

Graphical Models

Aarti Singh

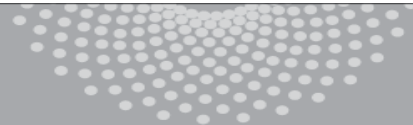
Slides Courtesy: Carlos Guestrin

Machine Learning 10-701/15-781

Apr 17, 2023



MACHINE LEARNING DEPARTMENT



Carnegie Mellon.
School of Computer Science

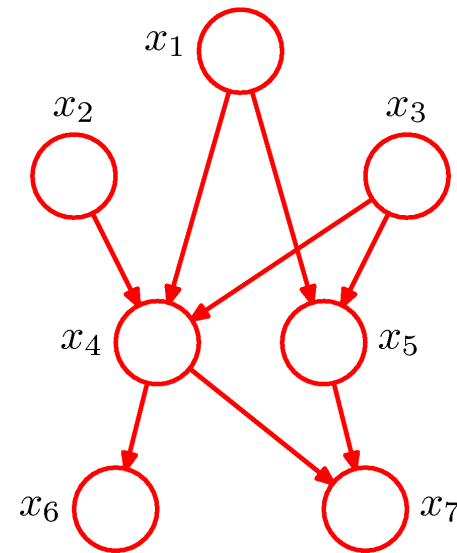
Directed – Bayesian Networks

- Compact representation for a joint probability distribution
- Bayes Net = Directed Acyclic Graph (DAG) + Conditional Probability Tables (CPTs)
- distribution factorizes according to graph

$$p(\mathbf{x}) = \prod_{k=1}^K p(x_k | \text{pa}_k)$$

≡ distribution satisfies **local Markov assumption**

x_k is independent of its non-descendants
given its parents pa_k



Independencies encoded by BN

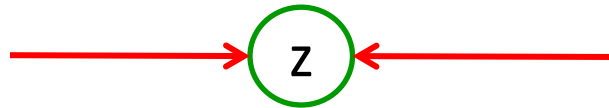
- Set of distributions that factorize according to the graph – \mathcal{F}
 \equiv satisfy local Markov assumption
- Set of distributions that respect conditional independencies implied by d-separation properties of graph – \mathcal{I}

D-separation

- A, B, C – non-intersecting set of nodes
- A is D-separated from B by C $\equiv A \perp B | C$
if all paths between nodes in A & B are “blocked”
i.e. path contains a node z such that either



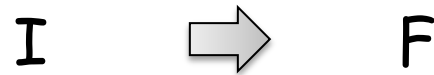
and z in C, OR



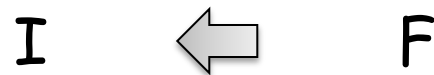
and neither z nor any of its descendants is in C.

Representation Theorem

- Set of distributions that factorize according to the graph - \mathcal{F}
- Set of distributions that respect conditional independencies implied by d-separation properties of graph - \mathcal{I}



Important because: **Given independencies of P can get BN structure G**



Important because: **Read independencies of P from BN structure G**

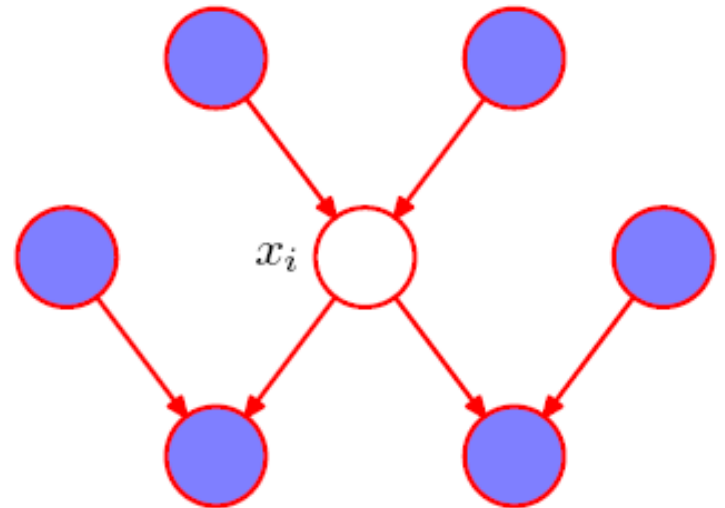
Markov Blanket

- Conditioning on the Markov Blanket, node i is independent of all other nodes.

$$p(\mathbf{x}_i | \mathbf{x}_{\{j \neq i\}}) = \frac{p(x_1, \dots, x_n)}{\sum_i p(x_1, \dots, x_n)} = \frac{\prod_k p(x_k | pa(x_k))}{\sum_i \prod_k p(x_k | pa(x_k))} = p(\mathbf{x}_i | \text{MB}(\mathbf{x}_i))$$

Only terms that remain are the ones which involve i

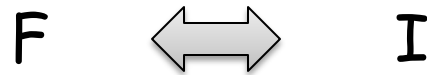
$$p(x_i | pa(x_i)) \quad p(x_k | pa(x_k) \ni i)$$



- Markov Blanket of node i - Set of parents, children and co-parents of node i

Directed – Bayesian Networks

- Graph encodes local independence assumptions (local Markov Assumptions)
- Other independence assumptions can be read off the graph using d-separation
- distribution factorizes according to graph \equiv distribution satisfies all independence assumptions found by d-separation



- Does the graph capture all independencies? Yes, for *almost all* distributions that factorize according to graph. More in 10-708

Topics in Graphical Models

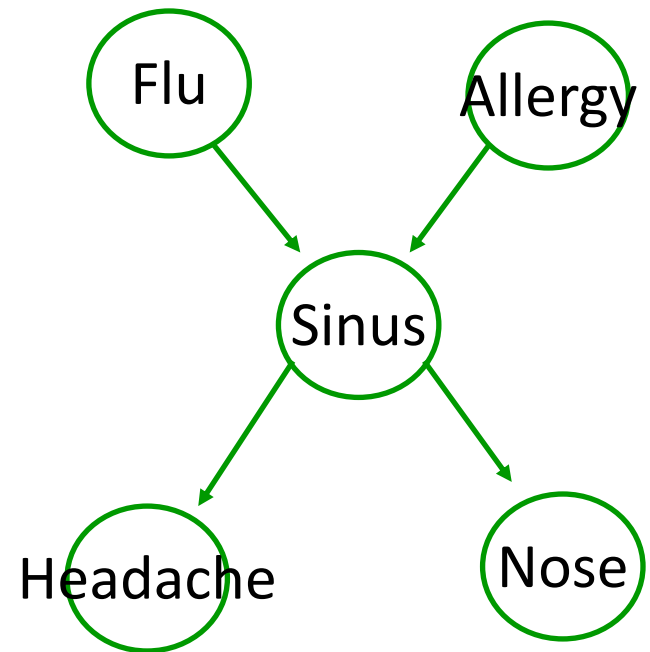
- Representation
 - Which joint probability distributions does a graphical model represent?
- Inference
 - How to answer questions about the joint probability distribution?
 - Marginal distribution of a node variable
 - Most likely assignment of node variables
- Learning
 - How to learn the parameters and structure of a graphical model?

Inference

- Possible queries:

1) Marginal distribution e.g. $P(S)$
Posterior distribution e.g. $P(F|H=1)$

2) Most likely assignment of nodes
 $\arg \max_{f,a,s,n} P(F=f,A=a,S=s,N=n|H=1)$



Inference

- Possible queries:

1) Marginal distribution e.g. $P(S)$

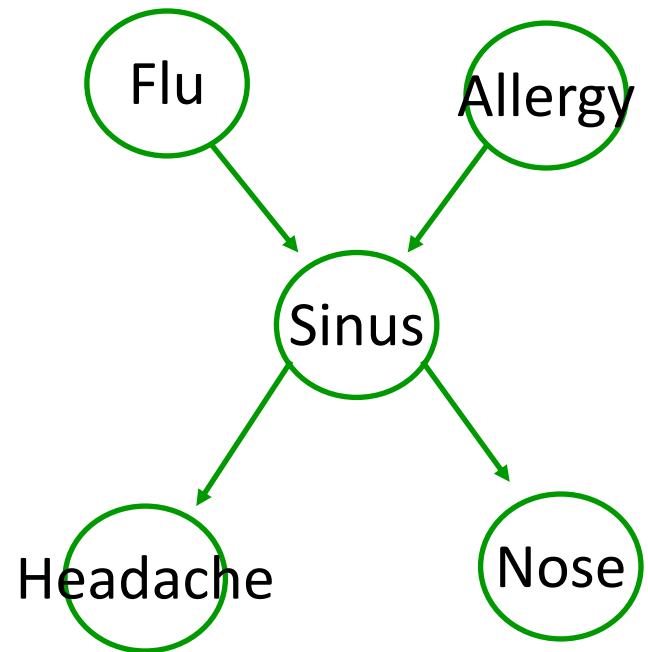
Posterior distribution e.g. $P(F|H=1)$

$P(F|H=1)$?

$$\begin{aligned} P(F|H=1) &= \frac{P(F, H=1)}{P(H=1)} \\ &= \frac{P(F, H=1)}{\sum_f P(F=f, H=1)} \end{aligned}$$

$$\propto P(F, H=1)$$

will focus on computing this, posterior will follow with only constant times more effort



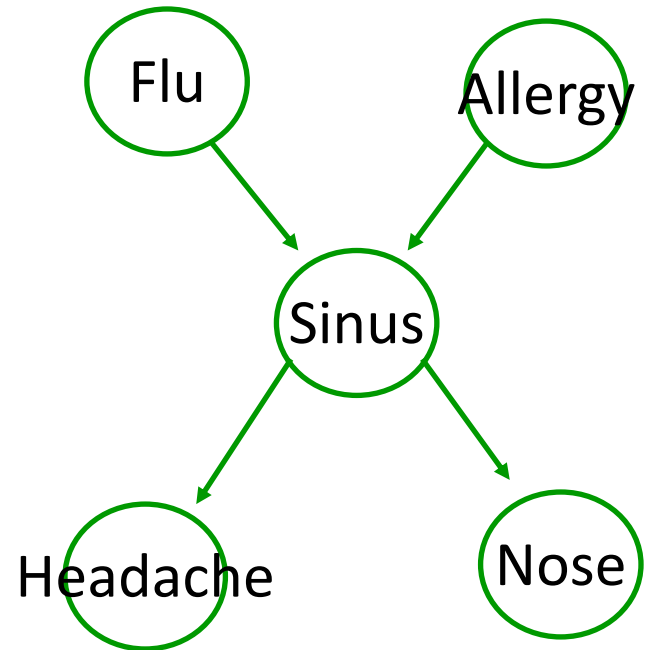
Marginalization

Need to marginalize over other vars

$$P(S) = \sum_{f,a,n,h} P(f,a,S,n,h)$$

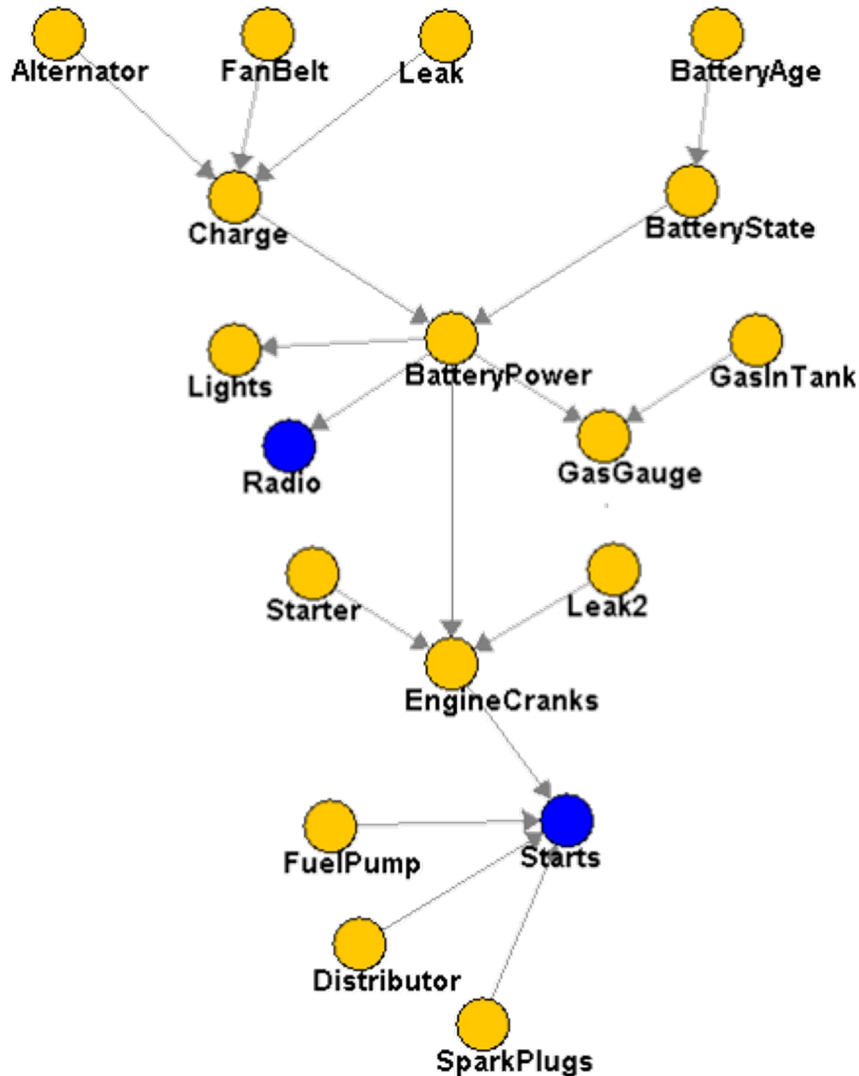
$$P(F,H=1) = \sum_{\underbrace{a,s,n}_{2^3 \text{ terms}}} P(F,a,s,n,H=1)$$

To marginalize out n binary variables,
need to sum over 2^n terms



*Inference seems exponential in number of variables!
Actually, inference in graphical models is NP-hard ☹️*

Bayesian Networks Example



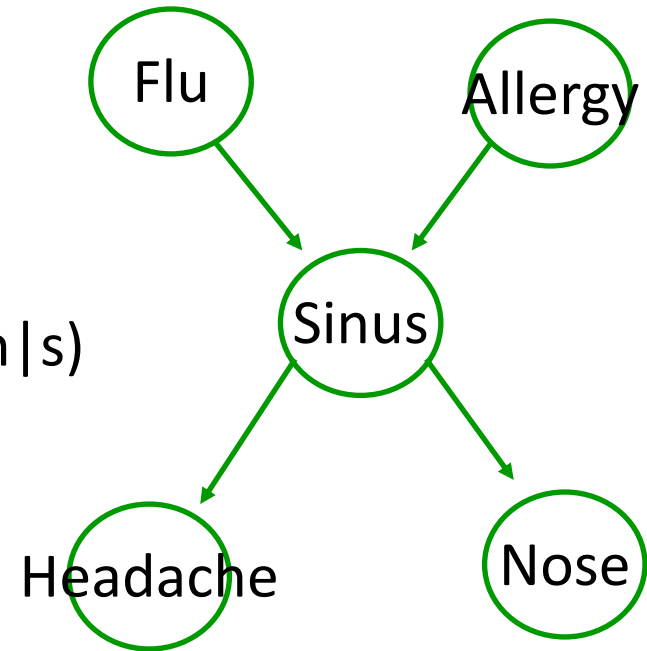
- 18 binary attributes
- Inference
 - $P(\text{BatteryAge} \mid \text{Starts}=f)$
- need to sum over 2^{16} terms!
- Not impressed?
 - HailFinder BN – more than $3^{54} = 58149737003040059690390169$ terms

Fast Probabilistic Inference

$$\begin{aligned} P(F, H=1) &= \sum_{a,s,n} P(F,a,s,n,H=1) \\ &= \sum_{a,s,n} P(F)P(a)P(s|F,a)P(n|s)P(H=1|s) \\ &= P(F) \sum_a P(a) \sum_s P(s|F,a)P(H=1|s) \sum_n P(n|s) \end{aligned}$$

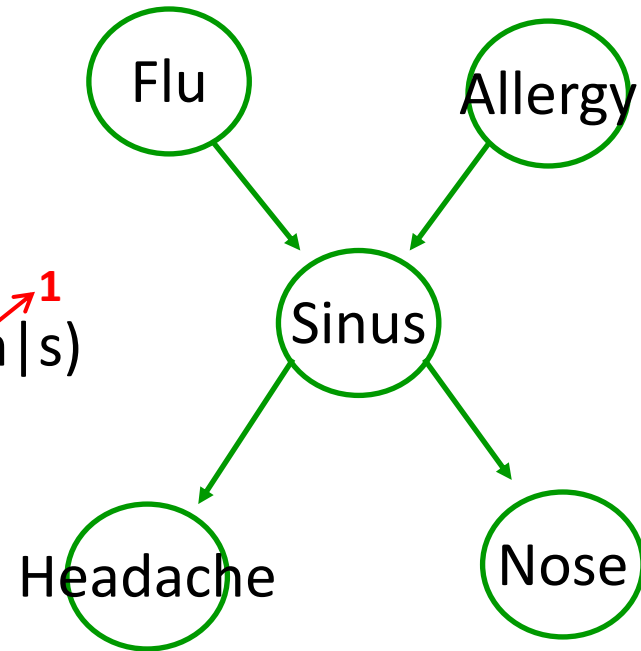
Push sums in as far as possible

Distributive property: $x_1z + x_2z = z(x_1+x_2)$
2 multiply 1 multiply



Fast Probabilistic Inference

$$\begin{aligned}
 P(F, H=1) &= \sum_{a,s,n} P(F,a,s,n,H=1) \\
 &= \sum_{a,s,n} P(F)P(a)P(s|F,a)P(n|s)P(H=1|s) \\
 &= P(F) \sum_a P(a) \sum_s P(s|F,a)P(H=1|s) \sum_n P(n|s) \\
 &= P(F) \sum_a P(a) \sum_s P(s|F,a)P(H=1|s) \\
 &= P(F) \sum_a P(a) g_1(F,a) \\
 &= P(F) g_2(F)
 \end{aligned}$$



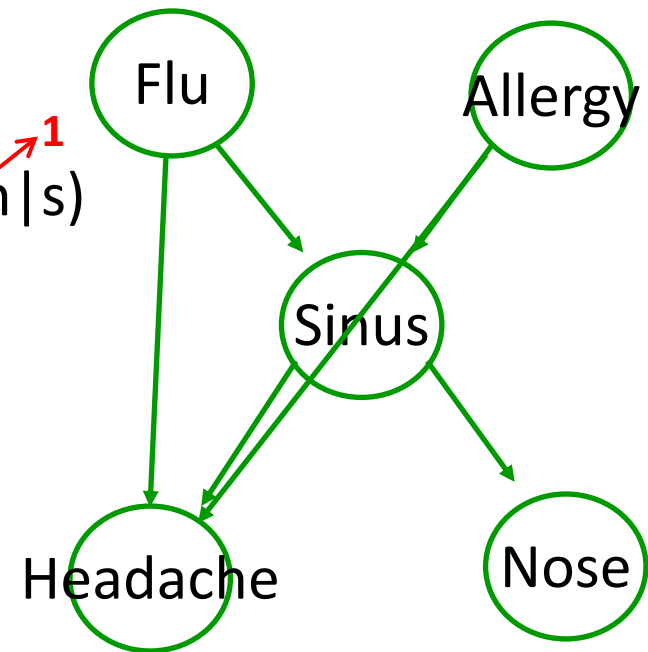
2^n vs. $n 2^k$ multiplies
 k - scope of (number of
 variables in) largest factor

(Potential for) exponential reduction in computation!

Fast Probabilistic Inference – Variable Elimination

$$\begin{aligned} P(F, H=1) &= \sum_{a,s,n} P(F)P(a)P(s|F,a)P(n|s)P(H=1|s) \\ &= P(F) \underbrace{\sum_a P(a) \sum_s P(s|F,a)P(H=1|s)}_{P(H=1|F)} \sum_n P(n|s) \end{aligned}$$

The diagram shows the simplification of the joint probability expression. The inner sum over s is labeled $P(H=1|F,a)$, and the outer sum over a is labeled $P(H=1|F)$. A red arrow points from the n index in the final sum to the number 1, indicating that the variable n is eliminated because its value is fixed to 1.

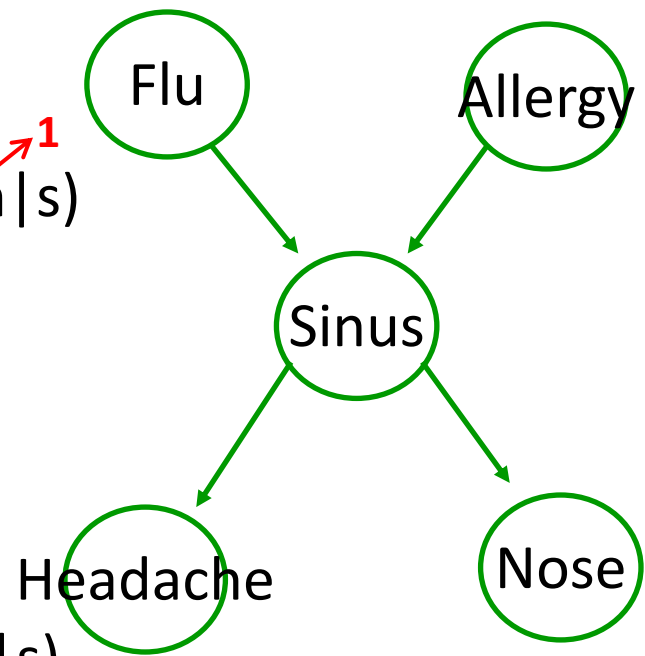


(Potential for) exponential reduction in computation!

Variable Elimination – Order can make a HUGE difference

$$\begin{aligned}
 P(F, H=1) &= \sum_{a,s,n} P(F)P(a)P(s|F,a)P(n|s)P(H=1|s) \\
 &= P(F) \sum_a P(a) \underbrace{\sum_s P(s|F,a)P(H=1|s)}_{P(H=1|F,a)} \sum_n P(n|s)
 \end{aligned}$$

$\underbrace{\hspace{10em}}_{P(H=1|F)}$

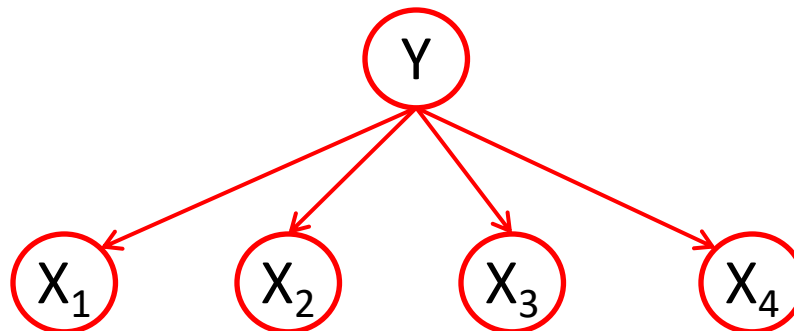


$$P(F, H=1) = P(F) \sum_a P(a) \underbrace{\sum_n \sum_s P(s|F,a)P(n|s)P(H=1|s)}_{g(F,a,n)}$$

3 - scope of largest factor

(Potential for) exponential reduction in computation!

Variable Elimination – Order can make a HUGE difference



$$P(X_1) = \sum_{Y, X_2, \dots, X_n} P(Y)P(X_1|Y) \prod_{i=2}^n P(X_i|Y)$$

$$= \sum_{Y, X_3, \dots, X_n} P(Y)P(X_1|Y) \prod_{i=3}^n P(X_i|Y) \underbrace{\sum_{X_2} P(X_2|Y)}_{g(Y)}$$

1 - scope of largest factor

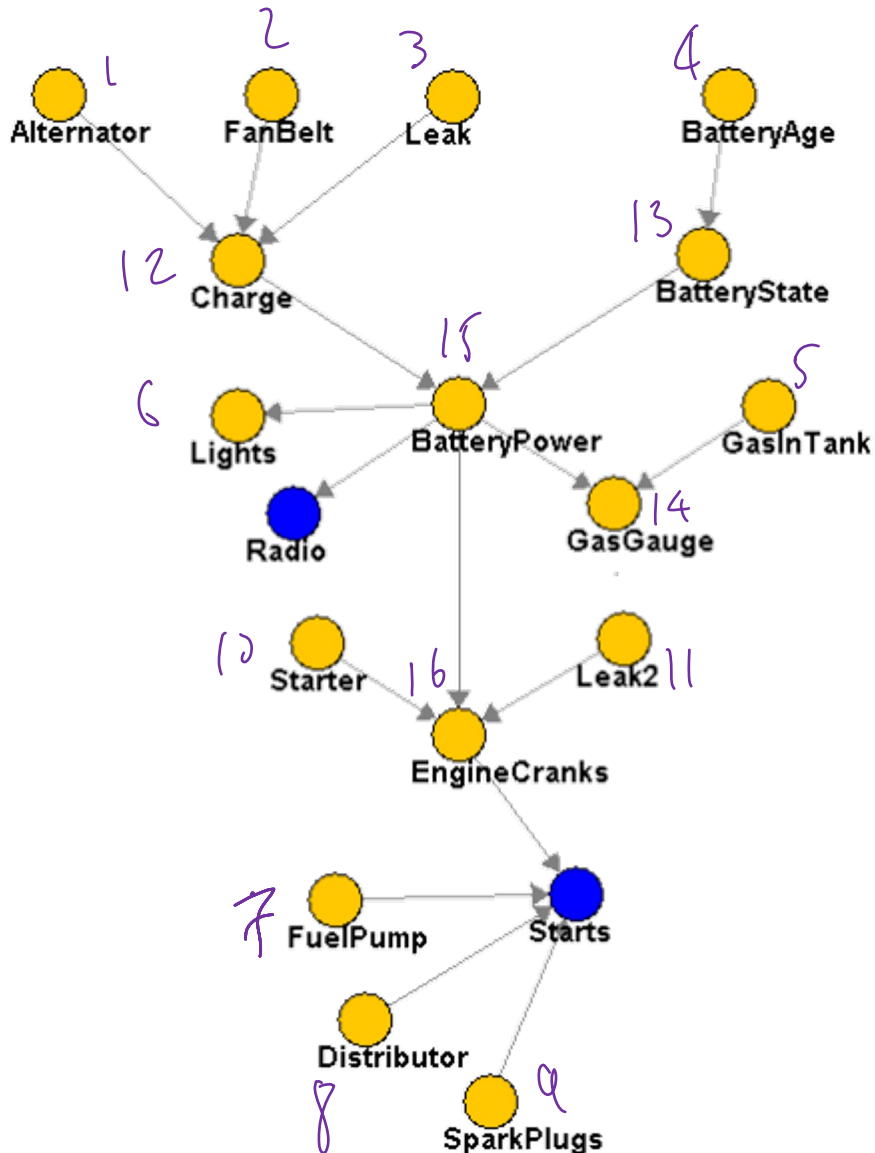
$$= \sum_{X_2, \dots, X_n} \underbrace{\sum_Y P(Y)P(X_1|Y) \prod_{i=2}^n P(X_i|Y)}_{g(X_1, X_2, \dots, X_n)}$$

n - scope of largest factor

Variable Elimination Algorithm

- Given BN – DAG and CPTs (initial factors – $p(x_i | pa_i)$ for $i=1, \dots, n$)
- Given Query $P(X|e) \equiv P(X,e)$ X – set of variables e - evidence
- Instantiate evidence e e.g. set $H=1$ **IMPORTANT!!!**
- Choose an ordering on the variables e.g., $X_{(1)}, \dots, X_{(n)}$
- For $i = 1$ to n , If $X_{(i)} \notin \{X, e\}$ (i.e. need to marginalize it out)
 - Collect factors g_1, \dots, g_k that include $X_{(i)}$
 - Generate a new factor by eliminating $X_{(i)}$ from these factors
$$g = \sum_{X_i} \prod_{j=1}^k g_j$$
 - Variable $X_{(i)}$ has been eliminated!
 - Remove g_1, \dots, g_k from set of factors but add g
- Normalize $P(X,e)$ to obtain $P(X|e)$

Complexity for (Poly)tree graphs



Variable elimination order:

- Consider undirected version (ignore edge directions)
- Start from “leaves” up
- find topological order
- eliminate variables in that order

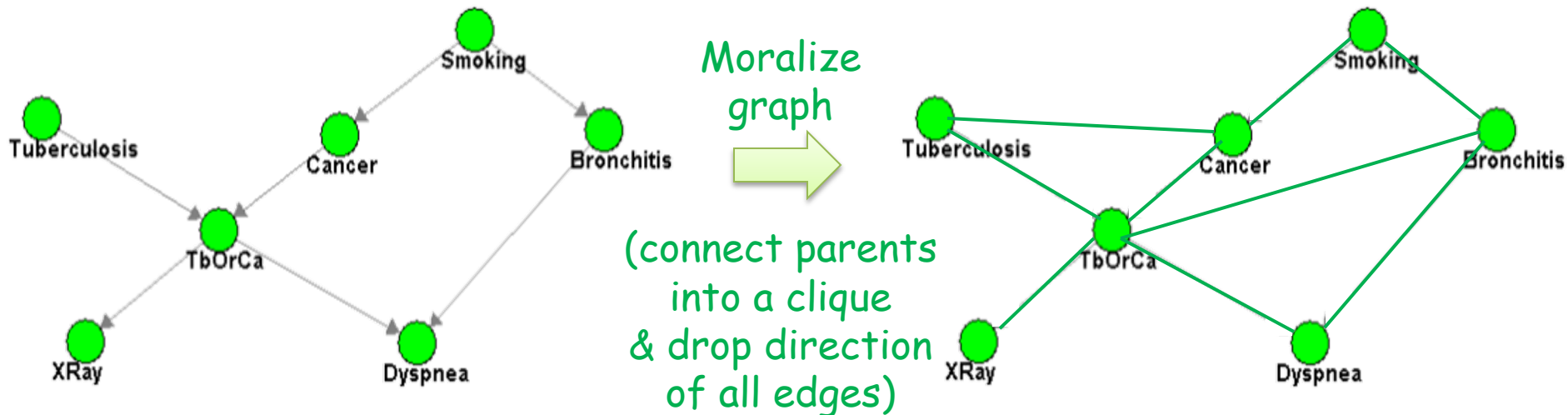
Does not create any factors bigger than original CPTs

For polytrees, inference is linear in # variables (vs. exponential in general)!

Complexity for graphs with loops

- Loop – undirected cycle

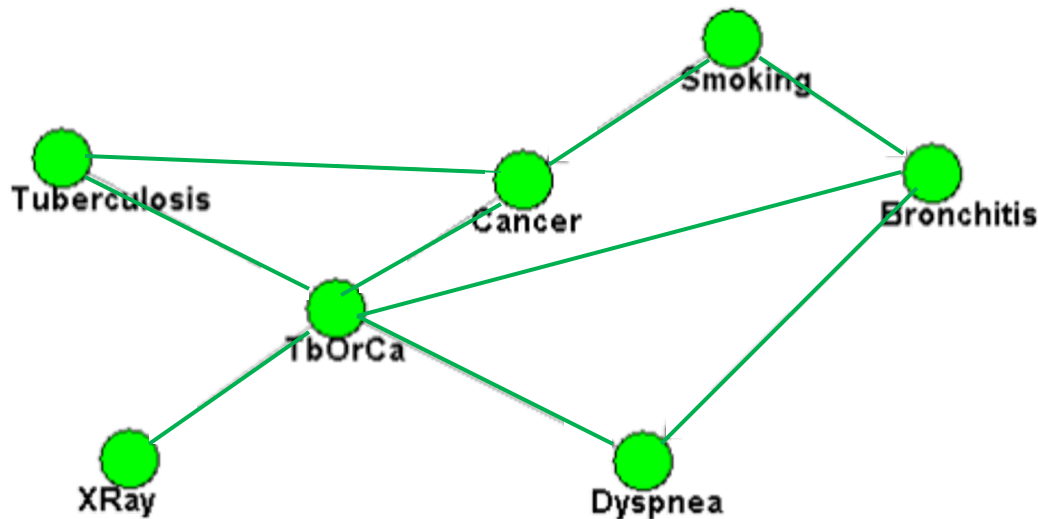
Linear in # variables but exponential in size of largest factor generated!



When you eliminate a variable, add edges between its neighbors

Complexity for graphs with loops

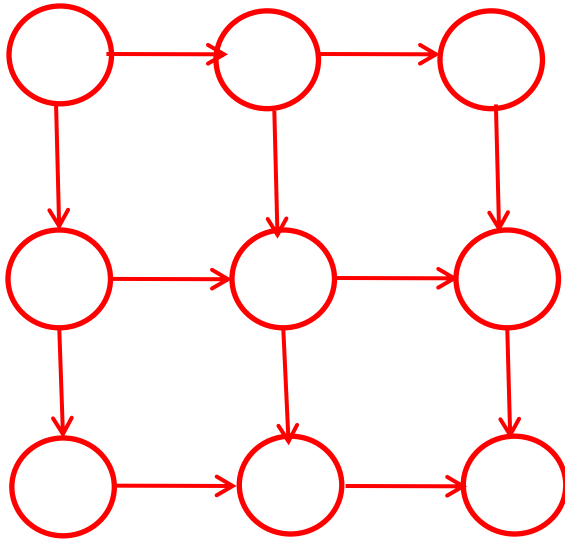
- Loop – undirected cycle



| Var eliminated | Factor generated |
|----------------|------------------|
| S | $g_1(C,B)$ |
| B | $g_2(C,O,D)$ |
| D | $g_3(C,O)$ |
| C | $g_4(T,O)$ |
| T | $g_5(O)$ |
| O | $g_6(X)$ |

Linear in # variables but exponential in size of largest factor generated \sim tree-width (max clique size-1) in resulting graph!

Example: Large tree-width with small number of parents



At most 2 parents per node, but tree width is $O(\sqrt{n})$

Compact representation \Rightarrow Easy inference ☹

Choosing an elimination order

- Choosing best order is NP-complete
 - Reduction from MAX-Clique
- Many good heuristics (some with guarantees)
- Ultimately, can't beat NP-hardness of inference
 - Even optimal order can lead to exponential variable elimination computation
- In practice
 - Variable elimination often very effective
 - Many (many many) approximate inference approaches available when variable elimination too expensive

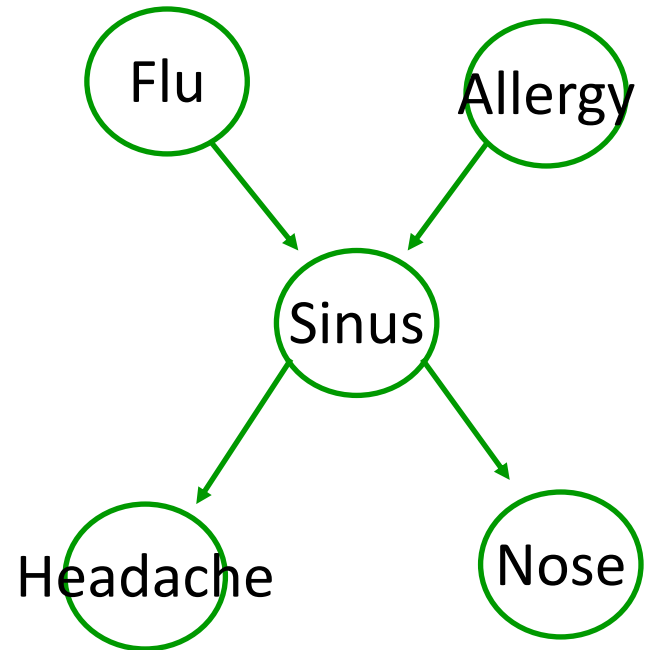
Inference

- Possible queries:
 - 2) Most likely assignment of nodes
$$\arg \max_{f,a,s,n} P(F=f, A=a, S=s, N=n | H=1)$$

Use Distributive property:

$$\max(x_1z, x_2z) = z \max(x_1, x_2)$$

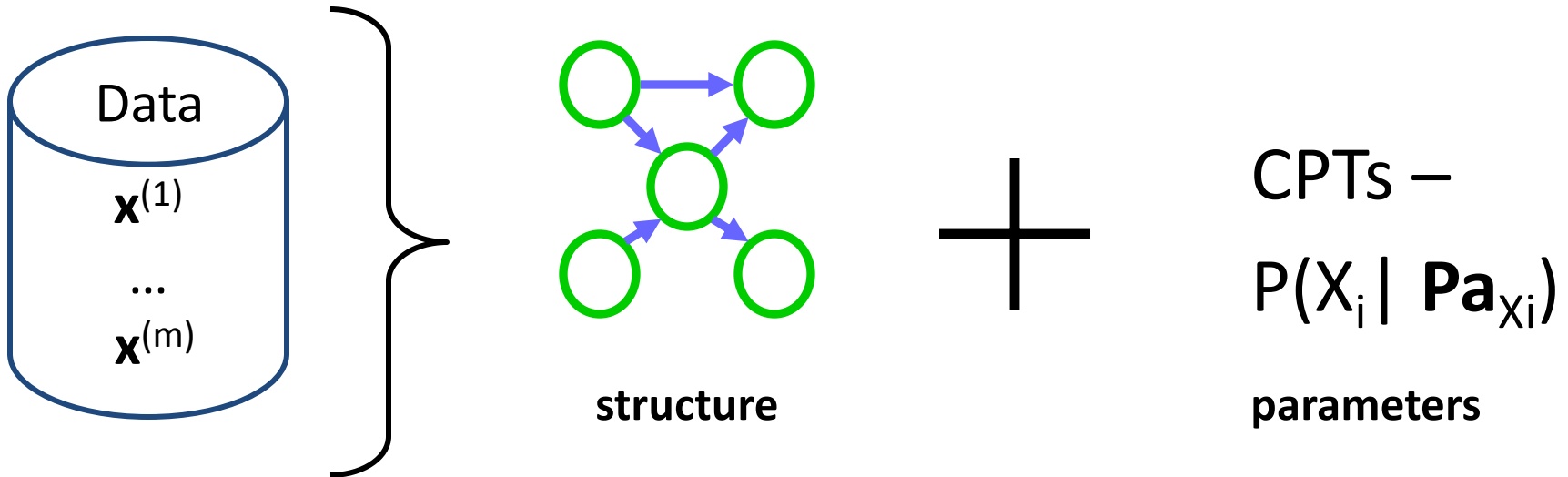
2 multiply 1 multiply



Topics in Graphical Models

- Representation
 - Which joint probability distributions does a graphical model represent?
- Inference
 - How to answer questions about the joint probability distribution?
 - Marginal distribution of a node variable
 - Most likely assignment of node variables
- Learning
 - How to learn the parameters and structure of a graphical model?

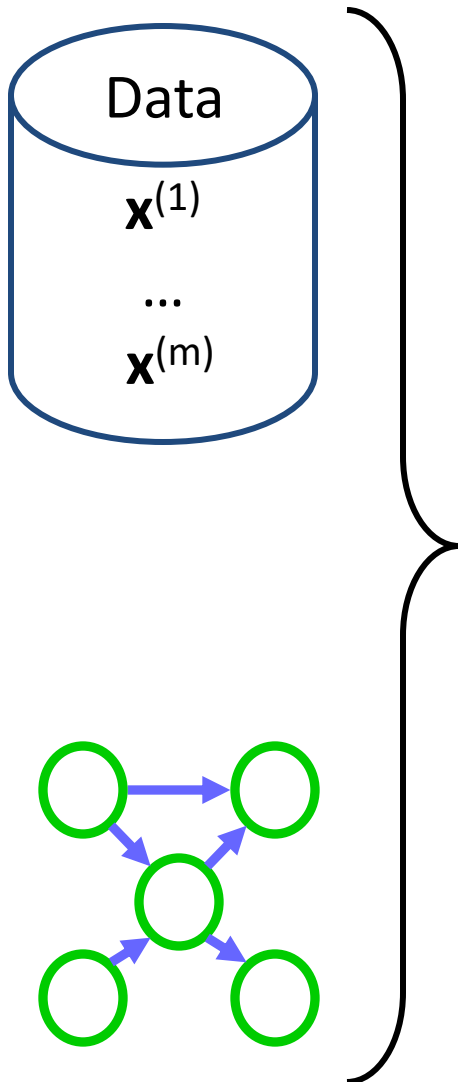
Learning



Given set of m independent samples (assignments of random variables),

find the best (most likely?) Bayes Net (graph Structure + CPTs)

Learning the CPTs (given structure)



For each discrete variable X_k

Compute MLE or MAP estimates for

$$p(x_k | \text{pa}_k)$$

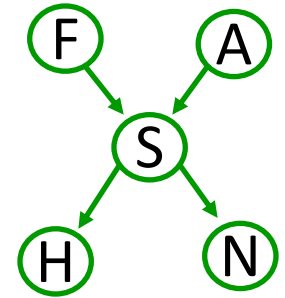
Recall

$$\text{MLE: } P(X_i = x_i | X_j = x_j) = \frac{\text{Count}(X_i = x_i, X_j = x_j)}{\text{Count}(X_j = x_j)}$$

MAP: Add psuedocounts

MLEs decouple for each CPT in Bayes Nets

- Given structure, log likelihood of data



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

$$= \log \prod_{j=1}^m P(f^{(j)}) P(a^{(j)}) P(s^{(j)} \mid f^{(j)}, a^{(j)}) P(h^{(j)} \mid s^{(j)}) P(n^{(j)} \mid s^{(j)})$$

$$= \sum_{j=1}^m [\log P(f^{(j)}) + \log P(a^{(j)}) + \log P(s^{(j)} \mid f^{(j)}, a^{(j)}) + \log P(h^{(j)} \mid s^{(j)}) + \log P(n^{(j)} \mid s^{(j)})]$$

$$= \underbrace{\sum_{j=1}^m \log P(f^{(j)})}_{\theta_F} + \underbrace{\sum_{j=1}^m \log P(a^{(j)})}_{\theta_A} + \underbrace{\sum_{j=1}^m \log P(s^{(j)} \mid f^{(j)}, a^{(j)})}_{\theta_{S|F,A}} +$$

Depends only on

θ_F

θ_A

$\theta_{S|F,A}$

$$+ \underbrace{\sum_{j=1}^m \log P(h^{(j)} \mid s^{(j)})}_{\theta_{H|S}} + \underbrace{\sum_{j=1}^m \log P(n^{(j)} \mid s^{(j)})}_{\theta_{N|S}}$$

$\theta_{H|S}$

$\theta_{N|S}$

Can compute MLEs of each parameter independently!

Information theoretic interpretation of MLE

$$\begin{aligned}\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G}) &= \sum_{j=1}^m \sum_{i=1}^n \log P \left(X_i = x_i^{(j)} \mid \mathbf{Pa}_{X_i} = \mathbf{x}_{\mathbf{Pa}_{X_i}}^{(j)} \right) \\ &= \sum_{i=1}^n \sum_{x_i} \sum_{\mathbf{x}_{\mathbf{Pa}_{X_i}}} \text{count}(X_i = x_i, \mathbf{Pa}_{X_i} = \mathbf{x}_{\mathbf{Pa}_{X_i}}) \log P \left(X_i = x_i \mid \mathbf{Pa}_{X_i} = \mathbf{x}_{\mathbf{Pa}_{X_i}} \right)\end{aligned}$$

Plugging in MLE estimates: ML score

$$\begin{aligned}\log \hat{P}(\mathcal{D} \mid \hat{\theta}_{\mathcal{G}}, \mathcal{G}) &= \sum_{j=1}^m \sum_{i=1}^n \log \hat{P} \left(x_i^{(j)} \mid \mathbf{x}_{\mathbf{Pa}_{X_i}}^{(j)} \right) \\ &= m \sum_{i=1}^n \sum_{x_i} \sum_{\mathbf{x}_{\mathbf{Pa}_{X_i}}} \hat{P}(x_i, \mathbf{x}_{\mathbf{Pa}_{X_i}}) \log \hat{P} \left(x_i \mid \mathbf{x}_{\mathbf{Pa}_{X_i}} \right)\end{aligned}$$

Reminds of entropy

Information theoretic interpretation of MLE

$$\log \hat{P}(\mathcal{D} | \hat{\theta}_{\mathcal{G}}, \mathcal{G}) = m \sum_{i=1}^n \sum_{x_i} \sum_{\mathbf{xPa}_{X_i}} \hat{P}(x_i, \mathbf{xPa}_{X_i}) \log \hat{P}(x_i | \mathbf{xPa}_{X_i})$$

$$= -m \sum_{i=1}^n \hat{H}(X_i | \mathbf{Pa}_{X_i})$$

$$= m \sum_{i=1}^n [\hat{I}(X_i, \mathbf{Pa}_{X_i}) - \underbrace{\hat{H}(X_i)}_{\text{Doesn't depend on graph structure } \mathcal{G}}]$$

Doesn't depend on graph structure \mathcal{G}

ML score for graph structure \mathcal{G}

$$\arg \max_{\mathcal{G}} \log \hat{P}(\mathcal{D} | \hat{\theta}_{\mathcal{G}}, \mathcal{G}) = \arg \max_{\mathcal{G}} \sum_{i=1}^n \hat{I}(X_i, \mathbf{Pa}_{X_i})$$

ML – Decomposable Score

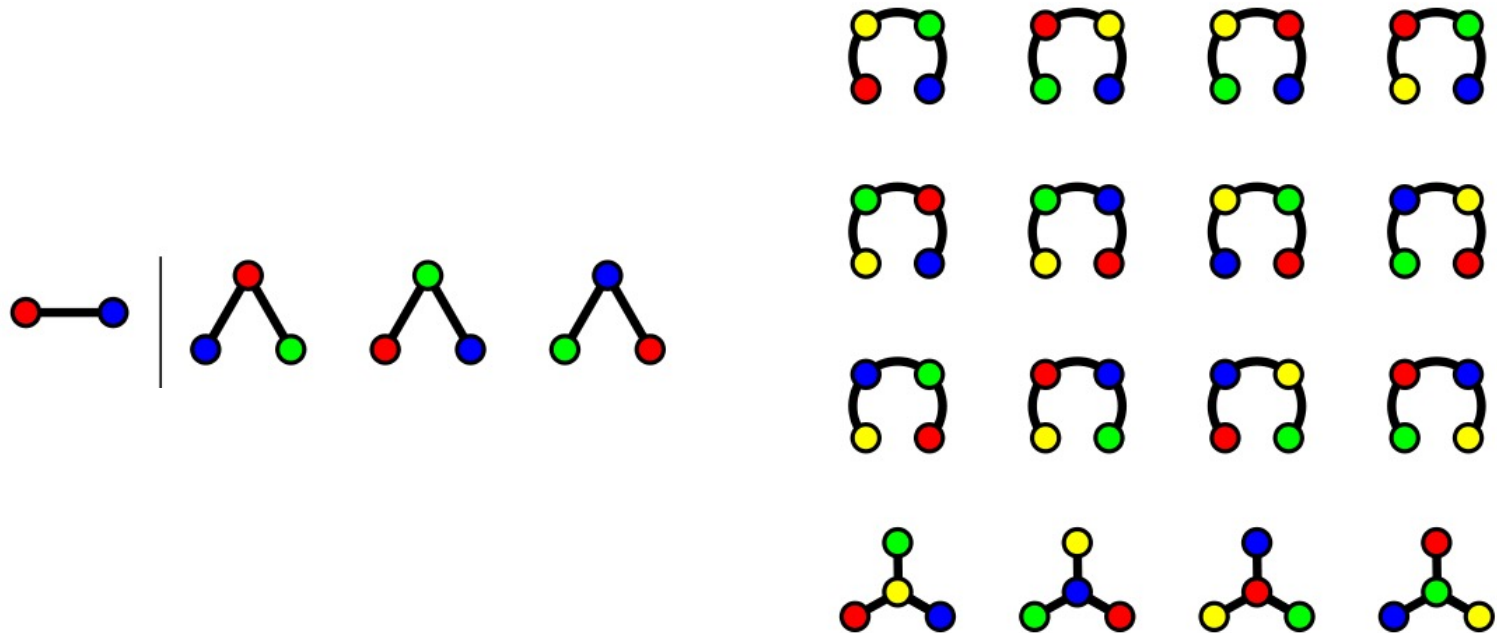
- Log data likelihood

$$\log \hat{P}(\mathcal{D} | \hat{\theta}_{\mathcal{G}}, \mathcal{G}) = m \sum_{i=1}^n [\hat{I}(X_i, \mathbf{Pa}_{X_i}) - \hat{H}(X_i)]$$

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

How many trees are there?

- Trees – every node has at most one parent
- n^{n-2} possible trees (Cayley's Theorem)

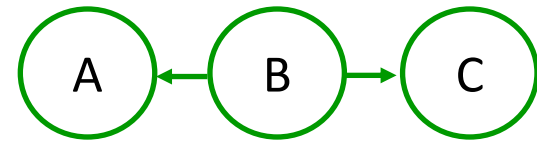
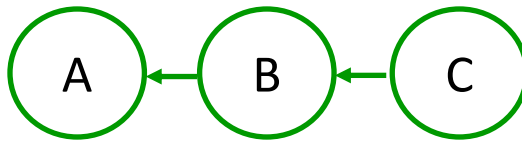
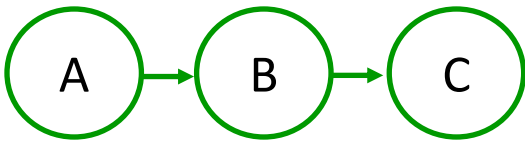


Nonetheless - Efficient optimal algorithm finds best tree!

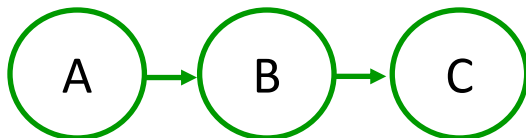
Scoring a tree

$$\arg \max_{\mathcal{G}} \log \hat{P}(\mathcal{D} \mid \hat{\theta}_{\mathcal{G}}, \mathcal{G}) = \arg \max_{\mathcal{G}} \sum_{i=1}^n \hat{I}(X_i, \mathbf{Pa}_{X_i})$$

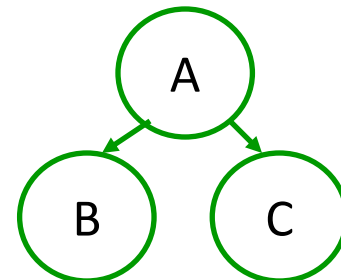
Equivalent Trees (same score): $I(A,B) + I(B,C)$



Score provides indication of structure:



$$I(A,B) + I(B,C)$$



$$I(A,B) + I(A,C)$$

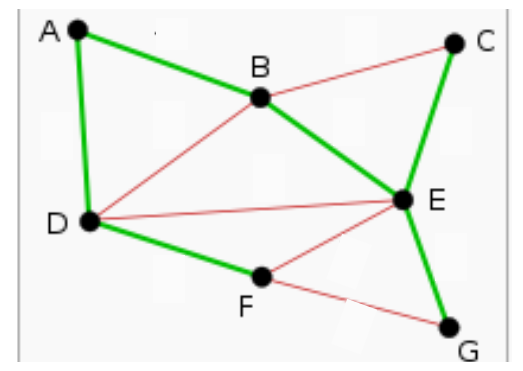
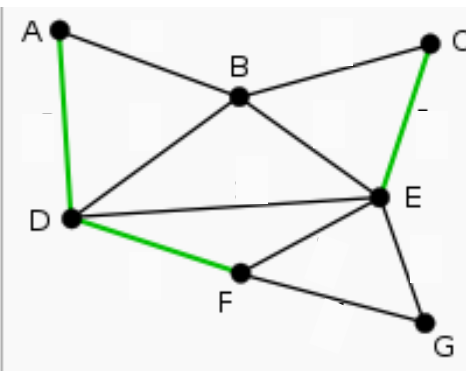
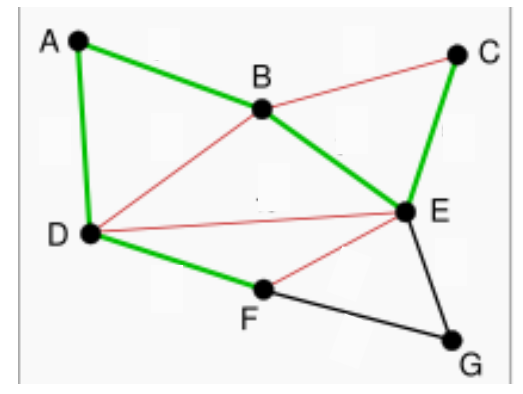
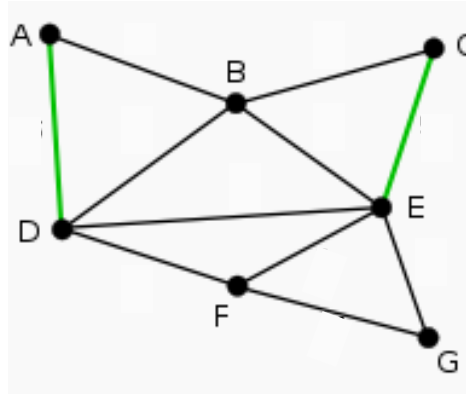
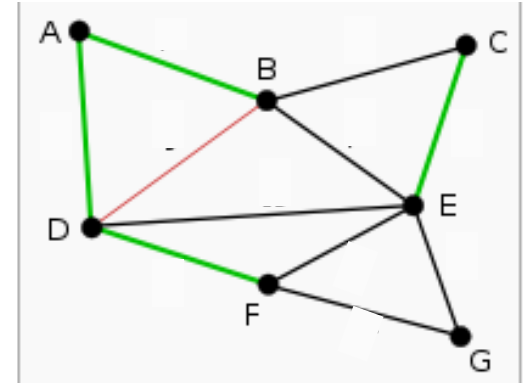
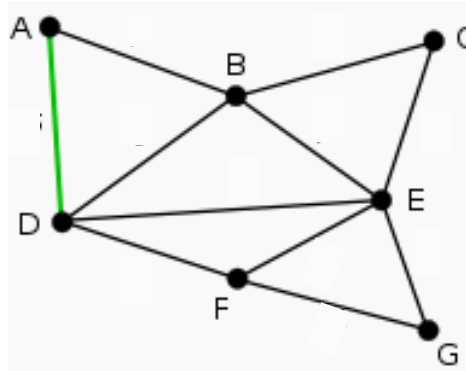
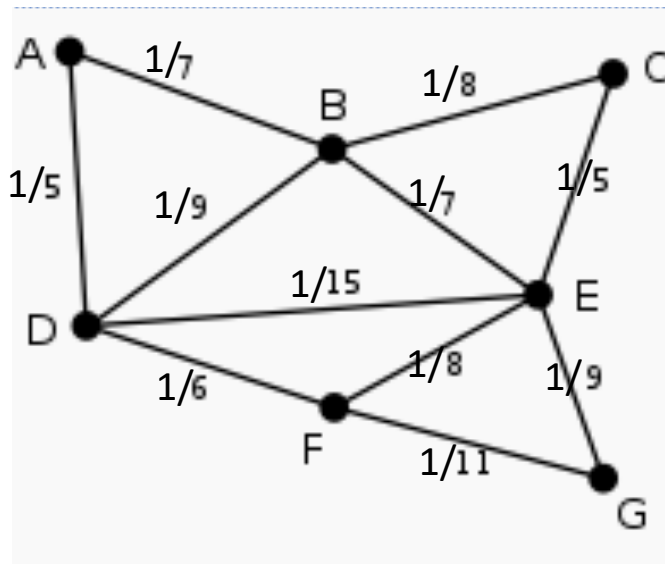
Chow-Liu algorithm

- For each pair of variables X_i, X_j
 - Compute empirical distribution: $\hat{P}(x_i, x_j) = \frac{\text{Count}(x_i, x_j)}{m}$
 - 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 $\hat{I}(X_i, X_j)$
- Optimal tree BN
 - Compute maximum weight spanning tree (e.g. Prim's, Kruskal's algorithm $O(n \log n)$)
 - Directions in BN: pick any node as root, breadth-first-search defines directions

Chow-Liu algorithm example



Scoring general graphical models

- Graph that maximizes ML score -> complete graph!

- Information never hurts

$$H(A|B) \geq H(A|B,C)$$

- Adding a parent always increases ML score

$$I(A,B,C) \geq I(A,B)$$

- The more edges, the fewer independence assumptions, the higher the likelihood of the data, but will overfit...

- Why does ML for trees work?

Restricted model space – tree graph

Regularizing

- Model selection
 - Use MDL (Minimum description length) score
 - BIC score (Bayesian Information criterion)

- Still NP –hard

Theorem: The problem of learning a BN structure with at most d parents is **NP-hard for any (fixed) $d > 1$** (Note: tree $d=1$)

- Mostly heuristic (exploit score decomposition)
- Chow-Liu: provides best tree approximation to any distribution.
- Start with Chow-Liu tree. Add, delete, invert edges. Evaluate BIC score

What you should know

- Learning BNs
 - Maximum likelihood or MAP learns parameters
 - ML score
 - Decomposable score
 - Information theoretic interpretation (Mutual information)
 - Best tree (Chow-Liu)
 - Other BNs, usually local search with BIC score