

# Bayesian Networks – Representation

Machine Learning – CSE546

Carlos Guestrin

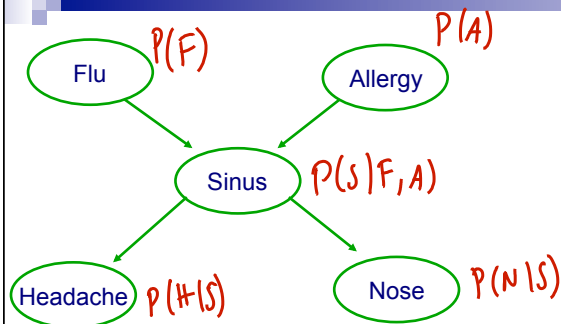
University of Washington

November 20, 2014

©Carlos Guestrin 2005-2014

1

## Factored joint distribution - Preview



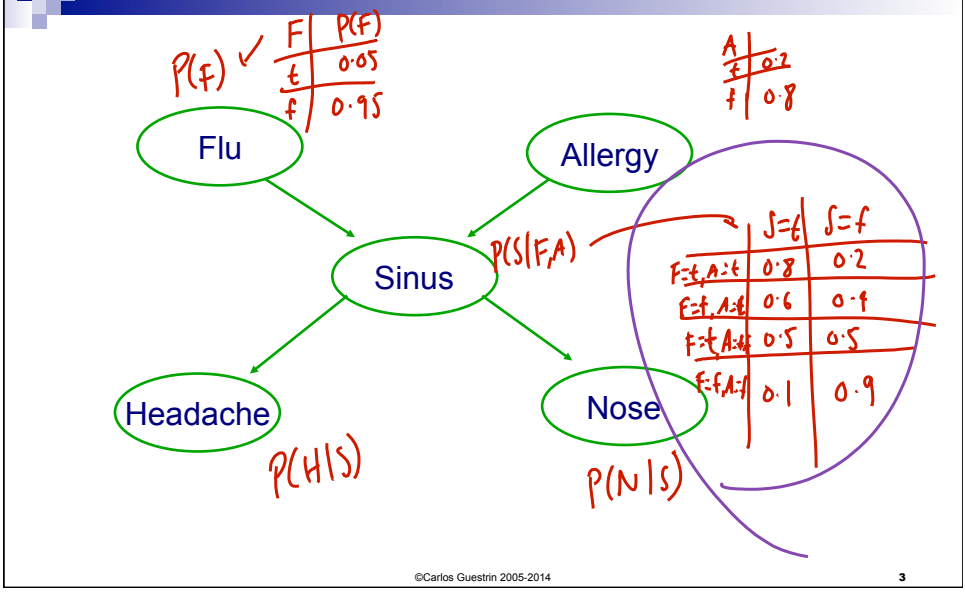
$$P(F, A, S, H, N) = P(F) P(A) P(S|F, A) P(H|S) P(N|S)$$

$2^5 - 1 = 31$  parameters

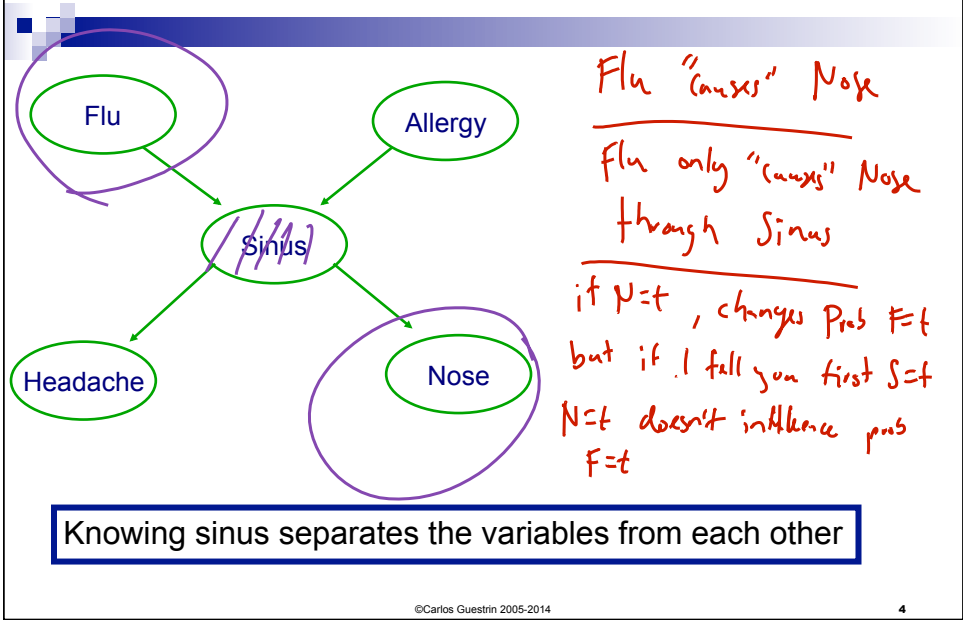
©Carlos Guestrin 2005-2014

2

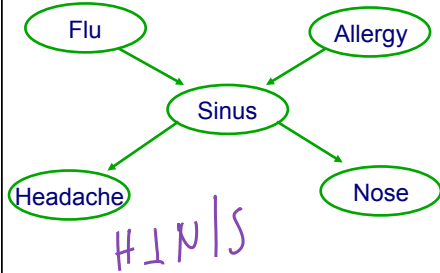
# What about probabilities? Conditional probability tables (CPTs)



# Key: Independence assumptions



# The independence assumption



**Local Markov Assumption:**  
 A variable X is independent of its non-descendants given its parents *and only its parents*

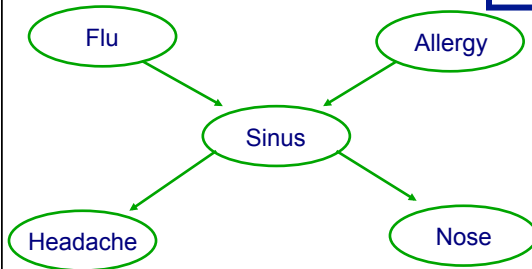
	F	A	S	H	N
non descendent	A	F	FA	FAU(S)	FAH
implies	F ⊥ A	A ⊥ F	S ⊥ FA / FA ⇒ nothing	H ⊥ FAN   S	N ⊥ FAH   S

©Carlos Guestrin 2005-2014

5

## Explaining away

**Local Markov Assumption:**  
 A variable X is independent of its non-descendants given its parents



$F \perp A$   
 $F \perp A | S ??$   
 don't know  
 $P(F=t | S=t) \neq P(F=t | S=t, A=t)$   
 $>$

©Carlos Guestrin 2005-2014

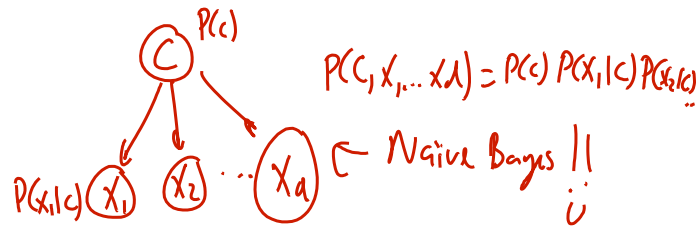
6

# Naïve Bayes revisited

$$\forall_i X_i \perp \{X_2, X_3, \dots, X_d\} \mid C$$

$$P(C, X_1, \dots, X_d) = P(C) \prod_i P(X_i \mid C)$$

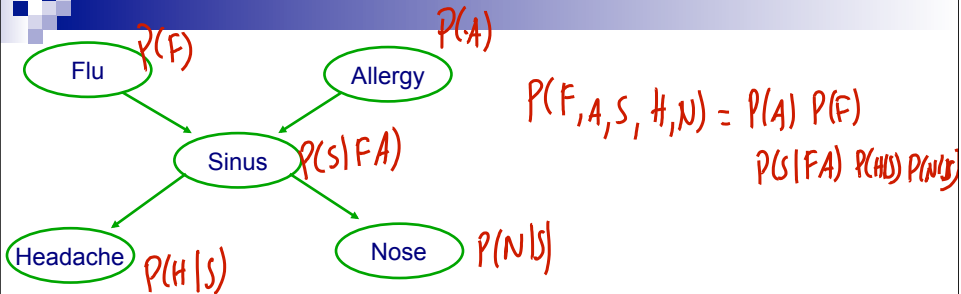
**Local Markov Assumption:**  
A variable X is independent of its non-descendants given its parents



©Carlos Guestrin 2005-2014

7

# Joint distribution



**Why can we decompose? Markov Assumption!**

©Carlos Guestrin 2005-2014

# The chain rule of probabilities

- $P(A,B) = P(A)P(B|A) = P(B)P(A|B)$

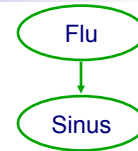
For any dist

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

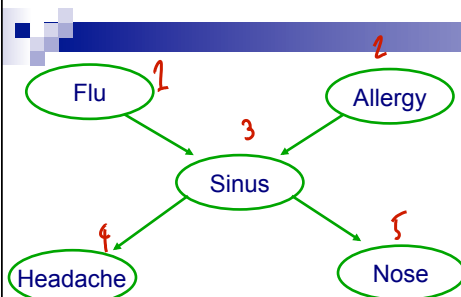
$$P(F,A,S) = P(F)P(A|F)P(S|F,A)$$

- More generally:

- $P(X_1, \dots, X_n) = P(X_1)P(X_2|X_1) \dots P(X_n|X_1, \dots, X_{n-1})$



# Chain rule & Joint distribution



## Local Markov Assumption:

A variable X is independent of its non-descendants given its parents

Proof by example! Holds for any BN

$$P(FASHN) = P(F)P(A|F)P(S|FA)P(H|FAS)P(N|FASH)$$

$$= P(F)P(A)P(S|FA)P(H|S)P(N|S) \quad !! \text{ good}$$

order on Hesi  
 $P(FASHN) = P(F)P(A|F)P(S|FA)P(H|FAS)P(N|FASH)$   
 would not get follow topological order

---


$$A \perp F \Rightarrow P(A|F) = P(A) \quad | \quad H \perp \{F,A\} | S \Rightarrow P(H|FAS) = P(H|S) \quad | \quad N \perp \{F,A,H\} | S \Rightarrow P(N|FASH) = P(N|S)$$

# The Representation Theorem – Joint Distribution to BN

BN:  Encodes independence assumptions

If conditional independencies in BN are subset of conditional independencies in  $P$

Obtain

Joint probability distribution:

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | \text{Pa}_{X_i})$$

# Two (trivial) special cases

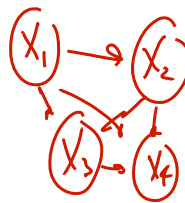
Edgeless graph



$X_i \perp\!\!\!\perp$  all others

fewest params  
high bias

Fully-connected graph



no independence

most params  
high variance

structure learning

# Bayesian Networks – (Structure) Learning

Machine Learning – CSE546

Carlos Guestrin

University of Washington

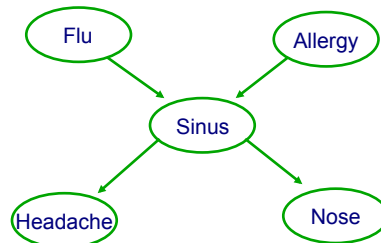
November 20, 2014

©Carlos Guestrin 2005-2014

13

## Review

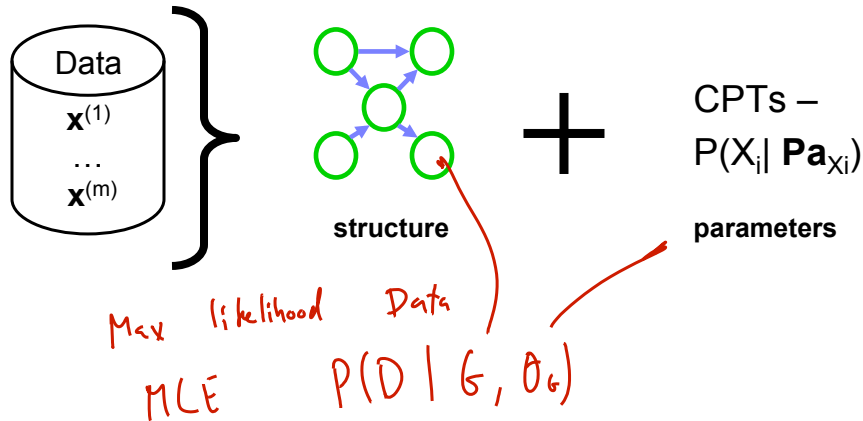
- Bayesian Networks
  - Compact representation for probability distributions
  - Exponential reduction in number of parameters
- Fast probabilistic inference
  - As shown in demo examples
  - Compute  $P(X|e)$
- Today
  - Learn BN structure



©Carlos Guestrin 2005-2014

14

# Learning Bayes nets

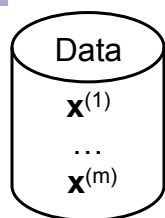


©Carlos Guestrin 2005-2014

15

# Learning the CPTs

*var with no parents*  
 $P(X_i = x_i) = \frac{\text{count}(X_i = x_i)}{\#m}$



For each discrete variable  $X_i$

$$P(S=t | A=t, F=f) \stackrel{\text{MLE}}{=} \frac{\text{count}(S=t, A=t, F=f)}{\text{count}(A=t, F=f)}$$

$$P(X_i = x_i | \text{Pa}_{X_i} = u) \stackrel{\text{MLE}}{=} \frac{\text{count}(X_i = x_i, \text{Pa}_{X_i} = u)}{\text{count}(\text{Pa}_{X_i} = u)}$$

*Subtlety: count(Pa\_{X\_i} = u) = 0 or very small*

*→ add smoothing / AKA regularization / AKA Bayesian prior*

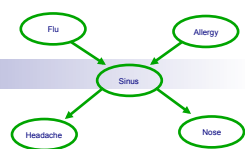
$$\text{MLE: } P(X_i = x_i | X_j = x_j) = \frac{\text{Count}(X_i = x_i, X_j = x_j)}{\text{Count}(X_j = x_j)}$$

©Carlos Guestrin 2005-2014

16



# Information-theoretic interpretation of maximum likelihood 1



Given structure, log likelihood of data:

$$\log P(\mathcal{D} | \theta_G, \mathcal{G}) \stackrel{\text{iid}}{=} \log \prod_{j=1}^m P(x_1^{(j)}, \dots, x_n^{(j)} | \theta_G)$$

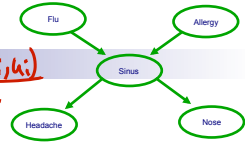
$$= \log \prod_{j=1}^m \prod_{i=1}^n P(x_i = x_i^{(j)} | \text{Pa}_{X_i} = u_i^{(j)})$$

$$= \sum_{j=1}^m \sum_{i=1}^n \log P(X_i = x_i^{(j)} | \text{Pa}_{X_i} = u_i^{(j)})$$

graph defines who parents are

$x_i^{(j)}$  ← data point (F=f, A=t, S=t, H=t, N=f)  
 ↙ variable F, A, S, H, N

# Information-theoretic interpretation of maximum likelihood 2



Given structure, log likelihood of data:

$$\log P(\mathcal{D} | \theta_G, \mathcal{G}) = \sum_{j=1}^m \sum_{i=1}^n \log P(X_i = x_i^{(j)} | \text{Pa}_{X_i} = \bar{x}^{(j)}[\text{Pa}_{X_i}])$$

$$= \sum_{i=1}^n \sum_{j=1}^m \log P(X_i = x_i^{(j)} | \text{Pa}_{X_i} = u_i^{(j)}) \leftarrow \sum_{j=1}^m \log P(h^{(j)} | s^{(j)})$$

$$= \sum_{i=1}^n \sum_{x_i \in \mathcal{X}_i} \sum_{u_i \in \text{Pa}_{X_i}} \text{Count}(X_i = x_i, \text{Pa}_{X_i} = u_i) \log \hat{P}(x_i | u_i)$$

$$= \text{Count}(H=t, S=t) \log P(H=t | S=t)$$

$$= m \sum_{i=1}^n \sum_{x_i} \sum_{u_i} \hat{P}(X_i = x_i, \text{Pa}_{X_i} = u_i) \log \hat{P}(x_i | u_i)$$

$$= \text{Count}(H=f, S=t) \log P(H=f | S=t)$$

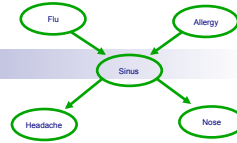
$$= -m \sum_{i=1}^n H(X_i | \text{Pa}_{X_i})$$

$$= \text{Count}(H=t, S=f) \log P(H=t | S=f)$$

$$= \text{Count}(H=f, S=f) \log P(H=f | S=f)$$

$$+ \text{Count}(H=t, S=f) \log P(H=t | S=f)$$

# Information-theoretic interpretation of maximum likelihood 3



■ Given structure, log likelihood of data:

$$\max_G \log \hat{P}(\mathcal{D} | \theta, \mathcal{G}) = m \sum_i \sum_{x_i, \text{Pa}_{x_i, \mathcal{G}}} \hat{P}(x_i, \text{Pa}_{x_i, \mathcal{G}}) \log \hat{P}(x_i | \text{Pa}_{x_i, \mathcal{G}})$$

$$\equiv \max_G -m \sum_{i=1}^n H(x_i | \text{Pa}_{x_i, \mathcal{G}}) = \min_G m \sum_{i=1}^n H(x_i | \text{Pa}_{x_i, \mathcal{G}})$$

$$\equiv \max_G \left[ \sum_{i=1}^n \hat{I}(x_i | \text{Pa}_{x_i, \mathcal{G}}) - m \sum_{i=1}^n H(x_i) \right]$$

Information Theoretic interpretation of MLE for  $\mathcal{G}$  doesn't depend on  $\mathcal{G}$

$\Rightarrow \max_G \equiv$  choosing parents with max mutual info with vars

in DTs  
 $I(A, B) = H(A) - H(A|B)$   
 if  $A \perp B$   
 $I(A, B) = 0$   
 if A very dependent with B,  $I(A, B)$  high  
 $\frac{\hat{I}(x_i, \text{Pa}_{x_i, \mathcal{G}})}{H(x_i) - H(x_i | \text{Pa}_{x_i, \mathcal{G}})}$

# Decomposable score

■ Log data likelihood

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

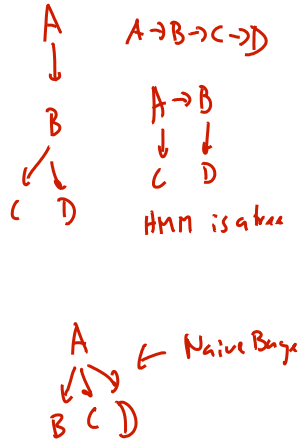
■ Decomposable score:

- Decomposes over families in BN (node and its parents)
- Will lead to significant computational efficiency!!!
- $\text{Score}(\mathcal{G} : \mathcal{D}) = \sum_{i=1}^n \text{FamScore}(X_i | \text{Pa}_{X_i} : \mathcal{D})$

$$\hat{I}(x_i | \text{Pa}_{x_i, \mathcal{G}})$$

# How many trees are there?

**Nonetheless – Efficient optimal algorithm finds best tree**

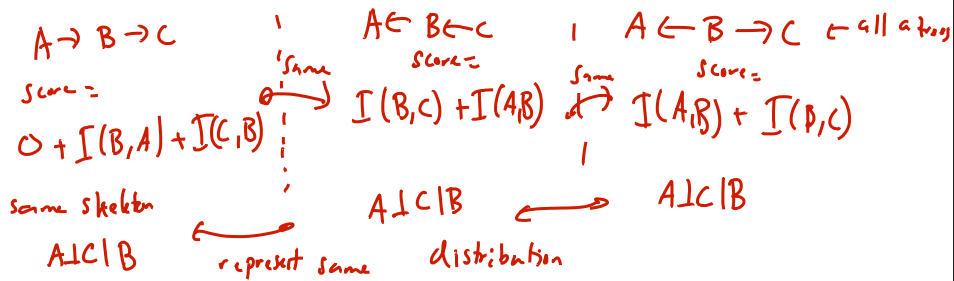


Every var has at most one parent  
 For  $n$  vars, how many possible trees?  
 $O(n \log n)$   
 $2$   
 exhaustive search is impossible

# Scoring a tree 1: equivalent trees

$$I(A, B) = I(B, A)$$

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

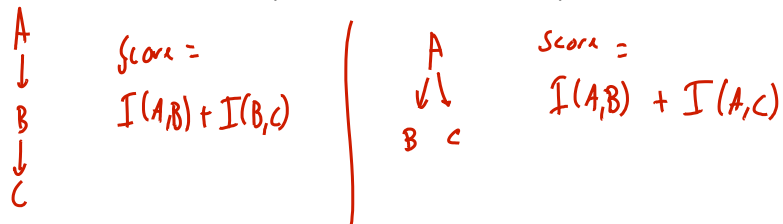


no tree for all edge directions:  

 not a tree, because B has 2 parents

## Scoring a tree 2: similar trees

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



in general for trees:  $\text{Score} = \max_{\mathcal{G}} \sum_{(i,j) \in E} I(x_i, x_j)$

©Carlos Guestrin 2005-2014

23

## Chow-Liu tree learning algorithm 1

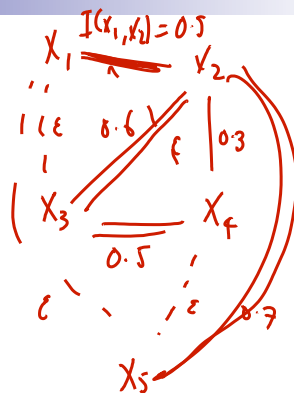
- For each pair of variables  $X_i, X_j$ 
  - Compute empirical distribution:

$$\hat{P}(x_i, x_j) = \frac{\text{Count}(x_i, x_j)}{m}$$

- 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  $\hat{I}(X_i, X_j)$



Run max spanning tree ← complexity is about  $O(E \log E)$   
 $O(n^2 \log n)$

©Carlos Guestrin 2005-2014

24

## Chow-Liu tree learning algorithm 2

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

### ■ Optimal tree BN

- Compute maximum weight spanning tree
- Directions in BN: pick any node as root, breadth-first-search defines directions

