

Bayes optimal classifier

Naïve Bayes

Machine Learning – CSE446

Carlos Guestrin

University of Washington

November 18, 2014

©Carlos Guestrin 2005-2014

1

Classification

■ Learn: $h: \mathbf{X} \mapsto Y$ ← *GPA, SF6 grade... hired, not hired*

- \mathbf{X} – features
- Y – target classes

■ Suppose you know $P(Y|\mathbf{X})$ exactly, how should you classify?

- Bayes optimal classifier:

$$\hat{y} = \underset{y}{\operatorname{argmax}} P(Y=y | X=x)$$

Logistic Regression $P(Y|x) = \frac{1}{1 + e^{-w \cdot x}}$

©Carlos Guestrin 2005-2014

2

Bayes Rule

use for classification

likelihood function

prior

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

posterior distribution

normalizer (partition function)

Which is shorthand for:

$$(\forall i, j) P(Y = y_i | X = x_j) = \frac{P(X = x_j | Y = y_i) P(Y = y_i)}{P(X = x_j)}$$

How hard is it to learn the optimal classifier?

	x_1	x_2	x_3	x_4	x_5	x_6	Y
Sky	Temp	Humid	Wind	Water	Forecast	EnjoySpt	
Sunny	Warm	Normal	Strong	Warm	Same	Yes	Yes
Sunny	Warm	High	Strong	Warm	Same	Yes	Yes
Rainy	Cold	High	Strong	Warm	Change	No	No
Sunny	Warm	High	Strong	Cool	Change	Yes	Yes

■ Data =

■ How do we represent these? How many parameters?

□ Prior, $P(Y)$:

■ Suppose Y is composed of k classes

$k-1$ params, because probs add up to 1

□ Likelihood, $P(\mathbf{X}|Y)$:

■ Suppose \mathbf{X} is composed of d binary features

$P(x=x | Y=y) \leftarrow$ for each $Y=y$, distribution over X

a lot, a lot \Rightarrow need a lot of data

$K(2^d - 1)$ params $P(S=\rightarrow, T=v, W=s, W=v, F=S | Y=no)$

■ Complex model \Rightarrow High variance with limited data!!!

Conditional Independence

$X \perp Y$ independent

$P(X, Y) = P(X)P(Y)$

- X is conditionally independent of Y given Z, if the probability distribution governing X is independent of the value of Y, given the value of Z
 $(\forall i, j, k) P(X = i | Y = j, Z = k) = P(X = i | Z = k)$

- e.g., $P(\overset{= \text{true}}{\text{Thunder}} | \overset{= \text{true}}{\text{Rain}}, \text{Lightning}) = P(\text{Thunder} | \text{Lightning})$
 $R \perp T ? \text{ No! But } R \perp T | L$

$P(T, R | L) = P(T | L) P(R | L)$

- Equivalent to:

$$P(X, Y | Z) = P(X | Z)P(Y | Z)$$

What if features are independent?

- Predict Thunder $P(T | LR) = P(T | L)$
- From two **conditionally Independent** features
 - Lightning
 - Rain

estimate $P(T | LR)$ ← $2^2(2-1) = 4$ *parameters*
 without independence

with independence only need $P(T | L)$ ← $2^1(2-1) = 2$ *parameters to learn*

The Naïve Bayes assumption

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

- Naïve Bayes assumption:
 - Features are independent given class:

$$P(X_1, X_2|Y) = P(X_1|X_2, Y)P(X_2|Y) \\ = P(X_1|Y)P(X_2|Y)$$

- More generally:

$$P(X_1 \dots X_d | Y) = \prod_{i=1}^d P(X_i | Y)$$

- How many parameters now?

- Suppose X is composed of d binary features

For each $P(X_i|Y) \leftarrow K(2-1)$ params \Rightarrow total: Kd params nice reduction!!

but we introduced bias

without assumption needed $K(2^d-1)$ params

The Naïve Bayes Classifier

dropped $P(x)$..

$$\underset{y}{\operatorname{argmax}} \frac{P(y)P(x|y)}{P(x)}$$

$$= \underset{y}{\operatorname{argmax}} P(y)P(x|y)$$

- Given:
 - Prior $P(Y)$
 - d conditionally independent features X given the class Y
 - For each X_i , we have likelihood $P(X_i|Y)$

$$P(Y = \text{hired}) = 0.03$$

$$P(Y = \text{not hired}) = 0.97$$

- Decision rule:

$$\hat{y}_{NB}^* = h_{NB}(x) = \underset{y}{\operatorname{argmax}} P(y)P(x_1, \dots, x_d | y) \\ = \underset{y}{\operatorname{argmax}} P(y) \prod_{i=1}^d P(x_i | y)$$

- If assumption holds, NB is optimal classifier!

MLE for the parameters of NB

- Given dataset
 - Count(A=a,B=b) == number of examples where A=a and B=b

- MLE for NB, simply:

- Prior: $P(Y=y) = \frac{\text{count}(Y=y)}{N}$ (MLE)

$P(Y=\text{hired})$

- Likelihood: $P(X_i=x_i|Y=y) = \frac{\text{count}(X_i=x_i, Y=y)}{\text{count}(Y=y)}$ (MLE)

$P(G=\text{high} | Y=\text{hired})$

$= \frac{\text{count}(G=\text{high}, Y=\text{hired})}{\text{count}(Y=\text{hired})}$ (MLE)

©Carlos Guestrin 2005-2014

9

Subtleties of NB classifier 1 – Violating the NB assumption

- ~~Usually~~ ^{always}, features are not conditionally independent:

$$P(X_1 \dots X_d | Y) \neq \prod_i P(X_i | Y)$$

- Actual probabilities $P(Y|\mathbf{X})$ often biased towards 0 or 1
- Nonetheless, NB is the single most used classifier out there
 - NB often performs well, even when assumption is violated
 - [Domingos & Pazzani '96] discuss some conditions for good performance

©Carlos Guestrin 2005-2014

10

Subtleties of NB classifier 2 – Insufficient training data

- What if you never see a training instance where $X_1=a$ when $Y=b$?

- e.g., $Y=\{\text{SpamEmail}\}$, $X_1=\{\text{'Enlargement'}\}$

- $P(X_1=a | Y=b) = 0$

- Thus, no matter what the values X_2, \dots, X_d take:

- $P(Y=b | X_1=a, X_2, \dots, X_d) = 0$

$P(y|x) \propto P(y) \prod_i P(x_i|y)$
 always zero if $x_i=a, y=b$, if use MLE

- “Solution”: smoothing/regularization
- Add “fake” counts, usually uniformly distributed
 - Equivalent to Bayesian Learning

Smooth Count ($X_i=x_i, Y=y$)
 $= \text{Count}(X_i=x_i, Y=y) + \alpha \text{Uniform}(X_i, Y)$
 $\alpha = \frac{1}{K R_i}$
 constant

©Carlos Guestrin 2005-2014

11

Text classification

- Classify e-mails
 - $Y = \{\text{Spam, NotSpam}\}$
- Classify news articles
 - $Y = \{\text{what is the topic of the article?}\}$
- Classify webpages
 - $Y = \{\text{Student, professor, project, ...}\}$
- What about the features \mathbf{X} ?
 - The text!

©Carlos Guestrin 2005-2014

12

Features X are entire document – X_i for i^{th} word in article

Article from rec.sport.hockey ← Y

Path: cantaloupe.srv.cs.cmu.edu!das-news.harvard.e
 From: xxx@yyy.zzz.edu (John Doe)
 Subject: Re: This year's biggest and worst (opinion)
 Date: 5 Apr 93 09:53:39 GMT

X_1 X_2 ...

I can only comment on the Kings, but the most obvious candidate for pleasant surprise is Alex Zhitnik. He came highly touted as a defensive defenseman, but he's clearly much more than that. Great skater and hard shot (though wish he were more accurate). In fact, he pretty much allowed the Kings to trade away that huge defensive liability Paul Coffey. Kelly Hrudey is only the biggest disappointment if you thought he was any good to begin with. But, at best, he's only a mediocre goaltender. A better choice would be Tomas Sandstrom, though not through any fault of his own, but because some things were decided Toronto decided

$P(Y|X)$

X

13

NB for Text classification

- $P(X|Y)$ is huge!!!
 - Article at least ^{10,000} 10,000 words, $X = \{X_1, \dots, X_{1000}\}$
 - X_i represents i^{th} word in document, i.e., the domain of X_i is entire vocabulary, e.g., Webster Dictionary (or more), 10,000 words, etc.
- NB assumption helps a lot!!!
 - $P(X_i=x_i|Y=y)$ is just the probability of observing word x_i in a document on topic y

$$\hat{y}_{NB} = \underset{y}{\arg \max} P(y) \prod_{i=1}^{LengthDoc} P(x_i|y)$$

what's prob of itl given topic
 $P(x_i = "hockey" | y: hockey)$

©Carlos Guestrin 2005-2014 14

Bag of words model

- Typical additional assumption – **Position in document doesn't matter**: $P(X_i=x_i|Y=y) = P(X_k=x_i|Y=y)$
 - “Bag of words” model – order of words on the page ignored
 - Sounds really silly, but often works very well!

$$P(y) \prod_{i=1}^{LengthDoc} P(x_i|y)$$

When the lecture is over, remember to wake up the person sitting next to you in the lecture room.

Bag of words model

- Typical additional assumption – **Position in document doesn't matter**: $P(X_i=x_i|Y=y) = P(X_k=x_i|Y=y)$
 - “Bag of words” model – order of words on the page ignored
 - Sounds really silly, but often works very well!

$$P(y) \prod_{i=1}^{LengthDoc} P(x_i|y)$$

in is lecture lecture next over person remember room
sitting the the the to to up wake when you

$$\operatorname{argmax}_y \log \left[P(y) \prod_i P(x_i|y) \right] = \operatorname{argmax}_y \log P(y) + \sum_i \log P(x_i|y)$$

Bag of Words Approach

word count



can't #
times each
word appears

can also
use N-grams...

= counting
sequence of
words

aardvark	0
about	2
all	2
Africa	1
apple	0
anxious	0
gas	1
...	
oil	1
...	
Zaire	0

©Carlos Guestrin 2005-2014

17

NB with Bag of Words for text classification

Learning phase:

□ Prior $P(Y)$

- Count how many documents you have from each topic (+ prior) *Smoothing*

□ $P(X_i|Y)$

- For each topic, count how many times you saw word in documents of this topic (+ prior)

Test phase:

□ For each document

- Use naïve Bayes decision rule

$$h_{NB}(x) = \arg \max_y P(y) \prod_{i=1}^{LengthDoc} P(x_i|y)$$

©Carlos Guestrin 2005-2014

18

Twenty News Groups results

Given 1000 training documents from each group
 Learn to classify new documents according to
 which newsgroup it came from

comp.graphics misc.forsale
 comp.os.ms-windows.misc rec.autos
 comp.sys.ibm.pc.hardware rec.motorcycles
 comp.sys.mac.hardware rec.sport.baseball
 comp.windows.x rec.sport.hockey

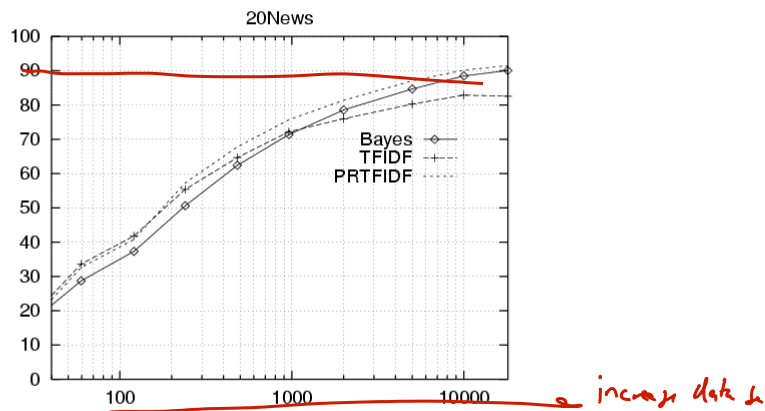
 alt.atheism sci.space
 soc.religion.christian sci.crypt
 talk.religion.misc sci.electronics
 talk.politics.mideast sci.med
 talk.politics.misc
 talk.politics.guns

Naive Bayes: 89% classification accuracy

©Carlos Guestrin 2005-2014

19

Learning curve for Twenty News Groups



Accuracy vs. Training set size (1/3 withheld for test)

©Carlos Guestrin 2005-2014

20

Bayesian Networks – Representation

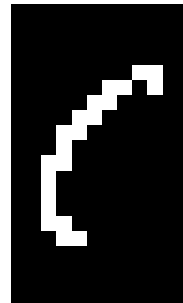
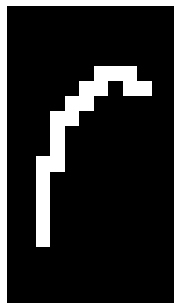
Machine Learning – CSE446
Carlos Guestrin
University of Washington

November 18, 2014

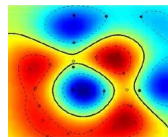
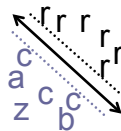
©Carlos Guestrin 2005-2014

21

Handwriting recognition



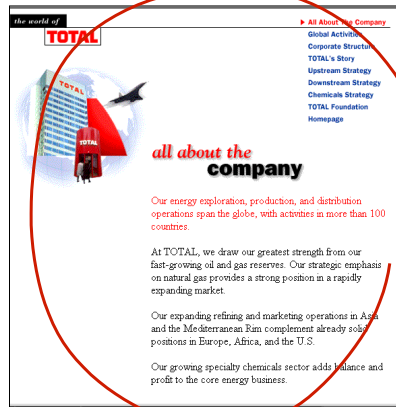
Character recognition, e.g., kernel SVMs



©Carlos Guestrin 2005-2014

22

Webpage classification



Company home page

VS

Personal home page

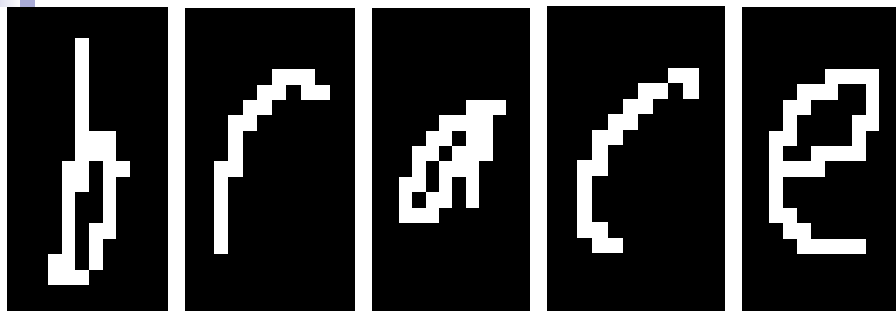
VS

University home page

VS

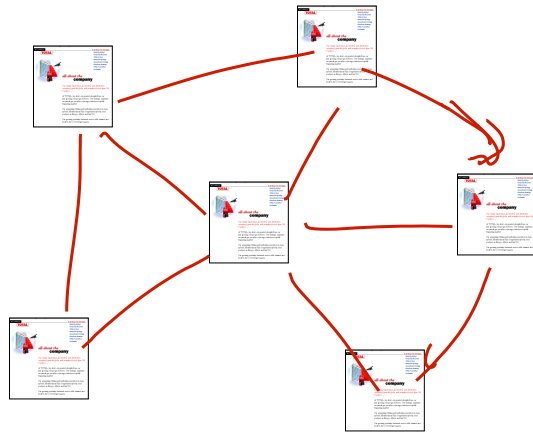
...

Handwriting recognition 2



"r" are more likely to come after a "b" than a "g"

Webpage classification 2



*Student's page
is likely to
point to prof page*

©Carlos Guestrin 2005-2014

25

Today – Bayesian networks

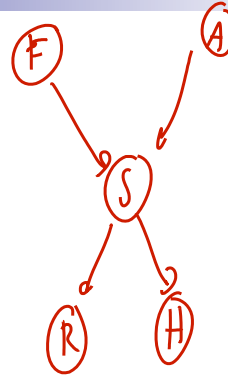
- One of the most exciting advancements in statistical AI in the last decades
- Generalizes naïve Bayes and logistic regression classifiers
- Compact representation for exponentially-large probability distributions
- Exploit conditional independencies

©Carlos Guestrin 2005-2014

26

Causal structure

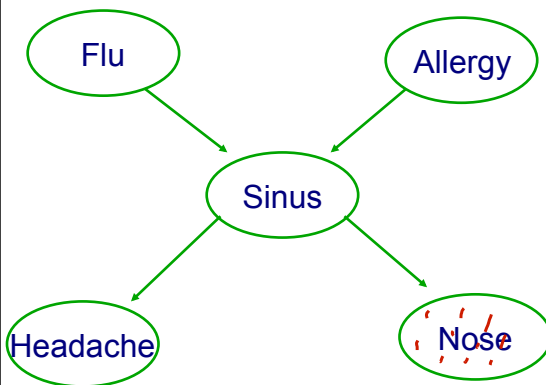
- Suppose we know the following:
 - The flu causes sinus inflammation
 - Allergies cause sinus inflammation
 - Sinus inflammation causes a runny nose
 - Sinus inflammation causes headaches
- How are these connected?



©Carlos Guestrin 2005-2014

27

Possible queries

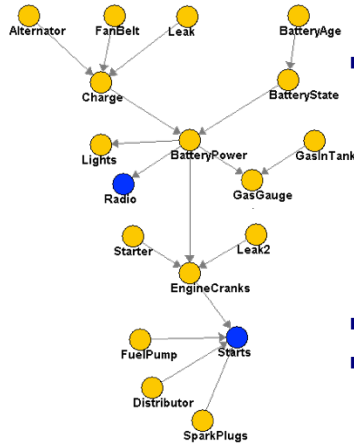


- Inference
 $P(F=t | N=t)$
- Most probable explanation
 $\max P(F,A | N=t)$
- Active data collection
what should I measure next?
 $H=t?$

©Carlos Guestrin 2005-2014

28

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

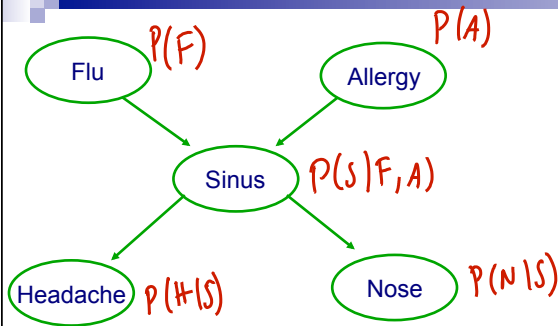
©Carlos Guestrin 2005-2014

29

exploit structure in graph to speed up computation

2^{18} possibilities

Factored joint distribution - Preview



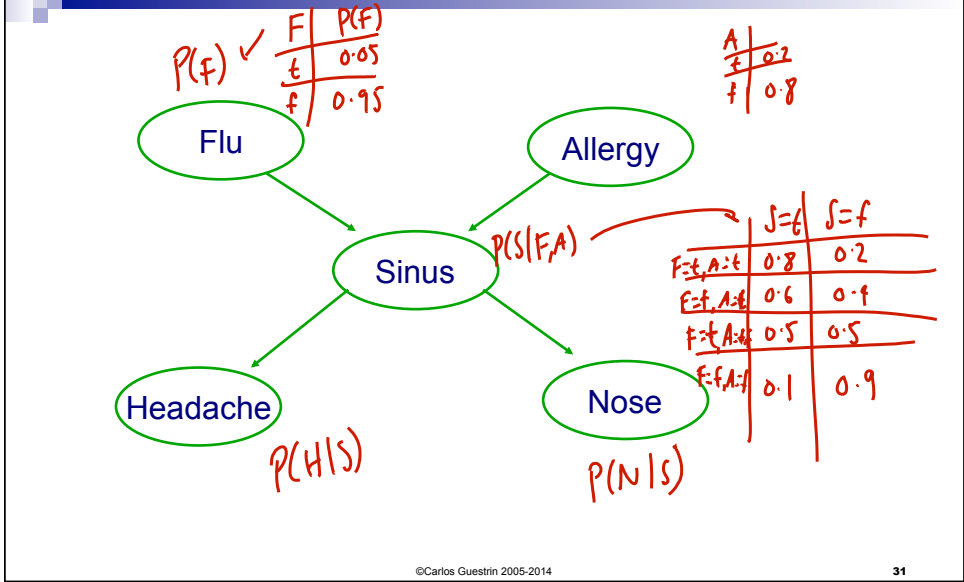
$$P(F, A, S, H, N) = P(F) P(A) P(S|F, A) P(H|S) P(N|S)$$

$2^5 - 1 = 31$ parameters

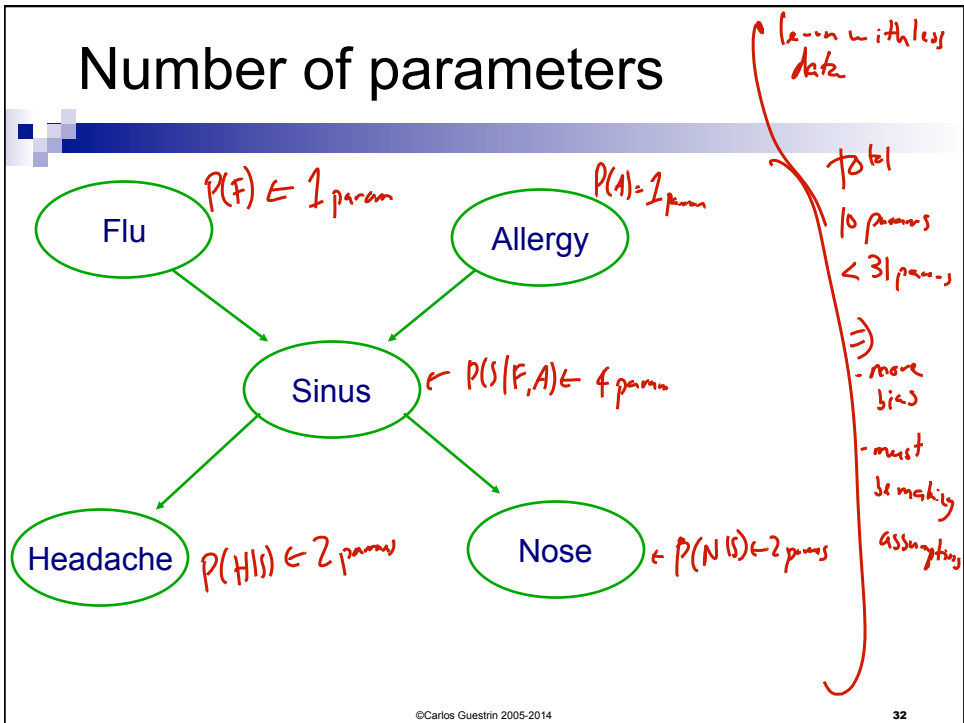
©Carlos Guestrin 2005-2014

30

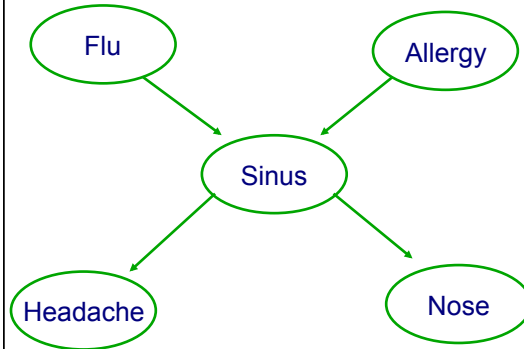
What about probabilities? Conditional probability tables (CPTs)



Number of parameters



Key: Independence assumptions



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

Knowing sinus separates the variables from each other

©Carlos Guestrin 2005-2014

33

(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 2005-2014

34

Marginally independent random variables

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

©Carlos Guestrin 2005-2014

35

Conditional independence

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

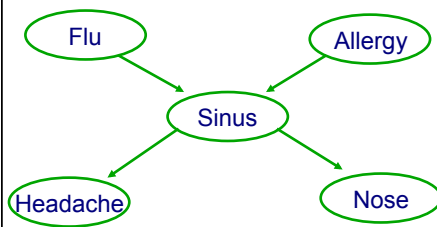
©Carlos Guestrin 2005-2014

36

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
 - $P(X, Y | Z) = P(X | Z) P(Y | Z)$

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
<i>non descends</i>	A	F	FA	FAU(S)	FAH
<i>implies</i>	F ⊥ A	A ⊥ F	S ⊥ FA / FA ⇒ nothing	H ⊥ {F, A, N} S	N ⊥ {F, A, H} S