

Readings:

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

Structure Learning in BNs 2: (the good,) the bad, the ugly

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

October 4th, 2006

1

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: $I(X_i, \text{Pa}_{X_i, \mathcal{G}}) = H(X_i) - H(X_i \mid \text{Pa}_{X_i, \mathcal{G}})$

$H(X \mid A) \geq H(X \mid A \cup Y)$

$I(X_i; \text{Pa}_{X_i, \mathcal{G}}) \leq I(X_i; \text{Pa}_{X_i, \mathcal{G}} \cup Y)$

\Rightarrow pick as many parents as possible

\Rightarrow fully connected graph if optimize

\Rightarrow Overfit!!!

likelihood

■ Adding a parent always increases score!!!

doesn't decrease

Bayesian score

prior $\rightarrow P(\theta_G, G) \stackrel{\text{chain rule}}{=} P(G) \cdot P(\theta_G | G)$
 simpler to express

Prior distributions:

Over structures

- hope for simpler structures
 - "regularize" for simpler structures

Over parameters of a structure

Posterior over structures given data:

$P(G|D) \propto P(D|G) \cdot P(G)$
 posterior likelihood term prior over graphs

$$P(D|G) = \int_{\theta_G} P(D|G, \theta_G) \cdot P(\theta_G | G) d\theta_G$$

prior over parameters

$$\log P(G | D) \propto \log P(G) + \log \int_{\theta_G} P(D | G, \theta_G) P(\theta_G | G) d\theta_G$$

10-708 - ©Carlos Guestrin 2006

3

Bayesian score and model complexity

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

Structure 1: X and Y independent

$$\log P(D|G) = \log \int_{\theta_x} P(D_x | \theta_x) \cdot P(\theta_x) d\theta_x + \log \int_{\theta_y} P(D_y | \theta_y) P(\theta_y) d\theta_y$$

Score doesn't depend on alpha

Structure 2: $X \rightarrow Y$

$$\log P(D|G) = \log \int_{\theta_x} P(D_x | \theta_x) P(\theta_x) d\theta_x$$

$$+ \log \int_{\theta_{y|x=t}} P(D_{y|x=t} | \theta_{y|x=t}) P(\theta_{y|x=t}) d\theta_{y|x=t}$$

Data points split between $P(Y=t|X=t)$ and $P(Y=t|X=f)$

For fixed M, only worth it for large α

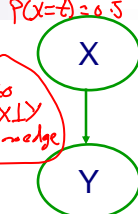
Because posterior over parameter will be more diffuse with less data

$\log P(D_{y|x=t} | G) \dots$ Dirichlet posterior

less data

likelihood $P(D|G)$ smaller

True model:



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

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

$$P(Y=f) = P(X=t) \cdot P(Y=f|X=t) + P(X=f) \cdot P(Y=f|X=f)$$

$$= 0.5$$

10-708 - ©Carlos Guestrin 2006

4

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 G and G' , then parameters have same prior
 - Finally, structure prior $P(G)$ satisfies **structure modularity**
 - Product of terms over families $P(G) \propto \prod_i P(E(X_i, \text{Par}_i))$
 - E.g., $P(G) \propto c^{|G|}$
 - Bayesian score decomposes along families!
- $\log P(G|D) \propto \sum_i \text{Score Fam}(X_i | \text{Par}_i : D)$
- Handwritten notes:*
 $\theta_{x|k} \perp \theta_x$
 $G \rightarrow G' = \dots$
 $\theta_{x_i} \perp \text{Par}_{x_i}$
 $|G| = \text{number edges}$
 $|G| < 1$
 \uparrow like fewer edges

10-708 - ©Carlos Guestrin 2006

5

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 $\text{Dim}(G)$ independent parameters, as $m \rightarrow \infty$:
- $$\log P(D | \mathcal{G}) \approx \log P(D | \mathcal{G}, \theta_{\mathcal{G}}) - \frac{\log m}{2} \text{Dim}(\mathcal{G}) + O(1)$$
- Handwritten notes:*
 $\log P(G|D) \propto \log P(G) + \log P(D|G)$
 \rightarrow likelihood score
 likes fully connected graph
 regularization
 penalizing
 likes simple graph

10-708 - ©Carlos Guestrin 2006

6

BIC approximation, a decomposable score

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

$\text{Dim}(\mathcal{G}) = \# \text{ params}$
 $= \sum_i \text{Dim}(P(X_i | \text{Pa}_{X_i})) = \sum_i (K_i - 1) \cdot K^{\text{Pa}_{X_i}}$
 if all vars take on K values exponentially unhappy with # parents

- Using information theoretic formulation:

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

$$= \sum_i \left[\underbrace{m \hat{I}(X_i, \text{Pa}_{X_i, \mathcal{G}}) - m \hat{H}(X_i)}_{\text{Score}_{\text{Fam BIC}}(X_i | \text{Pa}_{X_i} : D)} - \frac{\log m}{2} \text{Dim}(P(X_i | \text{Pa}_{X_i, \mathcal{G}})) \right]$$

10-708 - ©Carlos Guestrin 2006

7

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 \mathcal{G}^* , as $m \rightarrow \infty$, with probability 1
- \mathcal{G}^* maximizes the score
 - All structures **not I-equivalent** to \mathcal{G}^* have strictly lower score

- **Theorem:** BIC score is consistent ||| as $m \rightarrow \infty$

- **Corollary:** the Bayesian score is consistent with Dirichlet prior

- What about maximum likelihood score? true will max score
NO! - but fully connected graph, not I-equivalent also maximizes score

10-708 - ©Carlos Guestrin 2006

8

Priors for general graphs

- For finite datasets, prior is important!
- Prior over structure satisfying prior modularity
e.g. $P(G) \propto c^{|G|}$ $C \leq 1$
- What about prior over parameters, how do we represent it?
 - K2 prior: fix an α , $P(\theta_{X|PaX}) = \text{Dirichlet}(\alpha, \dots, \alpha)$
 - K2 is "inconsistent"

for more parents
prior is stronger ...
doesn't make sense ...

for any choice of parents $P(\theta_{X|PaX})$

no parents	$K\alpha$
one parent $U = \epsilon$ $U = f$	$K\alpha$ } $2K\alpha$
d parents $\forall i: PaX_i = K$	$K\alpha \cdot K^d$

K values of x_i if $U = PaX_i$
for all $u \in \text{Val}\{U\}$
equivalent sample size

10-708 - ©Carlos Guestrin 2006

9

BDe prior

- Remember that Dirichlet parameters analogous to "fictitious samples" (equivalent sample size)
- Pick a fictitious sample size m' (10)
- For each possible family, define a prior distribution $P'(X_i, Pa_{X_i})$
 - Represent with a BN
 - Usually independent (product of marginals) $P'(X_i, Pa_{X_i}) = P'(x_i) \prod_{x_j \in Pa_{X_i}} P'(x_j)$
- **BDe prior:**
 $P(\theta_{X_i|PaX_i=u}) = \text{Dirichlet}(m' P'(x_i=1, Pa_{X_i}=u), \dots)$
 $m' : P'(x_i=k, Pa_{X_i}=u)$
typically: $\forall i: P'(x_i) = \text{uniform}$
- Has "consistency property":

10-708 - ©Carlos Guestrin 2006

10

Announcements

30% of your grade

■ Project description is out on class website:

- Individual or groups of two only
- Suggested projects on the class website, or do something related to your research (preferable)
 - Must be something you started this semester
 - The semester goes really quickly, so be realistic (and ambitious ☺)

■ Project deliverables:

- one page proposal due next week (10/11)
- 5-page milestone report Nov. 1st
- Poster presentation on Dec. 1st, 3-6pm
- Write up, 8-pages, due Dec. 8th
- All write ups in NIPS format (see class website), page limits are strict

■ Objective:

- Explore and apply concepts in **probabilistic graphical models**
- Doing a fun project!

10-708 – ©Carlos Guestrin 2006

11

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

12

Chow-Liu for Bayesian score

- Edge weight $w_{X_j \rightarrow X_i}$ is advantage of adding X_j as parent for X_i

Bayesian
chow-Liu
→ could end
up with a forest

e.g. if all vars
independent,
then "regularization"
stop from adding edges

$\frac{m}{2} \log \frac{I(X_i, X_j)}{I(X_i)I(X_j)}$ X_1
 X_3

advantage of adding X_2 as parent for X_4
 X_4 advantage of
adding X_4 as
parent for X_2

BIC edge score: $w_{ji} = m \hat{I}(X_i, X_j) - \log \text{Dir}(\frac{1}{2} \text{Dir}(x_{ij}))$
if $w_{ji} \leq 0$, MST never pick edge.

- 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 undirected

10-708 – ©Carlos Guestrin 2006

13

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) $d \geq 2$

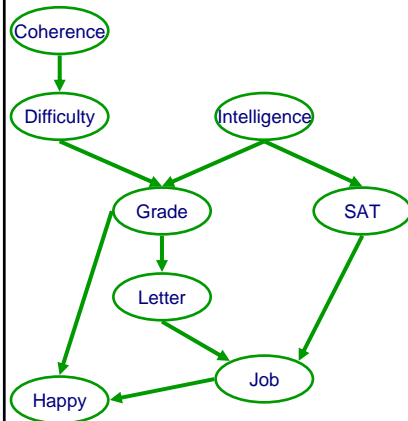
- Most structure learning approaches use heuristics

- Exploit score decomposition
- (Quickly) Describe two heuristics that exploit decomposition in different ways

10-708 – ©Carlos Guestrin 2006

14

Understanding score decomposition



10-708 - ©Carlos Guestrin 2006

15

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

16

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 $\mathbf{U} \subseteq \{X_1, \dots, X_{i-1}\} \text{ FamScore}(X_i | \mathbf{U} : D)$
- 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

17

Learn BN structure using local search

Starting from
Chow-Liu tree

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

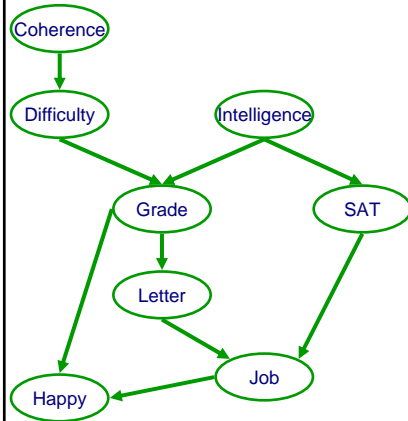
- Add edge
- Delete edge
- Invert edge

**Select using
favorite score**

10-708 – ©Carlos Guestrin 2006

18

Exploit score decomposition in local search



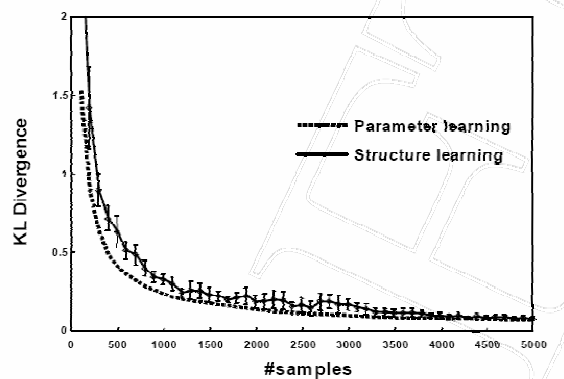
- Add edge and delete edge:
 - Only rescore one family!

- Reverse edge
 - Rescore only two families

10-708 — ©Carlos Guestrin 2006

19

Some experiments



Alarm network

10-708 — ©Carlos Guestrin 2006

20

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

21

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

22

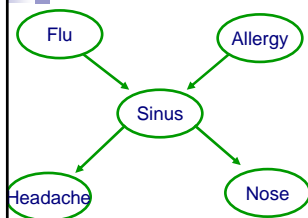
What you need to know about learning BN structures

- 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

10-708 – ©Carlos Guestrin 2006

23

Inference in graphical models: Typical queries 1

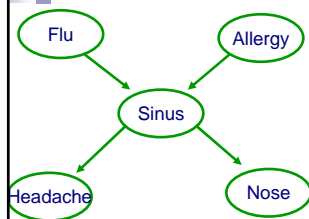


- Conditional probabilities
 - Distribution of some var(s). given evidence

10-708 – ©Carlos Guestrin 2006

24

Inference in graphical models: Typical queries 2 – Maximization



- Most probable explanation (MPE)
 - Most likely assignment to all hidden vars given evidence
- Maximum a posteriori (MAP)
 - Most likely assignment to some var(s) given evidence

10-708 – ©Carlos Guestrin 2006

25

Are MPE and MAP Consistent?



$P(S=t)=0.4$
 $P(S=f)=0.6$

$P(N|S)$

- Most probable explanation (MPE)
 - Most likely assignment to all hidden vars given evidence
- Maximum a posteriori (MAP)
 - Most likely assignment to some var(s) given evidence

10-708 – ©Carlos Guestrin 2006

26

Complexity of conditional probability queries 1

- How hard is it to compute $P(X|\mathbf{E}=\mathbf{e})$?

Reduction – 3-SAT

$$(\bar{X}_1 \vee X_2 \vee X_3) \wedge (\bar{X}_2 \vee X_3 \vee X_4) \wedge \dots$$

10-708 – ©Carlos Guestrin 2006

27

Complexity of conditional probability queries 2

- How hard is it to compute $P(X|\mathbf{E}=\mathbf{e})$?
 - At least NP-hard, but even harder!

10-708 – ©Carlos Guestrin 2006

28

Inference is #P-hard, hopeless?

- Exploit structure!
- Inference is hard in general, but easy for many (real-world relevant) BN structures

10-708 – ©Carlos Guestrin 2006

29

What about the maximization problems? First, most probable explanation (MPE)

- What's the complexity of MPE?

10-708 – ©Carlos Guestrin 2006

30

What about maximum a posteriori?

- At least, as hard as MPE!
- Actually, much harder!!! NP^{PP} -complete!

10-708 – ©Carlos Guestrin 2006

31

Can we exploit structure for maximization?

- For MPE
- For MAP

10-708 – ©Carlos Guestrin 2006

32

Exact inference is hard, what about approximate inference?

- Must define approximation criterion!
- Relative error of $\varepsilon > 0$
- Absolute error of $\varepsilon > 0$

10-708 – ©Carlos Guestrin 2006

33

Hardness of approximate inference

- Relative error of ε
- Absolute error of ε

10-708 – ©Carlos Guestrin 2006

34

What you need to know about inference

■ Types of queries

- ☐ probabilistic inference
- ☐ most probable explanation (MPE)
- ☐ maximum a posteriori (MAP)
 - MPE and MAP are truly different (don't give the same answer)

■ Hardness of inference

- ☐ Exact and approximate inference are NP-hard
- ☐ MPE is NP-complete
- ☐ MAP is much harder (NP^{PP} -complete)