



Bayesian Networks – Inference (cont.)

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

March 26th, 2007

©2005-2007 Carlos Guestrin

General probabilistic inference

■ Query:

$$P(\underline{X} \mid e)$$

Defn. cond. probs.

■ Using ~~Bayes rule~~:

$$P(X \mid e) = \frac{P(X, e)}{P(e)}$$

■ Normalization:

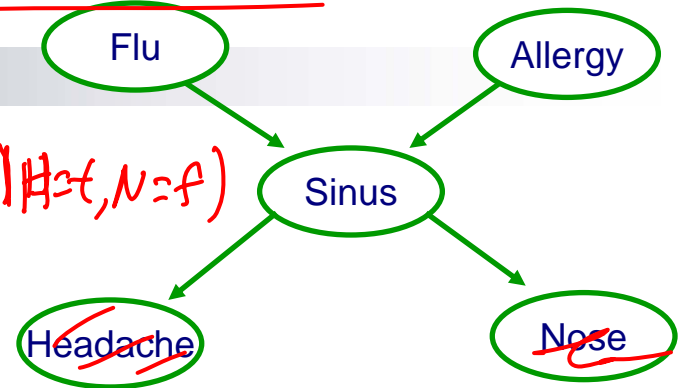
$$P(X \mid e) \propto P(X, e)$$

normalize to give answer

constant doesn't depend on X

compute

$$P(F=t \mid H=t, N=f)$$



$$P(\underline{F}, H=t, N=f)$$

normalize

F	
t	.3
f	.2

not normalized

$$P(F \mid H=t, N=f)$$

t	.6
f	.4

Marginalization



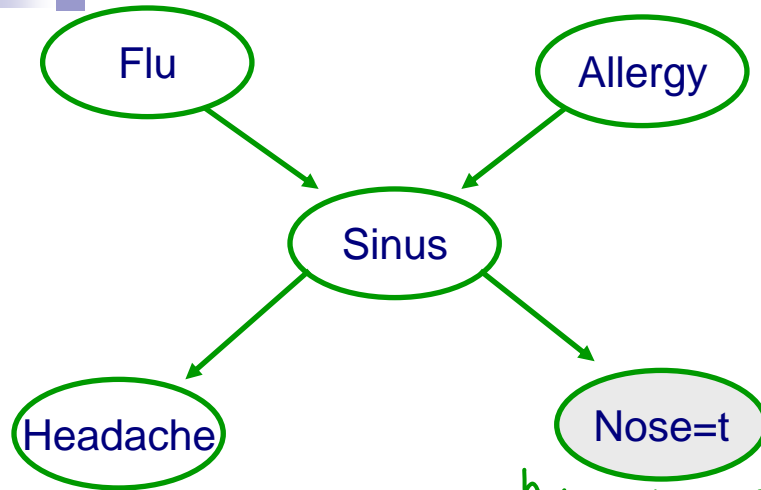
$$P(F, S, N) = P(F) \cdot P(S|F) \cdot P(N|S)$$

$$P(F=t, N=t) = P(F=t, S=t, N=t) + \underbrace{P(F=t, S=f, N=t)}$$

$$= P(F=t) \cdot P(S=t|F=t) \cdot P(N=t|S=t) + P(F=t) \cdot P(S=f|F=t) \cdot P(N=t|S=f)$$

↑
marginalize out S

Probabilistic inference example



$$P(F, N=t) \leftarrow \text{want}$$

know

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

$$P(F) \cdot P(A) \cdot P(S|FA) \cdot P(H|S) \cdot P(N=t|S)$$

how many terms adding? 2^3

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

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

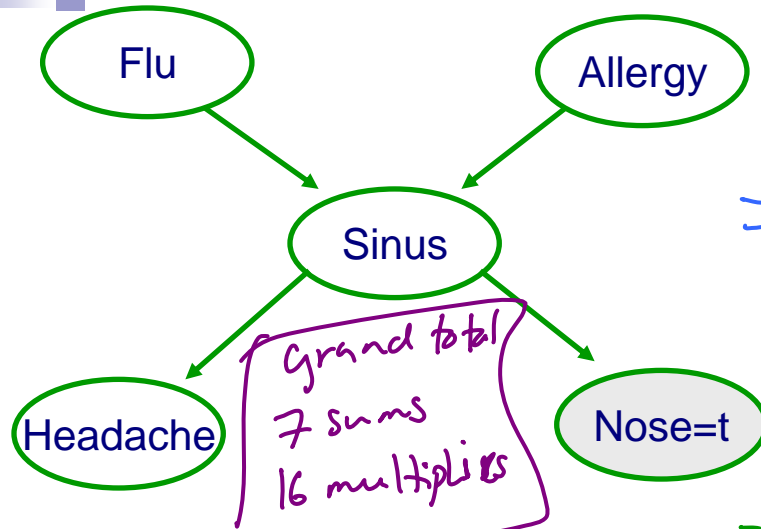
per value of F
7 sums
 $8 \times 4 = 32$ multiplies
grand total
17 sums
64 multiplies
 $P(N=t|S)$

Inference seems exponential in number of variables!
Actually, inference in graphical models is ^{in general} NP-hard ☹

Fast probabilistic inference

example – Variable elimination

eliminate (marginalize) vars one at a time



$$P(F, N=t) = \sum_a \sum_s \sum_h P(F) \cdot P(a) \cdot P(s|F, a) \cdot P(h|s) \cdot P(N=t|s)$$

$$= \sum_a \sum_s P(F) \cdot P(a) \cdot P(s|F, a) \cdot P(N=t|s)$$

$$\sum_h P(h|s)$$

1 sum
0 multiplies

special case

$$= \sum_a \sum_s P(F) \cdot P(a) \cdot P(s|F, a) \cdot P(N=t|s) \cdot 1$$

$$= \sum_a P(F) \cdot P(a) \sum_s P(s|F, a) \cdot P(N=t|s)$$

For each assignment of F & a

1 sum
2 multiplies

total
4 sums
8 multiplies

$$= \sum_a P(F) \cdot P(a) \cdot g_1(F, a) = P(F, N=t)$$

for each assignment of F: 1 sum
4 multiplies

total
2 sums
8 multiplies

(Potential for) Exponential reduction in computation!

Understanding variable elimination – Exploiting distributivity

$$a(b+c) = ab + ac$$

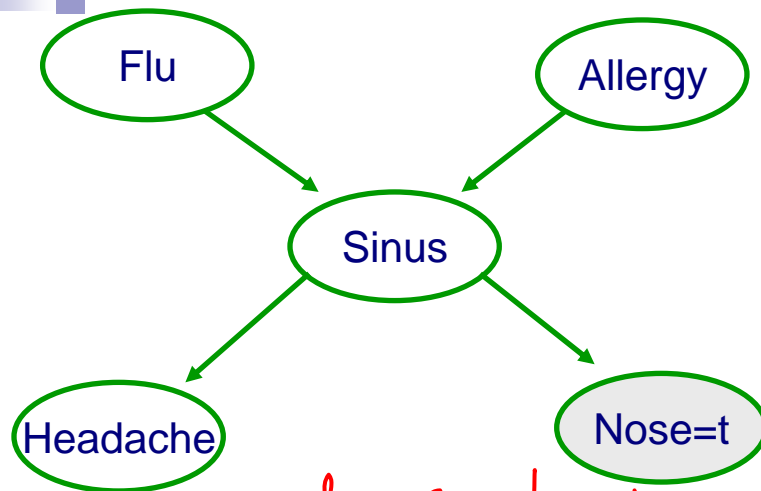


$$P(F, N=t) = P(F) \cdot P(S=t|F) \cdot P(N=t|S=t) + P(F) \cdot P(S=f|F) \cdot P(N=t|S=f)$$

$$= P(F) \left(P(S=t|F) \cdot P(N=t|S=t) + P(S=f|F) \cdot P(N=t|S=f) \right)$$

$g_1(F)$

Understanding variable elimination – Order can make a HUGE difference



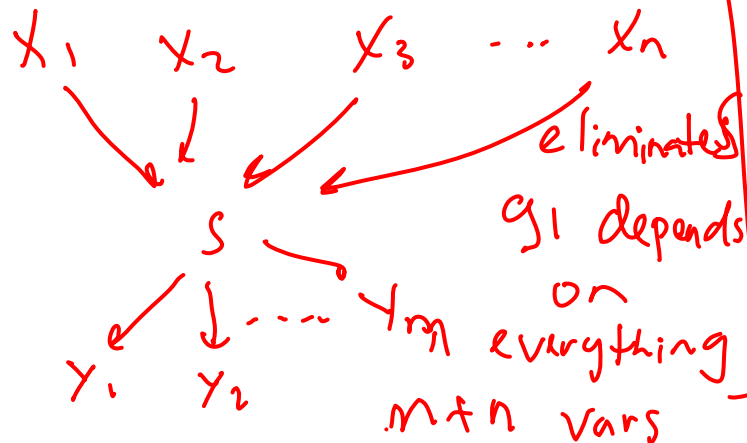
$$P(F, N=t) = \sum_a \sum_s \sum_h P(F) \cdot P(a) \cdot P(s|F, a) \cdot P(h|s) \cdot P(N=t|s)$$

$$= \sum_a \sum_h P(F) P(a)$$

increased g_1 by 1,

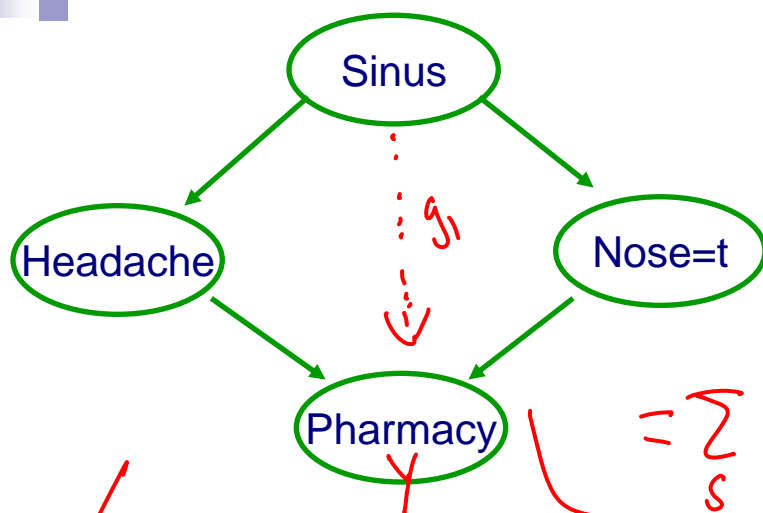
$$\sum_s P(s|F, a) \cdot P(h|s) \cdot P(N=t|s)$$

$g_1(F, a, h)$



exponentially larger.

Understanding variable elimination – Another example



$$P(Y|N=t)$$

$$P(Y, N=t) = \sum_s \sum_h P(s) \cdot P(h|s) \cdot P(N=t|s) \cdot P(Y|h, N=t)$$

$$= \sum_s P(s) P(N=t|s) \sum_h P(h|s) \cdot P(Y|h, N=t)$$

Y, S never appear together in a factor (in a CPT)

after eliminate H , distribution over S, N, Y

$$P(S, Y, N=t) = P(s) \cdot P(N=t|s) g_1(s, Y, N=t)$$

$$g_1(s, Y) \leftarrow S \text{ and } Y \text{ appear together} \\ \downarrow \\ g_1(s, Y, N=t)$$

V.E. generates distributions over smaller sets of variables

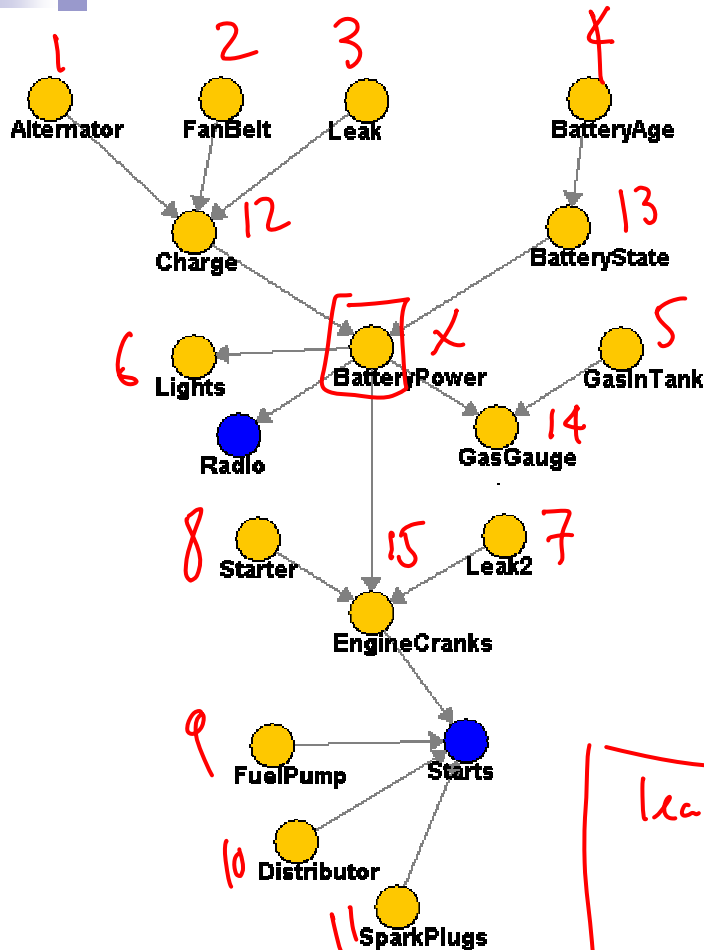
Variable elimination algorithm

- Given a BN and a query $P(X|e) \propto P(X,e)$
 - Instantiate evidence $e \leftarrow H=t, S=f$. **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 factors f_1, \dots, f_k that include X_i
 - Generate a new factor by eliminating X_i from these factors
- marginalize X_i* *multiply everything that depends on X_i*
- $$g = \sum_{X_i} \prod_{j=1}^k f_j$$
- 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



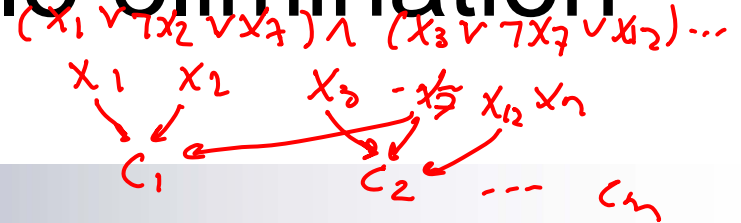
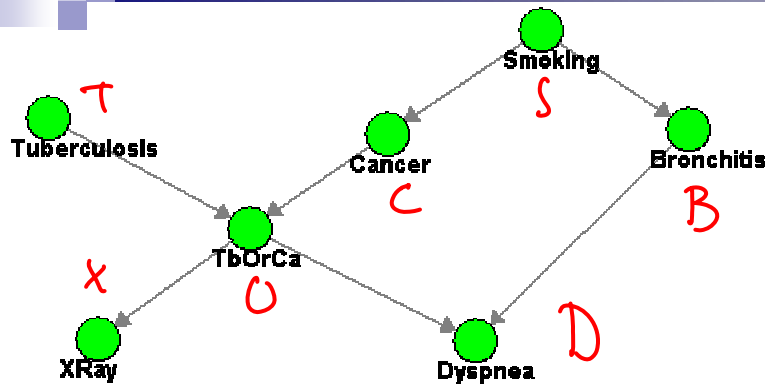
complexity not using VE
about $2^{15} \cdot c$

with VE
about $15 \cdot c$

leaves can be eliminated in
any order.

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

Complexity of variable elimination – Graphs with loops



$$\sum_B P(T) P(S) \cdot P(B|S) \cdot P(C|S) \cdot P(O|C,T) \cdot P(D|S,B) \cdot P(X|O)$$

$$g_i(S, O, D) =$$

K values

$$\sum_B P(B|S) \cdot P(D|O,B)$$

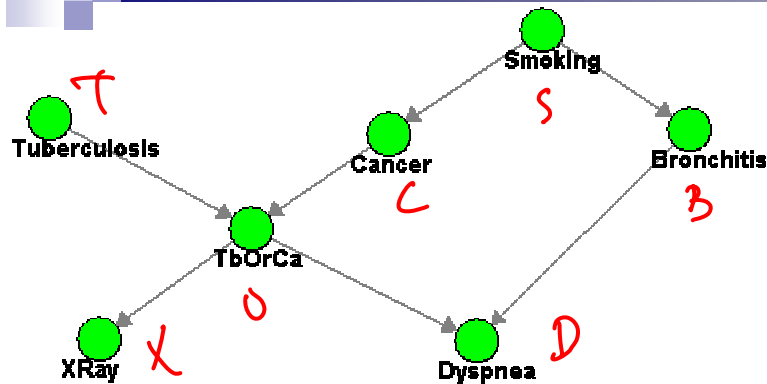
V.E. works on any graph! But, complexity depends on graph structure

table of k^3 values

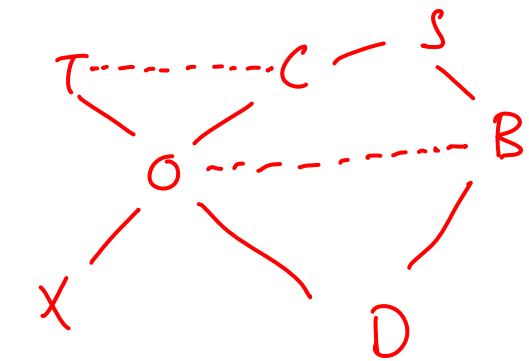
inference is #P-complete, even if nodes have small # of parents " \Rightarrow "
V.E. will generate large factors in some graphs

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

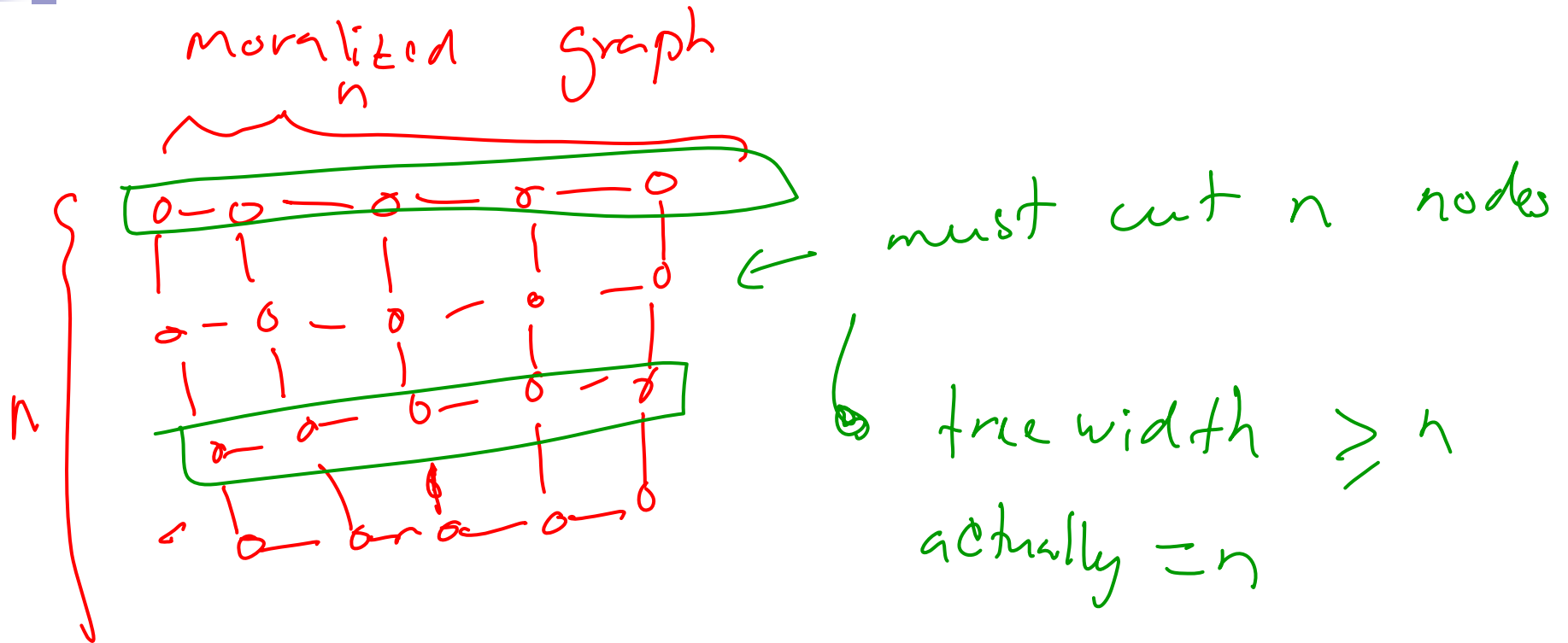


tree width of
this graph

↳ "max minimum
node cut between
pairs of vars"
1
none adjacent
sets vars

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



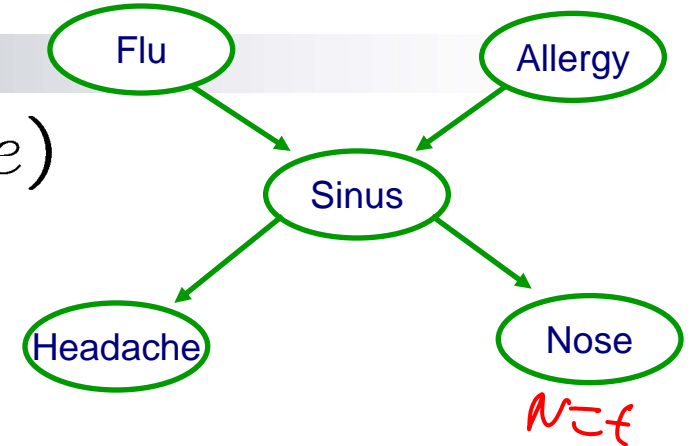
Compact representation \nRightarrow Easy inference ☹

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
- see readings* →

Most likely explanation (MLE)

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



■ Using Bayes rule:

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

constant

■ 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



Forward pass

$$\max_F \max_S P(f) P(s|f) \cdot P(N=t|s)$$

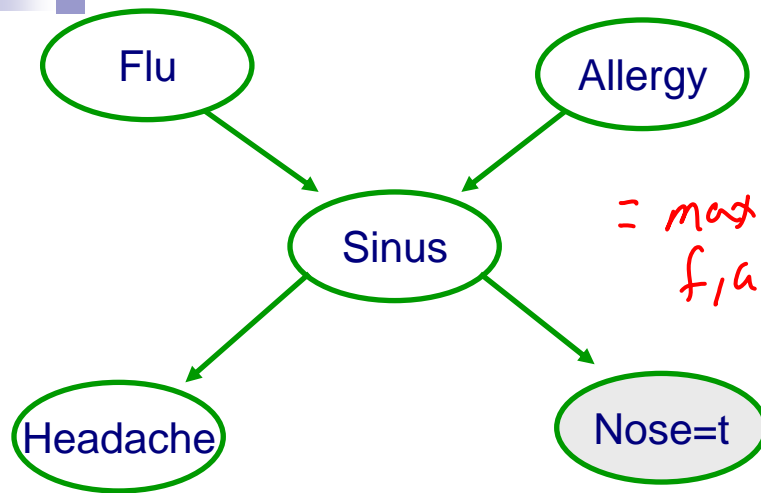
$$\max_f P(f) \underbrace{\max_s P(s|f) \cdot P(N=t|s)}_{g_1(f)}$$

$$\max_f P(f) \cdot g_1(f) = g_2 \leftarrow \text{constant} \leftarrow \text{prob. of best assignment}$$

Backward pass

$$f^* = \underset{f}{\operatorname{argmax}} P(f) \cdot g_1(f) \quad \Bigg| \quad s^* = \underset{s}{\operatorname{argmax}} P(s|f^*) \cdot P(N=t|s)$$

Example of variable elimination for MLE – Forward pass



$$\max_{f,a,s,h} P(f) \cdot P(a) \cdot P(s|f,a) \cdot P(h|s) \cdot P(n=t|s)$$

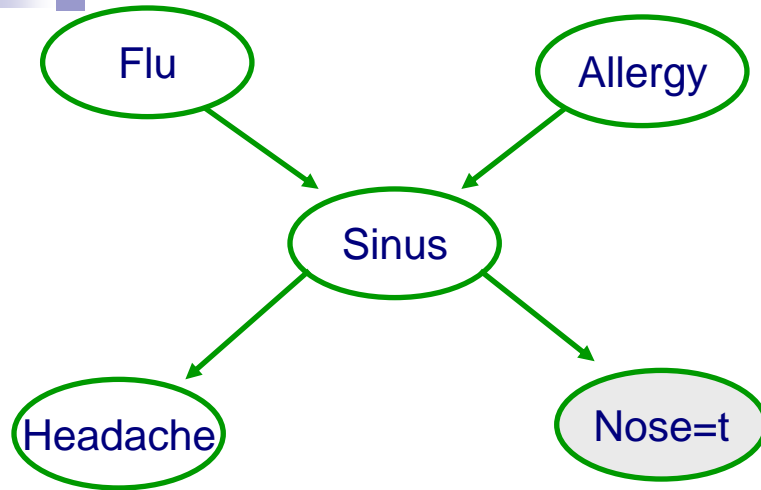
$$= \max_{f,a,s} P(f) \cdot P(a) \cdot P(s|f,a) \cdot P(n=t|s) \max_s P(h|s)$$

$$= \max_{f,a} P(f) \cdot P(a) \max_s P(s|f,a) \cdot g_1(s) \cdot P(n=t|s)$$

$$g_2(f,a)$$

then eliminate f,a

Example of variable elimination for MLE – Backward pass



$$(f^*, a^*) = \underset{f, a}{\operatorname{argmax}} P(f). P(a) g_2(f, a)$$

$$s^* = \underset{s}{\operatorname{argmax}} P(s|f^*, a^*) \cdot P(N=t|s) \cdot g_1(s)$$

$$h^* = \underset{h}{\operatorname{argmax}} P(h|s^*)$$

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

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

maximization (pointing to \max_{x_i}) *multiplication* (pointing to $\prod_{j=1}^k$)

- Variable X_i has been eliminated!

remember / cache g's

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”



HMMs

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

March 26th, 2007

©2005-2007 Carlos Guestrin

Adventures of our BN hero

- Compact representation for probability distributions
- Fast inference
- Fast learning

1. Naïve Bayes



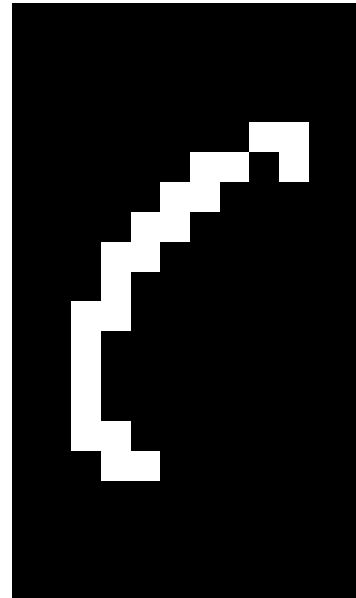
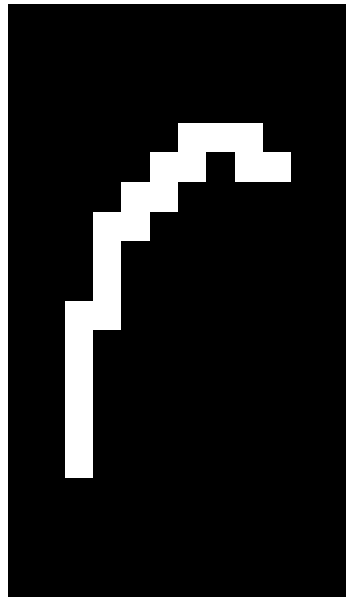
- But... Who are the most popular kids?

2 and 3.

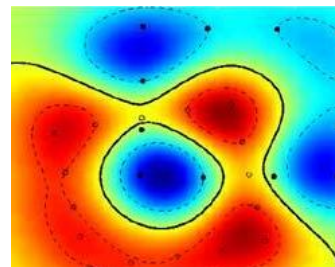
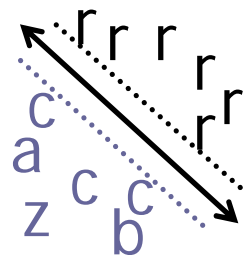
Hidden Markov models (HMMs)
Kalman Filters

→ a ~~Kalman~~ HMM with Gaussians
→ Kalman Filter with discrete dist.

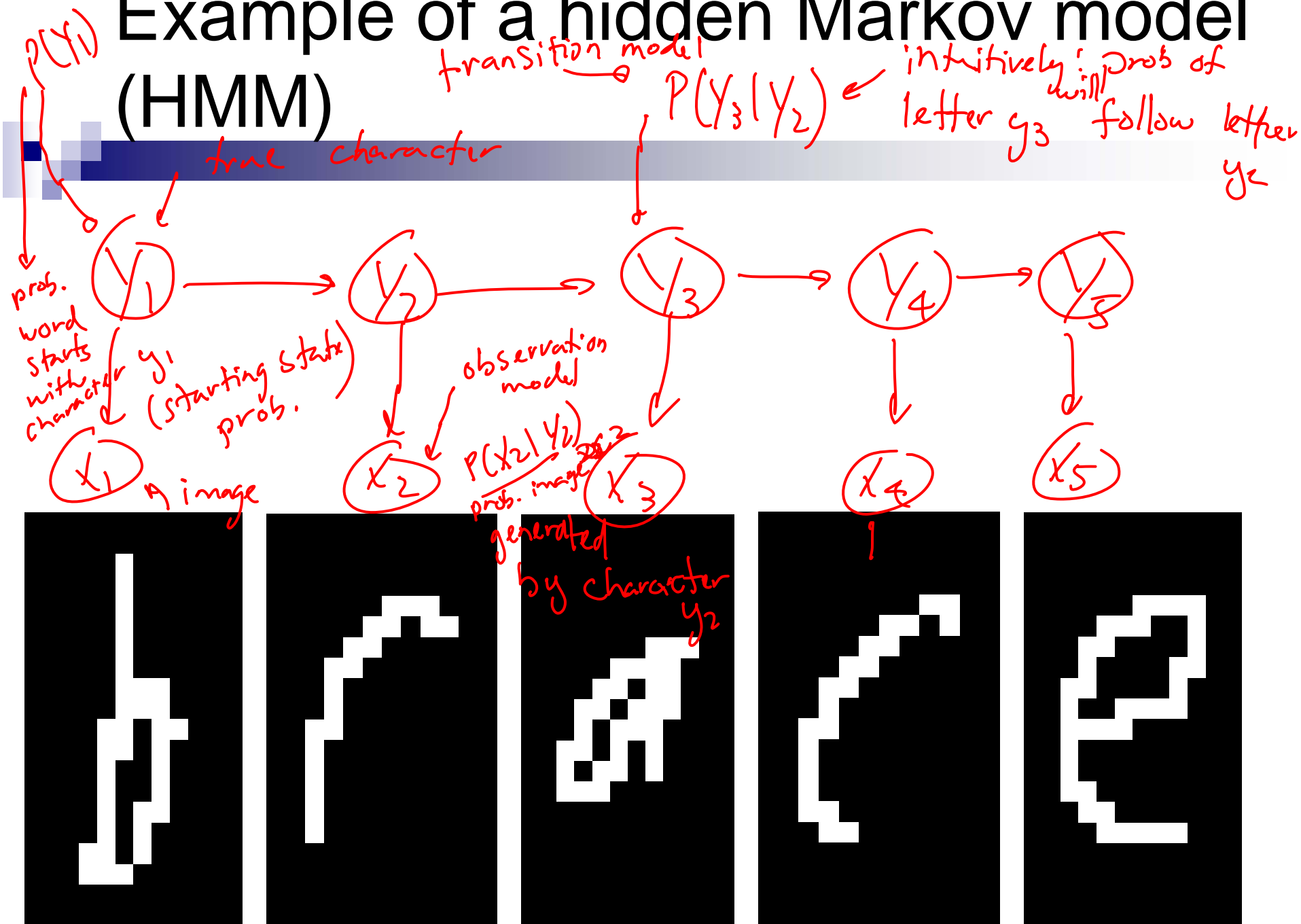
Handwriting recognition



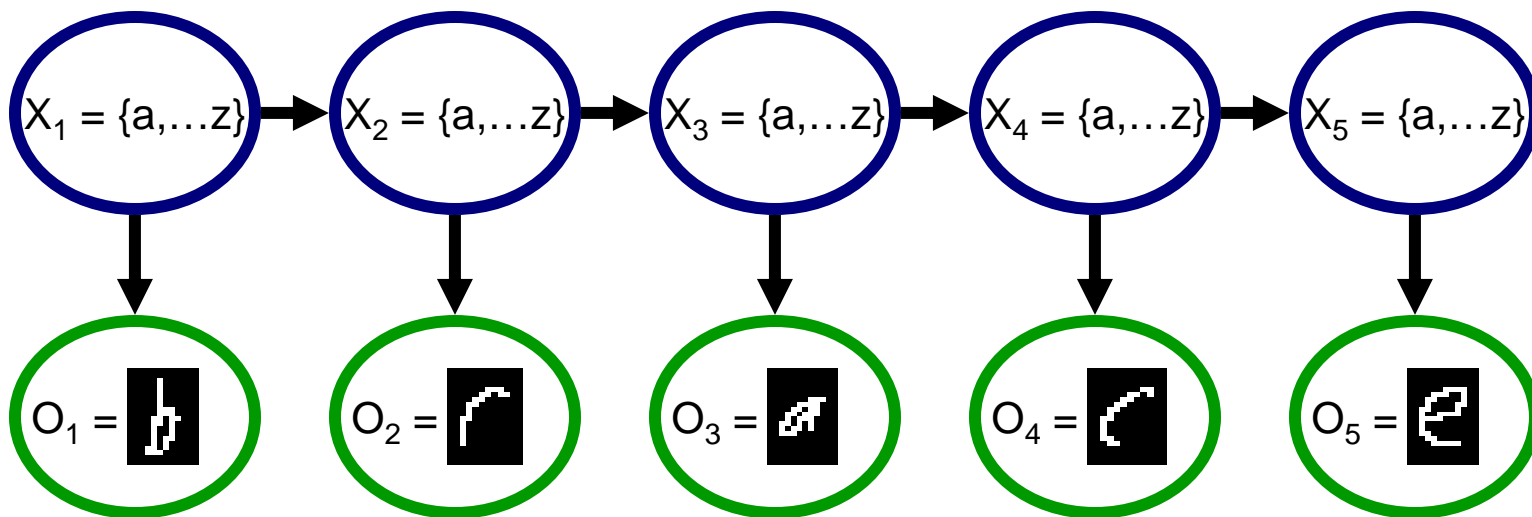
Character recognition, e.g., kernel SVMs



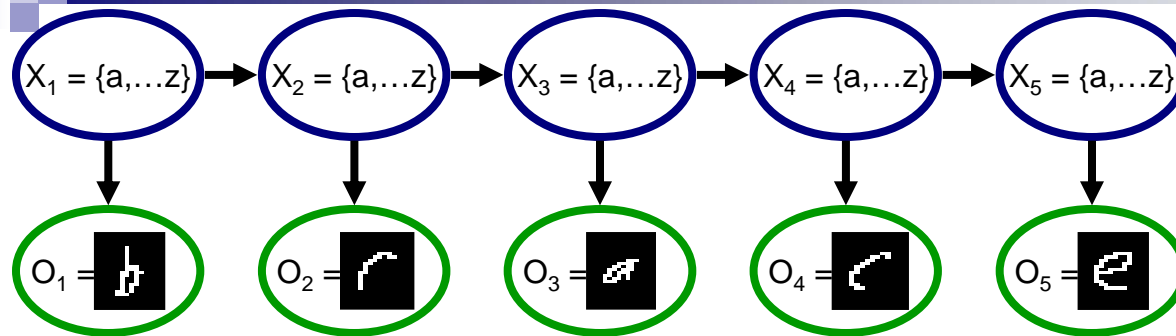
Example of a hidden Markov model (HMM)



Understanding the HMM Semantics



HMMs semantics: Details



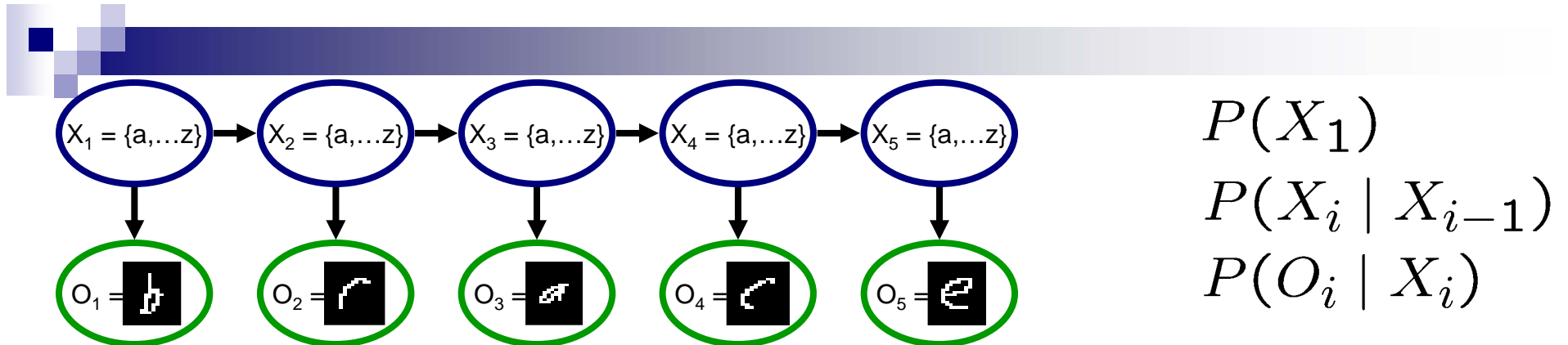
Just 3 distributions:

$$P(X_1)$$

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

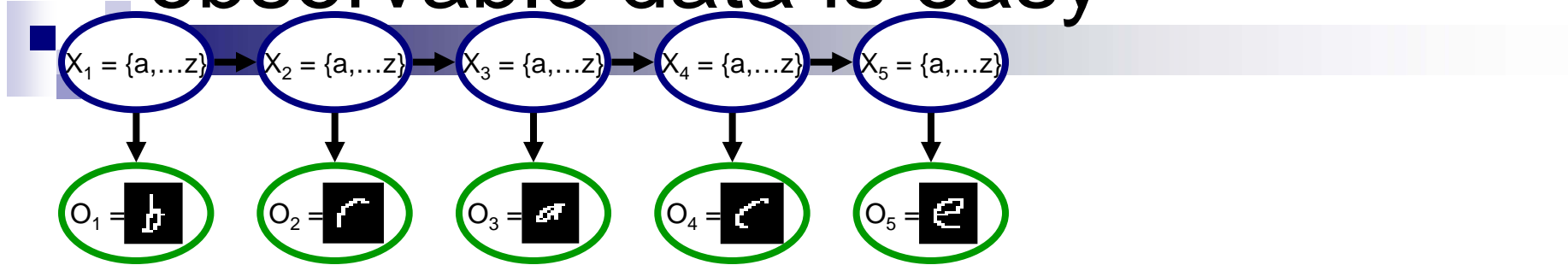
$$P(O_i \mid X_i)$$

HMMs semantics: Joint distribution



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

Learning HMMs from fully observable data is easy



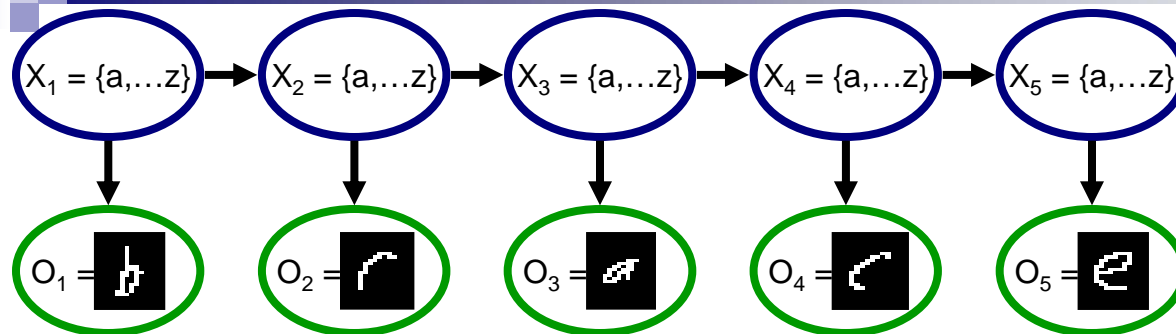
Learn 3 distributions:

$$P(X_1)$$

$$P(O_i \mid X_i)$$

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

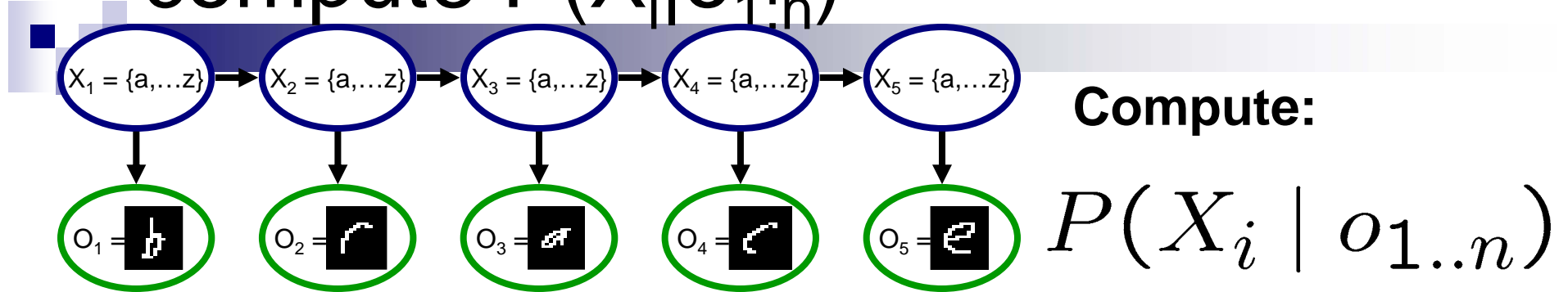
Possible inference tasks in an HMM



Marginal probability of a hidden variable:

Viterbi decoding – most likely trajectory for hidden vars:

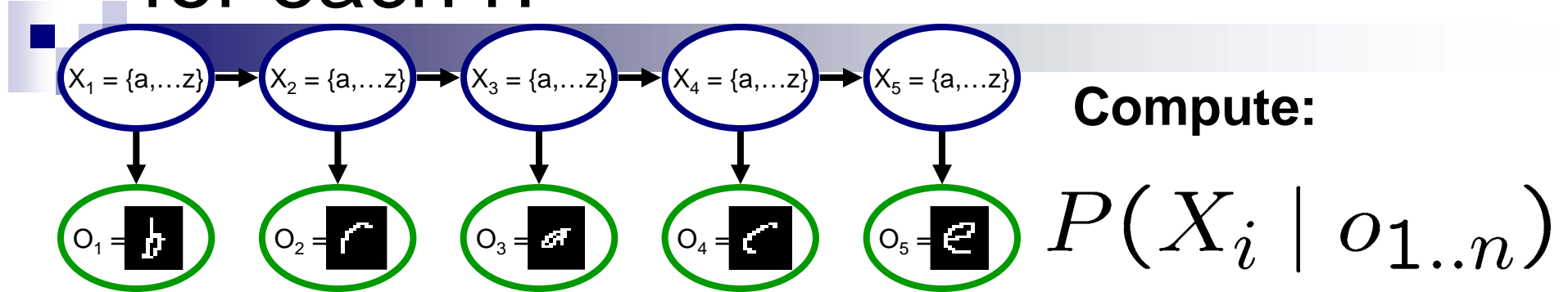
Using variable elimination to compute $P(X_i | o_{1:n})$



Variable elimination order?

Example:

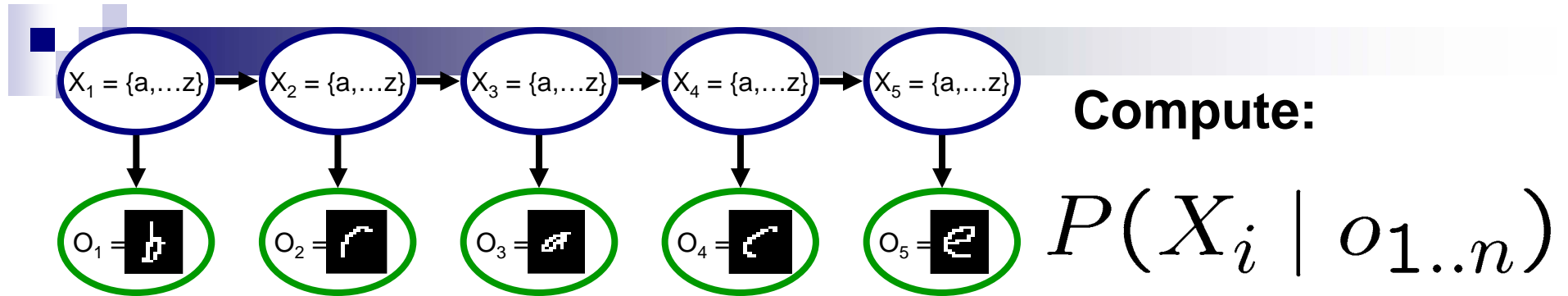
What if I want to compute $P(X_i | o_{1:n})$
for each i ?



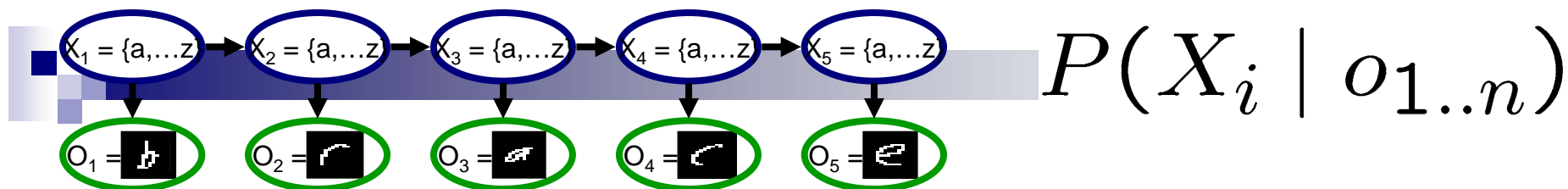
Variable elimination for each i ?

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

Reusing computation



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: $\beta_n(X_n) = 1$

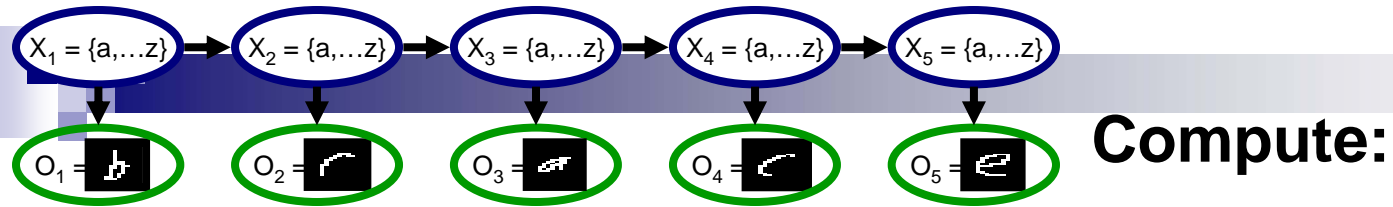
■ For $i = n-1$ to 1

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

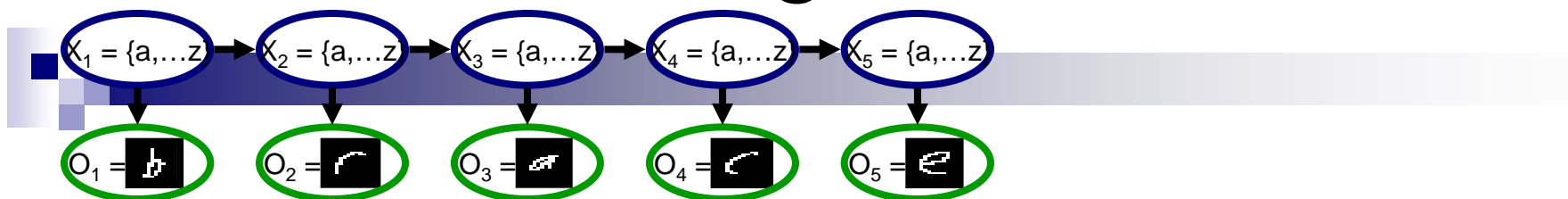
Most likely explanation



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 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

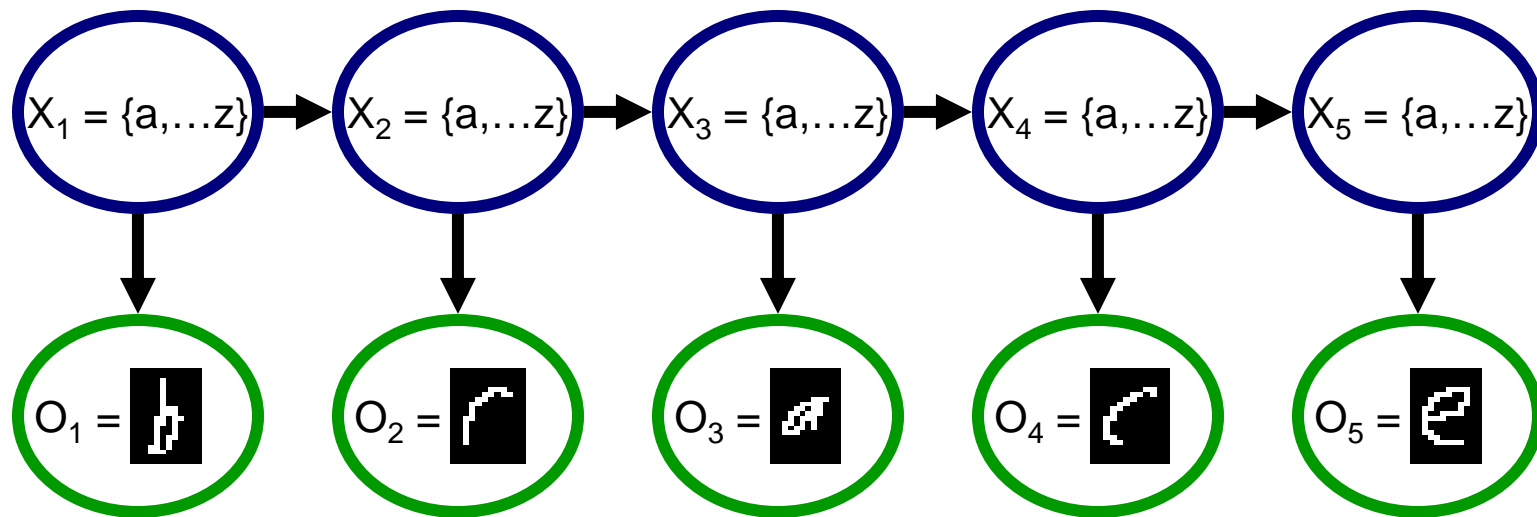
$$\alpha_i(X_i) = \max_{x_{i-1}} P(o_i | X_i) P(X_i | X_{i-1} = x_{i-1}) \alpha_{i-1}(x_{i-1})$$

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