# Structure Learning *finally*
## (The Good), The Bad, The Ugly

## Inference

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

October 13th, 2008

---

# 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!!!
  - Score($G : D$) = $\sum_i$ FamScore($X_i | \mathbf{Pa}_{X_i} : D$)

for MLE $\quad$ Fam Score $(X_i | Pa_{X_i} : D) = m\hat{I}(X_i; Pa_{X_i}) - m\hat{H}(X_i)$

# Structure learning for general graphs
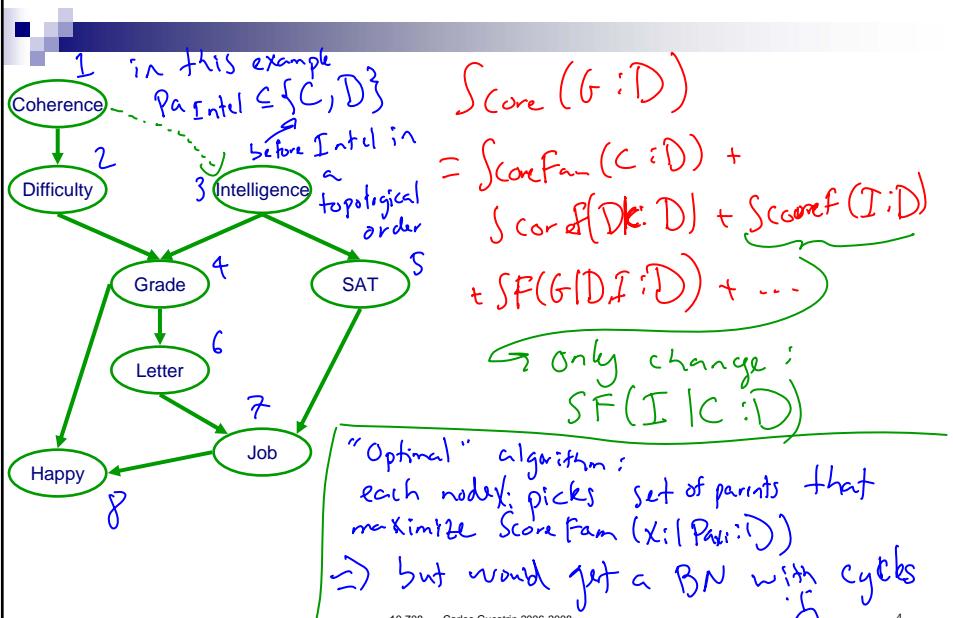
- In a tree, a node only has one parent

  *Chow-Liu*

- **Theorem**:
  - The problem of learning a BN structure with at most *d* parents is NP-hard for any (fixed) *d≥2*

- Most structure learning approaches use heuristics
  - Exploit score decomposition
  - (Quickly) Describe two heuristics that exploit decomposition in different ways

---

# Understanding score decomposition



1 in this example
Pa Intel ⊆ {C, D}
before Intel in a topological order

Coherence
2
Difficulty   3 Intelligence
4 Grade   SAT 5
6 Letter
7
Job
Happy
8

$Score(G : D)$
$= ScoreFam(C : D) +$
$Scoref(D|C : D) + Scoref(I : D)$
$+ SF(G|D,I : D) + \cdots$

→ only change :
$SF(I | C : D)$

"Optimal" algorithm :
each node $X_i$ picks set of parents that maximize $ScoreFam(X_i | Pa_{X_i} : D)$
⟹ but would get a BN with cycles

# Fixed variable order 1

*max number of parents = d*

- Pick a variable order ⟵
  - *just*
  - e.g., $X_1,\ldots,X_n$
- $X_i$ can only pick parents in $\{X_1,\ldots,X_{i-1}\}$
  
  $Pa_{X_i} \subseteq \{X_1,\ldots X_{i-1}\} \leftarrow n^d$
  - Any subset
  - Acyclicity guaranteed!
- Total score = sum score of each node

  *OPTIMAL BN, with d parents, consistent with order*

# Fixed variable order 2

- Fix max number of parents to k
- For each *i* in order
  - Pick $\mathbf{Pa}_{Xi} \subseteq \{X_1,\ldots,X_{i-1}\}$
    - Exhaustively search through all possible subsets
    - $\mathbf{Pa}_{Xi}$ is maximum $\mathbf{U} \subseteq \{X_1,\ldots,X_{i-1}\}$ FamScore$(X_i|\mathbf{U} : D)$

- Optimal BN for each order!!!
- Greedy search through space of orders:   1 2   4 5 3
  - 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

  *details in reading*

# Learn BN structure using local search

**Starting from Chow-Liu tree**



*result of fixed order search with K = d*

**Local search, possible moves:**
*Only if acyclic!!!* ✓
- Add edge
- Delete edge
- Invert edge
⋮

*if you are really eager, advanced search techniques like, tabu search beam ⋮ A\**
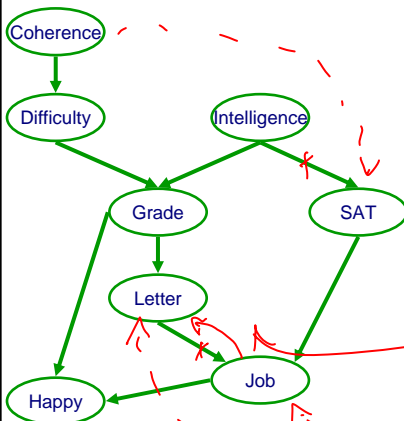
**Select using favorite score**

*BIC*

*or*

*Bayesian*

*or*

⋮

---

# Exploit score decomposition in local search



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

  *only rescore Score Fam (S|I:D)*

- **Reverse edge**
  - □ Rescore only two families

# Some experiments



KL Divergence vs #samples

- Parameter learning (dotted)
- Structure learning (solid)

*structure learn* (handwritten)

*Known "true" Structure* (handwritten)

Alarm network

---

# 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
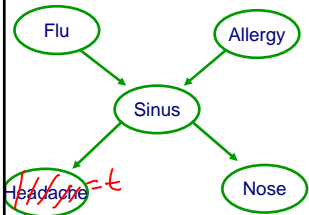
*e.g.: using decision trees  noisy-OR* (handwritten)

# 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

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

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

# 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$^{k+1}$))  2k+6
- Search techniques
  - Search through orders
  - Search through structures
- Bayesian model averaging

# Inference in graphical models: Typical queries 1

Flu   Allergy
  Sinus
Headache  Nose

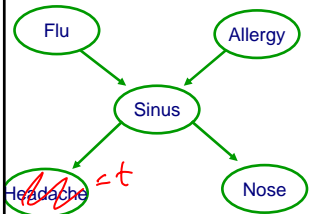Headache=t

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

$$P(A=t \mid H=t)$$

$$P(A=t \mid H=t) \propto P(A=t, H=t)$$

$$P(A=t, H=t) = \sum_s \sum_n \sum_f P(A=t, H=t, s, n, f)$$

# Inference in graphical models: Typical queries 2 – Maximization

some times called Maximum a posteriori (MAP)

Flu   Allergy
  Sinus
Headache  Nose

Headache=t

- Most probable explanation (MPE)
  - Most likely assignment to all hidden vars given evidence

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

- Maximum a posteriori (MAP)
  - Most likely assignment to some var(s) given evidence

$$\max_a P(A=a \mid H=t)$$

$$= \max_a \sum_s \sum_f \sum_n P(A=a, s, f, n \mid H=t)$$

# Are MPE and MAP Consistent?

Sinus → Nose

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

P(N|S) = 

|       | N=t | N=f |
|-------|-----|-----|
| S=t   | .9  | .1  |
| S=f   | .5  | .5  |

- **Most probable explanation (MPE)**
  - Most likely assignment to all hidden vars given evidence

  MPE: S=t, N=t

- **Maximum a posteriori (MAP)**
  - Most likely assignment to some var(s) given evidence

  $$\max_{s} P(S=s)$$

  ⇒ MAP(s) = S=f

---

# 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, andjunction trees
  - basic static and dynamic inference algorithms
  - parameter learning for Gaussian distributions, Chow Liu

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

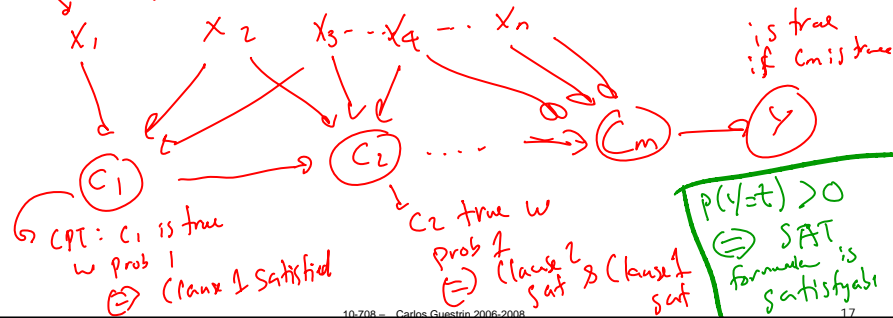# Complexity of conditional probability queries 1

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

Reduction – 3-SAT

$$(\overline{X}_1 \vee X_2 \vee X_3) \wedge (\overline{X}_2 \vee X_3 \vee X_4) \wedge ...$$



$E = \emptyset$

$P(X=\emptyset)$

want a satisfying assignment exists?

$C_1$

$C_2$ ... $C_m$

CPT: 50/50 Uniform

$X_1$ $X_2$ $X_3 \cdots X_4 \cdots X_n$

is true if $C_m$ is true

$C_1$ $C_2$ $\cdots$ $C_m$ $Y$

$C_1$ is true $w$ prob 1 $\Rightarrow$ Clause 1 Satisfied

$C_2$ true $w$ prob 1 $\Rightarrow$ Clause 2 sat & Clause 1 sat

$P(Y=t) > 0$ $\Leftrightarrow$ SAT formula is satisfiable

---

# Complexity of conditional probability queries 2

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

#P-complete problems: e.g., how many satisfying assignments a SAT formula has?

$X_1 \cdots X_n$

$C_1 \cdots C_m \to Y$

$2^n$ assignments, each has probability $\frac{1}{2^n}$

$$P(Y=t) = \frac{\# \text{ sat assignments}}{2^n}$$

# Inference is #P-complete, hopeless?

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

one key property → low tree width ✓ graphs

others, e.g., → CSI
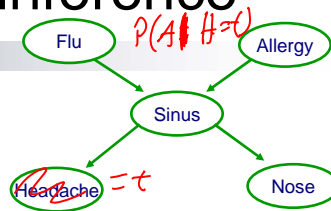↘ associative potentials

---

# Complexity for other inference questions

- Probabilistic inference   #P-complete
  - general graphs:
  - poly-trees and low tree-width:   polynomial

- Approximate probabilistic inference
  - Absolute error:   $|P(x) - \hat{P}(x)| \leq \varepsilon$   ← NP-hard for any $\varepsilon \leq 0.5$
  - Relative error:   $1 - \varepsilon \leq \dfrac{P(x)}{\hat{P}(x)} \leq 1 + \varepsilon$   ← NP-hard for any $\varepsilon > 0$

- Most probable explanation (MPE)
  - general graphs:   NP-Complete
  - poly-trees and low tree-width:   polynomial

- Maximum a posteriori (MAP)
  - general graphs:   $NP^{PP}$-Complete
  - poly-trees and low tree-width:   NP-hard

# 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

---

# General probabilistic inference

Flu          $P(A \mid H=t)$          Allergy

Sinus

Headache $=t$          Nose

- Query:   $P(X \mid e)$

- Using def. of cond. prob.:

$$P(X \mid e) = \frac{P(X, e)}{P(e)} \propto P(X,e) ;$$

compute
$\forall x$
$\Rightarrow P(X=x, \bar{t}=e)$

- Normalization:

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

normalize
$\begin{cases} P(A=t, H=t) = 0.2 \\ P(A=\bar{t}, H=t) = 0.1 \end{cases}$
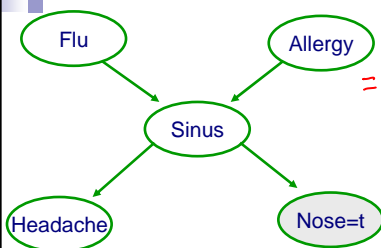
$P(A=t \mid H=t) = \frac{2}{3}$

# Marginalization

Flu → Sinus → Nose=t

$$P(F=t, N=t) = \sum_{S} P(F=t, S, N=t)$$

$$= P(F=t, S=t, N=t) + P(F=t, S=f, N=t)$$

---

# Probabilistic inference example

Flu     Allergy

Sinus

Headache          Nose=t

$$P(A|N=t) \propto P(A, N=t)$$

$$= \sum_{f} \sum_{s} \sum_{h} P(A, f, s, h, N=t)$$

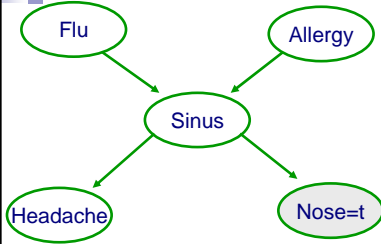$$= \sum_{f} \sum_{s} \sum_{h} P(f) P(A) P(s|f,A) P(h|s) P(N=t|s)$$

$$= \sum_{f} \sum_{s} P(f) P(A) P(s|f,A) P(N=t|s) \underbrace{\sum_{h} P(h|s)}_{1}$$

$$= \sum_{f} P(f) P(A) \underbrace{\sum_{s} P(s|f,A) P(N=t|s)}_{g_1(f,A)}$$

$$= P(A) \underbrace{\sum_{f} P(f) g_1(f,A)}_{g_2(A)} = P(A) g_2(A) = P(A, N=t)$$

**Inference seems exponential in number of variables!**

# Fast probabilistic inference example – Variable elimination



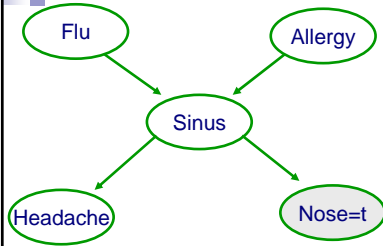(Potential for) Exponential reduction in computation!

# Understanding variable elimination – Exploiting distributivity
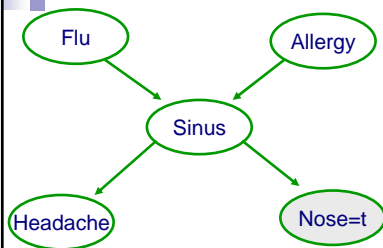
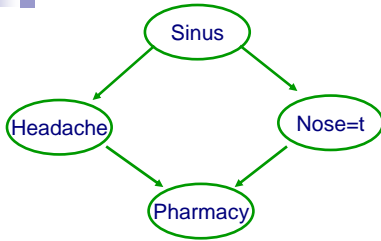# Understanding variable elimination – Order can make a HUGE difference

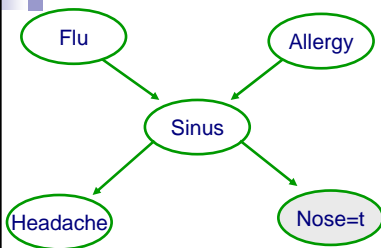# Understanding variable elimination – Intermediate results



Intermediate results are probability distributions

# Understanding variable elimination – Another example

Sinus

Headache    Nose=t

Pharmacy

# Pruning irrelevant variables

Flu    Allergy

Sinus

Headache    Nose=t

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