

Readings:

K&F: 17.3, 17.4, 17.5.1, 8.1, 12.1

Structure Learning (The Good), The Bad, The Ugly

Inference

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

October 13th, 2008

10-708 – ©Carlos Guestrin 2006-2008

1

Decomposable score

■ Log data likelihood

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

■ Decomposable score:

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

for MLE $\text{FamScore}(X_i \mid \mathbf{Pa}_{X_i} : D) = m \hat{I}(X_i \mid \mathbf{Pa}_{X_i}) - m \hat{H}(X_i)$

10-708 – ©Carlos Guestrin 2006-2008

2

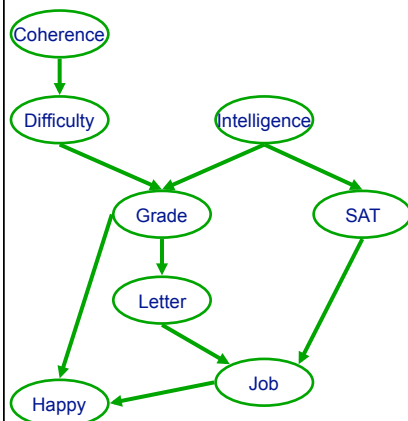
Structure learning for general graphs

- 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 2006-2008

3

Understanding score decomposition



10-708 – ©Carlos Guestrin 2006-2008

4

Fixed variable order 1

- Pick a variable order
 - 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 2006-2008

5

Fixed variable order 2

- Fix max number of parents to k
- For each i in order
 - 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 switched, only need to recompute score for this pair
 - $O(n)$ speed up per iteration

10-708 – ©Carlos Guestrin 2006-2008

6

Learn BN structure using local search



Starting from
Chow-Liu tree

Local search,
possible moves:

Only if acyclic!!!

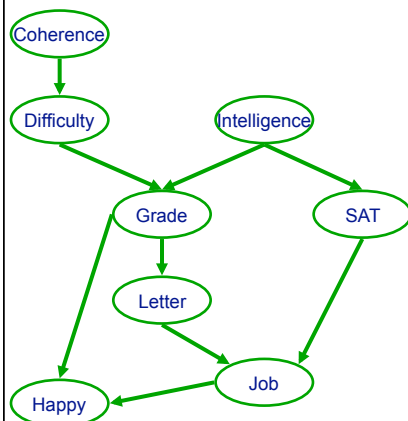
- Add edge
- Delete edge
- Invert edge

Select using
favorite score

10-708 – ©Carlos Guestrin 2006-2008

7

Exploit score decomposition in local search



■ Add edge and delete edge:

- Only rescore one family!

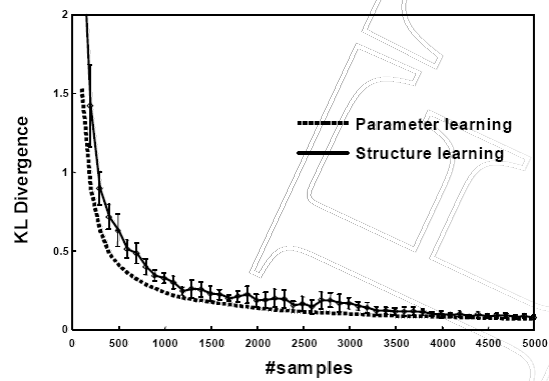
■ Reverse edge

- Rescore only two families

10-708 – ©Carlos Guestrin 2006-2008

8

Some experiments



Alarm network

10-708 – ©Carlos Guestrin 2006-2008

9

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 2006-2008

10

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 2006-2008

11

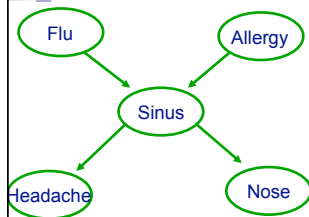
What you need to know about learning BN structures

- 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 2006-2008

12

Inference in graphical models: Typical queries 1



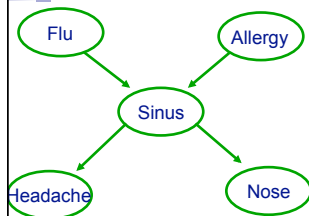
- Conditional probabilities

- Distribution of some var(s). given evidence

10-708 – ©Carlos Guestrin 2006-2008

13

Inference in graphical models: Typical queries 2 – Maximization



- Most probable explanation (MPE)

- Most likely assignment to all hidden vars given evidence

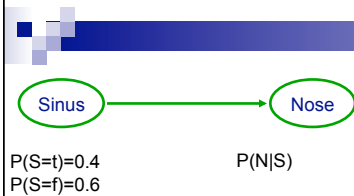
- Maximum a posteriori (MAP)

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

10-708 – ©Carlos Guestrin 2006-2008

14

Are MPE and MAP Consistent?



- 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 2006-2008

15

C++ Library

- Now available, join:

- ☐ <http://groups.google.com/group/10708-f08-code/>

- The library implements the following functionality:

- ☐ random variables, random processes, and linear algebra
- ☐ factorized distributions, such Gaussians, multinomial distributions, and mixtures
- ☐ graph structures and basic graph algorithms
- ☐ graphical models, including Bayesian networks, Markov networks, and junction trees
- ☐ basic static and dynamic inference algorithms
- ☐ parameter learning for Gaussian distributions, Chow Liu

- Fairly advanced C++ (not for everyone ☺)

10-708 – ©Carlos Guestrin 2006-2008

16

Complexity of conditional probability queries 1

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

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 2006-2008

17

Complexity of conditional probability queries 2

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

10-708 – ©Carlos Guestrin 2006-2008

18

Inference is #P-complete, hopeless?

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

10-708 – ©Carlos Guestrin 2006-2008

19

Complexity for other inference questions

- Probabilistic inference
 - general graphs:
 - poly-trees and low tree-width:
- Approximate probabilistic inference
 - Absolute error:
 - Relative error:
- Most probable explanation (MPE)
 - general graphs:
 - poly-trees and low tree-width:
- Maximum a posteriori (MAP)
 - general graphs:
 - poly-trees and low tree-width:

10-708 – ©Carlos Guestrin 2006-2008

20

Inference in BNs hopeless?

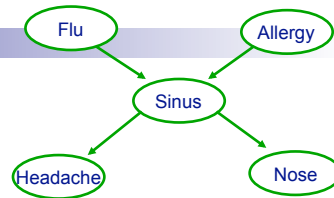
- In general, yes!
 - Even approximate!
- In practice
 - Exploit structure
 - Many effective approximation algorithms (some with guarantees)
- For now, we'll talk about exact inference
 - Approximate inference later this semester

10-708 – ©Carlos Guestrin 2006-2008

21

General probabilistic inference

- Query: $P(X | e)$



- Using def. of cond. prob.:

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

- Normalization:

$$P(X | e) \propto P(X, e)$$

10-708 – ©Carlos Guestrin 2006-2008

22

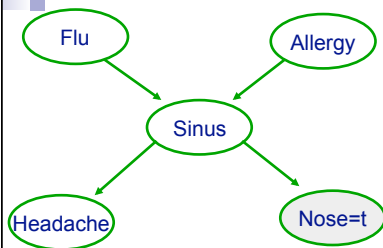
Marginalization



10-708 – ©Carlos Guestrin 2006-2008

23

Probabilistic inference example

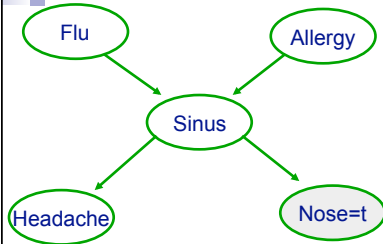


Inference seems exponential in number of variables!

10-708 – ©Carlos Guestrin 2006-2008

24

Fast probabilistic inference example – Variable elimination



(Potential for) Exponential reduction in computation!

10-708 – ©Carlos Guestrin 2006-2008

25

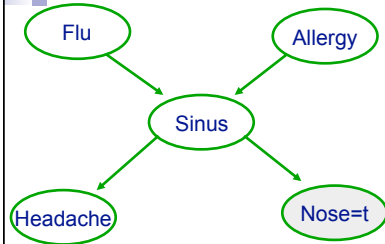
Understanding variable elimination – Exploiting distributivity



10-708 – ©Carlos Guestrin 2006-2008

26

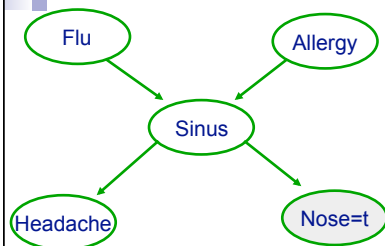
Understanding variable elimination – Order can make a HUGE difference



10-708 – ©Carlos Guestrin 2006-2008

27

Understanding variable elimination – Intermediate results

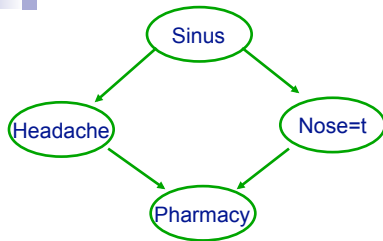


Intermediate results are probability distributions

10-708 – ©Carlos Guestrin 2006-2008

28

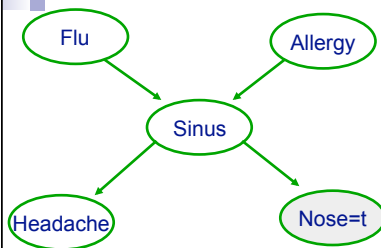
Understanding variable elimination – Another example



10-708 – ©Carlos Guestrin 2006-2008

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 2006-2008

30