

# Bayesian Networks – Representation

Machine Learning – CSE546

Carlos Guestrin

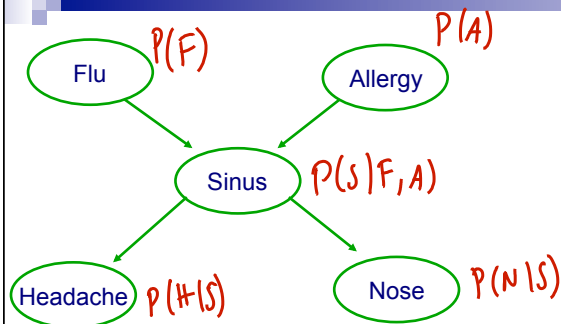
University of Washington

November 20, 2014

©Carlos Guestrin 2005-2014

1

## Factored joint distribution - Preview



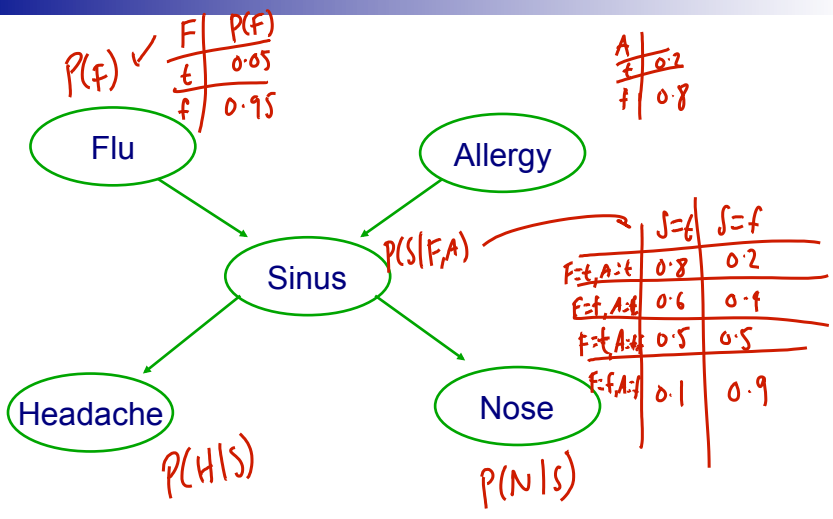
$$P(F,A,S,H,N) = P(F) P(A) P(S|F,A) P(H|S) P(N|S)$$

$2^5 - 1 = 31$  parameters

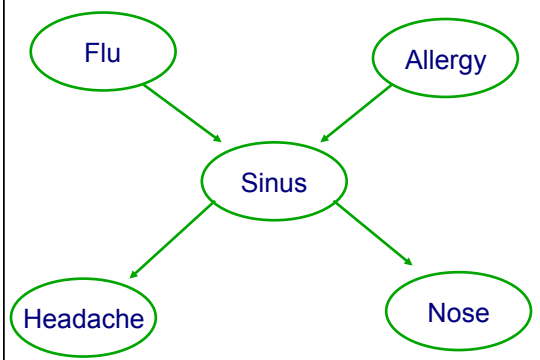
©Carlos Guestrin 2005-2014

2

# What about probabilities? Conditional probability tables (CPTs)



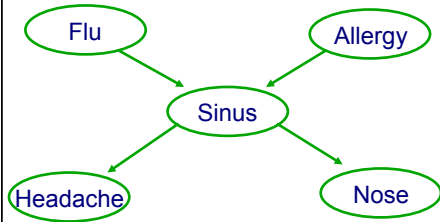
# Key: Independence assumptions



Flu "causes" Nose  
 Flu only "causes" Nose through Sinus  
if  $N=t$ , changes Prob  $F=t$  but if I tell you first  $S=t$   $N=t$  doesn't influence prob  $F=t$

Knowing sinus separates the variables from each other

# The independence assumption



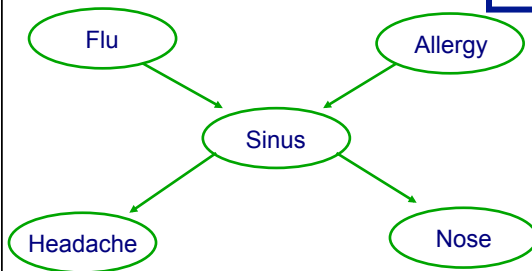
**Local Markov Assumption:**  
 A variable X is independent of its non-descendants given its parents *and only its parents*

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

©Carlos Guestrin 2005-2014

5

# Explaining away



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

©Carlos Guestrin 2005-2014

6

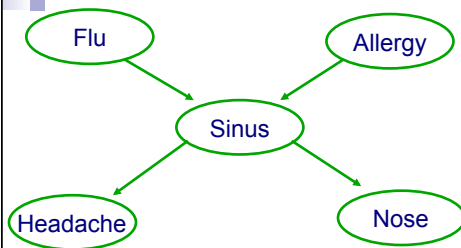
# Naïve Bayes revisited

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

©Carlos Guestrin 2005-2014

7

# Joint distribution

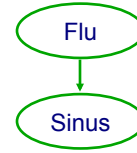


**Why can we decompose? Markov Assumption!**

©Carlos Guestrin 2005-2014

# The chain rule of probabilities

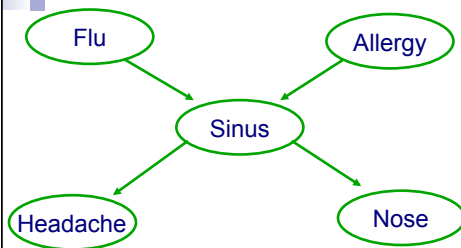
- $P(A,B) = P(A)P(B|A)$



- More generally:

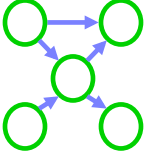
- $P(X_1, \dots, X_n) = P(X_1) P(X_2|X_1) \dots P(X_n|X_1, \dots, X_{n-1})$

# Chain rule & Joint distribution



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

## 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:**

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

©Carlos Guestrin 2005-2014

11

## Two (trivial) special cases

Edgeless graph

Fully-connected graph

©Carlos Guestrin 2005-2014

12

# Bayesian Networks – (Structure) Learning

Machine Learning – CSE546

Carlos Guestrin

University of Washington

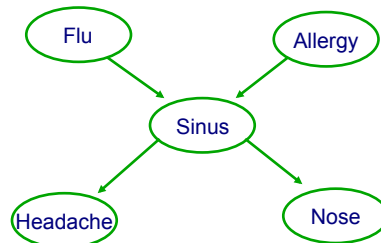
November 20, 2014

©Carlos Guestrin 2005-2014

13

## Review

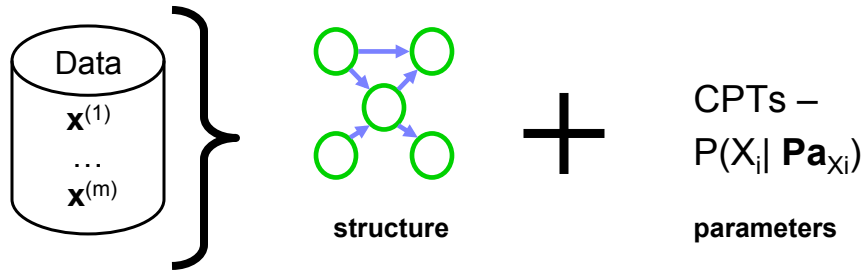
- Bayesian Networks
  - Compact representation for probability distributions
  - Exponential reduction in number of parameters
- Fast probabilistic inference
  - As shown in demo examples
  - Compute  $P(X|e)$
- Today
  - Learn BN structure



©Carlos Guestrin 2005-2014

14

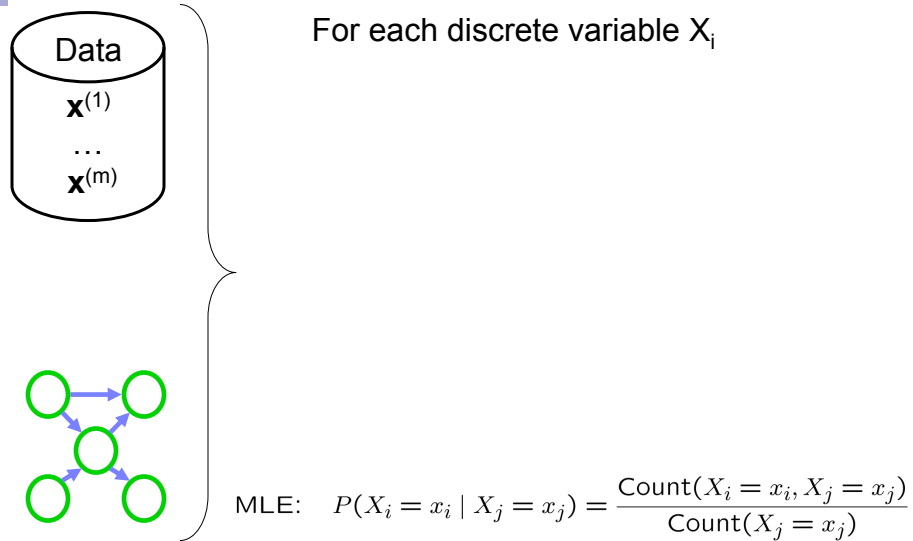
# Learning Bayes nets



©Carlos Guestrin 2005-2014

15

# Learning the CPTs



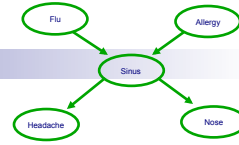
©Carlos Guestrin 2005-2014

16



## Information-theoretic interpretation of maximum likelihood 1

- Given structure, log likelihood of data:  
 $\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G})$



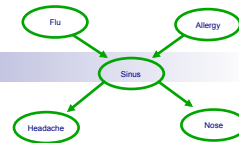
©Carlos Guestrin 2005-2014

17

## Information-theoretic interpretation of maximum likelihood 2

- Given structure, log likelihood of data:

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



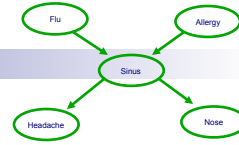
©Carlos Guestrin 2005-2014

18

## Information-theoretic interpretation of maximum likelihood 3

- Given structure, log likelihood of data:

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



## Decomposable score

- Log data likelihood

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

- Decomposable score:

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

## How many trees are there?

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

## Scoring a tree 1: equivalent trees

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

## Scoring a tree 2: similar trees

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

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

## 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

©Carlos Guestrin 2005-2014

25

## 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)  $d > 1$**
- Most structure learning approaches use heuristics
  - (Quickly) Describe the two simplest heuristic

©Carlos Guestrin 2005-2014

26

# Learn BN structure using local search

Starting from Chow-Liu tree

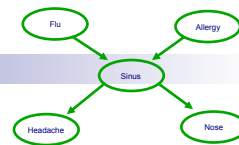
Local search, possible moves:

- Add edge
- Delete edge
- Invert edge

Score using BIC

# Learn Graphical Model Structure using LASSO

■ Graph structure is about selecting parents:



- If no independence assumptions, then CPTs depend on all parents:
- With independence assumptions, depend on key variables:
- One approach for structure learning, sparse logistic regression!

## What you need to know about learning BN structures

- Decomposable scores
  - Maximum likelihood
  - Information theoretic interpretation
- Best tree (Chow-Liu)
- Beyond tree-like models is NP-hard
- Use heuristics, such as:
  - Local search
  - LASSO