# Bayesian Networks – Representation

Machine Learning – CSE446

Carlos Guestrin

University of Washington
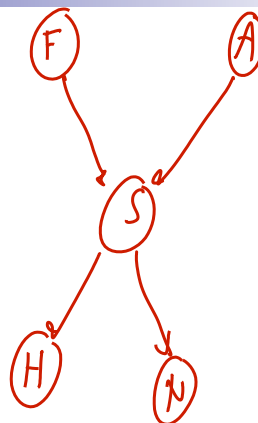
May 29, 2013
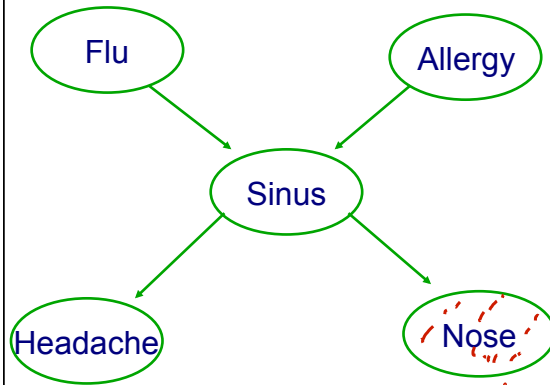
1

---

# 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?

2

---

1

# Possible queries

Flu · Allergy · Sinus · Headache · Nose

- **Inference**

$$P(F=t \mid N=t)$$

- **Most probable explanation**

$$\max_{f,a,s,h} P(f,s,a,h \mid N=t)$$

- **Active data collection**

$N=t$

What variable should I observe next?

$H=?, \quad S=?$

3

---

# Car starts BN

Alternator · FanBelt · Leak · BatteryAge · Charge · BatteryState · Lights · BatteryPower · GasInTank · Radio · GasGauge · Starter · Leak2 · EngineCranks · FuelPump · Starts · Distributor · SparkPlugs

- **18 binary attributes**

$2^{18}$ probabilities

- **Inference**
  - P(BatteryAge|Starts=f)

$$P(BA \mid S=f) = \sum_{a,f,b,\ \ell,\ c,\ \dots} P(a, fb, \ell, \dots, R=f, \dots, S=f)$$

$\rightarrow 2^{16}$

16

- $2^{16}$ terms, why so fast?
- **Not impressed?**
  - HailFinder BN – more than $3^{54} =$ 581497370030400596990390169 terms

4

2

# Factored joint distribution - Preview

Flu — $P(F)$

Allergy — $P(A)$

Sinus — $P(S|F,A)$

Headache — $P(H|S)$

Nose — $P(N|S)$

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

$2^5 = 32$ terms

31 params

5

---

# What about probabilities?
# Conditional probability tables (CPTs)

$P(F) \leftarrow$
$P(F=t) = 0.05$
$P(F=f) = 0.95$

$P(A) \rightarrow$

| | |
|---|---|
| t | 0.2 |
| f | 0.8 |

Flu

Allergy

Sinus

$P(S|FA) =$

| $P(S|FA)$ | $S=t$ | $S=f$ |
|---|---|---|
| $F=f, A=f$ | 0.1 | 0.9 |
| $F=t, A=f$ | 0.7 | 0.3 |
| $F=f, A=t$ | 0.4 | 0.6 |
| $F=t, A=t$ | 0.8 | 0.2 |

Headache — $P(H|S)$

Nose — $P(N|S)$

6

3

# Number of parameters



Flu — $P(F) \leftarrow 1$ param

Allergy — $P(A) \leftarrow 1$ param

Sinus — $4 \cdot (2-1)$ params A.K.A. 4 — $P(S|FA)$

Headache — $P(H|S) \leftarrow 2$ params

Nose — $P(N|S)$

Total: 10 params

$P(FASHN)$
$\to 31$ params
$10 < 31$
$\Rightarrow$
- more bias
- less flexible
- need less data to learn
- more accurate on smaller datasets

# Key: Independence assumptions



$F \perp N | S$

Flu only "causes" Nose through Sinus

if you tell N=t changes prob of F, but if I first tell you S=t, N doesn't affect prob. of Flu

**Knowing sinus separates the variables from each other**

# (Marginal) Independence

- Flu and Allergy are (marginally) independent

$$F \perp A$$

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

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

| $P(F)$ | | |
|---|---|---|
| Flu = t | .2 |
| Flu = f | .8 |

| $P(A)$ | | |
|---|---|---|
| Allergy = t | .4 |
| Allergy = f | .6 |

| $P(F, A)$ | Flu = t | Flu = f |
|---|---|---|
| Allergy = t | .4 x .2 = .008 | .4 x .8 |
| Allergy = f | .6 x .2 | .8 x .6 |

©Carlos Guestrin 2005-2013                    9

---

# Marginally independent random variables

- **Sets** of variables **X**, **Y** ⟶ entails
- X is independent of Y if
  - $P \vdash (X=x \perp Y=y)$, $\forall x \in Val(X)$, $y \in Val(Y)$

distribution $P(X=x, Y=y) = P(X=x) P(Y=y)$ $\forall x, \forall y$

- Shorthand:
  - **Marginal independence:** $P \vdash (X \perp Y)$

- **Proposition:** $P$ statisfies $(X \perp Y)$ if and only if
  - $P(X, Y) = P(X) P(Y)$
  
$P(X|Y) = P(X)$  } equivalent

©Carlos Guestrin 2005-2013                    10

5

# Conditional independence

- Flu and Headache are not (marginally) independent

$$P(H=t \mid F=t) \neq P(H=t)$$

- Flu and Headache are independent given Sinus infection

$$P(H=t \mid S=t) = P(H=t \mid S=t, F=t)$$

- More Generally:

$$X \perp Y \mid Z$$

$$P(X, Y \mid Z) = P(X \mid Z)\, P(Y \mid Z) \Big\} \text{ equivalent}$$

$$P(X \mid Y, Z) = P(X \mid Z)$$

11

---

# Conditionally independent random variables

- **Sets** of variables **X**, **Y**, **Z**
- X is independent of Y given Z if
  - $P \vdash (\mathbf{X}=\mathbf{x} \perp \mathbf{Y}=\mathbf{y} \mid \mathbf{Z}=\mathbf{z})$, $\forall \mathbf{x} \in \mathrm{Val}(\mathbf{X})$, $\mathbf{y} \in \mathrm{Val}(\mathbf{Y})$, $\mathbf{z} \in \mathrm{Val}(\mathbf{Z})$

$$P(X=x \mid Y=y, Z=z) = P(X=x \mid Z=z)$$

- Shorthand:
  - **Conditional independence:** $P \vdash (\mathbf{X} \perp \mathbf{Y} \mid \mathbf{Z})$
  - For $P \vdash (\mathbf{X} \perp \mathbf{Y} \mid \varnothing)$, write $P \vdash (\mathbf{X} \perp \mathbf{Y})$

- **Proposition:** $P$ statisfies $(\mathbf{X} \perp \mathbf{Y} \mid \mathbf{Z})$ if and only if
  - $P(\mathbf{X}, \mathbf{Y} \mid \mathbf{Z}) = P(\mathbf{X} \mid \mathbf{Z})\, P(\mathbf{Y} \mid \mathbf{Z})$

$$P(X \mid Y, Z) = P(X \mid Z) \Big\} \text{ equivalent}$$

12

# **The** independence assumption

Flu → Sinus ← Allergy
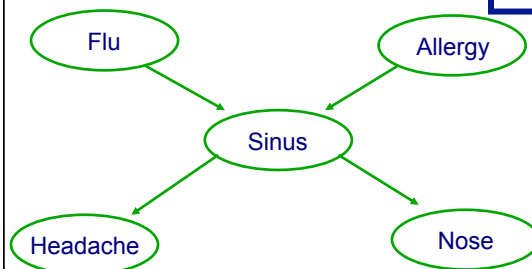Sinus → Headache
Sinus → Nose

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

| | F | A | S | H | N |
|---|---|---|---|---|---|
| non-desc. | A | F | FA | FAN | FAH |
| implies | F⊥A | A⊥F | S⊥FA\|FA ⇒ nothing | H⊥{F,A,N}\|S | N⊥{F,A,H}\|S |

©Carlos Guestrin 2005-2013    13

---

# Explaining away

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

Flu → Sinus ← Allergy
Sinus → Headache
Sinus → Nose

$F \perp A$

$F \perp A | S$ ??
———— don't know

$P(F=f | A=f, S=t) \neq P(F=t|S=t)$

No!!

Suppose $P(F=t|S=t)$ is high
but $P(F=t|S=t, A=t)$ is lower
because A=t explains away sinus infection
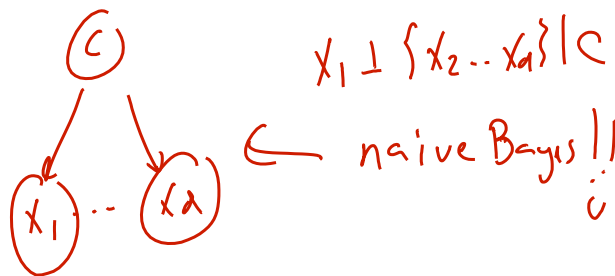
©Carlos Guestrin 2005-2013    14

# Naïve Bayes revisited

$$X_1 \perp \{X_2 \dots X_d\} \mid C$$
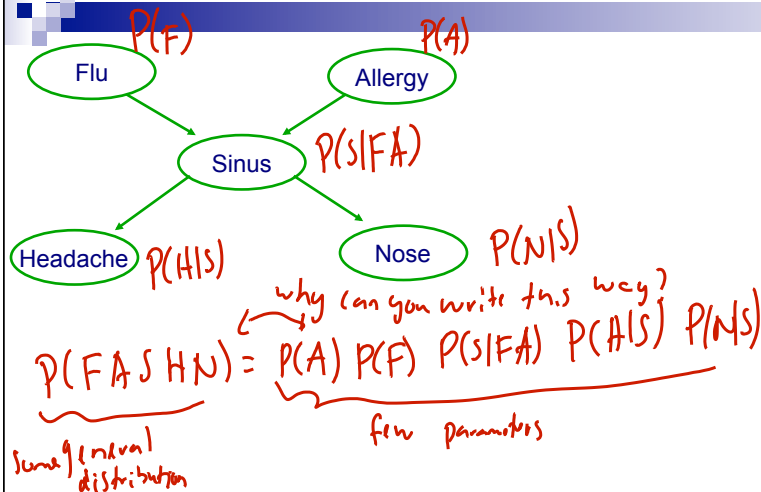
$$P(C, X_1, \dots X_d) = P(C) \prod P(X_i \mid C)$$

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

$$X_1 \perp \{X_2 \dots X_d\} \mid C$$

$\leftarrow$ naive Bayes !!

C → X₁ ... Xₐ

---

# Joint distribution

$P(F)$  Flu

$P(A)$  Allergy

$P(S\mid FA)$  Sinus

$P(H\mid S)$  Headache

$P(N\mid S)$  Nose

why can you write this way?

$$P(FASHN) = P(A)\, P(F)\, P(S\mid FA)\, P(H\mid S)\, P(N\mid S)$$

few parameters

Some general distribution

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

# The chain rule of probabilities

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

Flu → Sinus

For any dist:

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

for any ordering ✓

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

- More generally: any ordering over variables
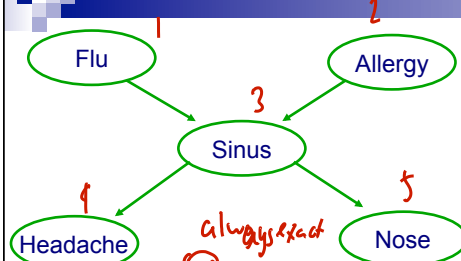  - $P(X_1,\ldots,X_n) = P(X_1) \, P(X_2|X_1) \, \ldots \, P(X_n|X_1,\ldots,X_{n-1})$

$P(X_3|X_2,X_1)$

17

---

# Chain rule & Joint distribution

1 Flu      2 Allergy

3 Sinus

4 Headache    always exact    5 Nose

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

proof by example: works for all BNs !!

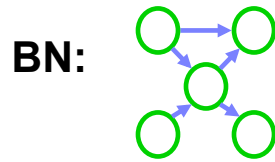$P(F,A,S,H,N) = P(F) \, P(A|F) \, P(S|F,A) \, P(H|SFA) \, P(N|SFAH)$

$P(A) \quad\quad P(H|S) \quad P(N|S)$

order of chain rule expansion matters a lot !! use topological order.

$A \perp F \Rightarrow P(A|F) = P(A)$

$H \perp \{FA\} | S$
$P(H|SFA) = P(H|S)$

$P(N|SFAH) = P(N|S)$
$N \perp \{FAH\} | S$

18

9

# The Representation Theorem – Joint Distribution to BN

**BN:**

Encodes independence assumptions *var indep of non-descend given its parents*

*there is also a converse*

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

Obtain

**Joint probability distribution:**

$$P(X_1, \ldots, X_n) = \prod_{i=1}^{n} P\left(X_i \mid \mathbf{Pa}_{X_i}\right)$$

*product over CPTs*

19

---

# Two (trivial) special cases

**Edgeless graph**

$X_1$  $X_2$  $X_3$

$X_4$  $X_5$

$X_i \perp \{everybody\} \mid \emptyset$

all vars independent

*structure learning*

**Fully-connected graph**

$X_1 \longrightarrow X_2 \longrightarrow X_3$

$X_4 \longrightarrow X_5$

$X_i \perp \{non\text{-}descendants\ and\ not\ parents\} \mid parents\ X_1 \ldots X_{i-1}$

$\emptyset$

no vars independents

20

# Bayesian Networks – (Structure) Learning

Machine Learning – CSE446

Carlos Guestrin

University of Washington

May 31, 2013

©Carlos Guestrin 2005-2013

21

---

# 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

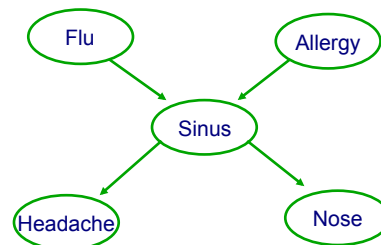Flu     Allergy

Sinus

Headache     Nose

©Carlos Guestrin 2005-2013

22

# Learning Bayes nets



Data

$\mathbf{x}^{(1)}$
…
$\mathbf{x}^{(m)}$

learn

CPTs –
$P(X_i | \mathbf{Pa}_{Xi})$

**structure**          **parameters**

maximizing       likelihood of data

$MLE$          $P(D | \theta_G, G)$

23

---

# Learning the CPTs

Data

$\mathbf{x}^{(1)}$
…
$\mathbf{x}^{(m)}$

For each discrete variable $X_i$

$P(S=t | A=t) \overset{MLE}{=} \dfrac{Count(S=t, A=t)}{Count(A=t)}$

given
we
know
parents

$P(X_i = x_i | Pa_{Xi} = u) \overset{mle}{=} \dfrac{Count(X=x_i, Pa_{Xi} = u)}{Count(u)}$

Small subtlety: $Count(u) = 0$, or small ...
need Smoothing / AKA regularization / AKA Bayesian learning

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

24

12

# Information-theoretic interpretation of maximum likelihood 1

*(handwritten: m data points, n variables)*

- Given structure, log likelihood of data:

*(handwritten: maximized over G)*

$$\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G}) \overset{iid}{=} \log \prod_{j=1}^{m} P(x_1^{(j)} \cdots x_n^{(j)} \mid \theta_G, G)$$

*(handwritten annotations: params, structure)*

$$= \log \prod_{j=1}^{m} \prod_{i=1}^{n} P\left(x_i^{(j)} \mid Pa_{X_i, G} = u_i^{(j)}\right)$$

$$= \sum_{j=1}^{m} \sum_{i=1}^{n} \log P\left(x_i^{(j)} \mid Pa_{X_i, G} = u_i^{(j)}\right)$$

*(handwritten: max over G)*

*(handwritten: $x_i^{(j)} \leftarrow$ data points (F=t, A=f, S=t ....))*

*(handwritten: $X_{i,j}$ — variables F, A, S, H, N)*

©Carlos Guestrin 2005-2013

25

---

# Information-theoretic interpretation of maximum likelihood 2

- Given structure, log likelihood of data:

*(handwritten: max over G, flip $u_i^{(j)}$)*

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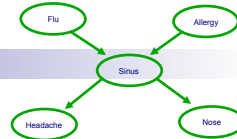
*(handwritten derivation:)*

$$= \sum_{i=1}^{n} \sum_{j=1}^{m} \log P\left(x_i^{(j)} \mid Pa_{X_i, G} = u_i^{(j)}\right) \quad \text{e.g.} \quad \sum_{j=1}^{m} \log P(h^{(j)} \mid s^{(j)})$$

$$= \sum_{i=1}^{n} \sum_{x_i} \sum_{u_i} Count\left(X_i = x_i, Pa_{X_i, G} = u_i\right) \log P(x_i \mid Pa_{X_i, G} = u_i)$$

$$= m \sum_{i=1}^{n} \sum_{x_i} \sum_{u_i} P(x_i, Pa_{X_i, G} = u_i) \log P(x_i \mid Pa_{X_i, G} = u_i)$$

$$\underbrace{\phantom{= m \sum_{i=1}^{n} \sum_{x_i} \sum_{u_i} P(x_i, Pa_{X_i, G} = u_i) \log P(x_i \mid Pa_{X_i, G} = u_i)}}_{-H(X_i \mid Pa_{X_i, G})}$$

$$P(x_i, u_i) \overset{MLE}{=} \frac{Count(X_i = x_i, Pa_{X_i, G} = u_i)}{m}$$

*(handwritten right column:)*

$$= \Big( Count(H=t, S=t) \log P(H=t \mid S=t)$$
$$+ Count(H=t, S=f) \log P(H=t \mid S=f)$$
$$+ Count(H=f, S=t) \log P(H=f \mid S=t)$$
$$+ Count(H=f, S=f) \log P(H=f \mid S=f) \Big)$$

©Carlos Guestrin 2005-2013

26

13

# Information-theoretic interpretation of maximum likelihood 3



- Given structure, log likelihood of data:

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

$\underset{\mathcal{G}}{max}$  $\underset{\mathcal{G}}{mal}$

$$= -m \sum_i H(X_i \mid Pa_{X_i}, G) \equiv \underset{G}{min}\ m \sum_{i=1}^{n} H(X_i \mid Pa_{X_i}, G)$$

$$\equiv \underset{G}{max}\ m \sum_{i=1}^{n} I(X_i, Pa_{X_i}, G) - m \sum_{i=1}^{n} H(X_i)$$

information theoretic criteria over G
⇒ max mutual info between vars & parents

just a constant WRT G

also measures how dependent vars are

For DTs:

$H(A|B)$

$= -\sum_{a,b} P(a,b) \log P(a|b)$

if $X_i$ is highly "correlated" with parents

$H(X_i | Pa_{X_i}, G)$ is low

⇒ good G

Mutual information

$= I(A, B)$

$= H(A) - H(A|B)$

---

# Decomposable score

- Log data likelihood

Constant

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

Example of a

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

e.g. $= \sum_{i=1}^{n} I(X_i, Pa_{X_i}, G)$

but there are many others