

Bayesian Networks – Representation

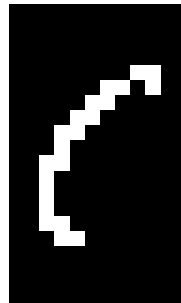
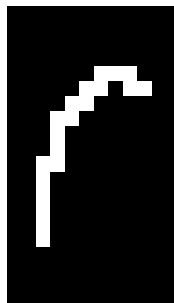
Machine Learning – CSE546
Carlos Guestrin
University of Washington

November 25, 2013

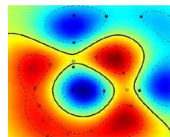
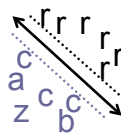
©Carlos Guestrin 2005-2013

1

Handwriting recognition



Character recognition, e.g., kernel SVMs



©Carlos Guestrin 2005-2013

2

Webpage classification



Company home page

VS

Personal home page

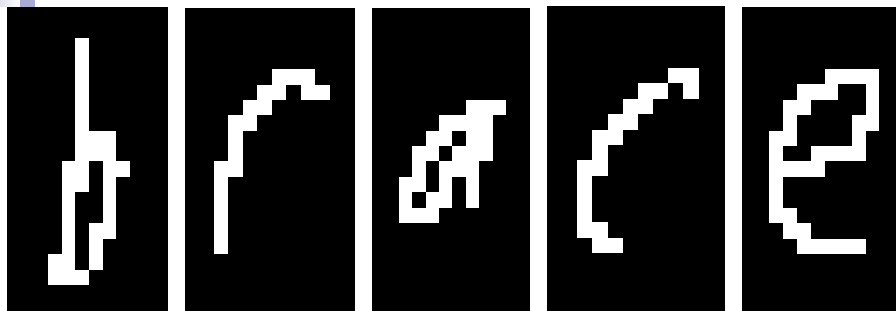
VS

University home page

VS

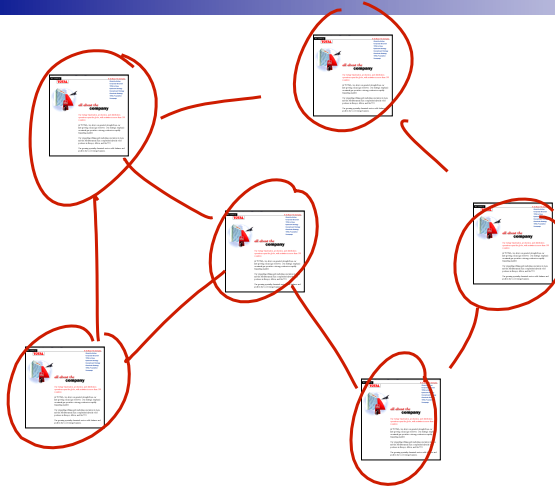
...

Handwriting recognition 2



○ — ○ — ○ — ○ — ○
"c" comes after an "a" much more
often than after a "e"

Webpage classification 2



©Carlos Guestrin 2005-2013

5

Today – Bayesian networks

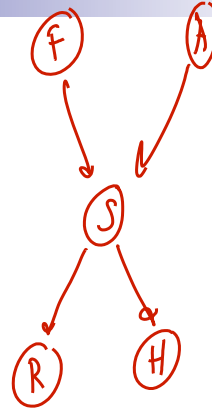
- One of the most exciting advancements in statistical AI in the last decades
- Generalizes naïve Bayes and logistic regression classifiers
- Compact representation for exponentially-large probability distributions
- Exploit conditional independencies

©Carlos Guestrin 2005-2013

6

Causal structure

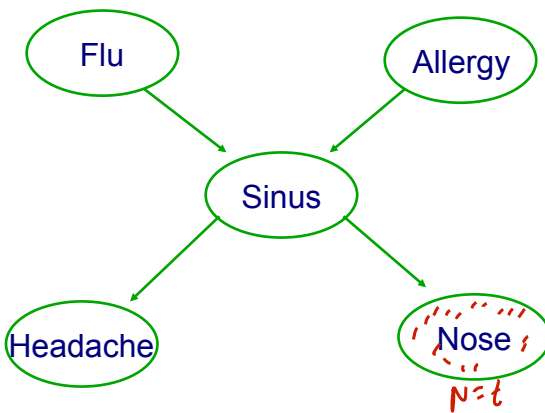
- Suppose we know the following:
 - The flu causes sinus inflammation
 - Allergies cause sinus inflammation
 - Sinus inflammation causes a runny nose
 - Sinus inflammation causes headaches
- How are these connected?



©Carlos Guestrin 2005-2013

7

Possible queries

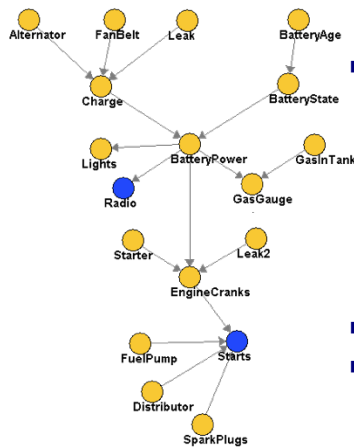


- Inference
 $P(F=t | N=t)$
- Most probable explanation
 $\max_{f,a,s,h} P(f,a,s,h | N=t)$
- Active data collection
given $N=t$, which test should I perform

©Carlos Guestrin 2005-2013

8

Car starts BN



- 18 binary attributes

2¹⁸ possibilities

- Inference

- P(BatteryAge|Starts=f)

$$P(BA | Starts=f) = \sum_{a,f,l,c,b,s,\dots} P(a,f,l,c,\dots, BA | Starts=f)$$

- 2¹⁶ terms, why so fast?

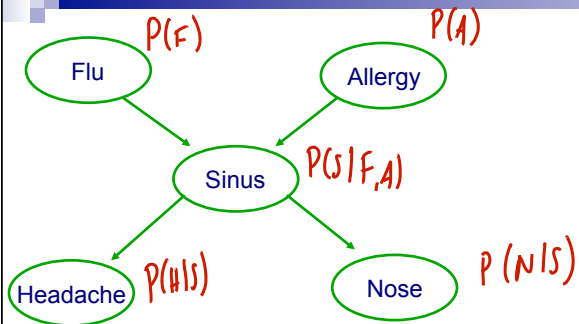
- Not impressed?

- HailFinder BN – more than 3⁵⁴ = 58149737003040059690390169 terms

©Carlos Guestrin 2005-2013

9

Factored joint distribution - Preview



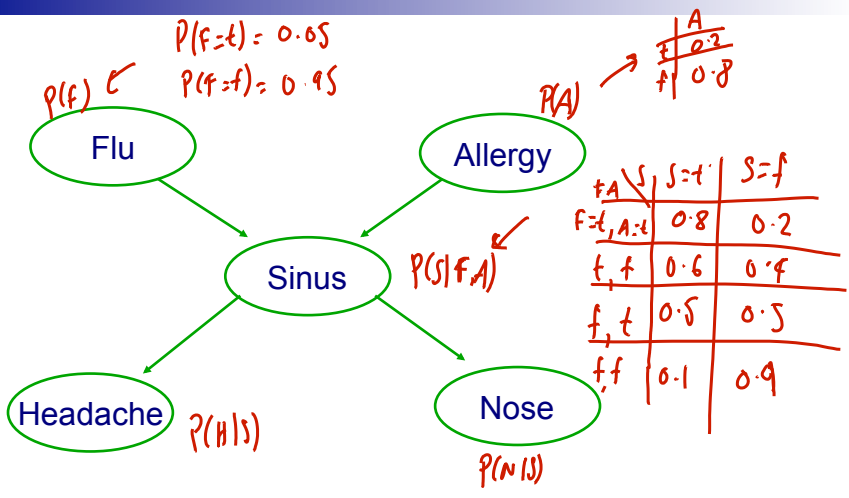
$$P(A, F, S, H, N) = P(F) P(A) P(S|F, A) P(H|S) P(N|S)$$

2⁵ = 32 terms *many factors*

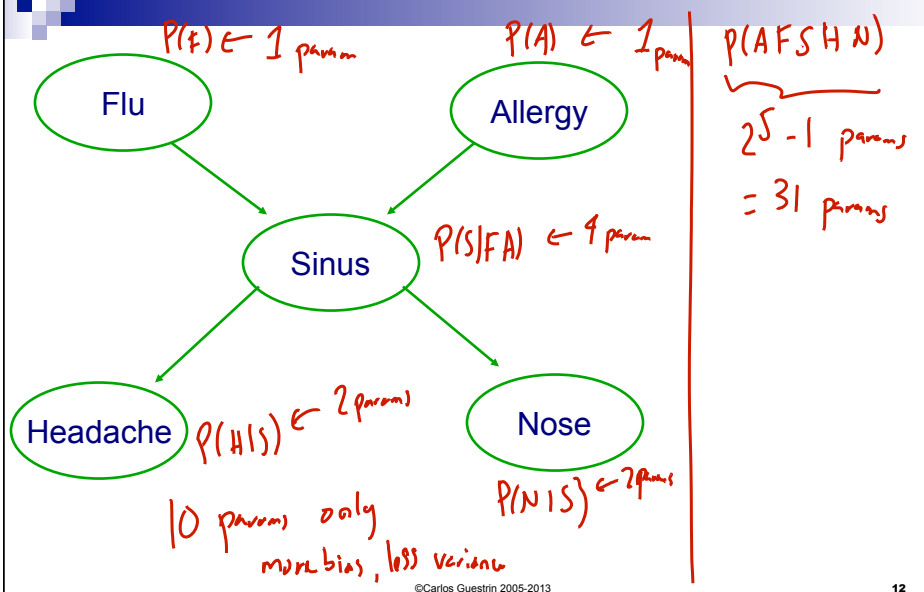
©Carlos Guestrin 2005-2013

10

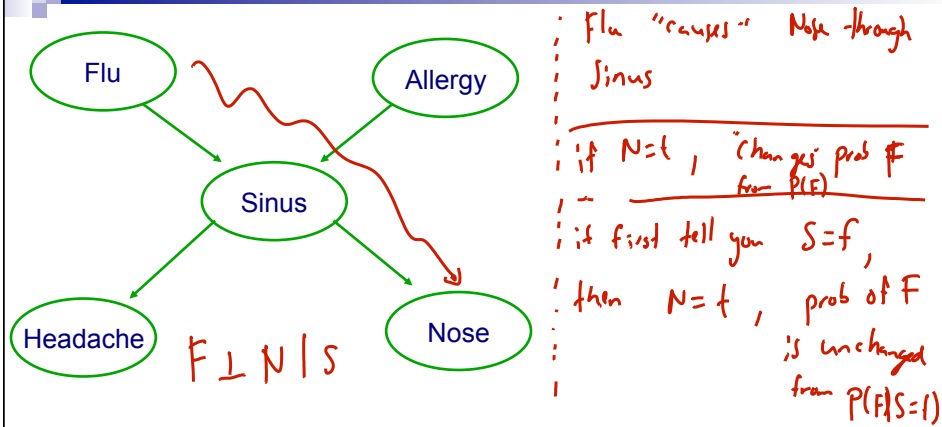
What about probabilities? Conditional probability tables (CPTs)



Number of parameters



Key: Independence assumptions



Knowing sinus separates the variables from each other

©Carlos Guestrin 2005-2013

13

(Marginal) Independence

- Flu and Allergy are (marginally) independent

$$F \perp A$$

$$\downarrow$$

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

OR

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

| | |
|---------|----|
| Flu = t | .2 |
| Flu = f | .8 |

| | |
|-------------|----|
| Allergy = t | .4 |
| Allergy = f | .6 |

| | Flu = t | Flu = f |
|-------------|-----------------------|----------------|
| Allergy = t | $.4 \times .2 = 0.08$ | $.4 \times .8$ |
| Allergy = f | $.6 \times .2$ | $.8 \times .6$ |

©Carlos Guestrin 2005-2013

14

Marginally independent random variables

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

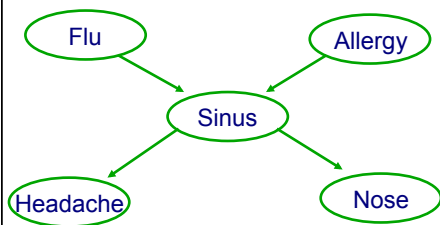
Conditional independence

- Flu and Headache are not (marginally) independent
 $P(H=t | F=t) \neq P(H=t) \quad F \not\perp H$
- Flu and Headache are independent given Sinus infection
 $F \perp H | S$
 $P(H=t | S=t) = P(H=t | S=t, F=t)$
- More Generally: $X \perp Y | Z$
 $P(X|Z) = P(X|Y, Z)$
 \Downarrow
 $P(X, Y|Z) = P(X|Z) P(Y|Z)$

Conditionally independent random variables

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

The independence assumption

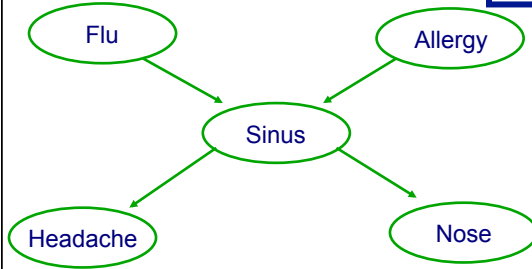


Local Markov Assumption:
 A variable X is independent of its non-descendants given its parents *(and only its parents)*

| | A | F | S | H | N |
|-----------------|-------------|-------------|---|------------------------------|------------------------------|
| non-descendants | F | A | FA | FAN | FAH |
| implies | $A \perp F$ | $F \perp A$ | $S \perp \{F, A\} \mid \emptyset$ <i>(nothing)</i> | $H \perp \{F, A, N\} \mid S$ | $N \perp \{F, A, H\} \mid S$ |

Explaining away

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



$F \perp A$
 $F \perp A | S$???
 No!

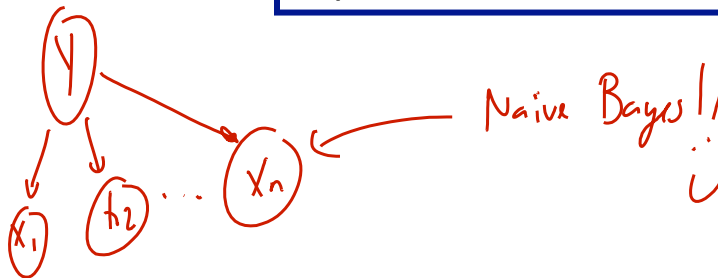
$P(F=t | S=t, A=f) \neq P(F=t | S=t)$
 No
 if it's not allergy, prob it's flu

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

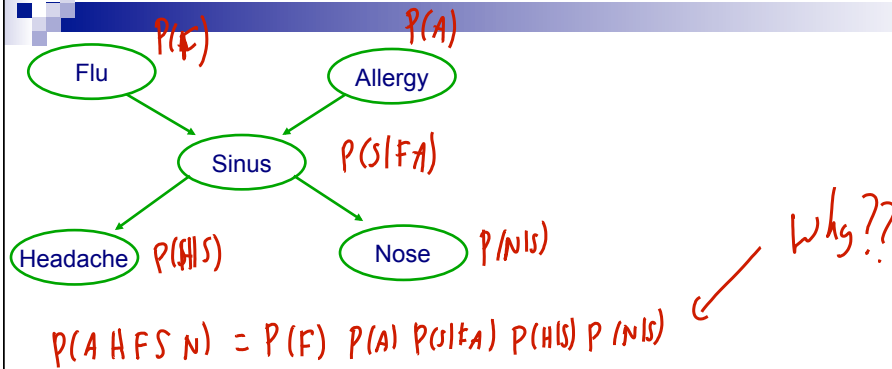
Naïve Bayes revisited

Variables: Y, X_1, \dots, X_n
 $X_1 \perp \{X_2, \dots, X_n\} | Y$
 \vdots

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

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

For any distribution:

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

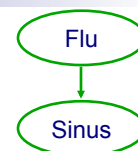
c.g. $P(F,S,H) = P(F) P(S|F) P(H|SF)$
 any distribution

- More generally:

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

always true

no indep assumption



Same as what Bayes Net says



Chain rule & Joint distribution

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

could have picked another ordering for chain rule, but for proof to work must pick any topological ordering

$$P(F A S H N) = P(F) \underbrace{P(A|F)}_{P(A)} \underbrace{P(S|F,A)}_{P(S)} \underbrace{P(H|SFA)}_{P(H|S)} \underbrace{P(N|HSFA)}_{P(N|S)}$$

| | | |
|---------------------------------------|--|--|
| $A \perp F \Rightarrow P(A F) = P(A)$ | $H \perp \{F, A\} S \Rightarrow P(H SFA) = P(H S)$ | $N \perp \{F, A, H\} S \Rightarrow P(N HSFA) = P(N S)$ |
|---------------------------------------|--|--|

©Carlos Guestrin 2005-2013 23

The Representation Theorem – Joint Distribution to BN

BN: **Encodes independence assumptions**

there is also a converse theorem

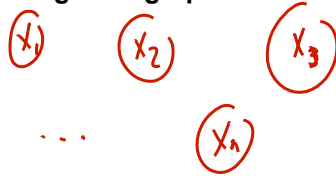
If conditional independencies in BN are subset of conditional independencies in P **Obtain** **Joint probability distribution:**

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

©Carlos Guestrin 2005-2013 24

Two (trivial) special cases

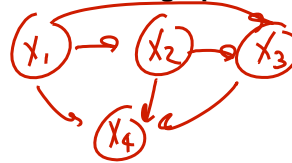
Edgeless graph



$x_i \perp \{ \text{everybody} \}$
else

all variables indep.
all the bins in the world

Fully-connected graph



$x_i \perp \{ \text{nobody} \} \mid \text{parents}$

no independence assumptions
 \Rightarrow represent anything
but exponentially many parameters
 \Rightarrow lot of variance (no bias)