

Bayesian Networks – Representation

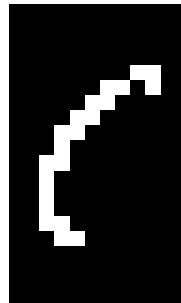
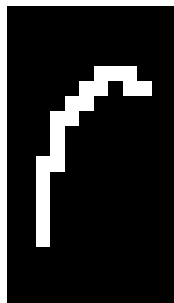
Machine Learning – CSE546
Carlos Guestrin
University of Washington

November 25, 2013

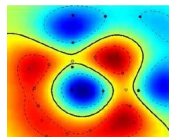
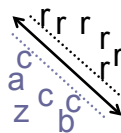
©Carlos Guestrin 2005-2013

1

Handwriting recognition



Character recognition, e.g., kernel SVMs



©Carlos Guestrin 2005-2013

2

Webpage classification



Company home page

VS

Personal home page

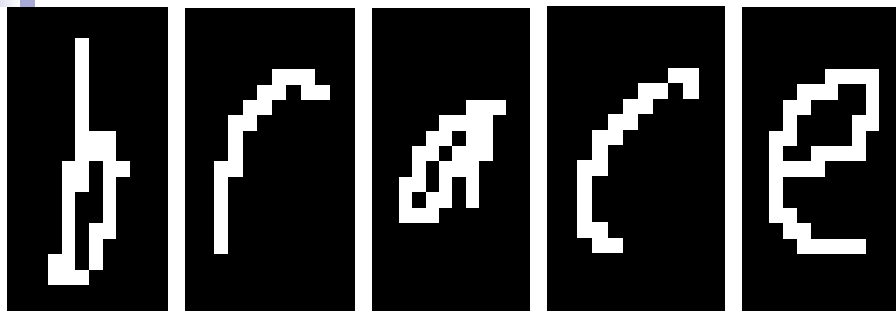
VS

University home page

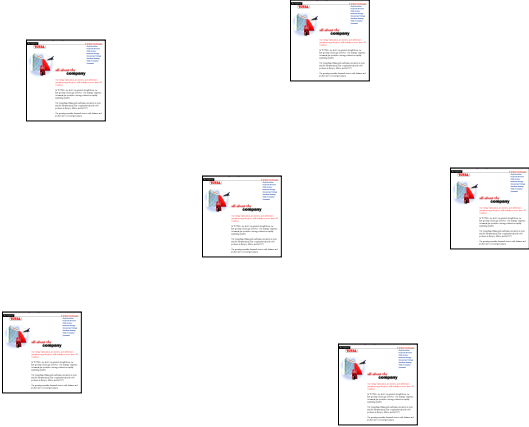
VS

...

Handwriting recognition 2



Webpage classification 2



©Carlos Guestrin 2005-2013 5

Today – Bayesian networks

- One of the most exciting advancements in statistical AI in the last decades
- Generalizes naïve Bayes and logistic regression classifiers
- Compact representation for exponentially-large probability distributions
- Exploit conditional independencies

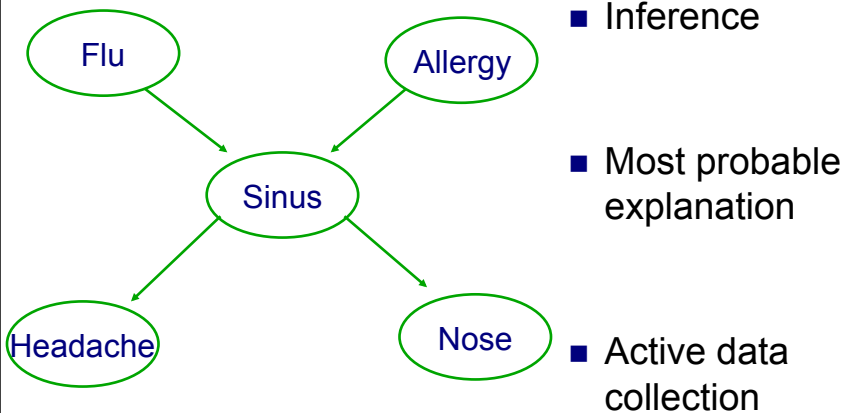
Causal structure

- Suppose we know the following:
 - The flu causes sinus inflammation
 - Allergies cause sinus inflammation
 - Sinus inflammation causes a runny nose
 - Sinus inflammation causes headaches
- How are these connected?

©Carlos Guestrin 2005-2013

7

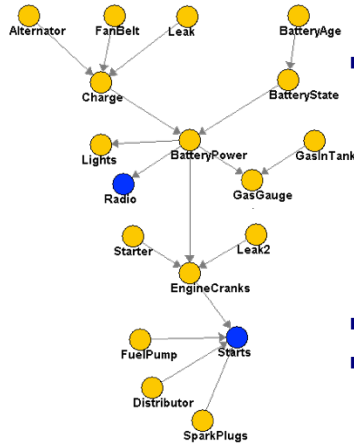
Possible queries



©Carlos Guestrin 2005-2013

8

Car starts BN



- 18 binary attributes

- Inference

- $P(\text{BatteryAge} | \text{Starts}=f)$

- 2^{16} terms, why so fast?

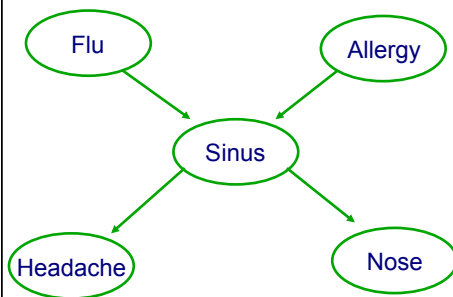
- Not impressed?

- HailFinder BN – more than $3^{54} = 58149737003040059690390169$ terms

©Carlos Guestrin 2005-2013

9

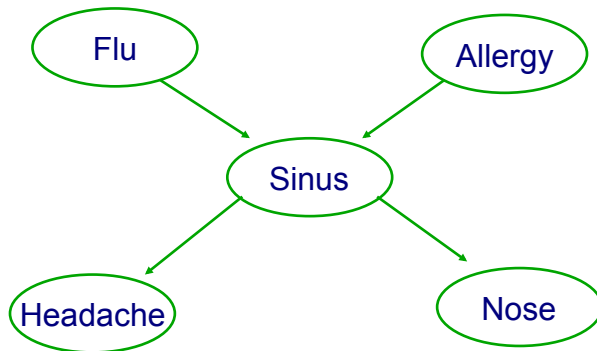
Factored joint distribution - Preview



©Carlos Guestrin 2005-2013

10

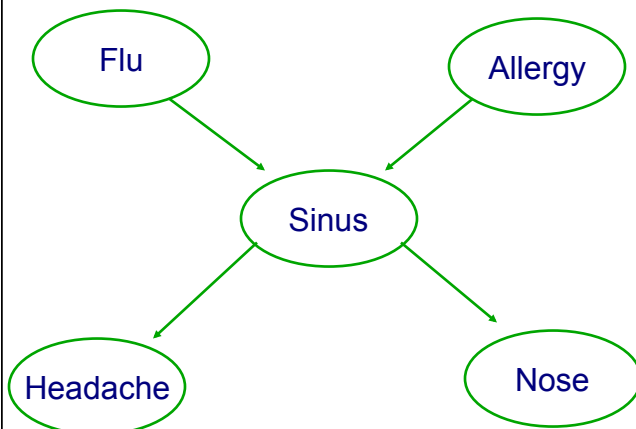
What about probabilities? Conditional probability tables (CPTs)



©Carlos Guestrin 2005-2013

11

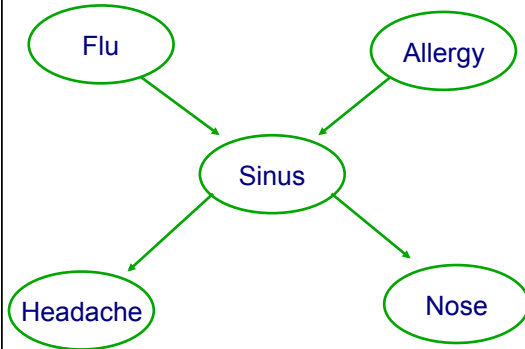
Number of parameters



©Carlos Guestrin 2005-2013

12

Key: Independence assumptions



Knowing sinus separates the variables from each other

©Carlos Guestrin 2005-2013

13

(Marginal) Independence

- Flu and Allergy are (marginally) independent

Flu = t	
Flu = f	

Allergy = t	
Allergy = f	

	Flu = t	Flu = f
Allergy = t		
Allergy = f		

©Carlos Guestrin 2005-2013

14

Marginally independent random variables

- **Sets** of variables \mathbf{X}, \mathbf{Y}
- \mathbf{X} is independent of \mathbf{Y} if
 - $P(\mathbf{X}=\mathbf{x} \perp \mathbf{Y}=\mathbf{y}), \forall \mathbf{x} \in \text{Val}(\mathbf{X}), \mathbf{y} \in \text{Val}(\mathbf{Y})$
- Shorthand:
 - **Marginal independence:** $P(\mathbf{X} \perp \mathbf{Y})$
- **Proposition:** P satisfies $(\mathbf{X} \perp \mathbf{Y})$ if and only if
 - $P(\mathbf{X}, \mathbf{Y}) = P(\mathbf{X}) P(\mathbf{Y})$

Conditional independence

- Flu and Headache are not (marginally) independent
- Flu and Headache are independent given Sinus infection
- More Generally:

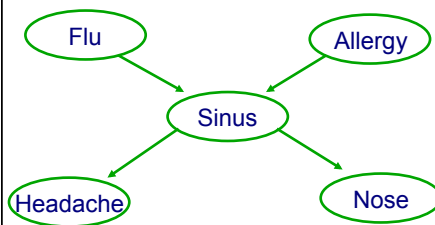
Conditionally independent random variables

- **Sets** of variables **X, Y, Z**
- X is independent of Y given Z if
 - $P F (X=x \perp Y=y | Z=z), \forall x \in \text{Val}(X), y \in \text{Val}(Y), z \in \text{Val}(Z)$
- Shorthand:
 - **Conditional independence:** $P F (X \perp Y | Z)$
 - For $P F (X \perp Y | \emptyset)$, write $P F (X \perp Y)$
- **Proposition:** P satisfies $(X \perp Y | Z)$ if and only if
 - $P(X, Y | Z) = P(X | Z) P(Y | Z)$

©Carlos Guestrin 2005-2013

17

The independence assumption



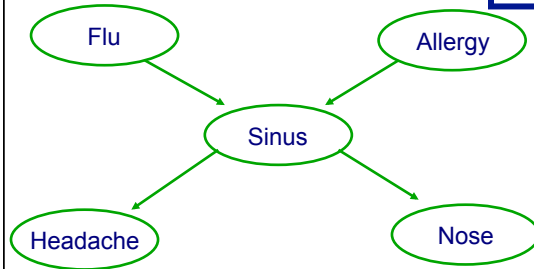
Local Markov Assumption:
A variable X is independent of its non-descendants given its parents

©Carlos Guestrin 2005-2013

18

Explaining away

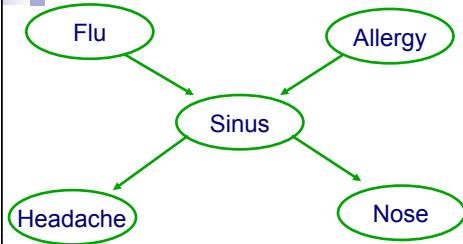
Local Markov Assumption:
A variable X is independent of its non-descendants given its parents



Naïve Bayes revisited

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

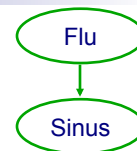
Joint distribution



Why can we decompose? Markov Assumption!

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

```

    graph TD
      Flu((Flu)) --> Sinus((Sinus))
      Allergy((Allergy)) --> Sinus((Sinus))
      Sinus((Sinus)) --> Headache((Headache))
      Sinus((Sinus)) --> Nose((Nose))
    
```

©Carlos Guestrin 2005-2013 23

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:**

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

©Carlos Guestrin 2005-2013 24

Two (trivial) special cases

Edgeless graph

Fully-connected
graph

©Carlos Guestrin 2005-2013

25

Bayesian Networks – (Structure) Learning

Machine Learning – CSE546

Carlos Guestrin

University of Washington

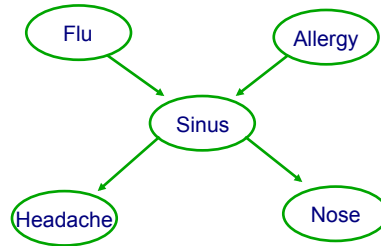
November 25, 2013

©Carlos Guestrin 2005-2013

26

Review

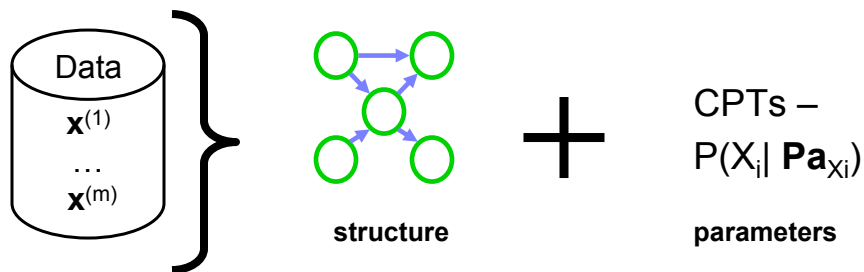
- Bayesian Networks
 - Compact representation for probability distributions
 - Exponential reduction in number of parameters
- Fast probabilistic inference
 - As shown in demo examples
 - Compute $P(X|e)$
- Today
 - Learn BN structure



©Carlos Guestrin 2005-2013

27

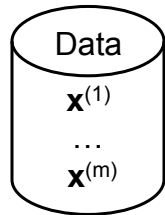
Learning Bayes nets



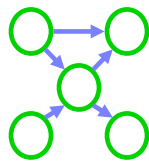
©Carlos Guestrin 2005-2013

28

Learning the CPTs



For each discrete variable X_i



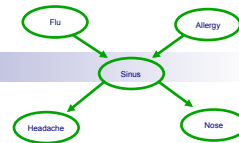
$$\text{MLE: } P(X_i = x_i | X_j = x_j) = \frac{\text{Count}(X_i = x_i, X_j = x_j)}{\text{Count}(X_j = x_j)}$$

©Carlos Guestrin 2005-2013

29

Information-theoretic interpretation of maximum likelihood 1

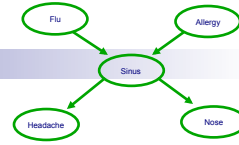
■ Given structure, log likelihood of data:
 $\log P(\mathcal{D} | \theta_{\mathcal{G}}, \mathcal{G})$



©Carlos Guestrin 2005-2013

30

Information-theoretic interpretation of maximum likelihood 2



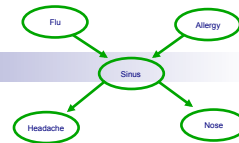
- Given structure, log likelihood of data:

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

©Carlos Guestrin 2005-2013

31

Information-theoretic interpretation of maximum likelihood 3



- Given structure, log likelihood of data:

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

©Carlos Guestrin 2005-2013

32

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

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

©Carlos Guestrin 2005-2013

37

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 2005-2013

38

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 2005-2013

39

Learn BN structure using local search

Starting from
Chow-Liu tree

Local search,
possible moves:

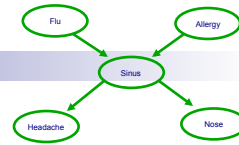
- Add edge
- Delete edge
- Invert edge

Score using BIC

©Carlos Guestrin 2005-2013

40

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!

©Carlos Guestrin 2005-2013

41

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

©Carlos Guestrin 2005-2013

42