Required Readings from Koller & Friedman:

Representation: 2.1, 2.2

Inference: 5.1, 6.1, 6.2, 6.7.1

Optional:

2.3, 5.2, 5.3, 6.3, 6.7.2

Bayesian Networks – Inference (cont.)

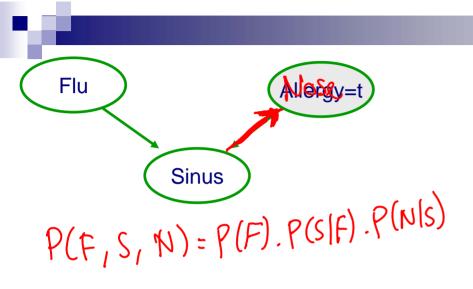
Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

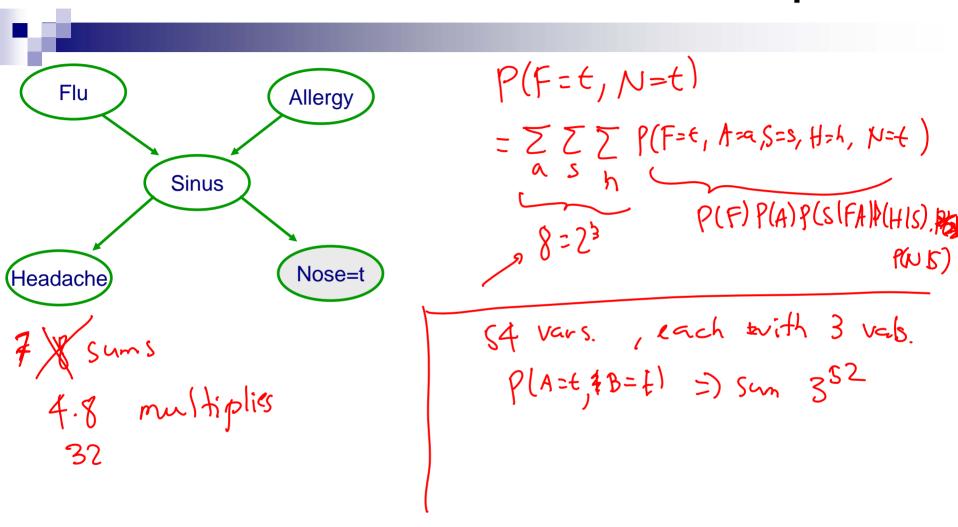
March 26th, 2006

Marginalization



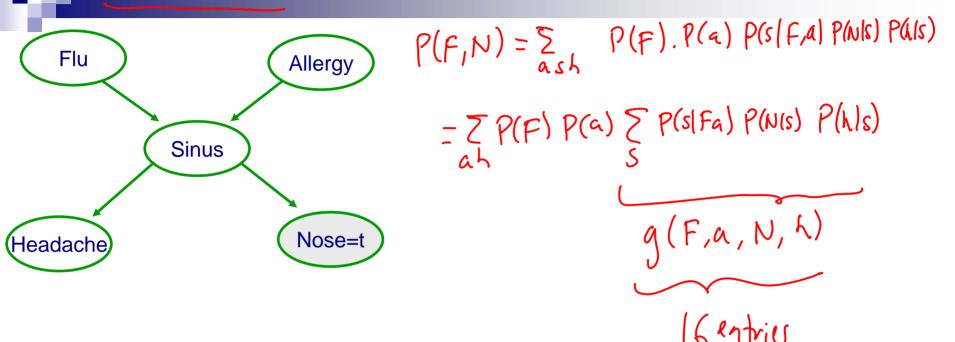
$$P(F = t, N = t)$$
= $Z P(F = t, S = s, N = t)$
= $P(F = t) . P(S = t | F = t) . P(N = t | S = t)$
+
 $P(F = t) . P(S = f | F = t) . P(N = t | S = t)$

Probabilistic inference example



Inference seems exponential in number of variables!

Understanding variable elimination – Order can make a HUGE difference



Variable elimination algorithm

- Instantiate evidence e fix values, N= E
- Prune non-ancestors of {X,e}

- Choose an ordering on variables, e.g., X₁, ..., X_n

- For i = 1 to n, If $X_i \notin \{X, e\}$ Solution Collect factors f_1, \ldots, f_k that include X_i Generate a new factor by eliminating X_i from these factors

$$f_{i}(A_{i}B), f_{2}(B_{i}C) = \prod_{j=1}^{k} f_{j}$$

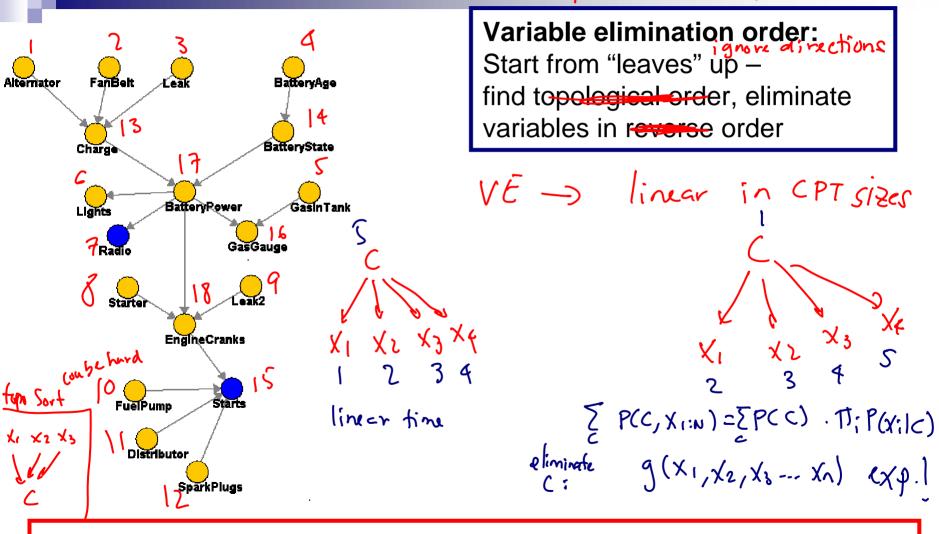
$$f(A_{j}B_{j}C) = f_{i} \cdot f_{2}$$

$$f(A_{j}B_{j}C) = f_{i} \cdot f_{2}$$

$$f(A_{j}C) = f_{i} \cdot f_{$$

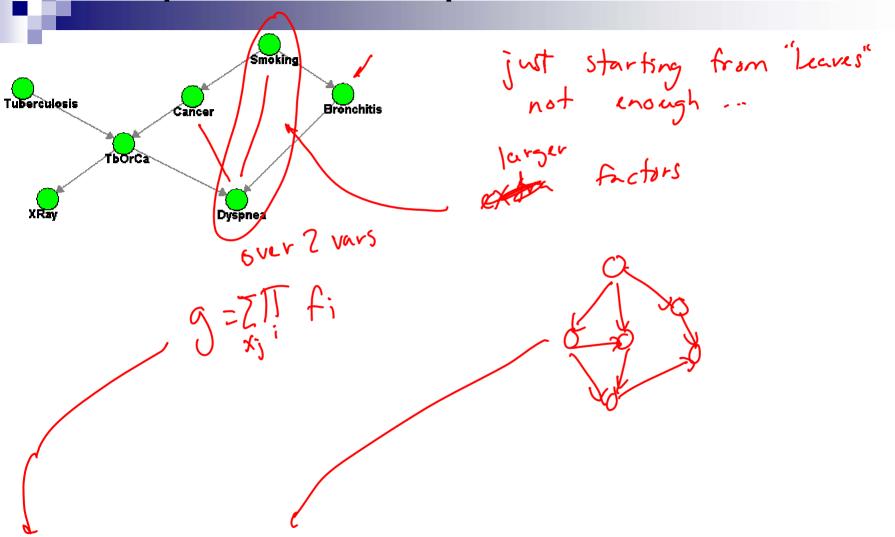
- Normalize P(X,e) to obtain P(X|e)

Complexity of variable elimination – (Poly)-tree graphs



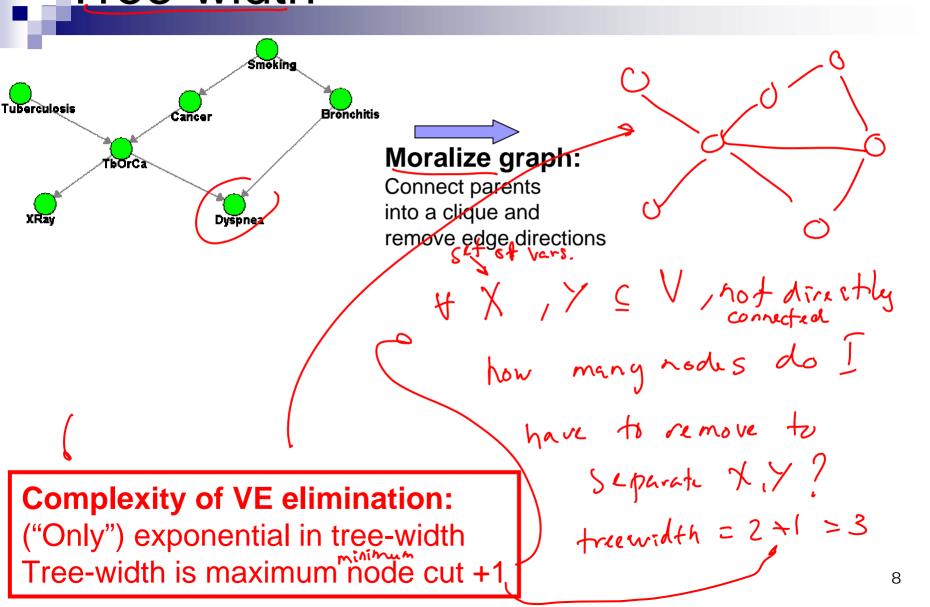
Linear in number of variables!!! (versus exponential)

Complexity of variable elimination – Graphs with loops

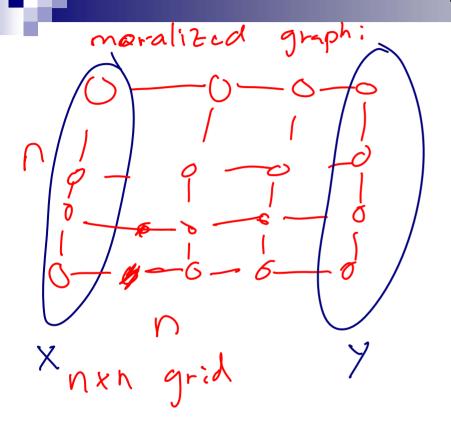


Exponential in number of variables in largest factor generated

Complexity of variable elimination – Tree-width



Example: Large tree-width with small number of parents

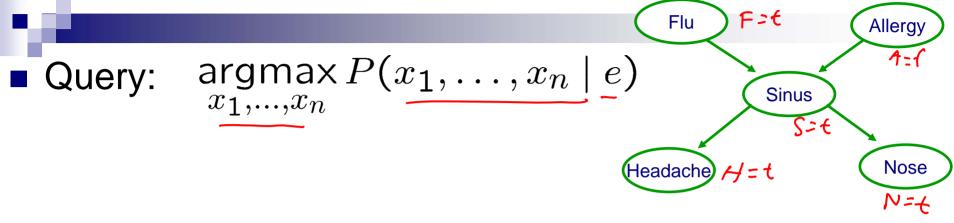


```
true width?
             def. of treewidth has at
even if vars have few parents, inflrence could be exponential in h
```

Choosing an elimination order

- Choosing best order is NP-complete
 - □ Reduction from MAX-Clique
- Many good heuristics (some with guarantees)
- Ultimately, can't beat NP-hardness of inference
 - □ Even optimal order can lead to exponential variable elimination computation
- In practice
 - □ Variable elimination often very effective
 - Many (many many) approximate inference approaches available when variable elimination too expensive

Most likely explanation (MLE)



Using Bayes rule:

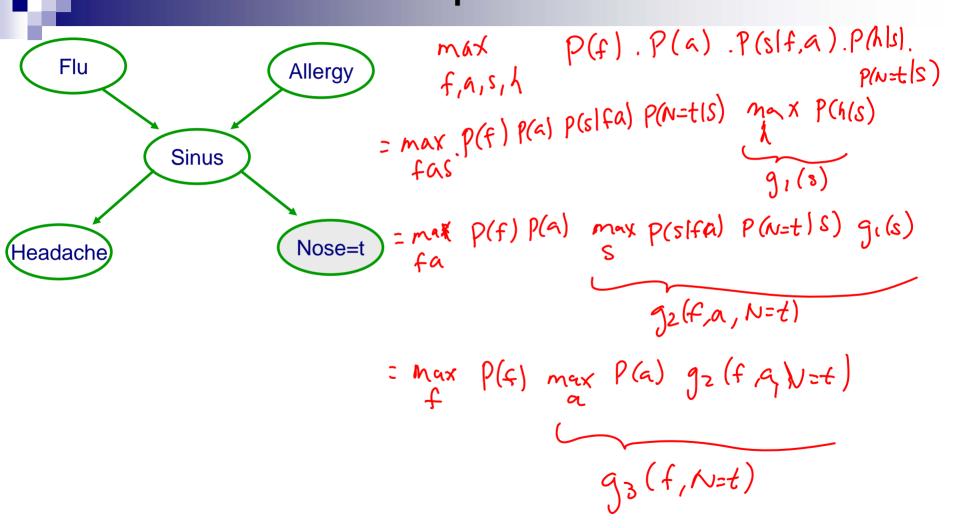
$$\underset{x_1,...,x_n}{\operatorname{argmax}} P(x_1,\ldots,x_n \mid e) = \underset{x_1,...,x_n}{\operatorname{argmax}} \frac{P(x_1,\ldots,x_n,e)}{P(e)}$$

Normalization irrelevant:

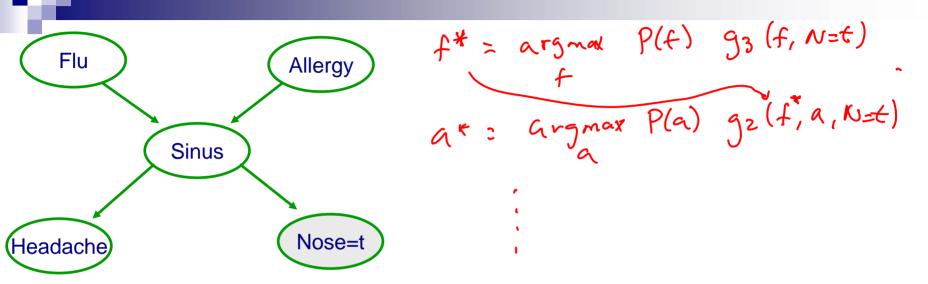
$$\underset{x_1,...,x_n}{\operatorname{argmax}} P(x_1,\ldots,x_n \mid e) = \underset{x_1,...,x_n}{\operatorname{argmax}} P(x_1,\ldots,x_n,e)$$

Max-marginalization

Example of variable elimination for MLE – Forward pass



Example of variable elimination for MLE – Backward pass



MLE Variable elimination algorithmForward pass

- Given a BN and a MLE query $\max_{x_1,...,x_n} P(x_1,...,x_n,\underline{e})$
- Instantiate evidence e
- Choose an ordering on variables, e.g., X₁, ..., X_n
- For i = 1 to n, If $X_i \notin \{e\}$
 - \square Collect factors $f_1, ..., f_k$ that include X_i
 - ☐ Generate a new factor by eliminating X_i from these factors

$$g = \max_{x_i} \prod_{j=1}^{k} f_j$$

□ Variable X_i has been eliminated!

MLE Variable elimination algorithm – Backward pass

- {x₁*,..., x_n*} will store maximizing assignment
- For i = n to 1, If $X_i \notin \{e\}$
 - \square Take factors $f_1, ..., f_k$ used when X_i was eliminated
 - □ Instantiate $f_1, ..., f_k$, with $\{x_{i+1}^*, ..., x_n^*\}$
 - Now each f_i depends only on X_i
 - □ Generate maximizing assignment for X_i:

$$x_i^* \in \operatorname{argmax} \prod_{j=1}^k f_j$$

What you need to know

- Bayesian networks
 - ☐ A useful compact **representation** for large probability distributions
- Inference to compute
 - □ Probability of X given evidence e
 - Most likely explanation (MLE) given evidence e
 - □ Inference is NP-hard
- Variable elimination algorithm
 - Efficient algorithm ("only" exponential in tree-width, not number of variables)
 - Elimination order is important!
 - □ Approximate inference necessary when tree-width to large
 - not covered this semester
 - Only difference between probabilistic inference and MLE is "sum" versus "max"

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

Machine Learning – 10701/15781
Carlos Guestrin
Carnegie Mellon University

March 26th, 2005

Adventures of our BN hero

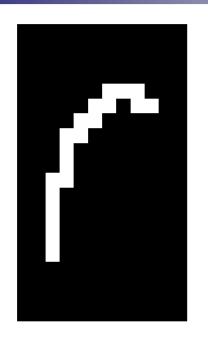
- Compact representation for 1. Naïve Bayes probability distributions
- Fast inference
- Fast learning
- But... Who are the most popular kids?

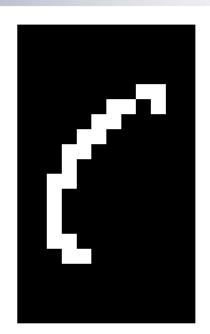
discrete vars.

2 and 3. Hidden Markov models (HMMs) Kalman Filters

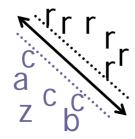
Gaussian Vars

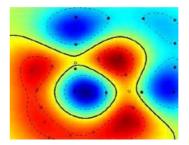
Handwriting recognition





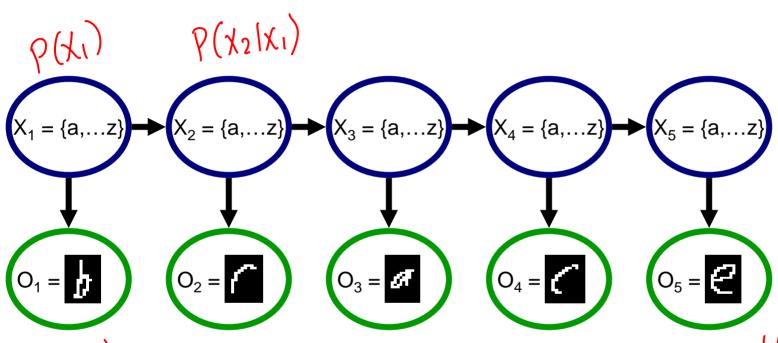
Character recognition, e.g., kernel SVMs





Example of a hidden Markov model (HMM)

Understanding the HMM Semantics



P(0,1X1)

one option is

Naive Bages

P(0;1X1) = IT P(0;1X1)

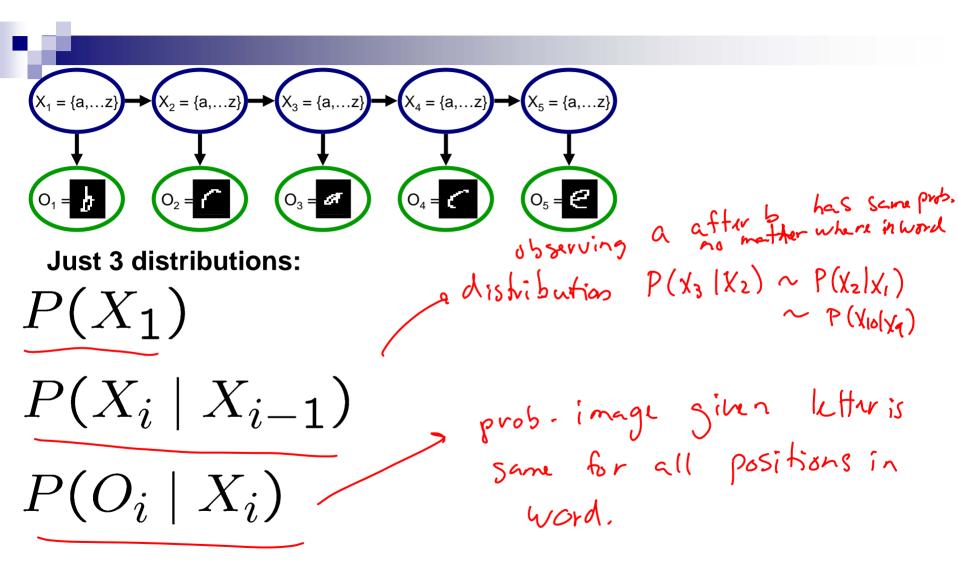
Sinpixels

Markov assumption: the future is indep- of past given present XI:E+ I XE+I:N |XE

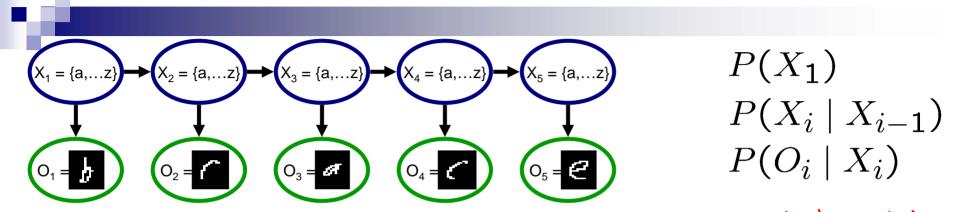
Ot I every Sody / Xt

Music

HMMs semantics: Details



HMMs semantics: Joint distribution



$$P(X_1, ..., X_n \mid o_1, ..., o_n) = P(X_{1:n} \mid o_{1:n})$$

$$\propto P(X_1)P(o_1 \mid X_1) \prod_{i=2}^n P(X_i \mid X_{i-1})P(o_i \mid X_i)$$

Learning HMMs from fully observable data is easy

$$X_1 = \{a, \dots z\} \longrightarrow X_2 = \{a, \dots z\} \longrightarrow (X_3 = \{a, \dots z\}) \longrightarrow (X_4 = \{a, \dots z\}) \longrightarrow (X_5 = \{a, \dots z\})$$

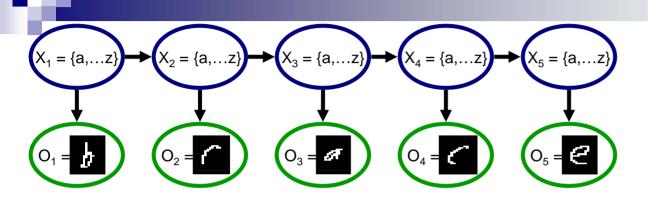
$$O_1 = \{c, \dots z\} \longrightarrow (O_3 = c, \dots z) \longrightarrow (O_4 = c, \dots z)$$

Learn 3 distributions:
$$P(X_1^{\circ}) = (\text{ount (# first letter a})$$
 galact training distributions of the letter was a whole letter was a possible of the letter was a count ($P(X_i^{\circ}|X_i^{\circ}) = (\text{ount (Pixel 12 was white, Xi=9})$ (ount (Xi=a))
$$P(X_i^{\circ}|X_i^{\circ}) = (\text{ount (a appears after b)})$$

$$= (\text{ount (# of b:s that are not ext)}$$

$$= (\text{ount (# of b:s that are not ext)}$$

Possible inference tasks in an HMM



Marginal probability of a hidden variable:

Viterbi decoding – most likely trajectory for hidden vars:

argmax
$$P(x_1, x_2, x_3, x_4, x_5 | 0_{1:5})$$

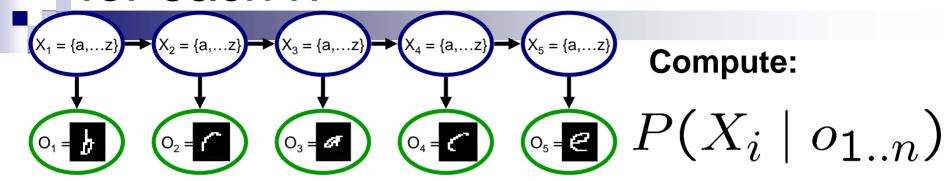
 $x_1 y_2 x_3 x_4 x_5$

Using variable elimination to compute $P(X_i|o_{1:n})$ $x_1 = \{a,...z\} \rightarrow (x_2 = \{a,...z\}) \rightarrow (x_3 = \{a,...z\}) \rightarrow (x_4 = \{a,...z\}) \rightarrow (x_5 = \{a,...z\})$ $x_1 = \{a,...z\} \rightarrow (x_2 = \{a,...z\}) \rightarrow (x_3 = \{a,...z\}) \rightarrow (x_4 = \{a,...z\}) \rightarrow (x_5 = \{a,...z\})$ $x_1 = \{a,...z\} \rightarrow (x_2 = \{a,...z\}) \rightarrow (x_3 = \{a,...z\}) \rightarrow (x_4 = \{a,...z\}) \rightarrow (x_5 = \{a,...z\}) \rightarrow$

Variable elimination order?

Example:

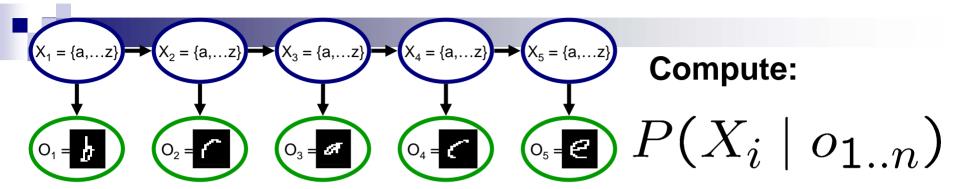
What if I want to compute $P(X_i|o_{1:n})$ for each i?



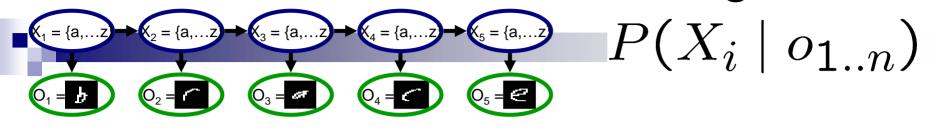
Variable elimination for each i?

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

Reusing computation



The forwards-backwards 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}

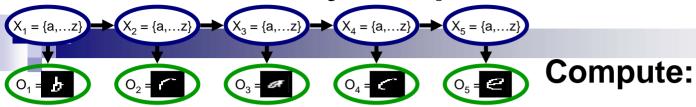
$$\alpha_i(X_i) = \sum_{x_{i-1}} P(o_i \mid X_i) P(X_i \mid X_{i-1} = x_{i-1}) \alpha_{i-1}(x_{i-1})$$

- 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: $P(X_i \mid o_{1..n}) = \alpha_i(X_i)\beta_i(X_i)$

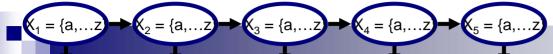
Most likely explanation



Variable elimination order?

Example:

The Viterbi algorithm













- Initialization: $\alpha_1(X_1) = P(X_1)P(o_1 \mid X_1)$
- \blacksquare For i = 2 to n
 - □ Generate a forwards factor by eliminating X_{i-1}

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

- Computing best explanation: $x_n^* = \operatorname{argmax} \alpha_n(x_n)$
- For i = n-1 to 1
 - □ Use argmax to get explanation:

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

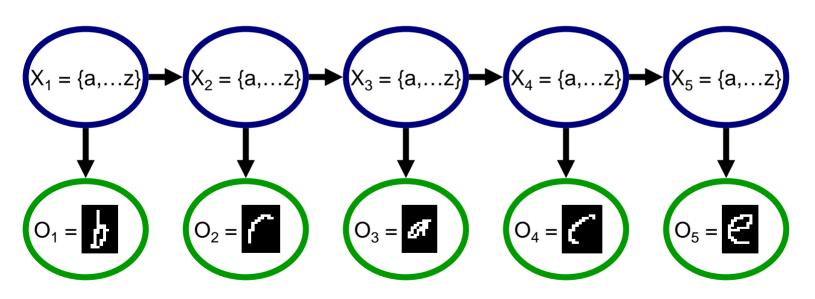
What you'll implement 1: multiplication

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

What you'll implement 2: max & argmax

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

Higher-order HMMs



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
 - \square Trick reduces inference from O(n²) to O(n)
 - □ Special case of BN