

Graphical Models

Aarti Singh

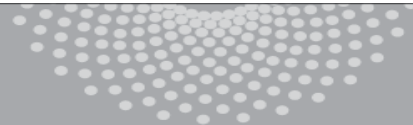
Slides Courtesy: Carlos Guestrin

Machine Learning 10-701/15-781

Apr 17, 2023



MACHINE LEARNING DEPARTMENT



Carnegie Mellon.
School of Computer Science

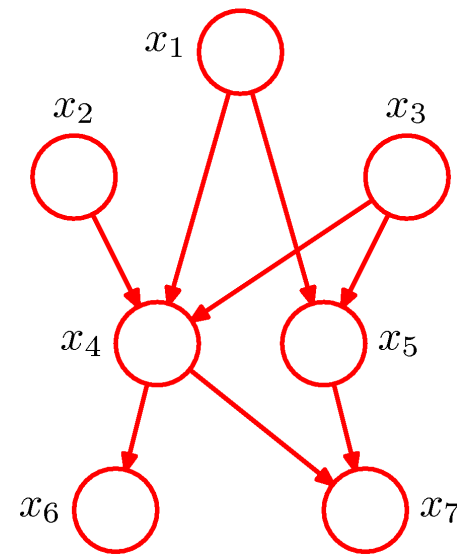
Directed – Bayesian Networks

- 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



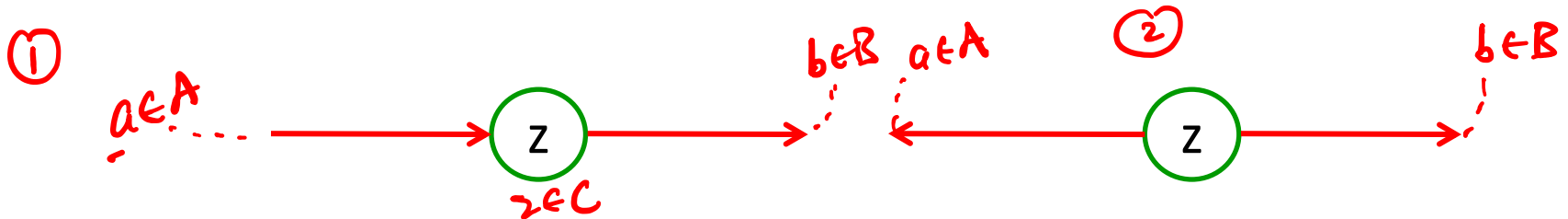
Independencies encoded by BN

- Set of distributions that factorize according to the graph – **F**
≡ satisfy local Markov assumption
- Set of distributions that respect conditional independencies implied by d-separation properties of graph – **I**

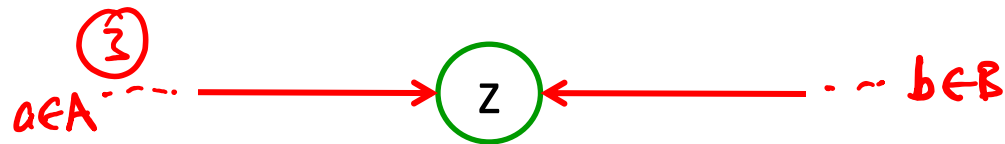
D-separation

sets of nodes
/ \

- 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



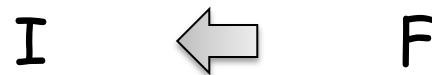
and neither z nor any of its descendants is in C.

Representation Theorem

- Set of distributions that factorize according to the graph - **F** ✓
- Set of distributions that respect conditional independencies implied by d-separation properties of graph - **I** ✓



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



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

$$P(B) = \sum_B P(A, B)$$

Markov Blanket

$$p(x_1, \dots, x_n) = \prod_k p(x_k | pa(x_k))$$

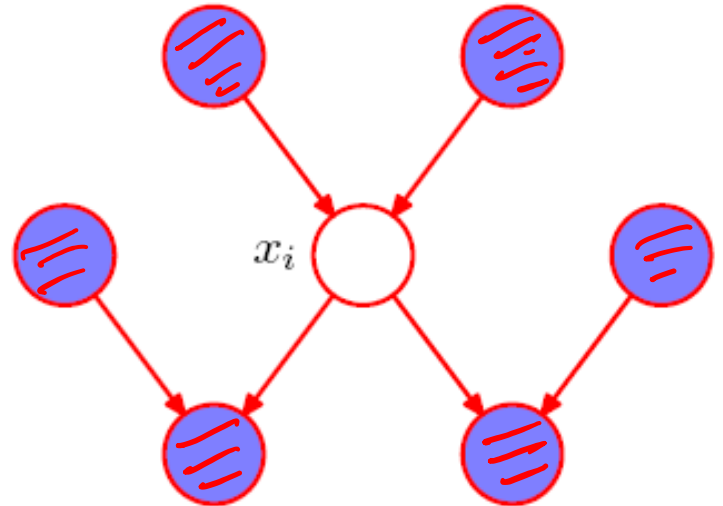
- Conditioning on the Markov Blanket, node i is independent of all other nodes.

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

$P(A|B) = \frac{\text{rest } P(A, B)}{P(B)} = \frac{P(A, B)}{\sum_b P(A, B)}$

Only terms that remain are the ones which involve i

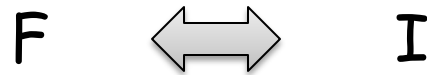
$$p(x_i | pa(x_i)) \quad p(x_k | pa(x_k) \ni i)$$



- 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

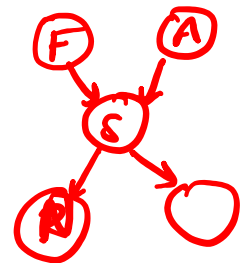
- How to answer questions about the joint probability distribution?

- Marginal distribution of a node variable
- Most likely assignment of node variables

$P(F), P(F, S)$

$P(F=1|S=1)$

$P(F, A, S, N, \dots)$



- Learning

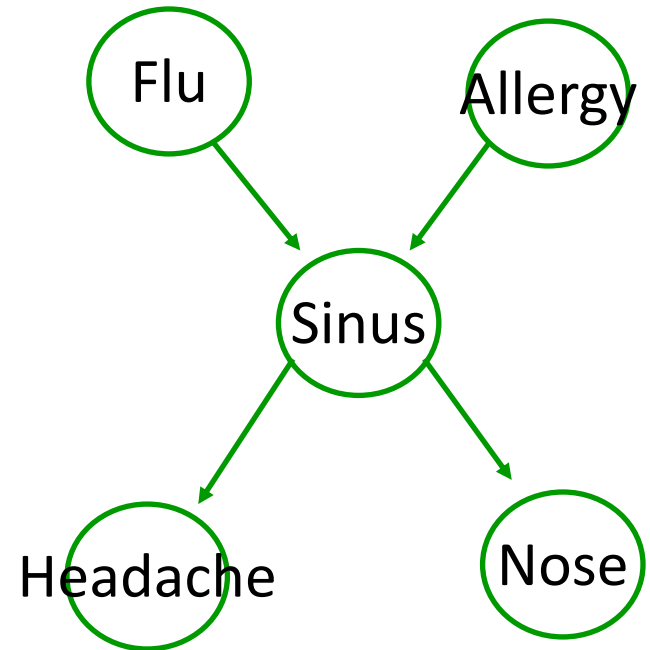
- How to learn the parameters and structure of a graphical model?

Inference

- Possible queries:

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

- 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. $P(F|H=1)$

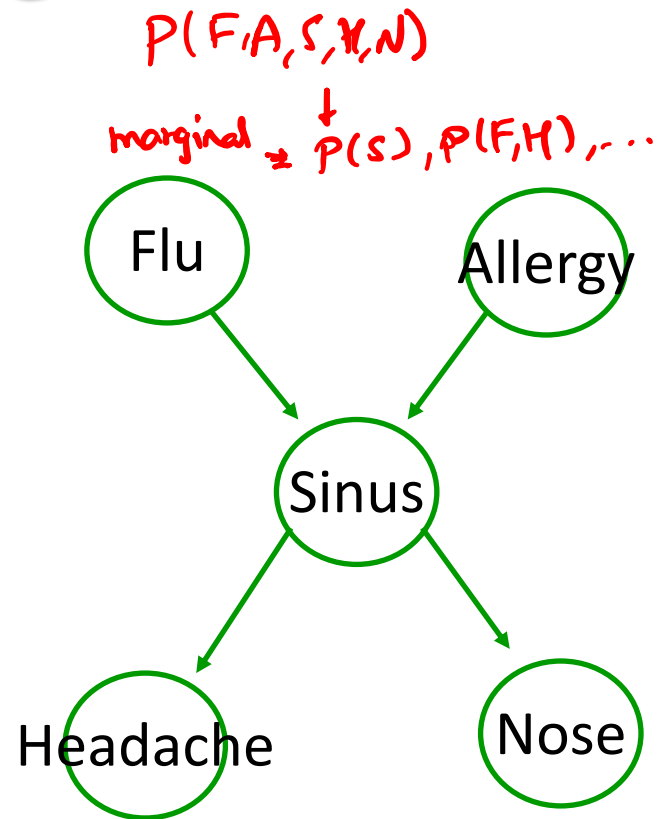
$P(F|H=1)$?

$$P(F|H=1) = \frac{P(F, H=1)}{P(H=1)}$$

$$= \frac{P(F, H=1)}{\sum_f P(F=f, H=1)}$$

$$\propto \underline{P(F, H=1)}$$

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



Marginalization

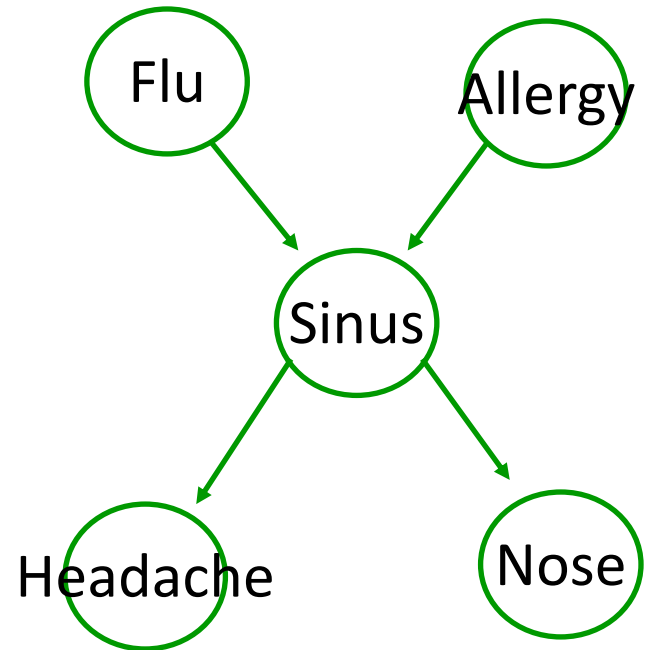
Need to marginalize over other vars

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

$$P(F,H=1) = \sum_{a,s,n} P(F,a,s,n,H=1)$$

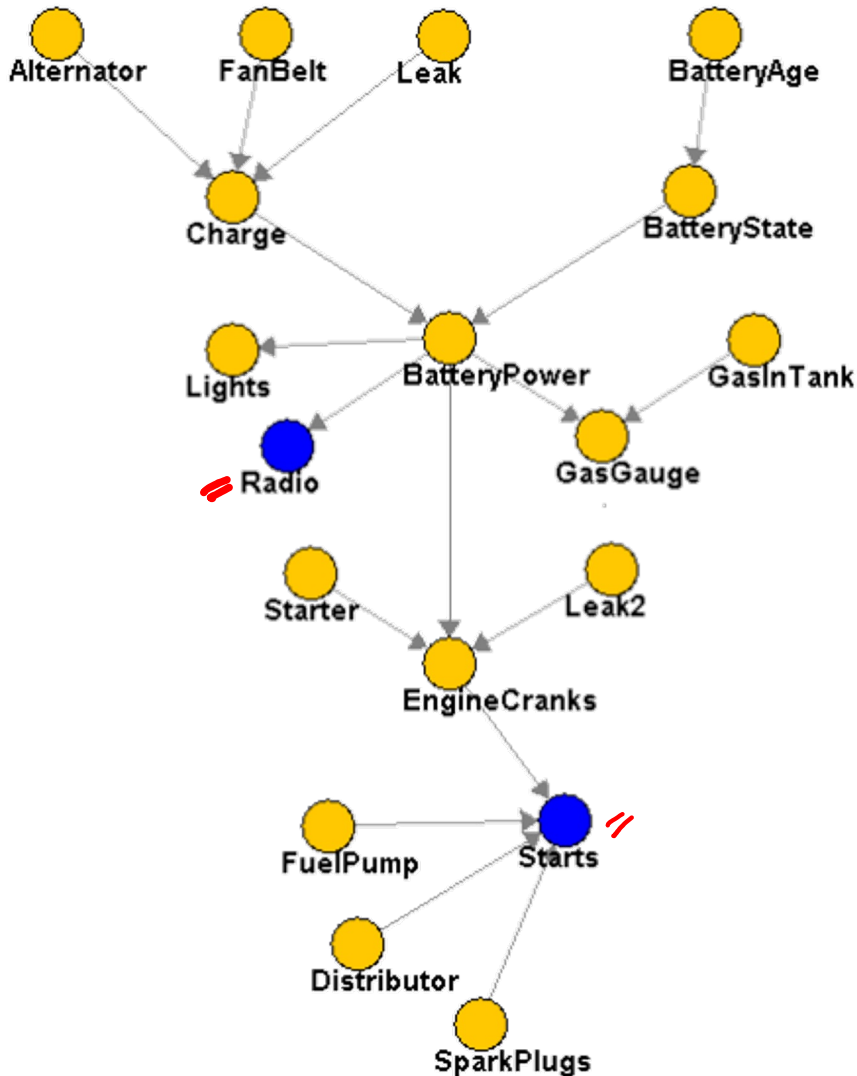
2^3 terms

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 ☹️

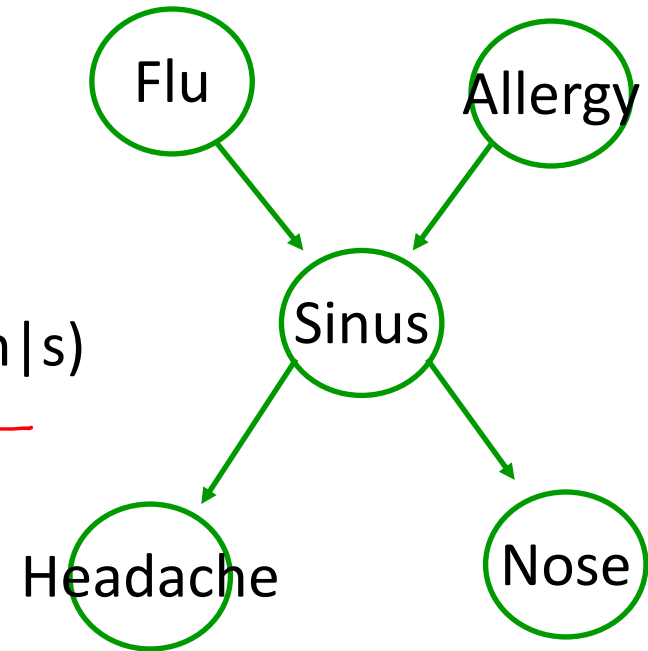
Bayesian Networks Example



- 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

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

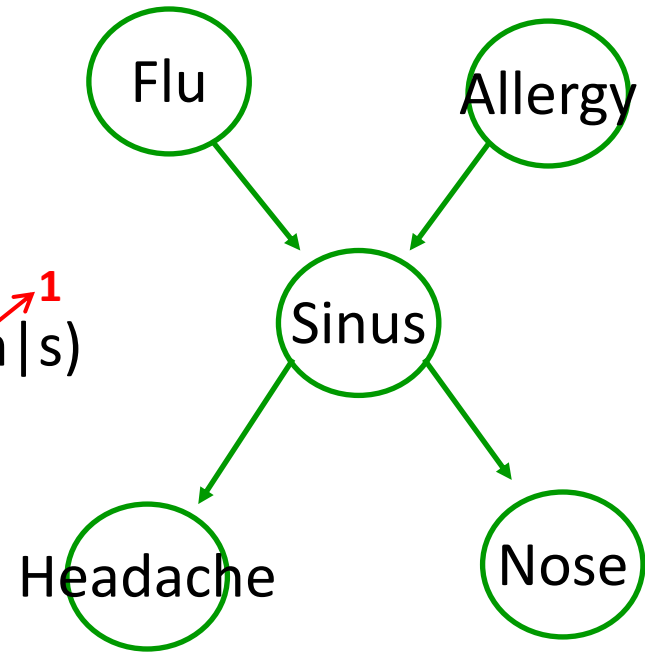


Push sums in as far as possible

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

Fast Probabilistic Inference

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



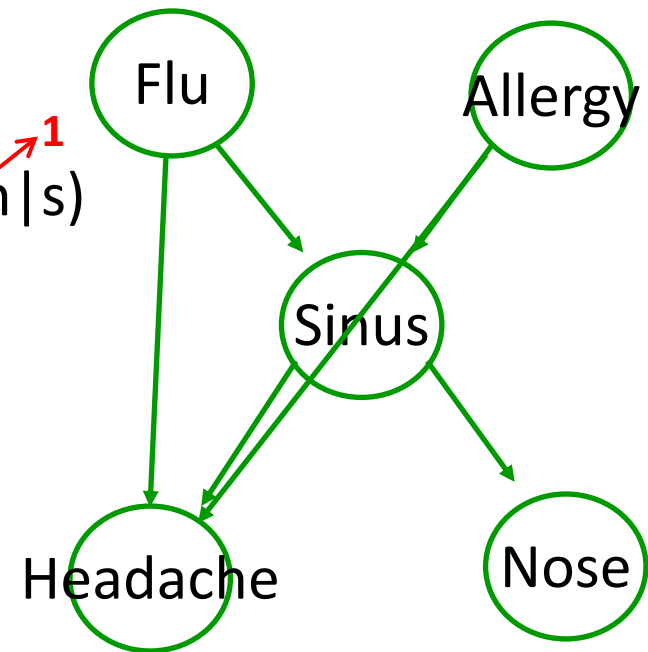
$\tilde{2}^n$ vs. $n \tilde{2}^k$ multiplies
 k - scope of (number of
 variables in) largest factor

(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. Red brackets and arrows highlight the elimination of variables. 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 a '1' above it, indicating that the sum over n is trivial (summing to 1).

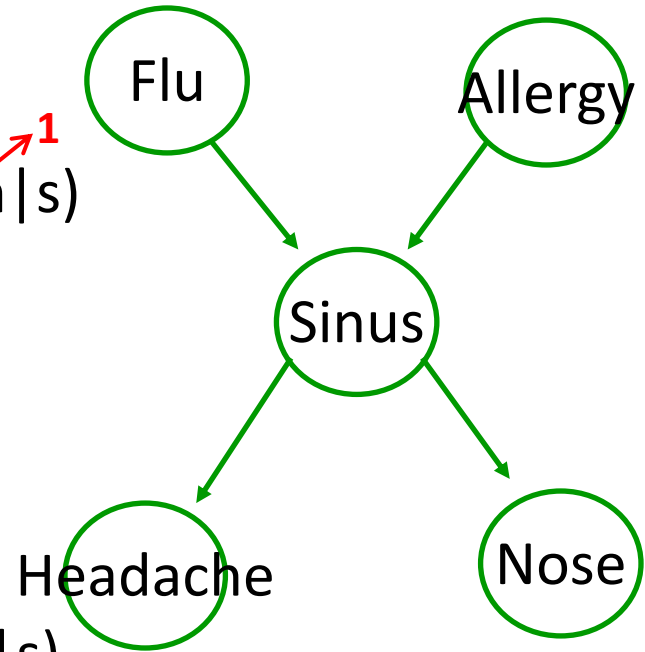


(Potential for) exponential reduction in computation!

Variable Elimination – Order can make a HUGE difference

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

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

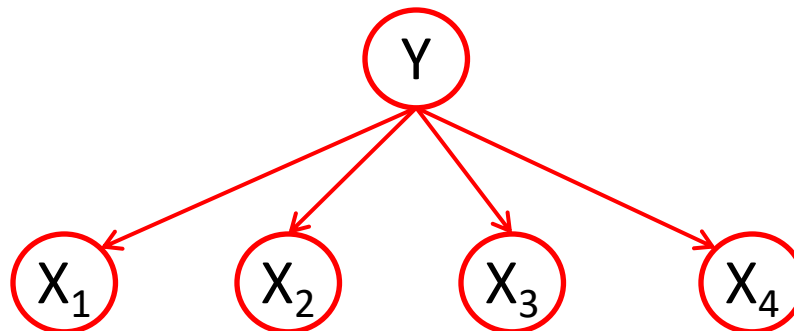


$$\checkmark \quad P(F, H=1) = P(F) \sum_a P(a) \underbrace{\sum_n \sum_s P(s|F,a)P(n|s)P(H=1|s)}_{g(F,a,n)}$$

3 - scope of largest factor

(Potential for) exponential reduction in computation!

Variable Elimination – Order can make a HUGE difference



Naive Bayes

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

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

1 - scope of largest factor

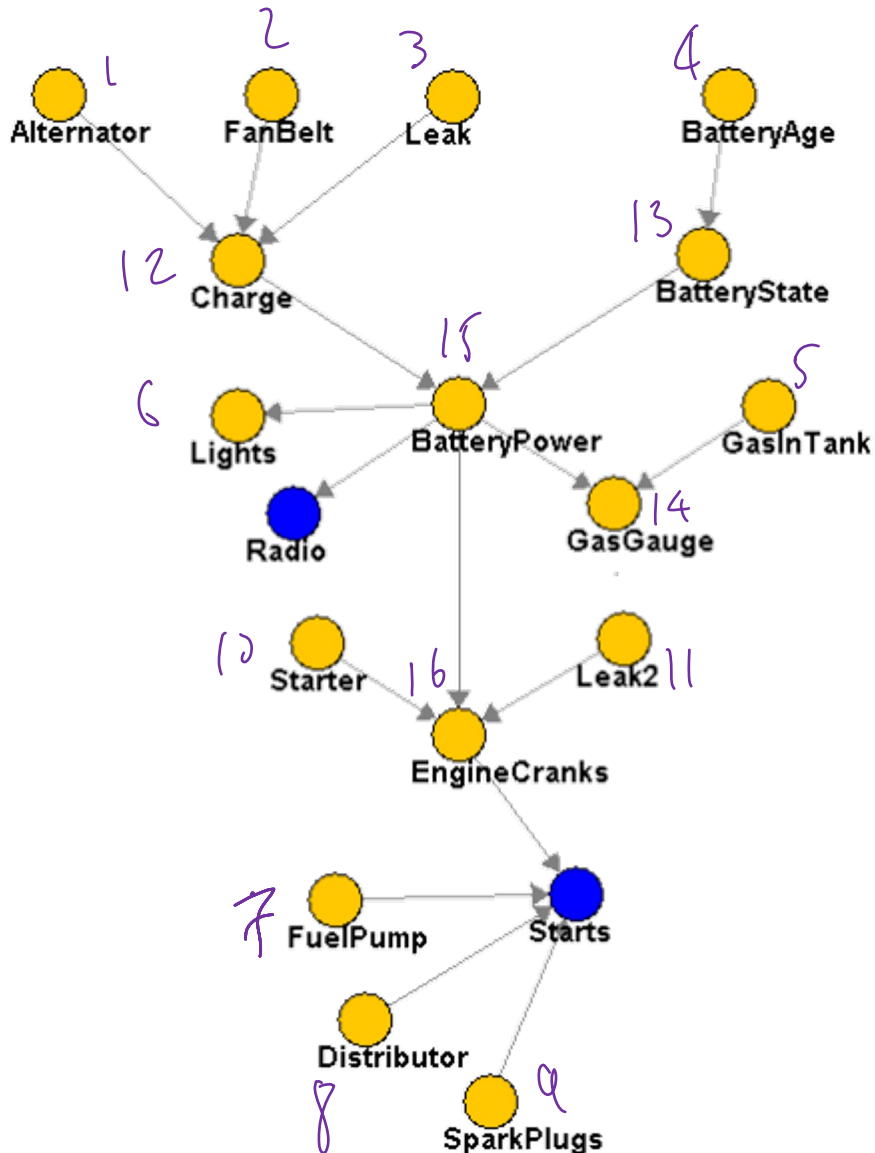
$$= \sum_{X_2, \dots, X_n} \underbrace{\sum_Y P(Y)P(X_1|Y) \prod_{i=2}^n P(X_i|Y)}_{g(X_1, X_2, \dots, X_n)}$$

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

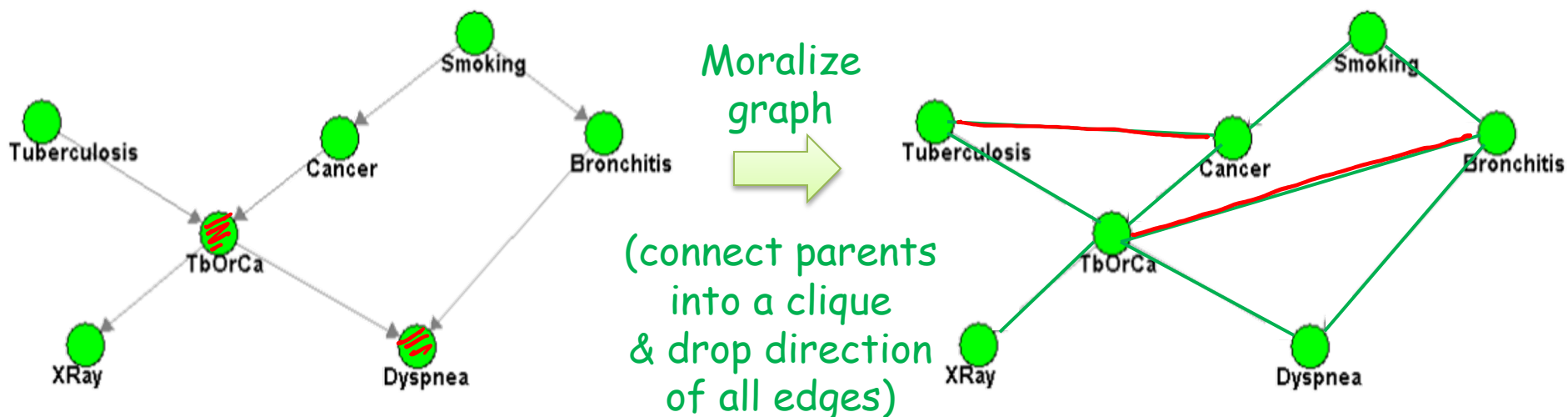
Does not create any factors bigger than original CPTs

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

Complexity for graphs with loops

- Loop – undirected cycle

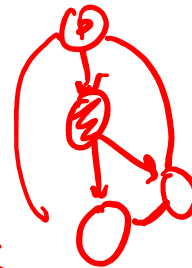
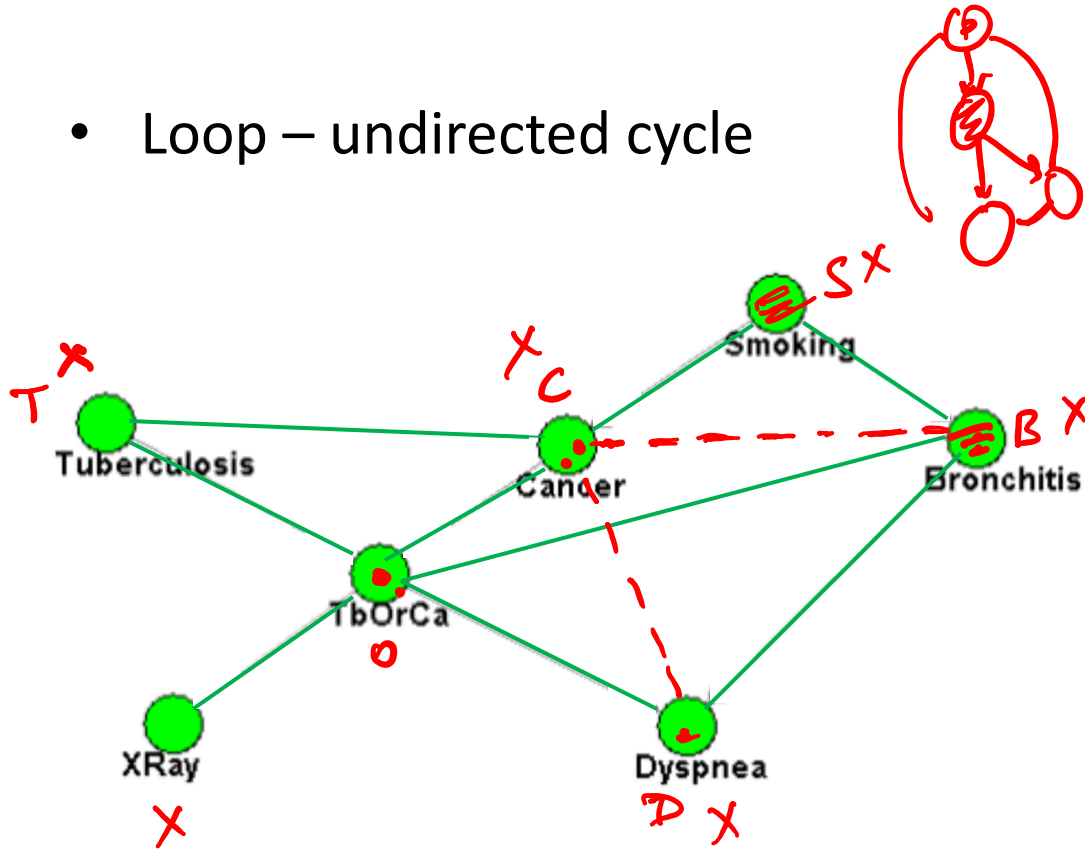
Linear in # variables but exponential in size of largest factor generated!



→ When you eliminate a variable, add edges between its neighbors

Complexity for graphs with loops

- Loop – undirected cycle



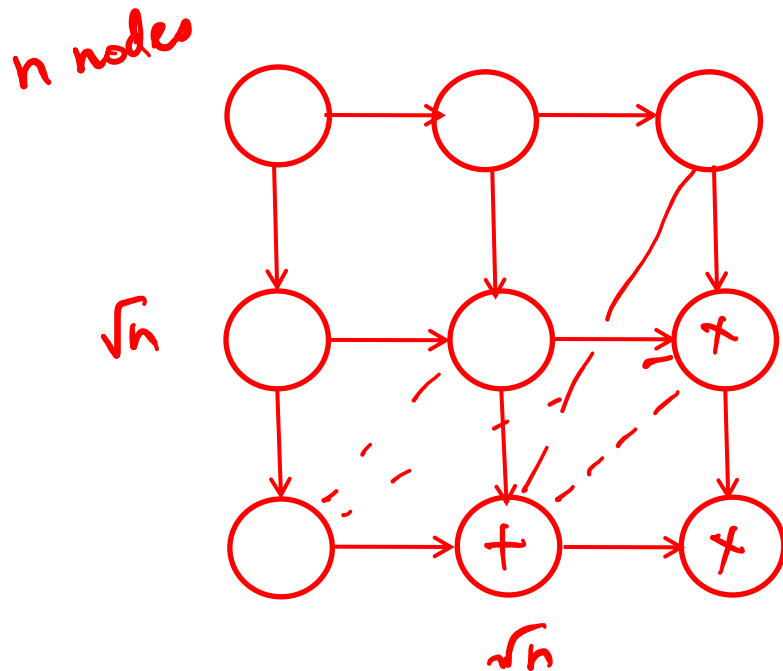
Var eliminated	Factor generated
S	$g_1(C, B)$
B	$g_2(C, O, D)$
D	$g_3(C, O)$
C	$g_4(T, O)$
T	$g_5(O)$
O	$g_6(X)$

T
 O
 C

$g_1(C, O)$
 $g_2(C, D)$
 $g_3(S, B, D)$

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

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

Inference

• Possible queries:

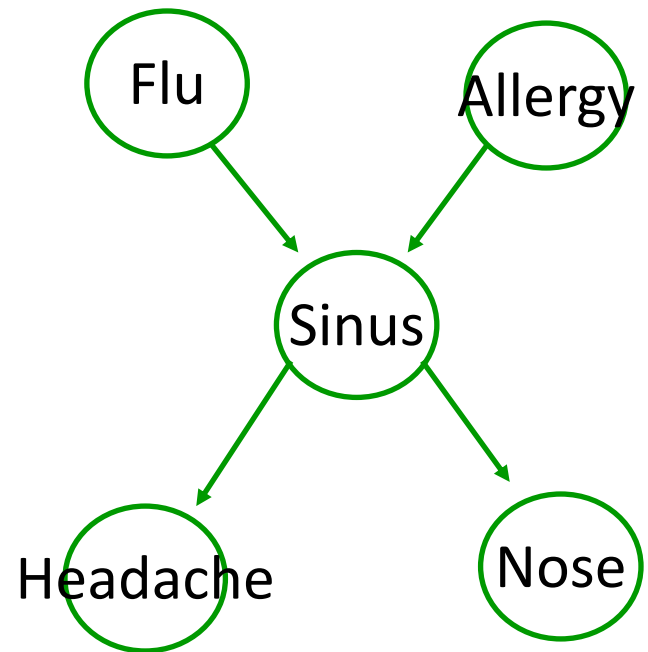
2) Most likely assignment of nodes

$$\arg \max_{f,a,s,n} P(F=f, A=a, S=s, N=n | H=1)$$

Use Distributive property:

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

2 multiply 1 multiply



Topics in Graphical Models

- Representation

- Which joint probability distributions does a graphical model represent?

- Inference

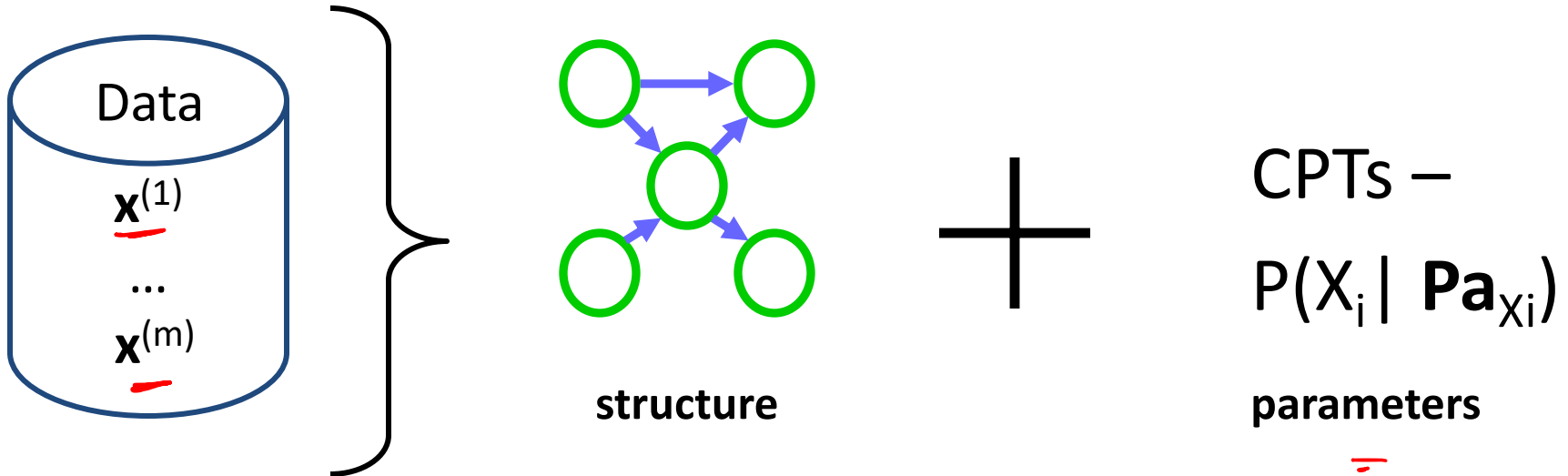
- How to answer questions about the joint probability distribution?

- Marginal distribution of a node variable
- Most likely assignment of node variables

- Learning

- How to learn the parameters and structure of a graphical model?

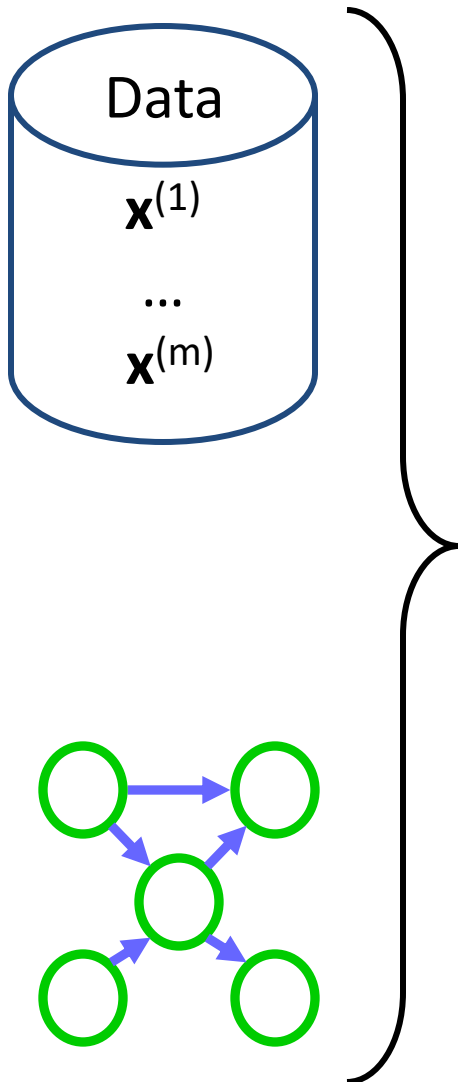
Learning



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

$$\underbrace{p(x_k | \text{pa}_k)}_{\text{graph}} \sim \text{Be}(\theta_{\text{pa}_k}^k)$$

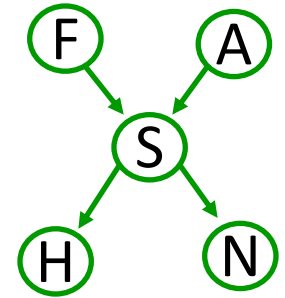
Recall

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

MAP: Add pseudocounts

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

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

examples / training data
nodes / variables

$$= \sum_{i=1}^n \sum_{x_i} \sum_{\mathbf{x}_{\mathbf{Pa}_{X_i}}} \text{count}(X_i = x_i, \mathbf{Pa}_{X_i} = \mathbf{x}_{\mathbf{Pa}_{X_i}}) \log P \left(X_i = x_i \mid \mathbf{Pa}_{X_i} = \mathbf{x}_{\mathbf{Pa}_{X_i}} \right)$$

$\approx m \hat{P}(x_i, \mathbf{x}_{\mathbf{Pa}_{X_i}})$

Plugging in MLE estimates: ML score *of a graphical model $\mathcal{G}_s = \mathcal{G}$*

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

$$= m \sum_{i=1}^n \sum_{x_i} \sum_{\mathbf{x}_{\mathbf{Pa}_{X_i}}} \hat{P}(x_i, \mathbf{x}_{\mathbf{Pa}_{X_i}}) \log \hat{P}(x_i \mid \mathbf{x}_{\mathbf{Pa}_{X_i}})$$

Reminds of entropy

Information theoretic interpretation of MLE

$$\begin{aligned}\log \hat{P}(\mathcal{D} \mid \hat{\theta}_{\mathcal{G}}, \mathcal{G}) &= m \sum_{i=1}^n \sum_{x_i} \sum_{\mathbf{xPa}_{X_i}} \hat{P}(x_i, \mathbf{xPa}_{X_i}) \log \hat{P}(x_i \mid \mathbf{xPa}_{X_i}) \\ &= -m \sum_{i=1}^n \hat{H}(X_i \mid \mathbf{Pa}_{X_i}) \\ &= m \sum_{i=1}^n [\hat{I}(X_i, \mathbf{Pa}_{X_i}) - \hat{H}(X_i)]\end{aligned}$$

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

ML score for graph structure \mathcal{G}

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