation
- Best tree (Chow-Liu)
- Best TAN
- Nearly best k-treewidth (in $O(N^{k+1})$)

10-708 - ©Carlos Guestrin 2006

19

Announcements

- Homework 2 out
 - Due Oct. 11th
- Project description out next week

10-708 - ©Carlos Guestrin 2006

20

Maximum likelihood score overfits!

$$\log \hat{P}(\mathcal{D} \mid \theta, \mathcal{G}) = M \sum_i \hat{I}(x_i, \text{Pa}_{x_i, \mathcal{G}}) - M \sum_i \hat{H}(X_i)$$

- Information never hurts:
- Adding a parent always increases score!!!

10-708 – ©Carlos Guestrin 2006

21

Bayesian score

- Prior distributions:
 - Over structures
 - Over parameters of a structure
- Posterior over structures given data:

$$\log P(\mathcal{G} \mid D) \propto \log P(\mathcal{G}) + \log \int_{\theta_{\mathcal{G}}} P(D \mid \mathcal{G}, \theta_{\mathcal{G}}) P(\theta_{\mathcal{G}} \mid \mathcal{G}) d\theta_{\mathcal{G}}$$

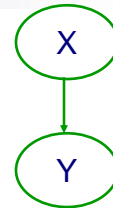
10-708 – ©Carlos Guestrin 2006

22

Bayesian score and model complexity

$$\log P(D | \mathcal{G}) = \log \int_{\theta_{\mathcal{G}}} P(D | \mathcal{G}, \theta_{\mathcal{G}}) P(\theta_{\mathcal{G}} | \mathcal{G}) d\theta_{\mathcal{G}}$$

True model:



- Structure 1: X and Y independent

- Score doesn't depend on alpha

- Structure 2: $X \rightarrow Y$

$$P(Y=t|X=t) = 0.5 + \alpha$$

$$P(Y=t|X=f) = 0.5 - \alpha$$

- Data points split between $P(Y=t|X=t)$ and $P(Y=t|X=f)$
- For fixed M, only worth it for large α
 - Because posterior over parameter will be more diffuse with less data

10-708 - ©Carlos Guestrin 2006

23

Bayesian, a decomposable score

$$\log P(D | \mathcal{G}) = \log \int_{\theta_{\mathcal{G}}} P(D | \mathcal{G}, \theta_{\mathcal{G}}) P(\theta_{\mathcal{G}} | \mathcal{G}) d\theta_{\mathcal{G}}$$

- As with last lecture, assume:

- Local and global parameter independence

- Also, prior satisfies **parameter modularity**:

- If X_i has same parents in \mathcal{G} and \mathcal{G}' , then parameters have same prior

- Finally, structure prior $P(\mathcal{G})$ satisfies **structure modularity**

- Product of terms over families
- E.g., $P(\mathcal{G}) \propto c^{|\mathcal{G}|}$

- Bayesian score decomposes along families!

10-708 - ©Carlos Guestrin 2006

24

BIC approximation of Bayesian score

- Bayesian has difficult integrals
- For Dirichlet prior, can use simple Bayes information criterion (BIC) approximation
 - In the limit, we can forget prior!
 - **Theorem:** for Dirichlet prior, and a BN with $\text{Dim}(\mathcal{G})$ independent parameters, as $M \rightarrow \infty$:

$$\log P(D | \mathcal{G}) = \log P(D | \mathcal{G}, \theta_{\mathcal{G}}) - \frac{\log M}{2} \text{Dim}(\mathcal{G}) + O(1)$$

10-708 – ©Carlos Guestrin 2006

25

BIC approximation, a decomposable score

- BIC: $\text{Score}_{\text{BIC}}(\mathcal{G} : D) = \log P(D | \mathcal{G}, \theta_{\mathcal{G}}) - \frac{\log M}{2} \text{Dim}(\mathcal{G})$

- Using information theoretic formulation:

$$\text{Score}_{\text{BIC}}(\mathcal{G} : D) = M \sum_i \hat{I}(x_i, \text{Pa}_{x_i, \mathcal{G}}) - M \sum_i \hat{H}(X_i) - \frac{\log M}{2} \sum_i \text{Dim}(P(X_i | \text{Pa}_{x_i, \mathcal{G}}))$$

10-708 – ©Carlos Guestrin 2006

26

Consistency of BIC and Bayesian scores

Consistency is limiting behavior, says nothing about finite sample size!!!

- A scoring function is **consistent** if, for true model G^* , as $M \rightarrow \infty$, with probability 1
 - G^* maximizes the score
 - All structures **not I-equivalent** to G^* have strictly lower score
- **Theorem:** BIC score is consistent
- **Corollary:** the Bayesian score is consistent
- What about maximum likelihood score?

10-708 - ©Carlos Guestrin 2006

27

Priors for general graphs

- For finite datasets, prior is important!
- Prior over structure satisfying prior modularity
- What about prior over parameters, how do we represent it?
 - *K2 prior:* fix an α , $P(\theta_{X|jPaXi}) = \text{Dirichlet}(\alpha, \dots, \alpha)$
 - K2 is "inconsistent"

10-708 - ©Carlos Guestrin 2006

28

BDe prior

- Remember that Dirichlet parameters analogous to “fictitious samples”
- Pick a fictitious sample size m'
- For each possible family, define a prior distribution $P(X_i, \mathbf{Pa}_{X_i})$
 - Represent with a BN
 - Usually independent (product of marginals)
- **BDe prior:**
- Has “consistency property”:

10-708 – ©Carlos Guestrin 2006

29

Score equivalence

- If G and G' are I-equivalent then they have same score
- **Theorem 1:** Maximum likelihood score and BIC score satisfy score equivalence
- **Theorem 2:**
 - If $P(G)$ assigns same prior to I-equivalent structures (e.g., edge counting)
 - and parameter prior is dirichlet
 - then **Bayesian score satisfies score equivalence if and only if** prior over parameters represented as a **BDe prior!!!!!!**

10-708 – ©Carlos Guestrin 2006

30

Chow-Liu for Bayesian score

- Edge weight $w_{X_j \rightarrow X_i}$ is advantage of adding X_j as parent for X_i
- Now have a directed graph, need directed spanning forest
 - Note that adding an edge can hurt Bayesian score – choose forest not tree
 - But, if score satisfies score equivalence, then $w_{X_j \rightarrow X_i} = w_{X_i \rightarrow X_j}$!
 - Simple maximum spanning forest algorithm works

10-708 – ©Carlos Guestrin 2006

31

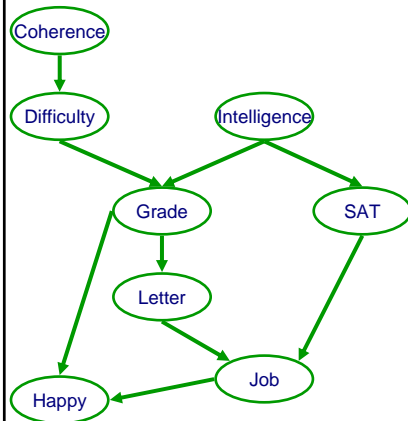
Structure learning for general graphs

- In a tree, a node only has one parent
- **Theorem:**
 - The problem of learning a BN structure with at most d parents is **NP-hard for any (fixed) $d \geq 2$**
- Most structure learning approaches use heuristics
 - Exploit score decomposition
 - (Quickly) Describe two heuristics that exploit decomposition in different ways

10-708 – ©Carlos Guestrin 2006

32

Understanding score decomposition



10-708 - ©Carlos Guestrin 2006

33

Fixed variable order 1

- Pick a variable order \prec
 - e.g., X_1, \dots, X_n
- X_i can only pick parents in $\{X_1, \dots, X_{i-1}\}$
 - Any subset
 - Acyclicity guaranteed!
- Total score = sum score of each node

10-708 - ©Carlos Guestrin 2006

34

Fixed variable order 2

- Fix max number of parents to k
- For each i in order \prec
 - Pick $\mathbf{Pa}_{X_i} \subseteq \{X_1, \dots, X_{i-1}\}$
 - Exhaustively search through all possible subsets
 - \mathbf{Pa}_{X_i} is maximum $\mathbf{U} \subseteq \{X_1, \dots, X_{i-1}\}$ FamScore($X_i | \mathbf{U} : D$)
- Optimal BN for each order!!!
- Greedy search through space of orders:
 - E.g., try switching pairs of variables in order
 - If neighboring vars in order are switch, only need to recompute score for this pair
 - $O(n)$ speed up per iteration
 - Local moves may be worse

10-708 – ©Carlos Guestrin 2006

35

Learn BN structure using local search

Starting from
Chow-Liu tree

Local search,
possible moves:
Only if acyclic!!!

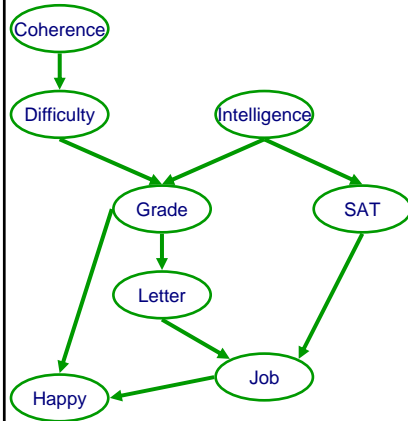
- Add edge
- Delete edge
- Invert edge

**Select using
favorite score**

10-708 – ©Carlos Guestrin 2006

36

Exploit score decomposition in local search



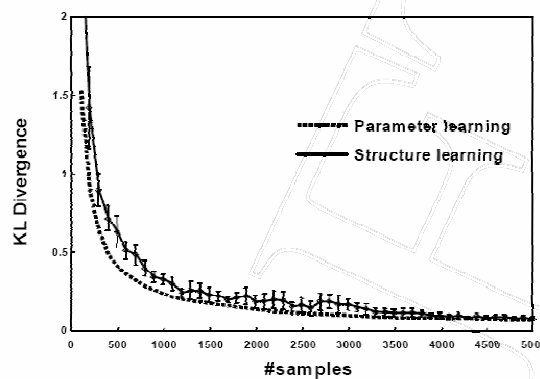
- Add edge and delete edge:
 - Only rescore one family!

- Reverse edge
 - Rescore only two families

10-708 — ©Carlos Guestrin 2006

37

Some experiments



Alarm network

10-708 — ©Carlos Guestrin 2006

38

Order search versus graph search

■ Order search advantages

- For fixed order, optimal BN – more “global” optimization
- Space of orders much smaller than space of graphs

■ Graph search advantages

- Not restricted to k parents
 - Especially if exploiting CPD structure, such as CSI
- Cheaper per iteration
- Finer moves within a graph

10-708 – ©Carlos Guestrin 2006

39

Bayesian model averaging

- So far, we have selected a single structure
- But, if you are really Bayesian, must average over structures

- Similar to averaging over parameters

$$\log P(D | \mathcal{G}) = \log \int_{\theta_{\mathcal{G}}} P(D | \mathcal{G}, \theta_{\mathcal{G}}) P(\theta_{\mathcal{G}} | \mathcal{G}) d\theta_{\mathcal{G}}$$

- Inference for structure averaging is very hard!!!
 - Clever tricks in reading

10-708 – ©Carlos Guestrin 2006

40

What you need to know about learning BN structures

- Decomposable scores
 - Data likelihood
 - Information theoretic interpretation
 - Bayesian
 - BIC approximation
- Priors
 - Structure and parameter assumptions
 - BDe if and only if score equivalence
- Best tree (Chow-Liu)
- Best TAN
- Nearly best k-treewidth (in $O(N^{k+1})$)
- Search techniques
 - Search through orders
 - Search through structures
- Bayesian model averaging