**Classic HMM tutorial – see class website:**
*L. R. Rabiner, "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition," Proc. of the IEEE, Vol.77, No.2, pp.257--286, 1989.

# HMMs (cont.)

Machine Learning – 10701/15781
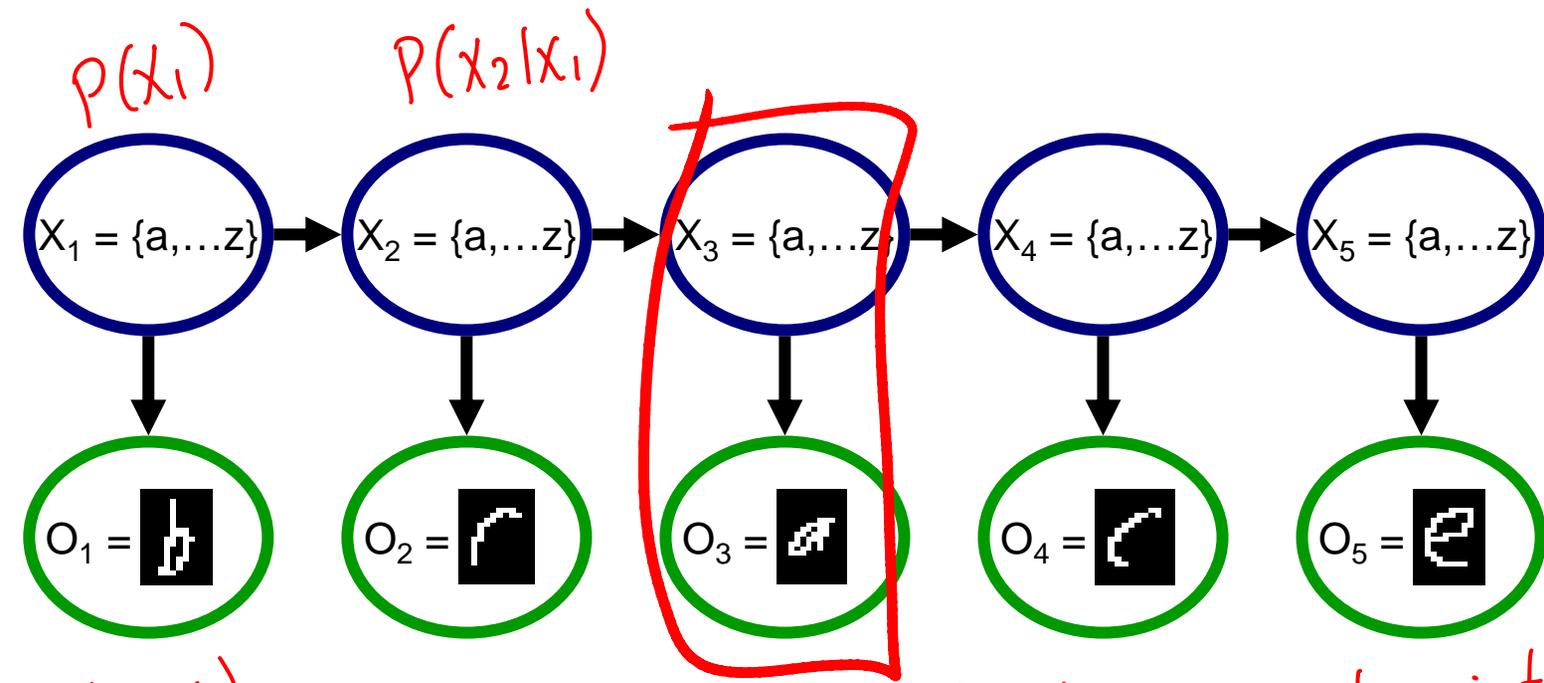
Carlos Guestrin

Carnegie Mellon University

March 29th, 2006

1

# Announcements

- This weeks recitation:
  - Go through several BNs topics, representation, inference, learning, in the context of an example $\rightarrow$ very useful for homework

# Understanding the HMM Semantics

$P(X_1)$

$P(X_2|X_1)$

$X_1 = \{a,\ldots z\}$ → $X_2 = \{a,\ldots z\}$ → $X_3 = \{a,\ldots z\}$ → $X_4 = \{a,\ldots z\}$ → $X_5 = \{a,\ldots z\}$

$O_1 =$ ▮  $O_2 =$ ▮  $O_3 =$ ▮  $O_4 =$ ▮  $O_5 =$ ▮

$P(O_1|X_1)$

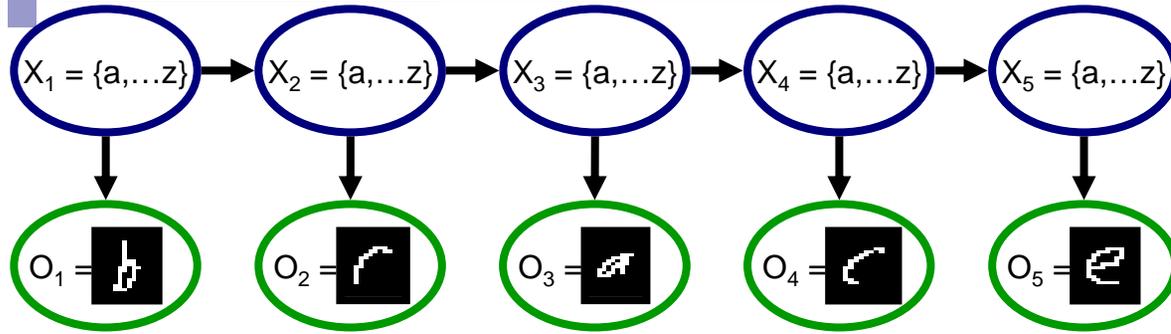one option is
Naive Bayes

$P(O_1|X_1) = \prod_{i \in pixels} P(O_i^i|X_1)$

Markov assumption : the future is
indep. of past given present
$X_{1:t-1} \perp X_{t+1:N} | X_t$

$O_t \perp everybody | X_t$

# HMMs semantics: Details

*basic!*



X₁ = {a,…z} → X₂ = {a,…z} → X₃ = {a,…z} → X₄ = {a,…z} → X₅ = {a,…z}

$X_1 = \{a,\ldots z\} \rightarrow X_2 = \{a,\ldots z\} \rightarrow X_3 = \{a,\ldots z\} \rightarrow X_4 = \{a,\ldots z\} \rightarrow X_5 = \{a,\ldots z\}$

$O_1 = \quad O_2 = \quad O_3 = \quad O_4 = \quad O_5 =$

**Just 3 distributions:**

$$P(X_1)$$

$$P(X_i \mid X_{i-1})$$

$$P(O_i \mid X_i)$$

*observing a distribution*

*a after b has same prob. no matter where in word*

$P(X_3 | X_2) \sim P(X_2 | X_1)$
$\sim P(X_{10} | X_9)$

*prob. image given letter is same for all positions in word.*

4

# Learning HMMs from fully observable data is easy

$X_1 = \{a,\dots z\}$ → $X_2 = \{a,\dots z\}$ → $X_3 = \{a,\dots z\}$ → $X_4 = \{a,\dots z\}$ → $X_5 = \{a,\dots z\}$

$O_1 =$    $O_2 =$    $O_3 =$    $O_4 =$    $O_5 =$

**Learn 3 distributions:**

$$P(X_1^{=a}) = \frac{\text{count}(\#\text{ first letter was }a)}{N = \text{data set size}}$$

*select training data where letter was a*

$$P(O_i^{=\text{pixel 17 is white}} \mid X_i^{=a}) = \frac{\text{count}(\text{pixel 17 was white}, X_i = a)}{\text{count}(X_i = a)}$$

*any position in word*

$$P(X_i^{=a} \mid X_{i-1}^{=b}) = \frac{\text{count}(a \text{ appears after } b)}{\text{count}(\# \text{ of } b\text{-s that are not at the end of the word})}$$

# Possible inference tasks in an HMM



**Marginal probability of a hidden variable:**

$$\rightarrow P(X_i \mid O_1 = \boxed{b}, O_2 \boxed{r}, O_3 = \boxed{a}, \ldots)$$

**Viterbi decoding – most likely trajectory for hidden vars:**

$$\underset{x_1 x_2 x_3 x_4 x_5}{\text{argmax}} \quad P(x_1, x_2, x_3, x_4, x_5 \mid O_{1:5})$$

# Using variable elimination to compute $P(X_i|o_{1:n})$



$X_1 = \{a,\ldots z\} \rightarrow X_2 = \{a,\ldots z\} \rightarrow X_3 = \{a,\ldots z\} \rightarrow X_4 = \{a,\ldots z\} \rightarrow X_5 = \{a,\ldots z\}$

$O_1 = \quad O_2 = \quad O_3 = \quad O_4 = \quad O_5 = $

**Compute:**

$$P(X_i \mid o_{1..n})$$

**Variable elimination order?**

$$1, 2, 3, 4 \quad \cdots \quad n$$

**Example:**

$$P(X_3 \mid O) = \sum_{x_1 x_2 x_4 x_5} P(x_1, x_2, X_3, x_4, x_5 \mid O)$$

$$\propto \sum_{x_1 x_2 x_4 x_5} P(x_1) P(O_1 \mid x_1) P(x_2 \mid x_1) P(O_2 \mid x_2) P(x_3 \mid x_2) \cdots$$

$$= \sum_{x_2 x_4 x_5} P(O_2 \mid x_2) P(x_3 \mid x_2) \cdots \underbrace{\sum_{x_1} P(x_1) P(O_1 \mid x_1) P(x_2 \mid x_1)}_{g_1(x_2, O_1)}$$

7

# What if I want to compute $P(X_i|o_{1:n})$ for each i?



X_1 = {a,...z} → X_2 = {a,...z} → X_3 = {a,...z} → X_4 = {a,...z} → X_5 = {a,...z}

$O_1 =$ □  $O_2 =$ □  $O_3 =$ □  $O_4 =$ □  $O_5 =$ □

**Compute:**

$$P(X_i \mid o_{1..n})$$

**Variable elimination for each i?**
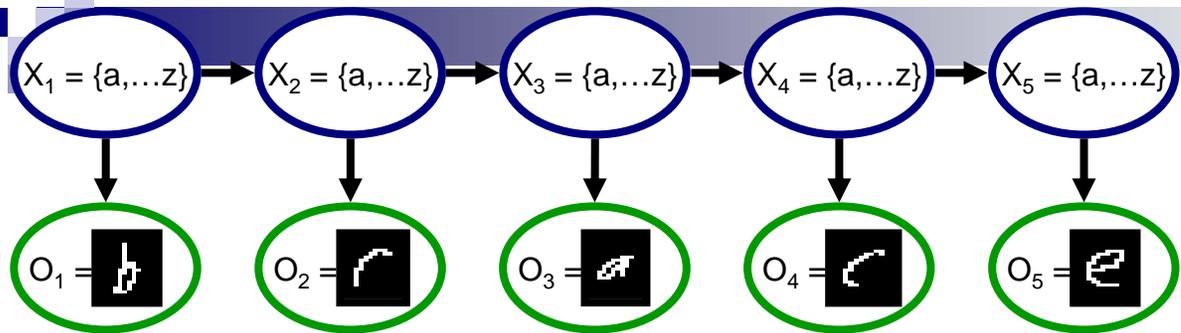
$P(X_1|0)$ , $P(X_2|0)$ , $P(X_3|0)$ ...

**Variable elimination for each i, what's the complexity?**

n letters    $P(X_i|0) \rightarrow O(n)$

$\rightarrow O(n^2)$    [can solve in $O(n)$

# Reusing computation



**Compute:**

$$P(X_i \mid o_{1..n})$$

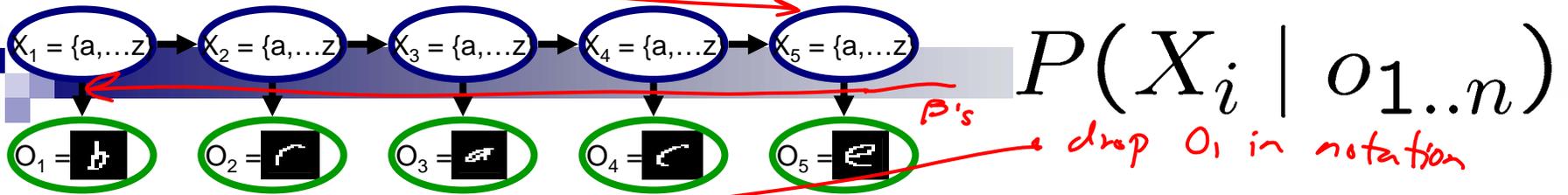$P(X_5 \mid O)$  $\rightarrow$  elimination order $1, 2, 3, 4$

$g_1(X_2, O_1)$ , $g_2(X_3, O_{1:2})$, $g_3(X_4, O_{1:3})$,

$g_4(X_5, O_{1:4})$

---

Same  g's  useful  for  $P(X_4 \mid O)$, order $1, 2, 3, 5$

eliminate 5  $\rightarrow$  use: $\sum_{X_5} P(O_5 \mid X_5) \cdot P(X_5 \mid X_4)$    computing $P(X_4 \mid O)$:

$\underbrace{\phantom{\sum_{X_5} P(O_5 \mid X_5) \cdot P(X_5 \mid X_4)}}_{g_5(X_4, O_5)}$

use $g_1, g_2, g_3, g_5$

# The forwards-backwards algorithm



$$P(X_i \mid o_{1..n})$$

*(handwritten annotations: α's, β's, drop $O_i$ in notation)*

- Initialization: $\alpha_1(X_1) = P(X_1)P(o_1 \mid X_1)$

- For i = 2 to n

  - ☐ Generate a forwards factor by eliminating $X_{i-1}$

  *(handwritten: sum out previous var prob obs, transition prob)*

  $$\alpha_i(X_i) = \sum_{x_{i-1}} P(o_i \mid X_i)P(X_i \mid X_{i-1} = x_{i-1})\alpha_{i-1}(x_{i-1})$$

- Initialization: $\beta_n(X_n) = 1$

- For i = n-1 to 1

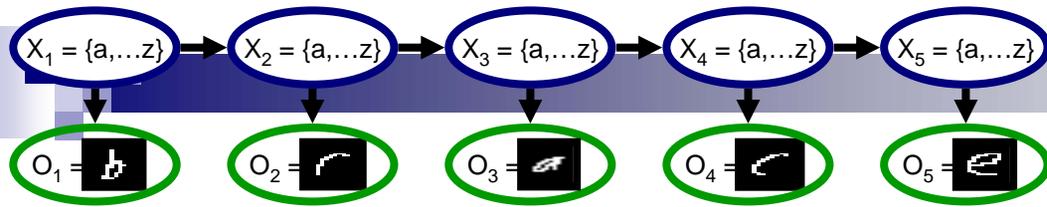  - ☐ Generate a backwards factor by eliminating $X_{i+1}$

  $$\beta_i(X_i) = \sum_{x_{i+1}} P(o_{i+1} \mid x_{i+1})P(x_{i+1} \mid X_i)\beta_{i+1}(x_{i+1})$$

- $\forall$ i, probability is: $\boxed{P(X_i \mid o_{1..n}) \propto \alpha_i(X_i)\beta_i(X_i)}$

*(handwritten annotations on right: $\alpha_n(X_n)$ normalized $= P(X_n \mid O_{1:n})$; $\beta_1(X_1)\alpha_1(X_1)$ normalized $= P(X_1 \mid O_{1:n})$; $\alpha_5(a)$, $\alpha_5(b)$, $\vdots$, $\alpha_5(z)$)*

10

# Most likely explanation

$\neq$ not equal to $\arg\max_{x_i} P(x_i | O_{1:5})$

$X_1 = \{a,\ldots z\}$ → $X_2 = \{a,\ldots z\}$ → $X_3 = \{a,\ldots z\}$ → $X_4 = \{a,\ldots z\}$ → $X_5 = \{a,\ldots z\}$

$O_1 = $   $O_2 = $   $O_3 = $   $O_4 = $   $O_5 = $ 

**Compute:** $\arg\max_{x_1 \cdots x_5} P(x_1, \ldots x_5 | O_{1:5})$

## Variable elimination order?

$1, 2, 3, 4, 5$

## Example:

$$P(x_1, \ldots, x_5 | O_{1:5})$$

$$\max_{x_1 x_2 x_3 x_4, x_5}$$

$$= \max_{x_2, \cdots x_5} P(x_3|x_2) P(O_2|x_2) P(O_3|x_3) P(x_4|x_3) \cdots$$

$$\max_{x_1} P(x_1) P(O_1|x_1) P(x_2|x_1)$$

$$\underbrace{\qquad\qquad}_{\alpha_1(x_2)}$$

$x_5^* = \arg\max_{x_5} \alpha_5(x_5)$ ; backwards

$x_4^* = \arg\max_{x_4} \alpha_4(x_4) \cdot P(x_5^*|x_4)$

$A, B$ binary

$P(A, B) =$

| B \\ A | t | f |
|---|---|---|
| t | 0.3 | 0 |
| f | 0.3 | 0.4 |

$\arg\max_{ab} P(a,b) = \begin{cases} a = f \\ b = f \end{cases}$

$P(a = t) = 0.6$

$P(c = f) = 0.4$

$\arg\max_a P(a) = (a = t)$

$\boxed{P(O_5 | x_5^*)}$ doesn't depend on $x_4$

# The Viterbi algorithm



- Initialization: $\alpha_1(X_1) = P(X_1)P(o_1 \mid X_1)$

- For i = 2 to n

  - Generate a forwards factor by eliminating $X_{i-1}$

$$\alpha_i(X_i) = \max_{x_{i-1}} P(o_i \mid X_i) P(X_i \mid X_{i-1} = x_{i-1}) \alpha_{i-1}(x_{i-1})$$

obs. prob.      transition      message

- Computing best explanation: $x_n^* = \arg\max_{x_n} \alpha_n(x_n)$

- For i = n-1 to 1

  - Use argmax to get explanation:

$$x_i^* = \arg\max_{x_i} P(x_{i+1}^* \mid x_i) \alpha_i(x_i)$$

best for next letter

12

# What you'll implement 1: multiplication

$g_2$

$g_1$

$$\alpha_i(X_i) = \max_{x_{i-1}} P(o_i \mid X_i) P(X_i \mid X_{i-1} = x_{i-1}) \alpha_{i-1}(x_{i-1})$$

$f_1(X_i) \cdot f_2(X_{i-1}, X_i)$ (factors)

def. new factor $g$, domain union $dom(f_1) \cup dom(f_2)$

$\forall x_i, x_{i-1}$

$g(X_{i-1} = x_{i-1}, X_i = x_i) = f_1(X_i = x_i) \cdot f_2(X_{i-1} = x_{i-1}, X_i = x_i)$

# What you'll implement 2: max & argmax

$$\alpha_i(X_i) = \max_{x_{i-1}} P(o_i \mid X_i)P(X_i \mid X_{i-1} = x_{i-1})\alpha_{i-1}(x_{i-1})$$

$$\max_{x_{i-1}} g(X_{i-1}, X_i)$$
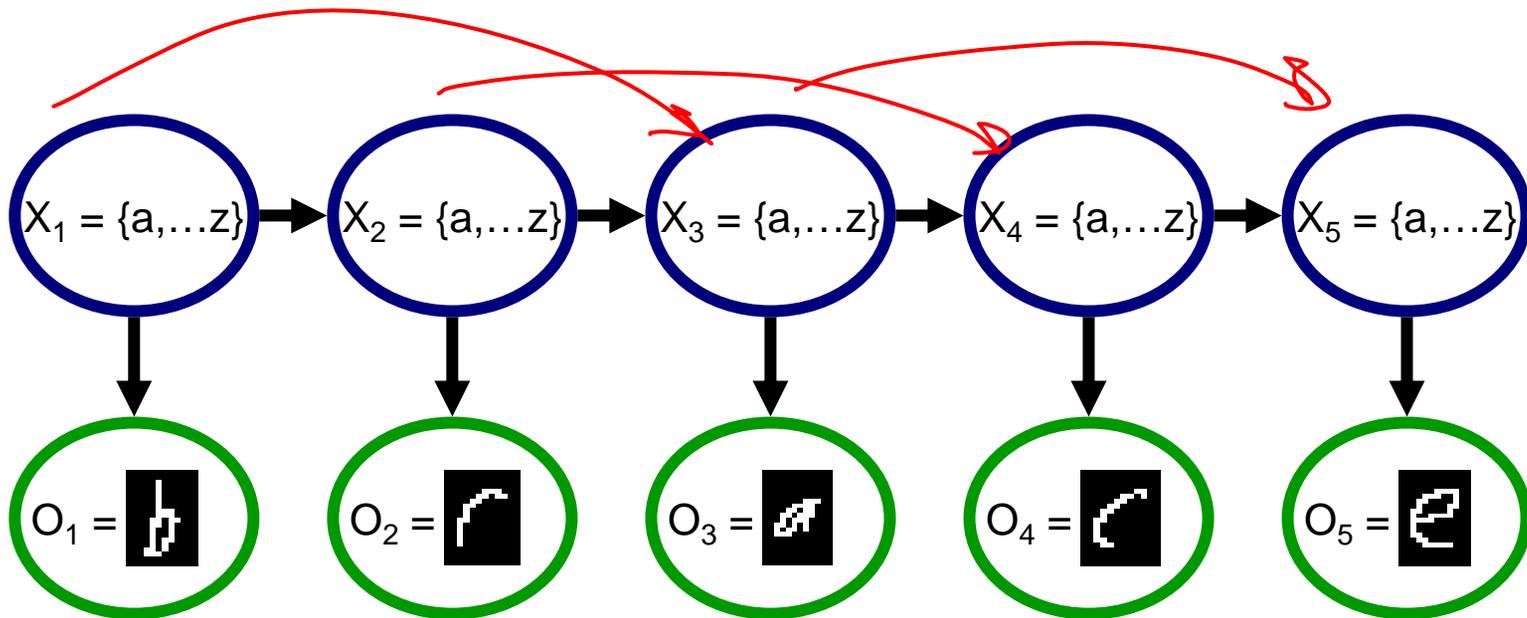
get function $h$, domain $dom(g) - \{X_{i-1}\}$

$$\forall x_i$$

$$h(X_i = x_i) = \max_{x_{i-1}} g(X_{i-1} = x_{i-1}, X_i = x_i)$$

# Higher-order HMMs

2nd order → depend on last 2 time steps



$X_1 = \{a,\ldots z\}$ → $X_2 = \{a,\ldots z\}$ → $X_3 = \{a,\ldots z\}$ → $X_4 = \{a,\ldots z\}$ → $X_5 = \{a,\ldots z\}$

$O_1 =$   $O_2 =$   $O_3 =$   $O_4 =$   $O_5 =$

**Add dependencies further back in time →**
**better representation, harder to learn**

# What you need to know

- Hidden Markov models (HMMs)
    - Very useful, very powerful!
    - Speech, OCR,…
    - Parameter sharing, only learn 3 distributions
    - Trick reduces inference from $O(n^2)$ to $O(n)$
    - Special case of BN

**Koller & Friedman Chapters (handed out):**
    **Chapter 11 (short)**
    **Chapter 12: 12.1, 12.2, 12.3 (covered in the beginning of semester)**
        **12.4 (Learning parameters for BNs)**
    **Chapter 13: 13.1, 13.3.1, 13.4.1, 13.4.3 (basic structure learning)**
**Learning BN tutorial (class website):**
    **ftp://ftp.research.microsoft.com/pub/tr/tr-95-06.pdf**
**TAN paper (class website):**
    **http://www.cs.huji.ac.il/~nir/Abstracts/FrGG1.html**

# Bayesian Networks – (Structure) Learning
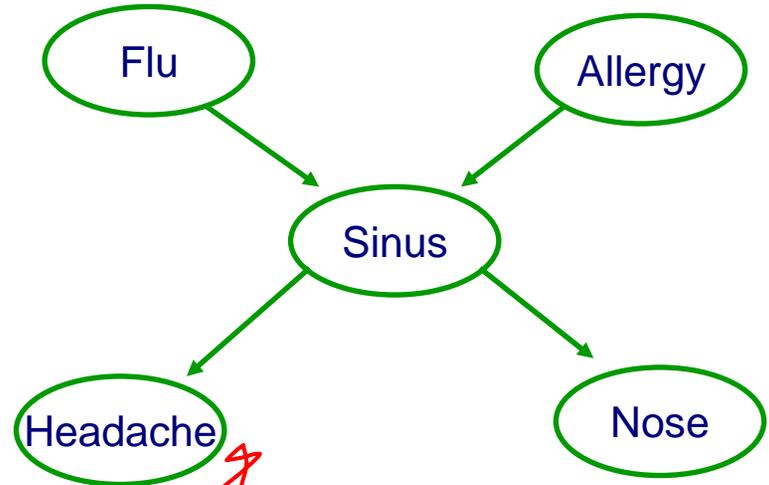
Machine Learning – 10701/15781

Carlos Guestrin

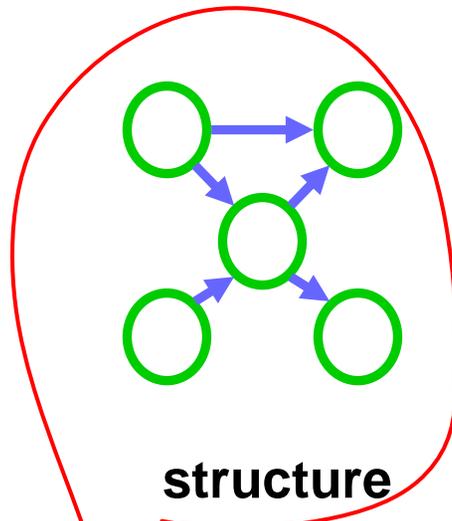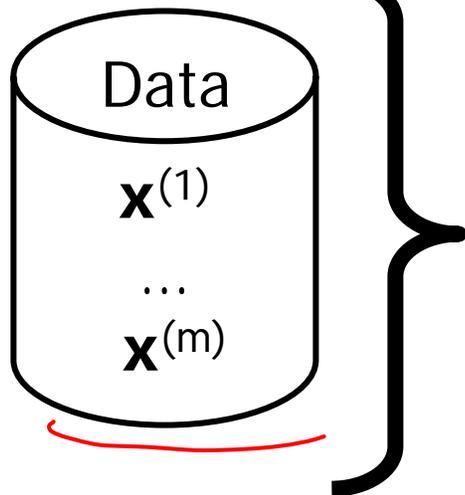Carnegie Mellon University

March 29th, 2006

# Review

- Bayesian Networks
  - Compact representation for probability distributions
  - Exponential reduction in number of parameters
- Fast probabilistic inference using variable elimination
  - Compute P(X|e)
  - Time exponential in tree-width, not number of variables
- Today
  - Learn BN structure

Flu

Allergy

Sinus

Headache

Nose

# Learning Bayes nets

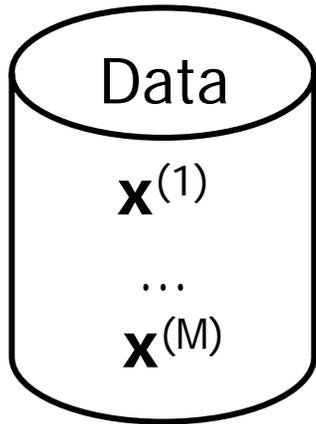| | Known structure | Unknown structure |
|---|---|---|
| Fully observable data | *easy !! :)* | → NP-hard (not always) we'll see next week |
| Missing data | hard, in 2 weeks | really hard, next semester |

A, B, C,
⟨A=a, B=?, C=c⟩



Data

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

…

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

**structure**

$+$

CPTs –
$P(X_i | \mathbf{Pa}_{Xi})$

**parameters**

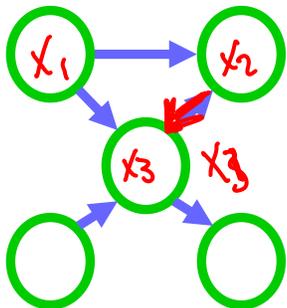# Learning the CPTs

Data

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

…

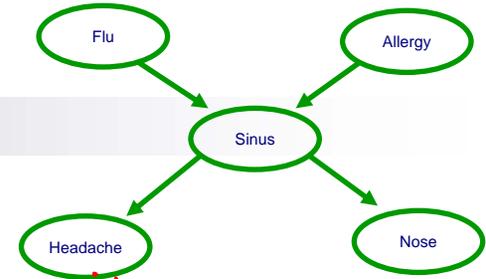$\mathbf{x}^{(M)}$

For each discrete variable $X_i$

$$P(X_3 \mid X_1, X_2)$$

$$P(X_3 = x_3 \mid X_1 = x_1, X_2 = x_2) = \frac{\text{Count}(X_3 = x_3, X_1 = x_1, X_2 = x_2)}{\text{Count}(X_1 = x_1, X_2 = x_2)}$$

**WHY?????????**

$X_1$  $X_2$

$X_3$  $X_3$

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

# Information-theoretic interpretation of maximum likelihood

- Given structure, log likelihood of data:

$$\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G})$$

$\arg\max_{\theta_{\mathcal{G}}}$

$$= \log \prod_j P(x^{(j)} \mid \theta, G)$$

$$x^{(s)} = \langle f=t, S=f, A=f, H=t, N=f \rangle$$
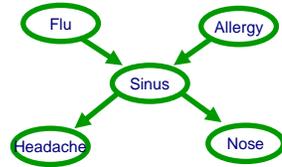
$$= \sum_j \log P(x^{(j)} \mid \theta, G)$$

$$= \sum_j \log P(f^{(j)}) \, P(a^{(j)}) \, P(s^{(j)} \mid f^{(j)}, a^{(j)}) \, P(h^{(j)} \mid s^{(j)}) \, P(n^{(j)} \mid s^{(j)})$$

$$= \left[ \sum_j \log P(f^{(j)}) \right] + \left[ \sum_j \log P(a^{(j)}) \right] \leftarrow \left[ \sum_j \log P(s^{(j)} \mid f^{(j)}, a^{(j)}) \right] + \cdots$$

$$\underbrace{\phantom{\left[ \sum_j \log P(f^{(j)}) \right]}}_{\text{learn } P(F)} \qquad \underbrace{\phantom{\left[ \sum_j \log P(a^{(j)}) \right]}}_{\text{learn } P(A)} \qquad \underbrace{\phantom{\left[ \sum_j \log P(s^{(j)} \mid f^{(j)}, a^{(j)}) \right]}}_{\text{learn } P(S \mid F, A)}$$

log likelihood decomposes
with BN structure

# Maximum likelihood (ML) for learning BN structure

**Possible structures**

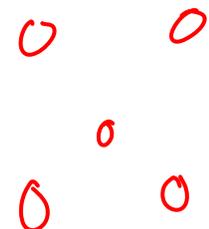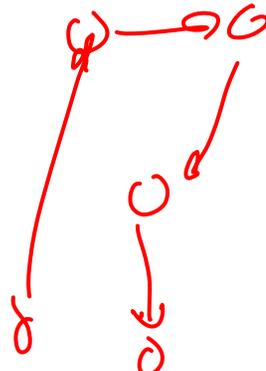Learn parameters using ML

**Score structure**

$$\log P(D \mid \theta_G, G)$$

Flu    Allergy

Sinus

Headache    Nose

**Data**

$$<x_1^{(1)},\ldots,x_n^{(1)}>$$
$$\ldots$$
$$<x_1^{(M)},\ldots,x_n^{(M)}>$$

$-10000$

wins

$-12000$

$-50000$

# How many trees are there?

with n variables?

$$\binom{n}{k} \geq \left(\frac{n}{k}\right)^k$$

A
B
E
C D

how many undirected trees

n vers $\rightarrow$ $\frac{n}{2}(n-1)$ possible edges

choosing n-1

A
E
B
C D

$$\binom{\frac{n}{2}(n-1)}{n-1} \geq \left(\frac{\frac{n}{2}n-1}{n-1}\right)^{n-1} = \left(\frac{n}{2}\right)^{n-1}$$

sorry... I counted    A

B—C
D—E

go back to
d tree slide

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

# Information-theoretic interpretation of maximum likelihood 2

Flu     Allergy

Sinus

Headache     Nose

■ Given structure, log likelihood of data:

$$\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G}) = \sum_{i=1}^{n} \sum_{j=1}^{M} \log P(x_i^{(j)} \mid \mathbf{Pa}_{x_i, \mathcal{G}}^{(j)})$$

learning for each CPT

$s = t, \ f = t, \ a = t$
$= 50 \ times$

$$\sum_{j=1}^{M} \log P(x_i^{(j)} \mid Pa_{x_i, G}^{(j)}) = M \sum_{x_i, Pa_{x_i, G}} \frac{Count(x_i, Pa_{x_i, G})}{M} \cdot \log \hat{P}(x_i \mid Pa_{x_i, G})$$

$$Count(a, b) = \hat{P}(a, b)$$

$$= M \sum_{x_i, Pa_{x_i, G}} \hat{P}(x_i, Pa_{x_i, G}) \ \log \ \hat{P}(x_i \mid Pa_{x_i, G})$$

Conditional entropy
$H(A \mid B)$
$= -\sum_{a, b} P(a, b) \log P(a \mid b)$

$$= -M \hat{H}(X_i \mid Pa_{x_i, G})$$

# Information-theoretic interpretation of maximum likelihood 3

- Given structure, log likelihood of data:

$$\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G}) = M \sum_{i=1}^{n} \sum_{x_i, \mathbf{Pa}_{x_i, \mathcal{G}}} \hat{P}(x_i, \mathbf{Pa}_{x_i, \mathcal{G}}) \log \hat{P}(x_i \mid \mathbf{Pa}_{x_i, \mathcal{G}})$$

# Mutual information $\rightarrow$ Independence tests

- Statistically difficult task!
- Intuitive approach: **Mutual information**

$$I(X_i, X_j) = \sum_{x_i, x_j} P(x_i, x_j) \log \frac{P(x_i, x_j)}{P(x_i) P(x_j)}$$

- Mutual information and independence:
  - □ $X_i$ and $X_j$ independent if and only if $I(X_i, X_j) = 0$

- Conditional mutual information:

# Decomposable score

- **Log data likelihood**

$$\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G}) = M \sum_i \hat{I}(X_i, \mathbf{Pa}_{X_i, \mathcal{G}}) - M \sum_i \hat{H}(X_i)$$

# Scoring a tree 1: equivalent trees

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

# Scoring a tree 2: similar trees

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

# 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 $X_1, \ldots, X_n$
  - Edge (i,j) gets weight $\hat{I}(X_i, X_j)$

# Chow-Liu tree learning algorithm 2

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

- **Optimal tree BN**
  - □ Compute maximum weight spanning tree
  - □ Directions in BN: pick any node as root, breadth-first-search defines directions

# Can we extend Chow-Liu 1

- Tree augmented naïve Bayes (TAN)
[Friedman et al. '97]

  □ Naïve Bayes model overcounts, because correlation between features not considered

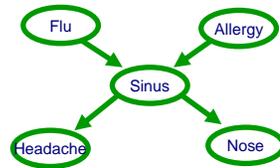  □ Same as Chow-Liu, but score edges with:

$$\hat{I}(X_i, X_j \mid C) = \sum_{c, x_i, x_j} \hat{P}(c, x_i, x_j) \log \frac{\hat{P}(x_i, x_j \mid c)}{\hat{P}(x_i \mid c)\hat{P}(x_j \mid c)}$$

# Can we extend Chow-Liu 2

- (Approximately learning) models
  with tree-width up to *k*
  - □ [Narasimhan & Bilmes '04]
  - □ But, $O(n^{k+1})$…

# Scoring general graphical models – Model selection problem

**What's the best structure?**

Flu → Sinus ← Allergy
Sinus → Headache
Sinus → Nose

**Data**

$<x\_1\textasciicircum\{(1)\},\ldots,x\_n\textasciicircum\{(1)\}>$

…

$<x\_1\textasciicircum\{(m)\},\ldots,x\_n\textasciicircum\{(m)\}>$

**The more edges, the fewer independence assumptions, the higher the likelihood of the data, but will overfit…**

# Maximum likelihood overfits!

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

- Information never hurts:

- Adding a parent always increases score!!!

# Bayesian score avoids overfitting

- Given a structure, distribution over parameters

$$\log P(D \mid \mathcal{G}) = \log \int_{\theta_{\mathcal{G}}} P(D \mid \mathcal{G}, \theta_{\mathcal{G}}) P(\theta_{\mathcal{G}} \mid \mathcal{G}) d\theta_{\mathcal{G}}$$

- Difficult integral: use Bayes information criterion (BIC) approximation (equivalent as M$\to \infty$)

$$\log P(D \mid \mathcal{G}) \approx \log P(D \mid \mathcal{G}, \theta_{\mathcal{G}}) - \frac{\text{NumberParams}(\mathcal{G})}{2} \log M + \mathcal{O}(1)$$

- Note: regularize with MDL score
- Best BN under BIC still NP-hard

# How many graphs are there?

$$\sum_{k=1}^{n} \binom{n}{k} = 2^n - 1$$

# 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) $d \geq 2$

- Most structure learning approaches use heuristics
  - ☐ Exploit score decomposition
  - ☐ (Quickly) Describe two heuristics that exploit decomposition in different ways

# Learn BN structure using local search

**Starting from Chow-Liu tree**

**Local search,** possible moves:
- Add edge
- Delete edge
- Invert edge

**Score using BIC**

# What you need to know about learning BNs

- **Learning BNs**
  - Maximum likelihood or MAP learns parameters
  - Decomposable score
  - Best tree (Chow-Liu)
  - Best TAN
  - Other BNs, usually local search with BIC score