

**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.)

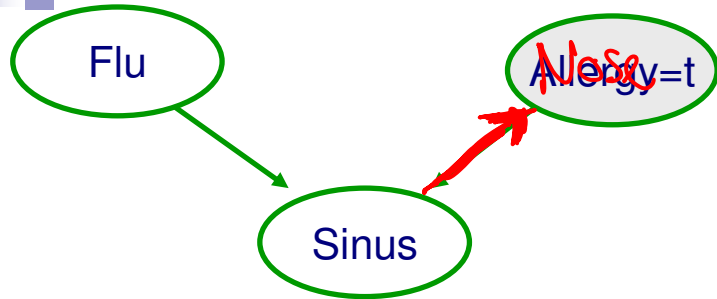
Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

March 26<sup>th</sup>, 2006

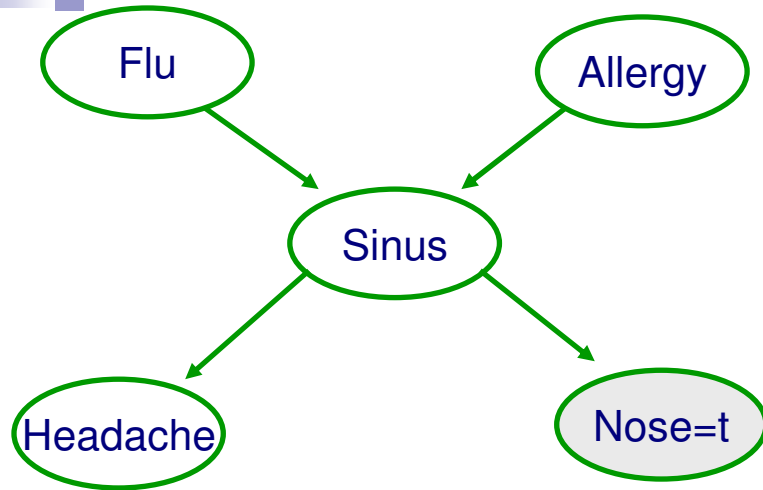
# Marginalization



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

$$\begin{aligned} P(F=t, N=t) &= \sum_s P(F=t, S=s, N=t) \\ &= P(F=t) \cdot P(S=t|F=t) \cdot P(N=t|S=t) \\ &\quad + \\ &\quad P(F=t) \cdot P(S=f|F=t) \cdot P(N=t|S=f) \end{aligned}$$

# Probabilistic inference example



$$P(F=t, N=t)$$

$$= \sum_a \sum_s \sum_h P(F=t, A=a, S=s, H=h, N=t)$$

$\underbrace{\sum_a \sum_s \sum_h}_{8=2^3}$ 
 $\underbrace{P(F)P(A)P(S|FA)P(H|S)P(N|S)}_{P(N|S)}$

~~7~~ sums

4.8 multiplies

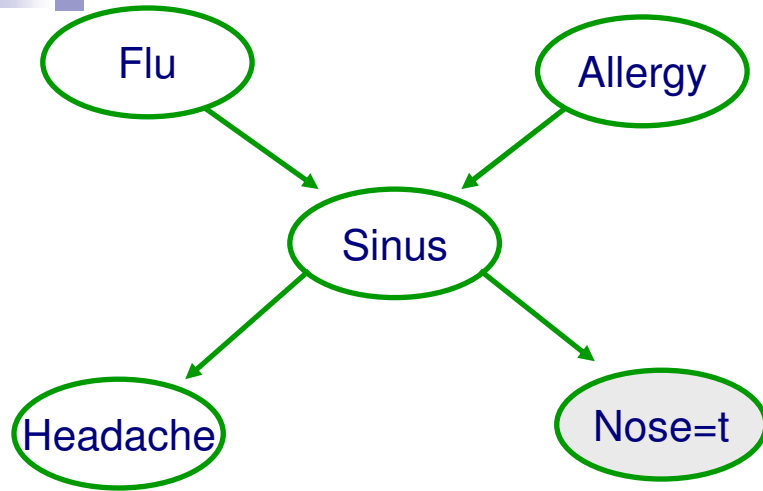
32

54 vars. , each with 3 vals.

$$P(A=t, B=t) \Rightarrow \text{sum } 3^{52}$$

**Inference seems exponential in number of variables!**

# Understanding variable elimination – Order can make a HUGE difference



$$P(F, N) = \sum_{a, s, h} P(F) \cdot P(a) P(s|F, a) P(N|s) P(h|s)$$

$$= \sum_{a, h} P(F) P(a) \sum_s P(s|F, a) P(N|s) P(h|s)$$

$$\underbrace{\hspace{10em}}_{g(F, a, N, h)}$$

16 entries

# Variable elimination algorithm

- Given a BN and a query  $P(X|e) \propto P(X,e)$
- Instantiate evidence  $e$  *fix values,  $N=e$*
- Prune non-ancestors of  $\{X,e\}$  IMPORTANT!!!
- Choose an ordering on variables, e.g.,  $X_1, \dots, X_n$
- For  $i = 1$  to  $n$ , If  $X_i \notin \{X,e\}$ 
  - Collect *CPTs, or g.s* factors  $f_1, \dots, f_k$  that include  $X_i$
  - Generate a new factor by eliminating  $X_i$  from these factors
 

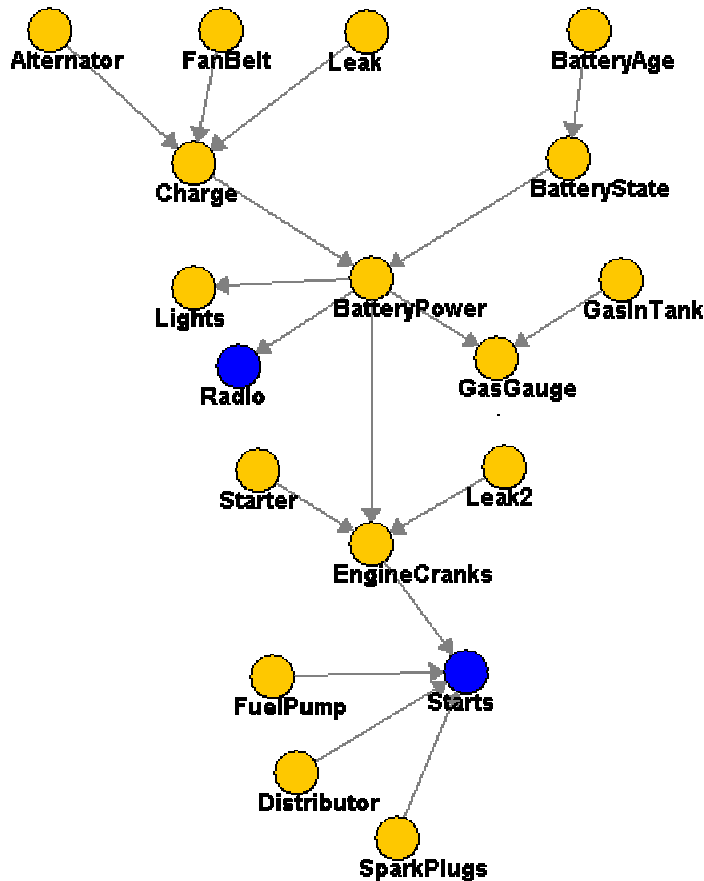
$$\underline{g} = \sum_{X_i} \prod_{j=1}^k \underline{f_j}$$

↙
 $\sum_s P(s|a,F) \cdot g(N,s)$
  - Variable  $X_i$  has been eliminated!
- Normalize  $P(X,e)$  to obtain  $P(X|e)$

# Complexity of variable elimination – (Poly)-tree graphs

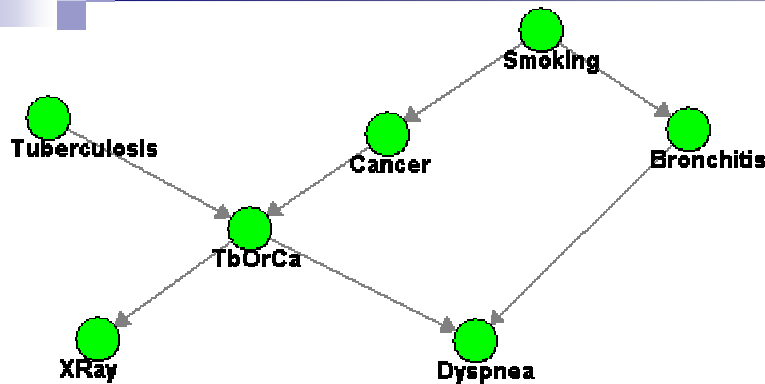
## Variable elimination order:

Start from “leaves” up –  
find topological order, eliminate  
variables in reverse order



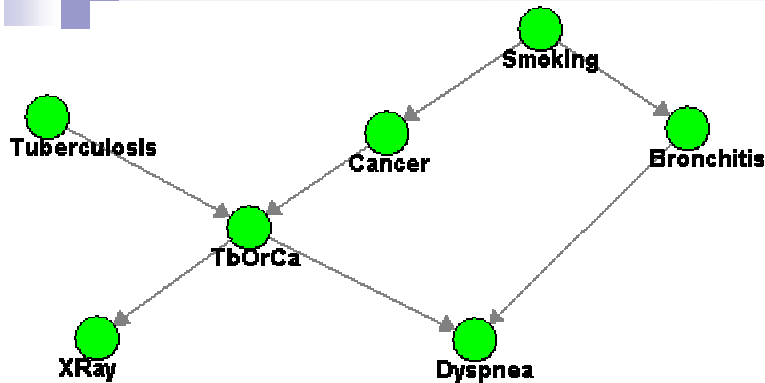
**Linear in number of variables!!! (versus exponential)**

# Complexity of variable elimination – Graphs with loops



**Exponential in number of variables in largest factor generated**

# Complexity of variable elimination – Tree-width



**Moralize graph:**  
Connect parents  
into a clique and  
remove edge directions

**Complexity of VE elimination:**  
("Only") exponential in tree-width  
Tree-width is maximum node cut + 1



# Example: Large tree-width with small number of parents

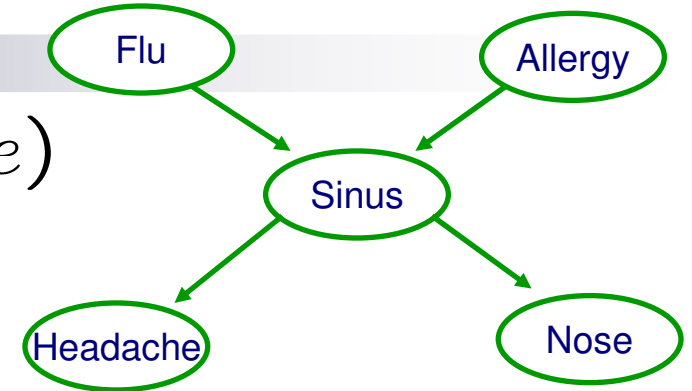
Compact representation  $\nRightarrow$  Easy inference 😞

# 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)

- Query:  $\operatorname{argmax}_{x_1, \dots, x_n} P(x_1, \dots, x_n \mid e)$



- Using Bayes rule:

$$\operatorname{argmax}_{x_1, \dots, x_n} P(x_1, \dots, x_n \mid e) = \operatorname{argmax}_{x_1, \dots, x_n} \frac{P(x_1, \dots, x_n, e)}{P(e)}$$

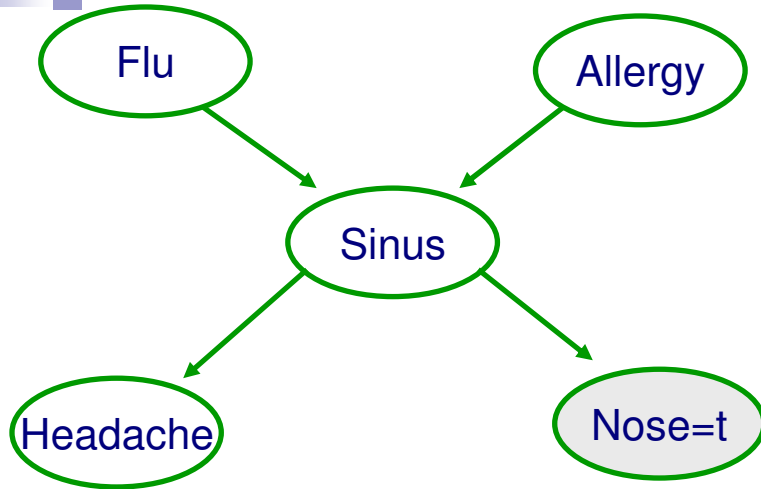
- Normalization irrelevant:

$$\operatorname{argmax}_{x_1, \dots, x_n} P(x_1, \dots, x_n \mid e) = \operatorname{argmax}_{x_1, \dots, x_n} P(x_1, \dots, x_n, e)$$

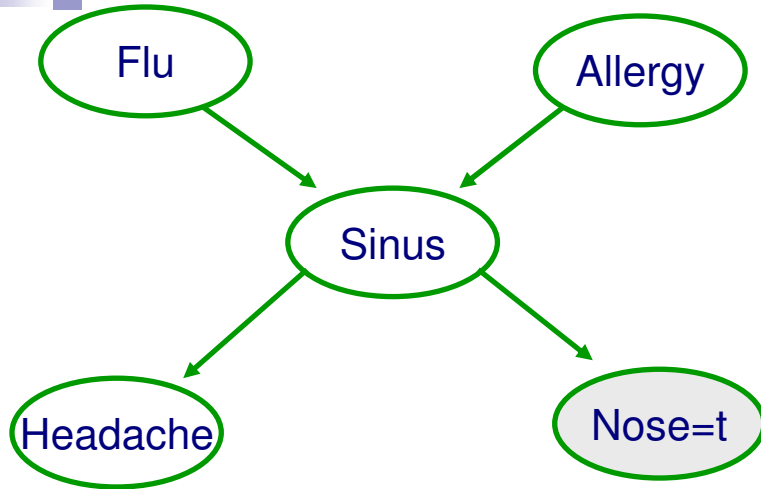
# Max-marginalization



# Example of variable elimination for MLE – Forward pass



# Example of variable elimination for MLE – Backward pass



# MLE Variable elimination algorithm

## – Forward pass

- Given a BN and a MLE query  $\max_{x_1, \dots, x_n} P(x_1, \dots, x_n, e)$
- Instantiate evidence  $e$
- Choose an ordering on variables, e.g.,  $X_1, \dots, X_n$
- For  $i = 1$  to  $n$ , If  $X_i \notin \{e\}$ 
  - Collect factors  $f_1, \dots, f_k$  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

- $\{x_1^*, \dots, x_n^*\}$  will store maximizing assignment
- For  $i = n$  to  $1$ , If  $X_i \notin \{e\}$ 
  - Take factors  $f_1, \dots, f_k$  used when  $X_i$  was eliminated
  - Instantiate  $f_1, \dots, f_k$ , with  $\{x_{i+1}^*, \dots, x_n^*\}$ 
    - Now each  $f_j$  depends only on  $X_i$
  - Generate maximizing assignment for  $X_i$ :

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



# 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 too 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

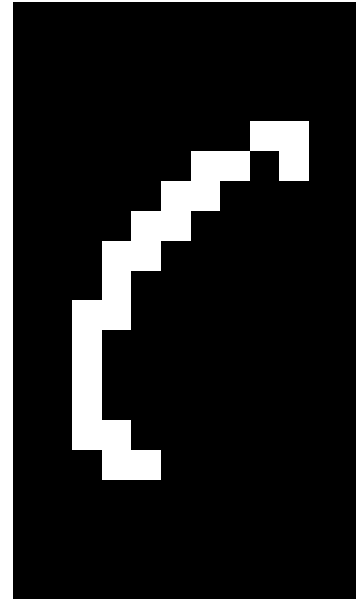
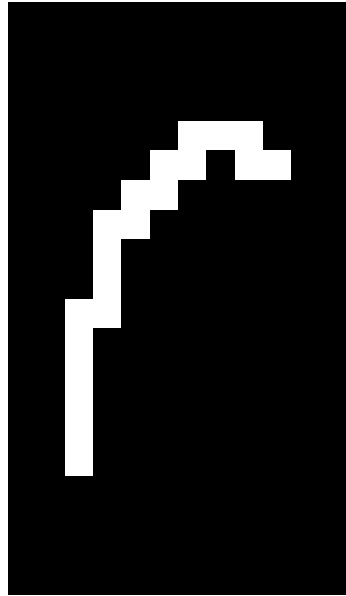
Carnegie Mellon University

March 26<sup>th</sup>, 2005

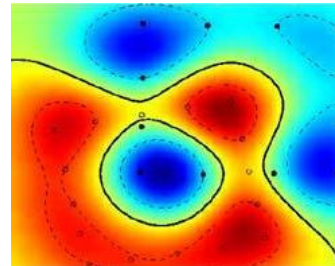
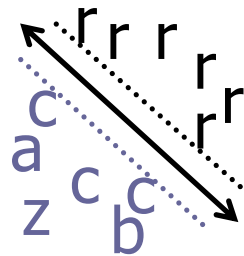
# Adventures of our BN hero

- Compact representation for probability distributions
  - Fast inference
  - Fast learning
  - But... Who are the most popular kids?
- 1. Naïve Bayes**
- 2 and 3.**  
**Hidden Markov models (HMMs)**  
**Kalman Filters**

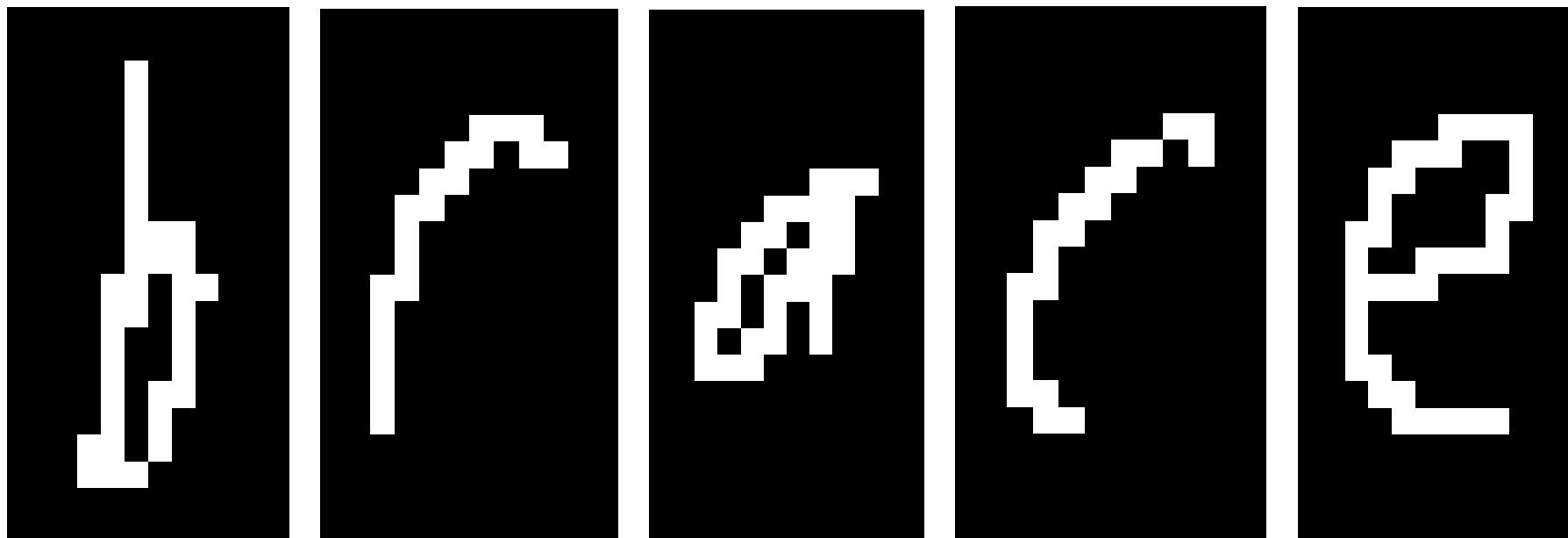
# Handwriting recognition



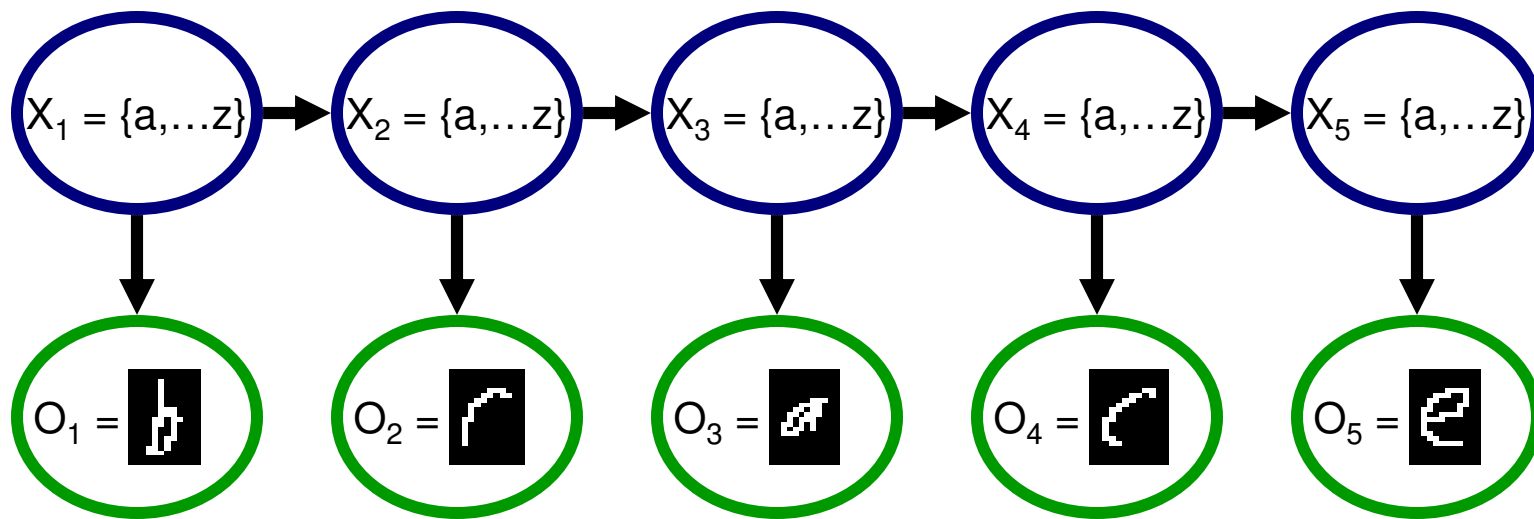
Character recognition, e.g., kernel SVMs



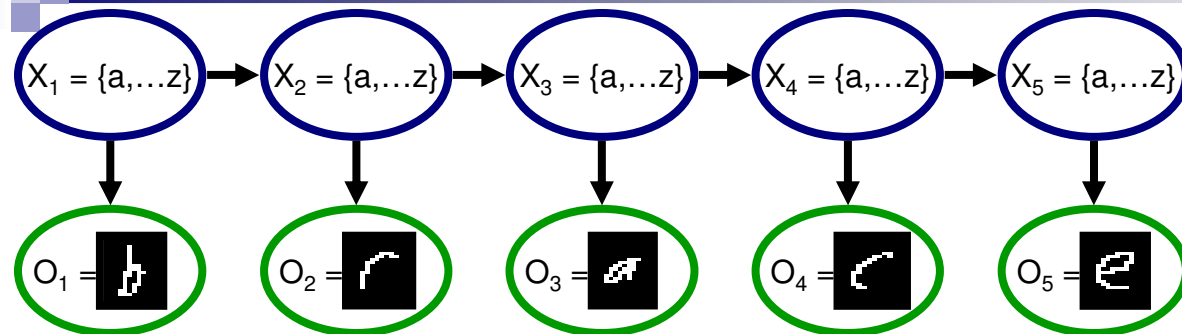
# Example of a hidden Markov model (HMM)



# Understanding the HMM Semantics



# HMMs semantics: Details



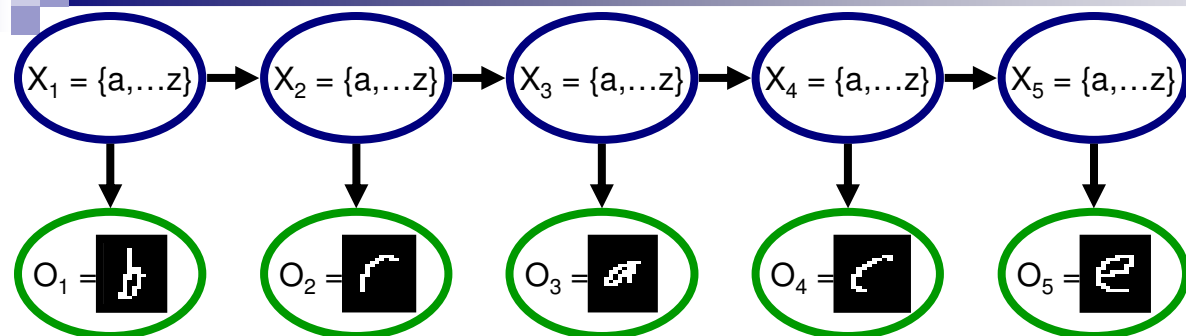
**Just 3 distributions:**

$$P(X_1)$$

$$P(X_i | X_{i-1})$$

$$P(O_i | X_i)$$

# HMMs semantics: Joint distribution



$$P(X_1)$$

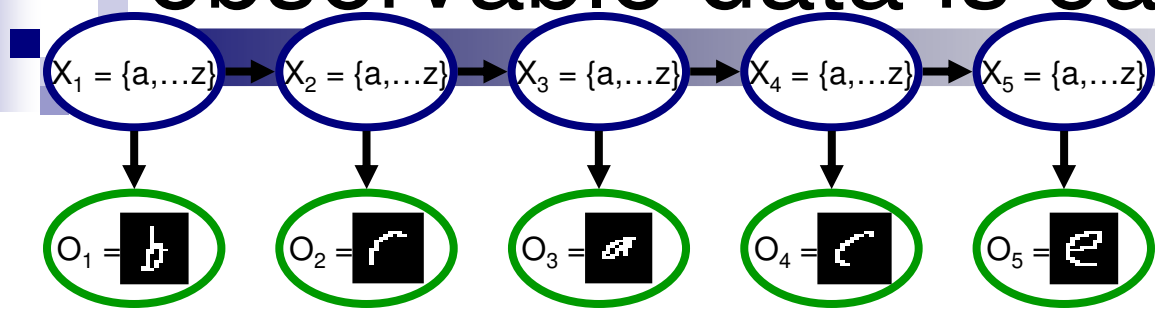
$$P(X_i | X_{i-1})$$

$$P(O_i | X_i)$$

$$P(X_1, \dots, X_n | o_1, \dots, o_n) = P(X_{1:n} | o_{1:n}) \\ \propto P(X_1)P(o_1 | X_1) \prod_{i=2}^n P(X_i | X_{i-1})P(o_i | X_i)$$



# Learning HMMs from fully observable data is easy



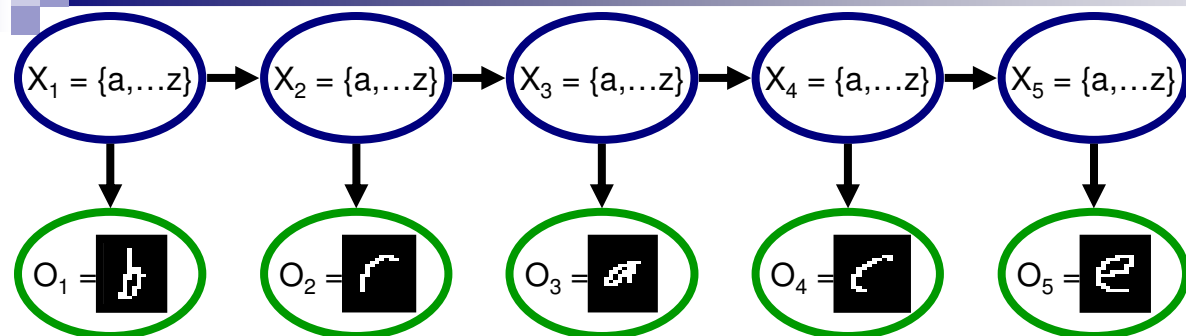
Learn 3 distributions:

$$P(X_1)$$

$$P(O_i | X_i)$$

$$P(X_i | X_{i-1})$$

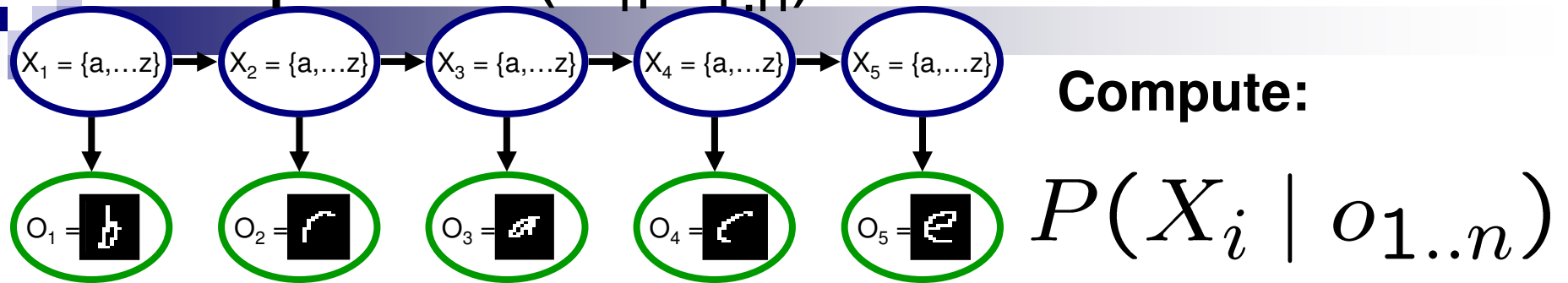
# Possible inference tasks in an HMM



**Marginal probability of a hidden variable:**

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

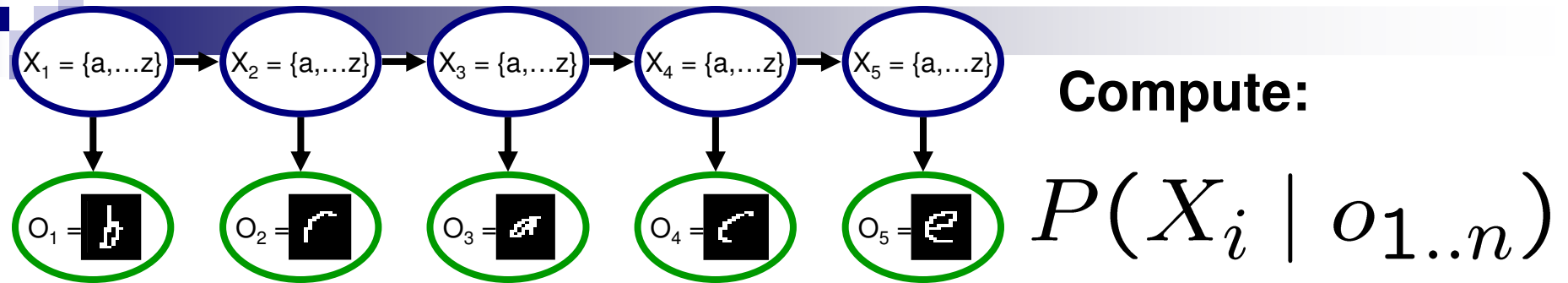
# Using variable elimination to compute $P(X_i | o_{1:n})$



**Variable elimination order?**

**Example:**

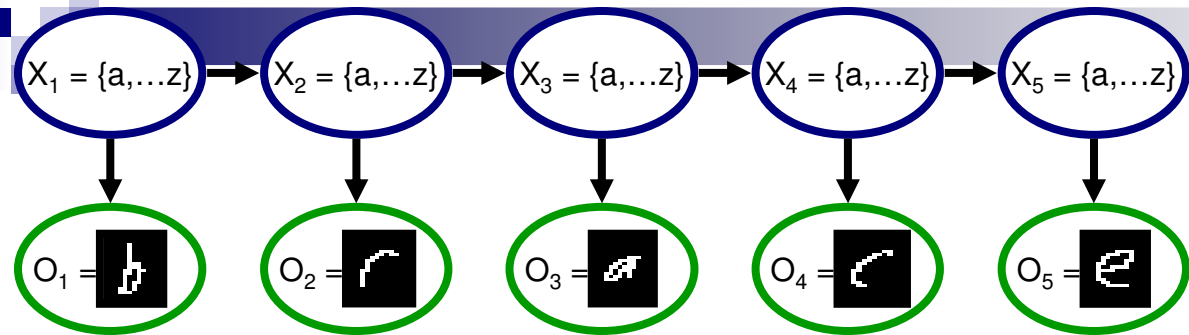
What if I want to compute  $P(X_i | o_{1:n})$   
for each  $i$ ?



**Variable elimination for each  $i$ ?**

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

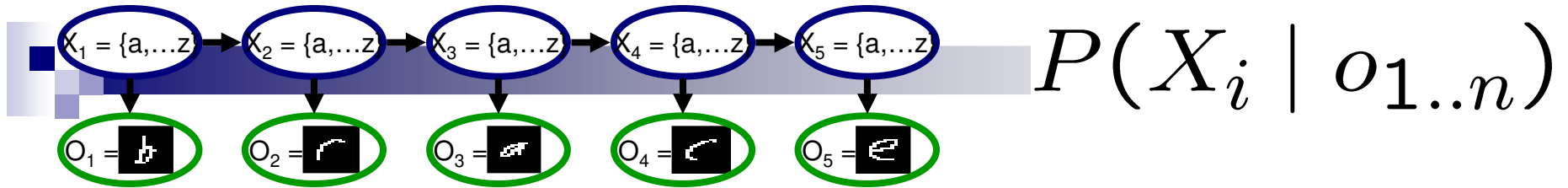
# Reusing computation



**Compute:**

$$P(X_i | o_{1..n})$$

# The forwards-backwards algorithm



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

- For  $i = 2$  to  $n$

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

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

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

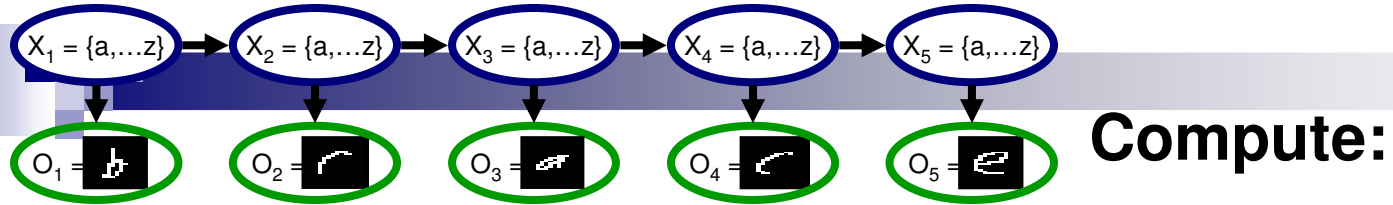
- For  $i = n-1$  to  $1$

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

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

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

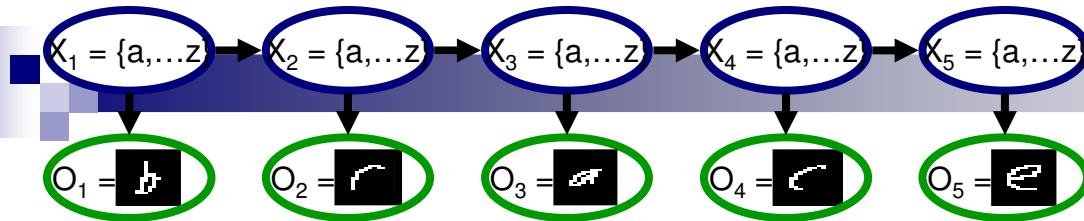
# Most likely explanation



**Variable elimination order?**

**Example:**

# The Viterbi algorithm



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

- For  $i = 2$  to  $n$

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

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

- Computing best explanation:  $x_n^* = \operatorname{argmax}_{x_n} \alpha_n(x_n)$

- For  $i = n-1$  to  $1$

- Use  $\operatorname{argmax}$  to get explanation:

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



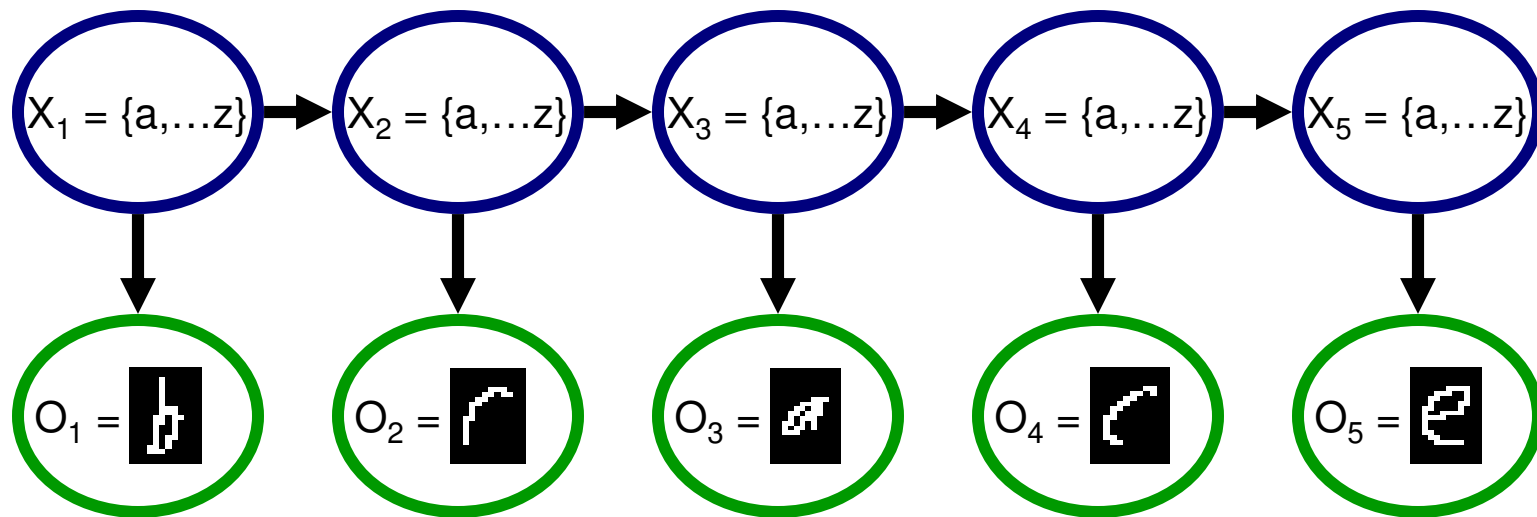
# What you'll implement 1: multiplication

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

# What you'll implement 2: max & argmax

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

# 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
  - Trick reduces inference from  $O(n^2)$  to  $O(n)$
  - Special case of BN