

Readings:

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

K&F: 7 (overview of inference)

K&F: 8.1, 8.2 (Variable Elimination)

Structure Learning in BNs 3:

(the good, the bad,) and, finally, the ugly

Inference

we now get to use the BNs!!

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

October 6th, 2006

1

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 $P(\mathcal{G}) \propto \prod P(E(X_i, \text{Par}_{X_i}))$
- E.g., $P(\mathcal{G}) \propto c^{|\mathcal{G}|}$

■ Bayesian score decomposes along families!

$$\log P(\mathcal{G} | D) \propto \sum_i \text{Score Fam}(X_i | \text{Par}_{X_i}, D)$$

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

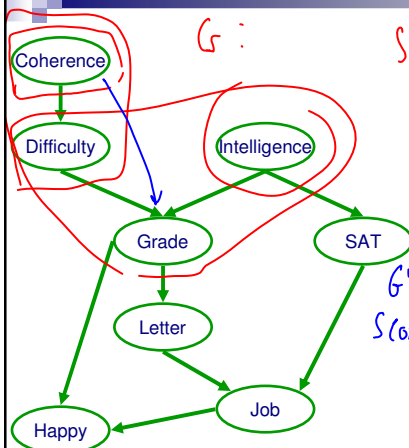
- Exploit score decomposition
- (Quickly) Describe two heuristics that exploit decomposition in different ways

10-708 - ©Carlos Guestrin 2006

3

Understanding score decomposition

$$\text{Score}_{\text{Fam}} = \text{SF}$$



$$\begin{aligned} \text{Score}(G:D) &= \sum_i \text{Score}_{\text{Fam}}(X_i | \text{Pa}_{X_i, G}:D) \\ &= \text{SF}(C|\emptyset:D) + \text{SF}(D|C:D) + \text{SF}(I|\emptyset:D) \\ &\quad + \text{SF}(G|D,I:D) + \dots \end{aligned}$$

G' Suppose C becomes parent of G

$$\begin{aligned} \text{Score}(G:D) - \text{Score}(G':D) \\ &= \text{SF}(G|D,I:D) - \text{SF}(G|D,I,C:D) \end{aligned}$$

10-708 - ©Carlos Guestrin 2006

4

Fixed variable order 1 F, A, S, N, H

- 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
- $G: X_i$ picked parents Pa_{X_i}
- $Score(G; D) = \sum_i SF(X_i | Pa_{X_i}; D)$
- DAG: $1 \dots i-1$
- $X_1 \rightarrow X_2$
 $X_3 \rightarrow X_2$
 $X_3 \rightarrow X_4$
 $X_4 \rightarrow X_2$
- Selecting parents of X_5
any subset of X_1, \dots, X_4
 \Rightarrow still have a DAG
- $Score(G; D) \Rightarrow$ indep. choice
 for each node
 is OPTIMAL!!
 max choice parents for each X_i

10-708 - ©Carlos Guestrin 2006

5

Fixed variable order 2

FASHN good order
HFANS bad order

- Fix max number of parents to k $\sum_{j=1}^k \binom{i-1}{j} \leq \frac{2^{i-1}}{2^{i-1}-1}$
 - For each i in order \prec
 - Pick $Pa_{X_i} \subseteq \{X_1, \dots, X_{i-1}\}$ complexity =
 - Exhaustively search through all possible subsets
 - Pa_{X_i} is maximum $U \subseteq \{X_1, \dots, X_{i-1}\}$ FamScore $(X_i | U; D)$
 - 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 switched, only need to recompute score for this pair
 - $O(n)$ speed up per iteration
- $\begin{matrix} F & A & S & H & N \\ \uparrow & & & & \\ F & S & A & H & N \end{matrix}$

10-708 - ©Carlos Guestrin 2006

6

Learn BN structure using local search

Starting from
Chow-Liu tree

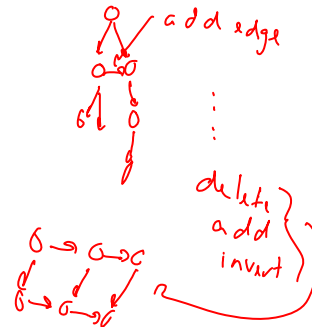


Local search,
possible moves:

- Only if acyclic!!!
- Add edge
- Delete edge
- Invert edge

Computed
locally/efficiently
because of
Score decomposition
Fam Score

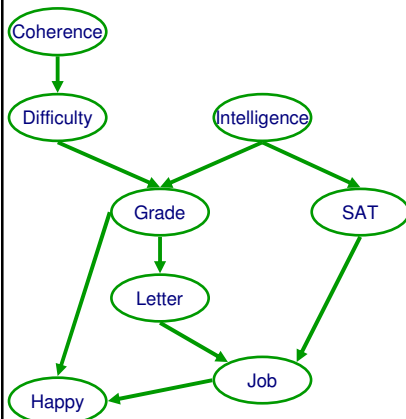
Select using
favorite score



10-708 - ©Carlos Guestrin 2006

7

Exploit score decomposition in local search



■ Add edge and delete edge:

- Only rescore one family!

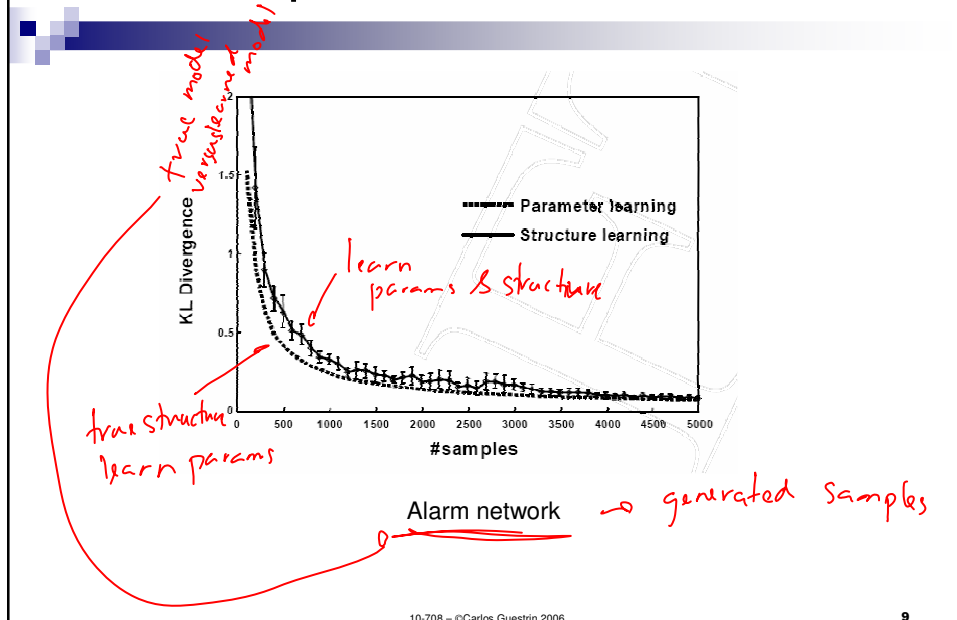
■ Reverse edge

- Rescore only two families

10-708 - ©Carlos Guestrin 2006

8

Some experiments



Order search versus graph search

- Order search advantages *random*
 - For fixed order, optimal BN – more “global” optimization
 - Space of orders much smaller than space of graphs
 $n!$ $2^{(2n^2)}$
- Graph search advantages
 - Not restricted to k parents
 - Especially if exploiting CPD structure, such as CSI
 - Cheaper per iteration
 - Finer moves within a graph

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

$P(\mathcal{G} | D) \rightarrow$ probability for each graph
- Inference for structure averaging is very hard!!!
 - Clever tricks in reading

10-708 – ©Carlos Guestrin 2006

11

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

12

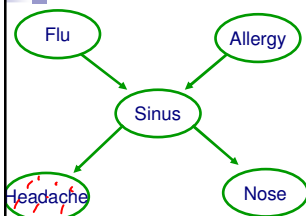
Announcements

- Don't forget project proposals due this Wednesday
- Special recitation on advanced topic:
 - Ajit Singh on Optimal Structure Learning
 - On Monday Oct 9, 5:30-7:00pm in Wean Hall 4615A

10-708 - ©Carlos Guestrin 2006

13

Inference in graphical models: Typical queries 1



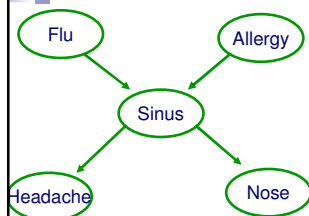
- Conditional probabilities
 - Distribution of some var(s). given evidence

evidence = t
$$P(F|e) = P(F|H=t)$$

10-708 - ©Carlos Guestrin 2006

14

Inference in graphical models: Typical queries 2 – Maximization



■ Most probable explanation (MPE)

- Most likely assignment to all hidden vars given evidence

$$\arg \max_{f, a, s, n} P(F=f, A=a, S=s, N=n | H=t)$$

■ Maximum a posteriori (MAP)

- Most likely assignment to some var(s) given evidence

$$\arg \max_f P(F=f | H=t)$$

10-708 – ©Carlos Guestrin 2006

15

Are MPE and MAP Consistent?



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

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

$$P(N|S) =$$

	t	f
t	0.9	0.1
f	0.5	0.5

■ Most probable explanation (MPE)

- Most likely assignment to all hidden vars given evidence

$$\arg \max_{s, n} P(S=s, N=n) = (S=t, N=t)$$

$$P(S=s, N=n) = P(S=s) \cdot P(N=n | S=s)$$

■ Maximum a posteriori (MAP)

- Most likely assignment to some var(s) given evidence

$$\arg \max_s P(s) = (S=f)$$

don't agree

10-708 – ©Carlos Guestrin 2006

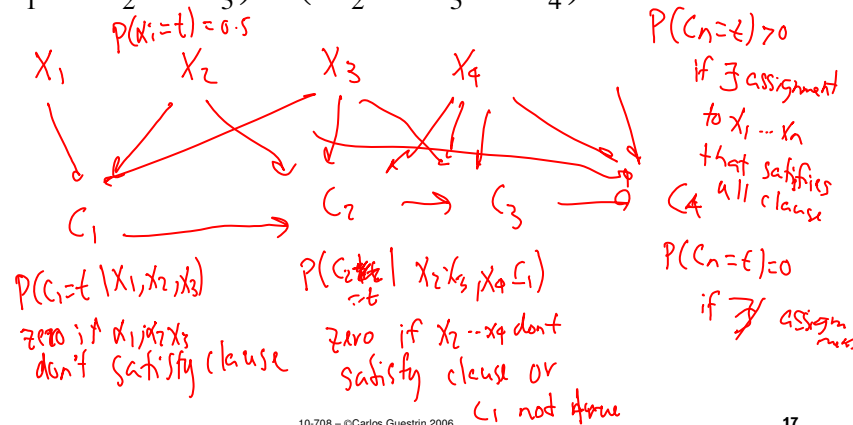
16

Complexity of conditional probability queries 1

- How hard is it to compute $P(X|E=e)$? NP-hard

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

17

Complexity of conditional probability queries 2

- How hard is it to compute $P(X|E=e)$? #P complete

- At least NP-hard, but even harder!

$P(C_n=t) = \alpha$
 have m literals $X_1 \dots X_m$
 $\text{prob}(X_1=t, X_2=f, \dots, X_m=t) = \frac{1}{2^m}$
 each satisfying assignment gives $C_n=t$ $\frac{1}{2^m}$ prob
 if $P(C_n=t) = \alpha \Rightarrow$ 3-SAT has $\alpha 2^m$ satisfying assignments
 \Rightarrow count # sat. assignments
 \Rightarrow #P-hard

10-708 – ©Carlos Guestrin 2006

18

Inference is ~~#P-hard~~^{Complete}, hopeless?

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

Structure "simple" \Rightarrow inference "easy"

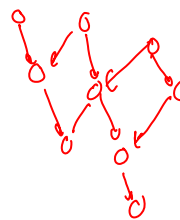
trees



poly-trees



low-treewidth graphs



10-708 - ©Carlos Guestrin 2006

19

Complexity for other inference questions

- Probabilistic inference
 - general graphs: ~~#P-complete~~
 - poly-trees and low tree-width: easy
- Approximate probabilistic inference ^{within $\epsilon > 0$}
 - Absolute error: $|\hat{P} - P| \leq \epsilon$
 - Relative error: $\frac{|\hat{P} - P|}{|P|} \leq \epsilon$

$\left. \begin{array}{l} \text{Absolute error} \\ \text{Relative error} \end{array} \right\} \text{NP-hard} \leftarrow \begin{array}{l} \epsilon < 0.5 \\ \text{any } \epsilon > 0 \end{array}$
- Most probable explanation (MPE)
 - general graphs: NP-complete
 - poly-trees and low tree-width: easy
- Maximum a posteriori (MAP)
 - general graphs: NP^{PP}-complete
 - poly-trees and low tree-width: NP-hard

10-708 - ©Carlos Guestrin 2006

20

Inference in BNs hopeless?

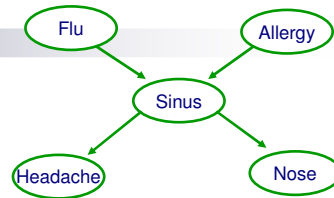
- In general, yes!
 - Even approximate!
- In practice
 - Exploit structure
 - Many effective approximation algorithms (some with guarantees)
- For now, we'll talk about exact inference
 - Approximate inference later this semester

10-708 - ©Carlos Guestrin 2006

21

General probabilistic inference

■ Query: $P(X | e)$



■ Using def. of cond. prob.:

$$P(X | e) = \frac{P(X, e)}{P(e)}$$

■ Normalization:

$$P(\underline{X} | e) \propto \underline{P(X, e)}$$

$$P(X=t, e) = 0.4$$

$$P(X=f, e) = 0.1$$

normalize

$$P(X=t | e) = \frac{0.4}{0.5} = 0.8$$

10-708 - ©Carlos Guestrin 2006

22

Marginalization

I know

$$P(F, S, N) = P(F) \cdot P(S|F) \cdot P(N|S)$$



marginalization

$$P(S, N=t) = \sum_f P(F=f, S, N=t)$$

$$= \sum_f P(F=f) \cdot P(S|F=f) \cdot P(N=t|S)$$

$$P(S=t|N=t)$$

first focus on

$$\begin{cases} P(S=t, N=t) \\ P(S=f, N=t) \end{cases}$$

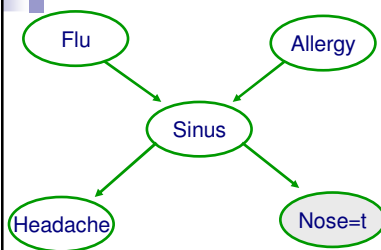
$$\rightarrow P(S, N=t)$$

↑
computing table
for each value
of S

10-708 - ©Carlos Guestrin 2006

23

Probabilistic inference example

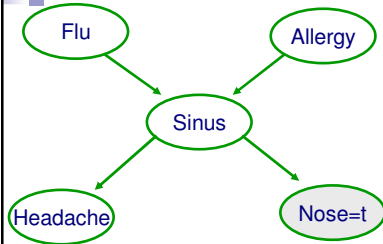


Inference seems exponential in number of variables!

10-708 - ©Carlos Guestrin 2006

24

Fast probabilistic inference example – Variable elimination

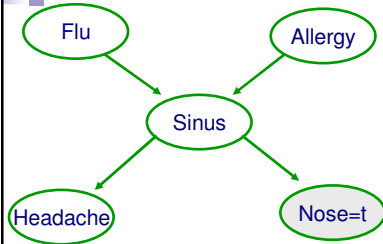


(Potential for) Exponential reduction in computation!

Understanding variable elimination – Exploiting distributivity



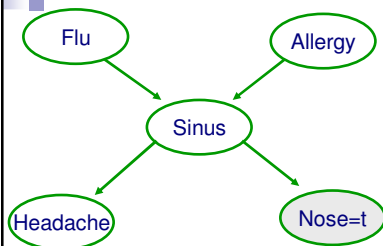
Understanding variable elimination – Order can make a HUGE difference



10-708 – ©Carlos Guestrin 2006

27

Understanding variable elimination – Intermediate results

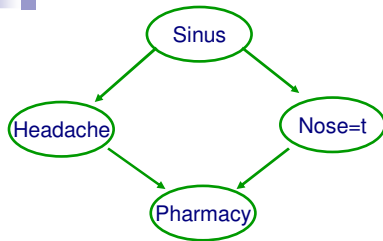


Intermediate results are probability distributions

10-708 – ©Carlos Guestrin 2006

28

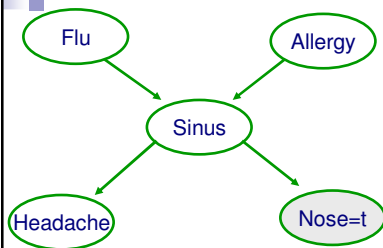
Understanding variable elimination – Another example



10-708 – ©Carlos Guestrin 2006

29

Pruning irrelevant variables



Prune all non-ancestors of query variables
More generally: Prune all nodes not on active trail between evidence and query vars

10-708 – ©Carlos Guestrin 2006

30

Variable elimination algorithm

- Given a BN and a query $P(X|e) \propto P(X,e)$
- Instantiate evidence \mathbf{e}
- Prune non-active vars for $\{X, \mathbf{e}\}$
- Choose an ordering on variables, e.g., X_1, \dots, X_n
- Initial *factors* $\{f_1, \dots, f_n\}$: $f_i = P(X_i | \mathbf{Pa}_{X_i})$ (CPT for X_i)
- For $i = 1$ to n , If $X_i \notin \{X, \mathbf{E}\}$
 - Collect factors f_1, \dots, f_k that include X_i
 - Generate a new factor by eliminating X_i from these factors

IMPORTANT!!!

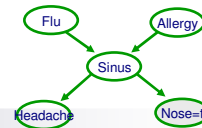
$$g = \sum_{X_i} \prod_{j=1}^k f_j$$

- Variable X_i has been eliminated!
- Normalize $P(X, \mathbf{e})$ to obtain $P(X|\mathbf{e})$

10-708 – ©Carlos Guestrin 2006

31

Operations on factors



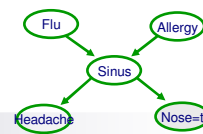
$$g = \sum_{X_i} \prod_{j=1}^k f_j$$

Multiplication:

10-708 – ©Carlos Guestrin 2006

32

Operations on factors



$$g = \sum_{X_i} \prod_{j=1}^k f_j$$

Marginalization:

10.708 – ©Carlos Guestrin 2006

33

Complexity of VE – First analysis

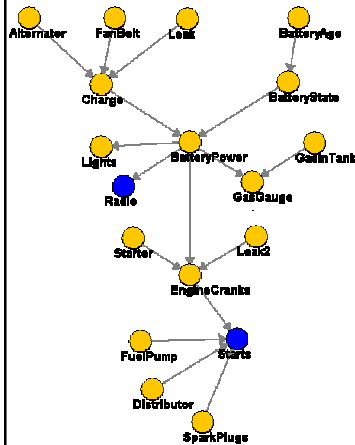
- Number of multiplications:

- Number of additions:

10.708 – ©Carlos Guestrin 2006

34

Complexity of variable elimination – (Poly)-tree graphs



Variable elimination order:

Start from “leaves” inwards:

- Start from skeleton!
- Choose a “root”, any node
- Find topological order for root
- Eliminate variables in reverse order

Linear in CPT sizes!!! (versus exponential)

35

What you need to know about inference thus far

- 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 (NPP-complete)
- Variable elimination algorithm
 - Eliminate a variable:
 - Combine factors that include this var into single factor
 - Marginalize var from new factor
 - Efficient algorithm (“only” exponential in induced-width, not number of variables)
 - If you hear: “Exact inference only efficient in tree graphical models”
 - You say: “No!!! Any graph with low induced width”
 - And then you say: “And even some with very large induced-width” (next week with context-specific independence)
- Elimination order is important!
 - NP-complete problem
 - Many good heuristics

10-708 – ©Carlos Guestrin 2006

36