

Probability and Structure in Natural Language Processing

Noah Smith, Carnegie Mellon University

2012 International Summer School in
Language and Speech Technologies

Quick Recap

- Yesterday:
 - Bayesian networks and some formal properties
 - Markov networks and some formal properties
 - Exact marginal inference using variable elimination
 - Sum-product version
 - Beginnings of the max-product version

Variable Elimination for Conditional Probabilities

Input: Graphical model on \mathbf{V} , set of query variables \mathbf{Y} , evidence $\mathbf{X} = \mathbf{x}$

Output: factor φ and scalar α

1. Φ = factors in the model
2. Reduce factors in Φ by $\mathbf{X} = \mathbf{x}$
3. Choose variable ordering on $\mathbf{Z} = \mathbf{V} \setminus \mathbf{Y} \setminus \mathbf{X}$
4. $\varphi = \text{Variable-Elimination}(\Phi, \mathbf{Z})$
5. $\alpha = \sum_{\mathbf{z} \in \text{Val}(\mathbf{Z})} \varphi(\mathbf{z})$
6. Return φ, α

From Marginals to MAP

- Replace factor marginalization steps with *maximization*.
 - Add bookkeeping to keep track of the maximizing values.
- Add a traceback at the end to recover the solution.
- This is analogous to the connection between the forward algorithm and the Viterbi algorithm.
 - Ordering challenge is the same.

Factor Maximization

- Given \mathbf{X} and Y ($Y \notin \mathbf{X}$), we can turn a factor $\phi(\mathbf{X}, Y)$ into a factor $\psi(\mathbf{X})$ via maximization:

$$\psi(\mathbf{X}) = \max_Y \phi(\mathbf{X}, Y)$$

- We can refer to this new factor by $\max_Y \phi$.

Factor Maximization

- Given \mathbf{X} and Y ($Y \notin \mathbf{X}$), we can turn a factor $\phi(\mathbf{X}, Y)$ into a factor $\psi(\mathbf{X})$ via maximization:

$$\psi(\mathbf{X}) = \max_Y \phi(\mathbf{X}, Y)$$

A	B	C	$\phi(A, B, C)$
0	0	0	0.9
0	0	1	0.3
0	1	0	1.1
0	1	1	1.7
1	0	0	0.4
1	0	1	0.7
1	1	0	1.1
1	1	1	0.2



“maximizing out” B

A	C	$\psi(A, C)$	
0	0	1.1	B=1
0	1	1.7	B=1
1	0	1.1	B=1
1	1	0.7	B=0

Distributive Property

- A useful property we exploited in variable elimination:

$$X \notin \text{Scope}(\phi_1) \Rightarrow \sum_X (\phi_1 \cdot \phi_2) = \phi_1 \cdot \sum_X \phi_2$$

- Under the same conditions, factor multiplication distributes over max, too:

$$\max_X (\phi_1 \cdot \phi_2) = \phi_1 \cdot \max_X \phi_2$$

Traceback

Input: Sequence of factors with associated variables: $(\psi_{Z_1}, \dots, \psi_{Z_k})$

Output: \mathbf{z}^*

- Each ψ_Z is a factor with scope including Z and variables eliminated *after* Z .
- Work backwards from $i = k$ to 1:
 - Let $z_i = \arg \max_z \psi_{Z_i}(z, z_{i+1}, z_{i+2}, \dots, z_k)$
- Return \mathbf{z}

About the Traceback

- No extra (asymptotic) expense.
 - Linear traversal over the intermediate factors.
- The factor operations for both sum-product VE and max-product VE can be generalized.
 - Example: get the K most likely assignments

Eliminating One Variable (Max-Product Version with Bookkeeping)

Input: Set of factors Φ , variable Z to eliminate

Output: new set of factors Ψ

1. Let $\Phi' = \{\varphi \in \Phi \mid Z \in \text{Scope}(\varphi)\}$
2. Let $\Psi = \{\varphi \in \Phi \mid Z \notin \text{Scope}(\varphi)\}$
3. Let τ be $\max_Z \prod_{\varphi \in \Phi'} \varphi$
 - Let ψ be $\prod_{\varphi \in \Phi'} \varphi$ (bookkeeping)
4. Return $\Psi \cup \{\tau\}, \psi$

Variable Elimination

(Max-Product Version with Decoding)

Input: Set of factors Φ , ordered list of variables Z to eliminate

Output: new factor

1. For each $Z_i \in Z$ (in order):
 - Let $(\Phi, \psi_{Z_i}) = \text{Eliminate-One}(\Phi, Z_i)$
2. Return $\prod_{\varphi \in \Phi} \varphi, \text{Traceback}(\{\psi_{Z_i}\})$

Variable Elimination Tips

- Any ordering will be correct.
- Most orderings will be too expensive.
- There are heuristics for choosing an ordering (you are welcome to find them and test them out).

(Rocket Science: True MAP)

- Evidence: $\mathbf{X} = \mathbf{x}$
- Query: \mathbf{Y}
- Other variables: $\mathbf{Z} = \mathbf{V} \setminus \mathbf{X} \setminus \mathbf{Y}$

$$\begin{aligned} \mathbf{y}^* &= \arg \max_{\mathbf{y} \in \text{Val}(\mathbf{Y})} P(\mathbf{Y} = \mathbf{y} \mid \mathbf{X} = \mathbf{x}) \\ &= \arg \max_{\mathbf{y} \in \text{Val}(\mathbf{Y})} \sum_{\mathbf{z} \in \text{Val}(\mathbf{Z})} P(\mathbf{Y} = \mathbf{y}, \mathbf{Z} = \mathbf{z} \mid \mathbf{X} = \mathbf{x}) \end{aligned}$$

- First, marginalize out \mathbf{Z} , then do MAP inference over \mathbf{Y} given $\mathbf{X} = \mathbf{x}$
- This is not usually attempted in NLP, with some exceptions.

Parting Shots

- You will probably never implement the general variable elimination algorithm.
- You will rarely use exact inference.
- There is value in understanding the *problem* that approximation methods are trying to solve, and what an *exact* (if intractable) solution would look like!

Lecture 3:

Structures and Decoding

Two Meanings of “Structure”

- Yesterday: structure of a **graph** for modeling a collection of random variables together.
- Today: **linguistic** structure.
 - Sequence labelings (POS, IOB chunkings, ...)
 - Parse trees (phrase-structure, dependency, ...)
 - Alignments (word, phrase, tree, ...)
 - Predicate-argument structures
 - Text-to-text (translation, paraphrase, answers, ...)

A Useful Abstraction?

- I think so.
- Brings out commonalities:
 - Modeling formalisms (e.g., linear models with features)
 - Learning algorithms (lectures 4-6)
 - Generic inference algorithms
- Permits sharing across a wider space of problems.
- Disadvantage: hides engineering details.

Familiar Example: Hidden Markov Models

Hidden Markov Model

- **X** and **Y** are both sequences of symbols
 - **X** is a sequence from the vocabulary Σ
 - **Y** is a sequence from the state space Λ

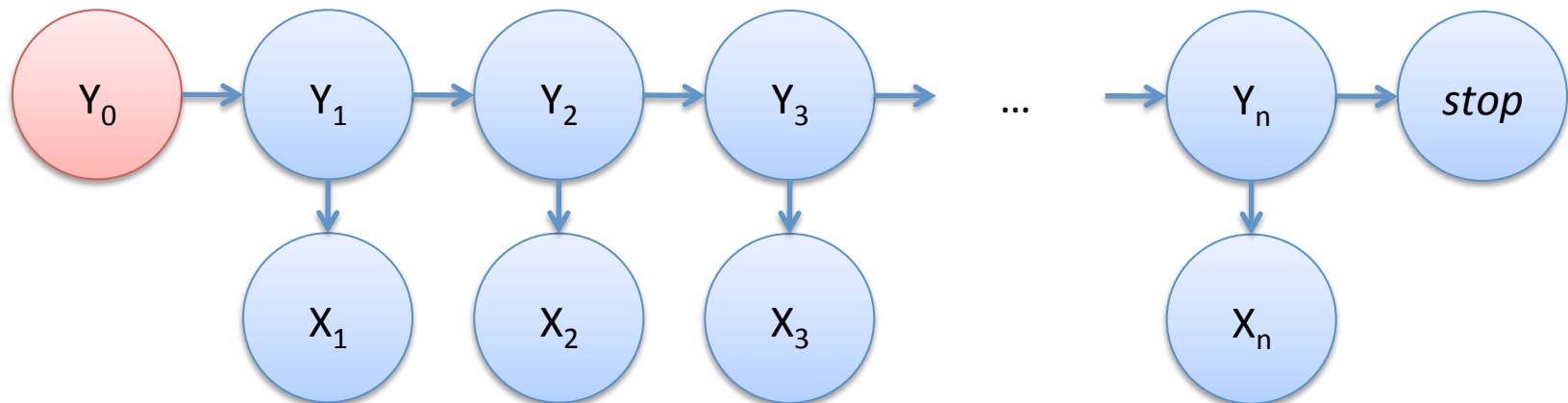
$$p(\mathbf{X} = \mathbf{x}, \mathbf{Y} = \mathbf{y}) = \left(\prod_{i=1}^n p(x_i | y_i) p(y_i | y_{i-1}) \right) p(\text{stop} | y_n)$$

- Parameters:
 - Transitions $p(y' | y)$
 - including $p(\text{stop} | y)$, $p(y | \text{start})$
 - Emissions $p(x | y)$

Hidden Markov Model

- The joint model's independence assumptions are easy to capture with a Bayesian network.

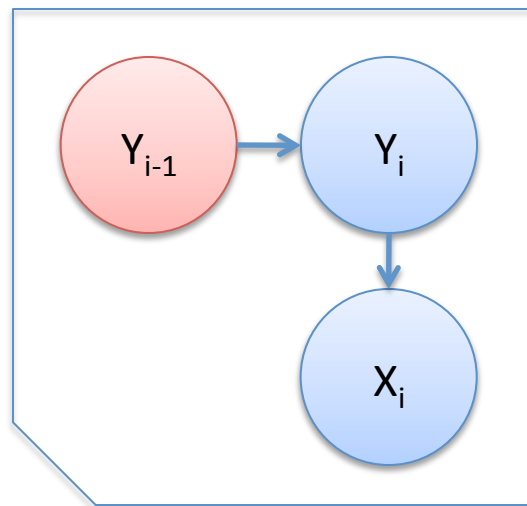
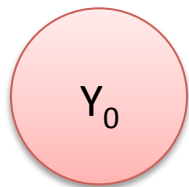
$$p(\mathbf{X} = \mathbf{x}, \mathbf{Y} = \mathbf{y}) = \left(\prod_{i=1}^n p(x_i | y_i) p(y_i | y_{i-1}) \right) p(\text{stop} | y_n)$$



Hidden Markov Model

- The joint model instantiates **dynamic Bayesian networks**.

$$p(\mathbf{X} = \mathbf{x}, \mathbf{Y} = \mathbf{y}) = \left(\prod_{i=1}^n p(x_i | y_i) p(y_i | y_{i-1}) \right) p(stop | y_n)$$

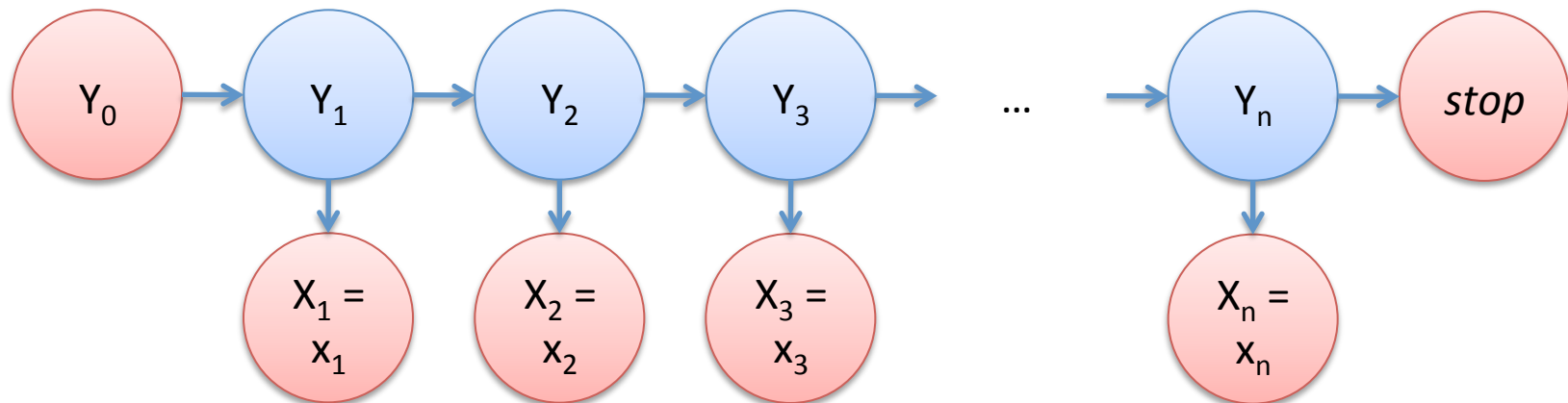


template that
gets copied as
many times as
needed

Hidden Markov Model

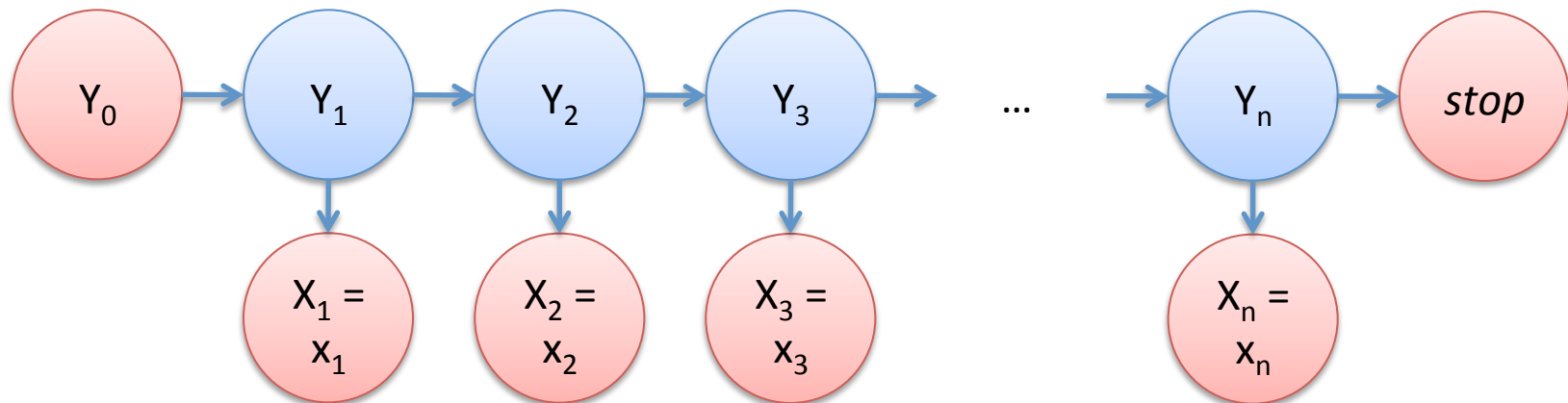
- Given \mathbf{X}' 's value as evidence, the dynamic part becomes unnecessary, since we know n .

$$p(\mathbf{X} = \mathbf{x}, \mathbf{Y} = \mathbf{y}) = \left(\prod_{i=1}^n p(x_i | y_i) p(y_i | y_{i-1}) \right) p(\text{stop} | y_n)$$



Hidden Markov Model

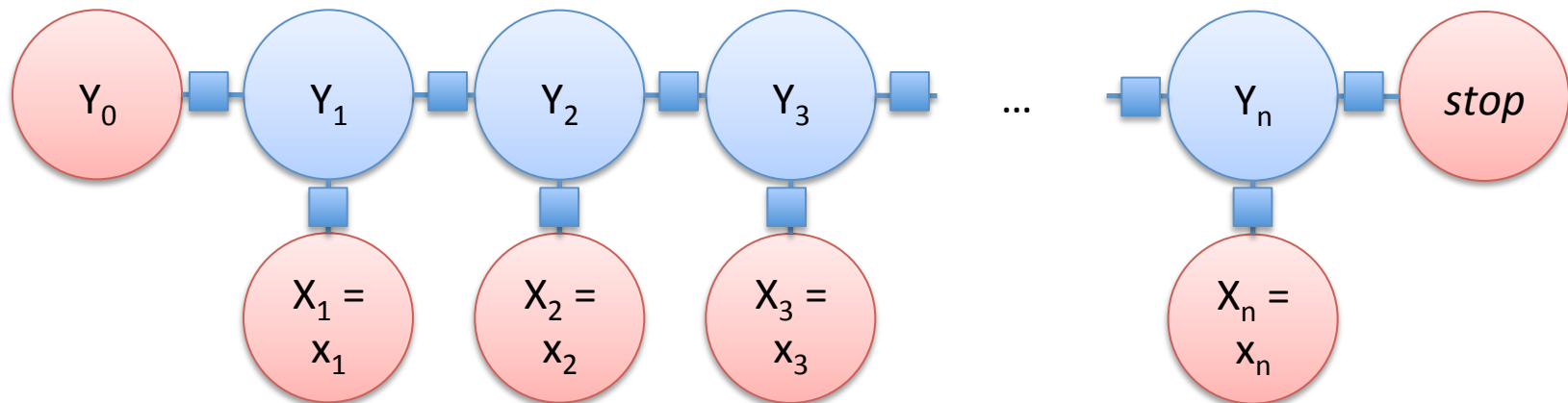
- The usual inference problem is to find the most probable value of \mathbf{Y} given $\mathbf{X} = \mathbf{x}$.



Hidden Markov Model

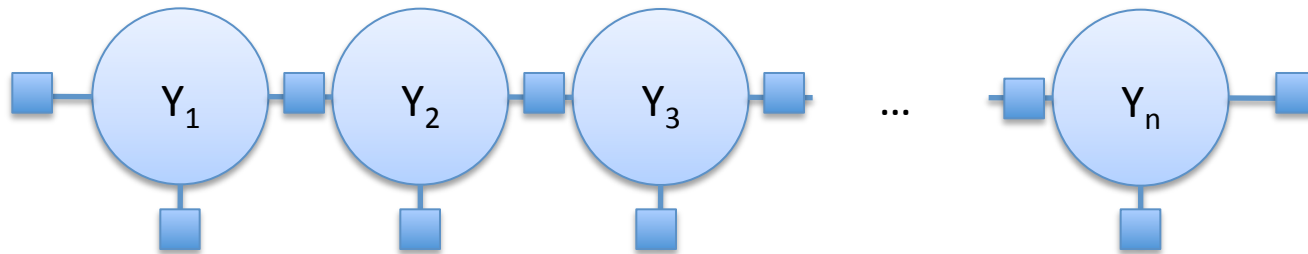
- The usual inference problem is to find the most probable value of \mathbf{Y} given $\mathbf{X} = \mathbf{x}$.

- Factor graph:



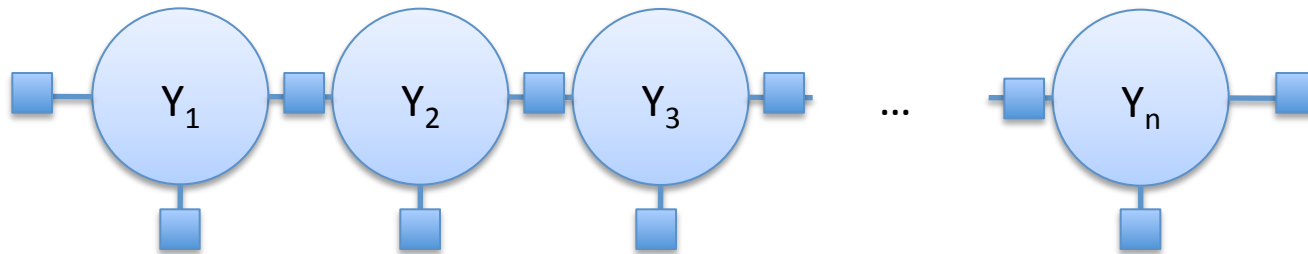
Hidden Markov Model

- The usual inference problem is to find the most probable value of \mathbf{Y} given $\mathbf{X} = \mathbf{x}$.
- Factor graph after reducing factors to respect evidence:



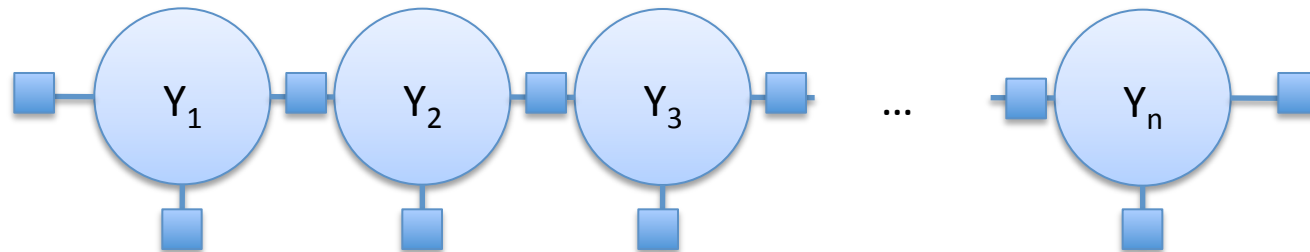
Hidden Markov Model

- The usual inference problem is to find the most probable value of \mathbf{Y} given $\mathbf{X} = \mathbf{x}$.
- Clever ordering should be apparent!



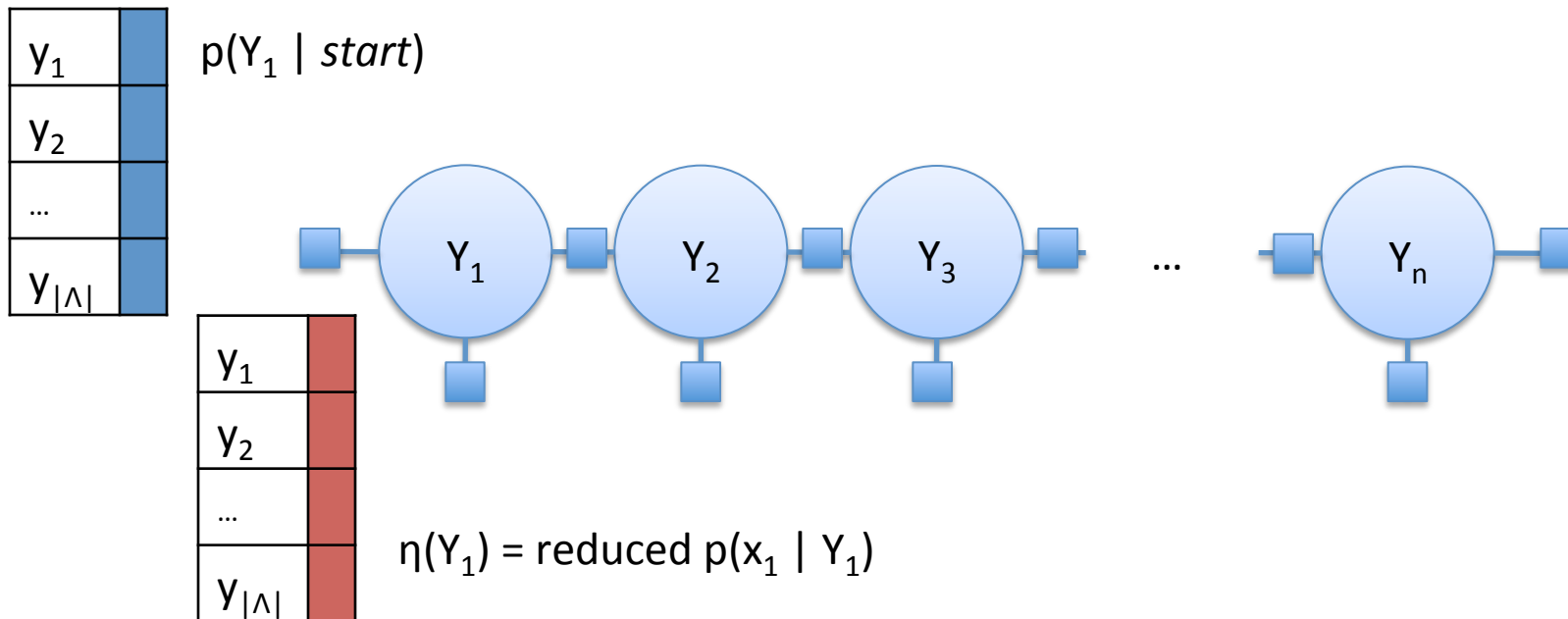
Hidden Markov Model

- When we eliminate Y_1 , we take a product of three relevant factors.
 - $p(Y_1 \mid \text{start})$
 - $\eta(Y_1) = \text{reduced } p(x_1 \mid Y_1)$
 - $p(Y_2 \mid Y_1)$



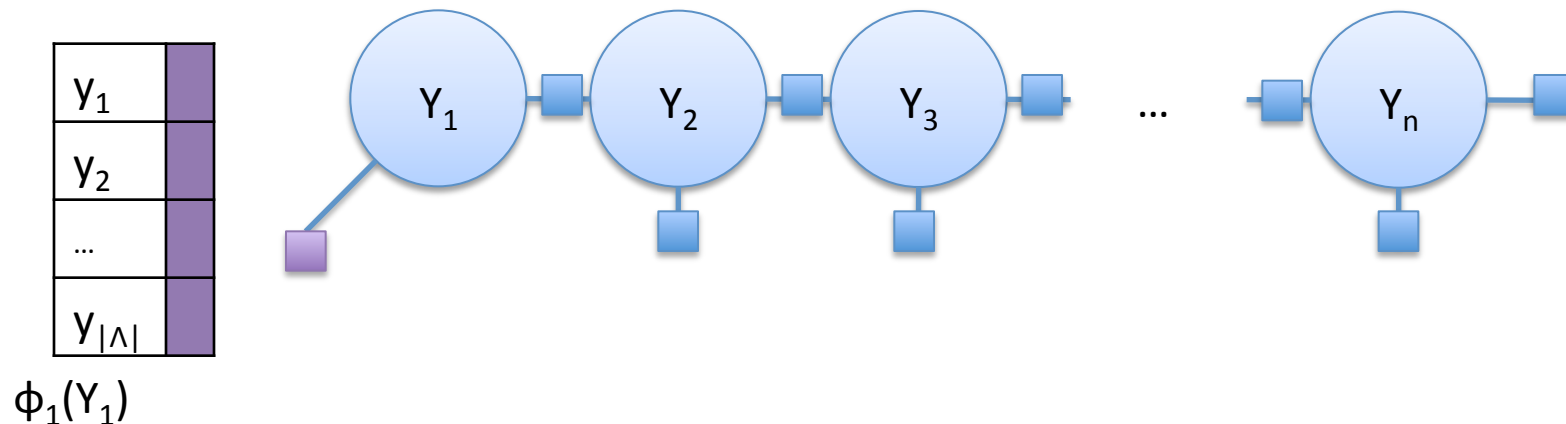
Hidden Markov Model

- When we eliminate Y_1 , we first take a product of two factors that only involve Y_1 .



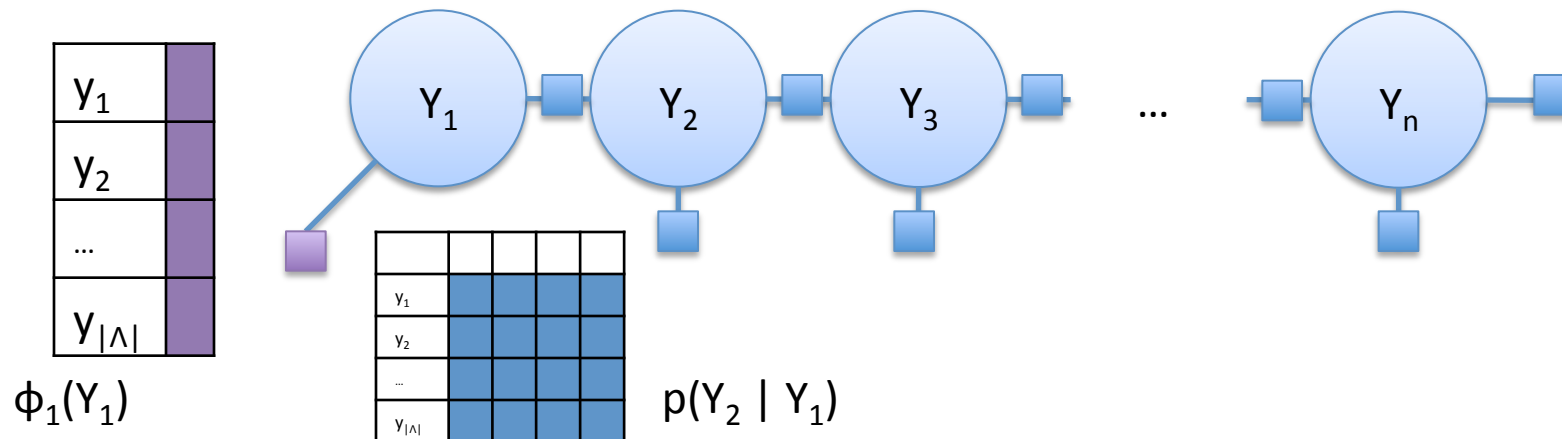
Hidden Markov Model

- When we eliminate Y_1 , we first take a product of two factors that only involve Y_1 .
- This is the Viterbi probability vector for Y_1 .



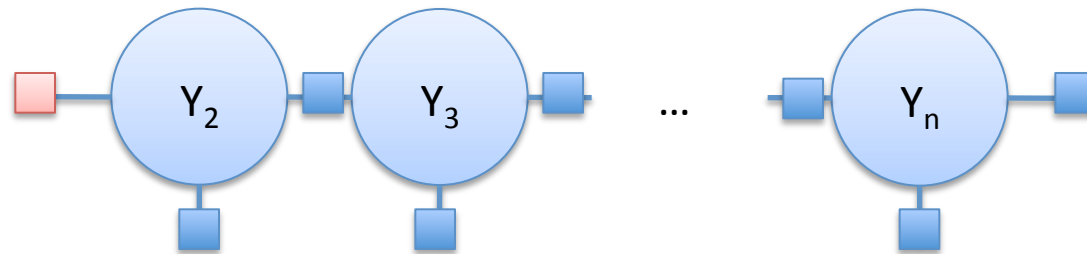
Hidden Markov Model

- When we eliminate Y_1 , we first take a product of two factors that only involve Y_1 .
- This is the Viterbi probability vector for Y_1 .
- Eliminating Y_1 equates to solving the Viterbi probabilities for Y_2 .



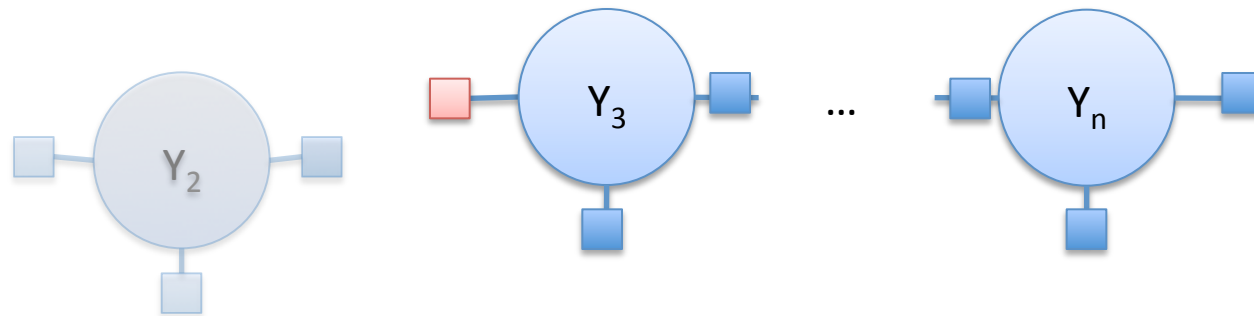
Hidden Markov Model

- Product of all factors involving Y_1 , then reduce.
 - $\phi_2(Y_2) = \max_{y \in \text{Val}(Y_1)} (\phi_1(y) \times p(Y_2 \mid y))$
 - This factor holds Viterbi probabilities for Y_2 .



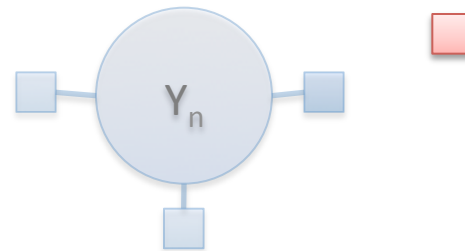
Hidden Markov Model

- When we eliminate Y_2 , we take a product of the analogous two relevant factors.
- Then reduce.
 - $\phi_3(Y_3) = \max_{y \in \text{Val}(Y_2)} (\phi_2(y) \times p(Y_3 \mid y))$



Hidden Markov Model

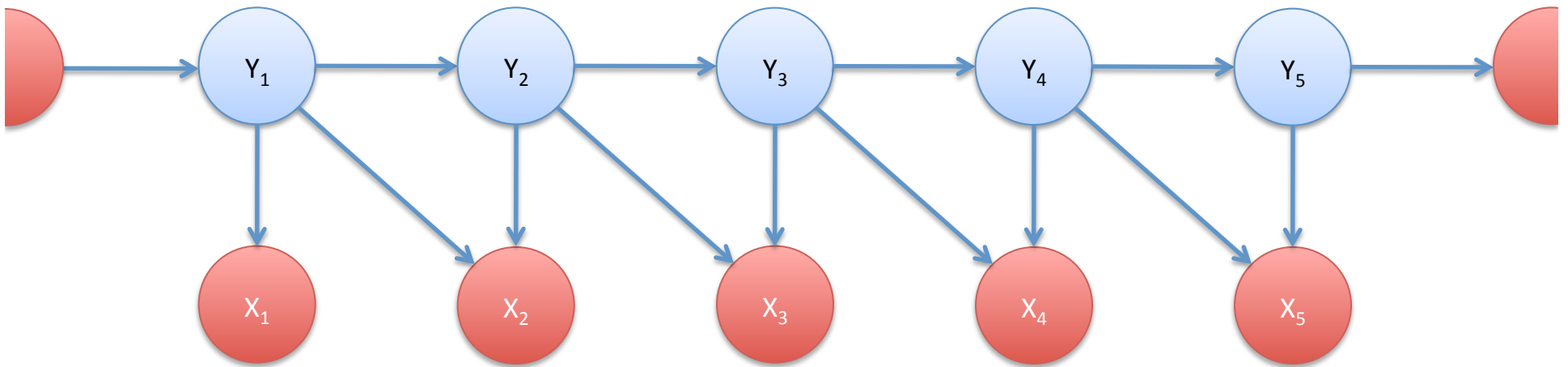
- At the end, we have one final factor with one row, ϕ_{n+1} .
- This is the score of the best sequence.
- Use backtrack to recover values.



Why Think This Way?

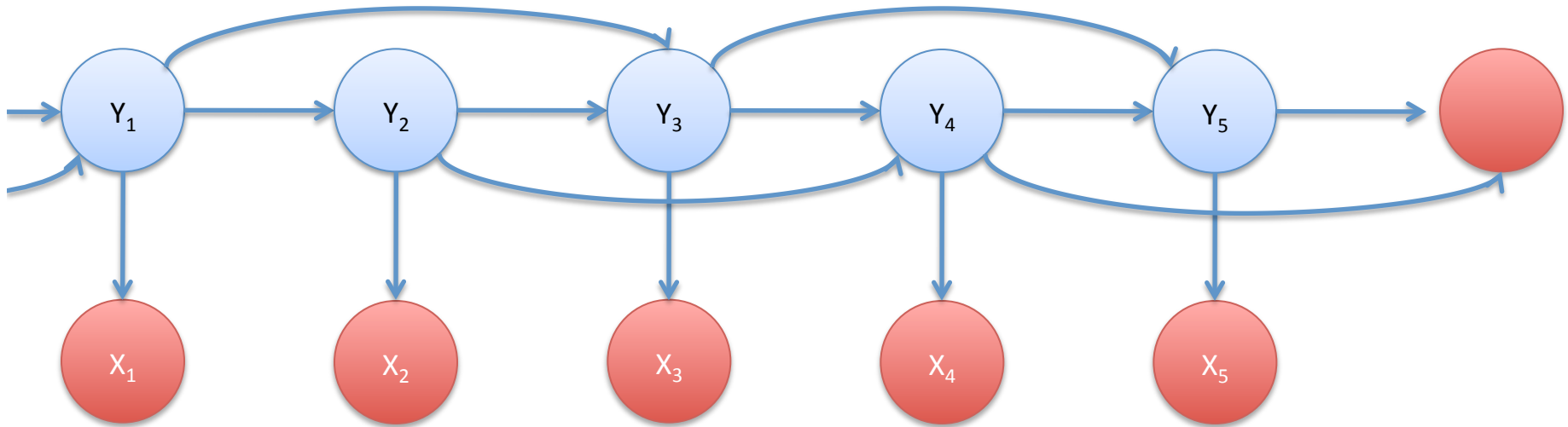
- Easy to see how to generalize HMMs.
 - More evidence
 - More factors
 - More hidden structure
 - More dependencies
- Probabilistic interpretation of factors is *not* central to finding the “best” \mathbf{Y} ...
 - Many factors are not conditional probability tables.

Generalization Example 1



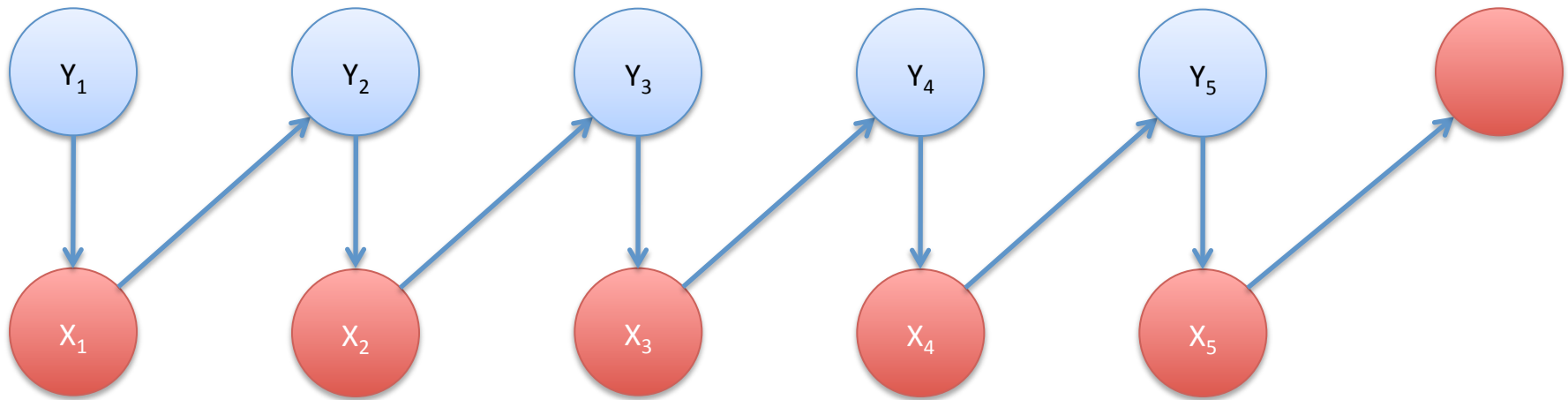
- Each word also depends on previous state.

Generalization Example 2



- “Trigram” HMM

Generalization Example 3



- Aggregate bigram model (Saul and Pereira, 1997)

General Decoding Problem

- Two structured random variables, \mathbf{X} and \mathbf{Y} .
 - Sometimes described as *collections* of random variables.
- “Decode” observed value $\mathbf{X} = \mathbf{x}$ into some value of \mathbf{Y} .
- Usually, we seek to maximize some score.
 - E.g., MAP inference from yesterday.

Linear Models

- Define a feature vector function \mathbf{g} that maps (\mathbf{x}, \mathbf{y}) pairs into d-dimensional real space.
- Score is linear in $\mathbf{g}(\mathbf{x}, \mathbf{y})$.

$$\begin{aligned} \text{score}(\mathbf{x}, \mathbf{y}) &= \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}) \\ \mathbf{y}^* &= \arg \max_{\mathbf{y} \in \mathcal{Y}_{\mathbf{x}}} \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}) \end{aligned}$$

- Results:
 - **decoding** seeks \mathbf{y} to maximize the score.
 - **learning** seeks \mathbf{w} to ... do something we'll talk about later.
- Extremely general!

Generic Noisy Channel as Linear Model

$$\begin{aligned}\hat{\mathbf{y}} &= \arg \max_{\mathbf{y}} \log (p(\mathbf{y}) \cdot p(\mathbf{x} \mid \mathbf{y})) \\ &= \arg \max_{\mathbf{y}} \log p(\mathbf{y}) + \log p(\mathbf{x} \mid \mathbf{y}) \\ &= \arg \max_{\mathbf{y}} w_{\mathbf{y}} + w_{\mathbf{x}|\mathbf{y}} \\ &= \arg \max_{\mathbf{y}} \mathbf{w}^{\top} \mathbf{g}(\mathbf{x}, \mathbf{y})\end{aligned}$$

- Of course, the two probability terms are typically composed of “smaller” factors; each can be understood as an exponentiated weight.

Max Ent Models as Linear Models

$$\begin{aligned}\hat{\boldsymbol{y}} &= \arg \max_{\boldsymbol{y}} \log p(\boldsymbol{y} \mid \boldsymbol{x}) \\ &= \arg \max_{\boldsymbol{y}} \log \frac{\exp \boldsymbol{w}^\top \boldsymbol{g}(\boldsymbol{x}, \boldsymbol{y})}{z(\boldsymbol{x})} \\ &= \arg \max_{\boldsymbol{y}} \boldsymbol{w}^\top \boldsymbol{g}(\boldsymbol{x}, \boldsymbol{y}) - \log z(\boldsymbol{x}) \\ &= \arg \max_{\boldsymbol{y}} \boldsymbol{w}^\top \boldsymbol{g}(\boldsymbol{x}, \boldsymbol{y})\end{aligned}$$

HMMs as Linear Models

$$\begin{aligned}\hat{\mathbf{y}} &= \arg \max_{\mathbf{y}} \log p(\mathbf{x}, \mathbf{y}) \\&= \arg \max_{\mathbf{y}} \left(\sum_{i=1}^n \log p(x_i \mid y_i) + \log p(y_i \mid y_{i-1}) \right) + \log p(\text{stop} \mid y_n) \\&= \arg \max_{\mathbf{y}} \left(\sum_{i=1}^n w_{y_i \downarrow x_i} + w_{y_{i-1} \rightarrow y_i} \right) + w_{y_n \rightarrow \text{stop}} \\&= \arg \max_{\mathbf{y}} \sum_{y, x} w_{y \downarrow x} \text{freq}(y \downarrow x; \mathbf{y}, \mathbf{x}) + \sum_{y, y'} w_{y \rightarrow y'} \text{freq}(y \rightarrow y'; \mathbf{y}) \\&= \arg \max_{\mathbf{y}} \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y})\end{aligned}$$

Running Example

		1	2	3	4	5	6	7	8	9	10
x	=	Britain	sent	warships	across	the	English	Channel	Monday	to	rescue
y	=	B	O	O	O	O	B	I	B	O	O
y'	=	O	O	O	O	O	B	I	B	O	O

		11	12	13	14	15	16	17	18	19	20
		Britons	stranded	by	Eyjafjallajökull	's	volcanic	ash	cloud	.	
		B	O	O	B	O	O	O	O	O	O
		B	O	O	B	O	O	O	O	O	O

- IOB sequence labeling, here applied to NER
- Often solved with HMMs, CRFs, M³Ns ...

feature function $g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$		$g(\mathbf{x}, \mathbf{y})$	$g(\mathbf{x}, \mathbf{y}')$
<i>bias:</i>	count of i s.t. $y_i = B$	5	4
	count of i s.t. $y_i = I$	1	1
	count of i s.t. $y_i = O$	14	15
<i>lexical:</i>	count of i s.t. $x_i = \textit{Britain}$ and $y_i = B$	1	0
	count of i s.t. $x_i = \textit{Britain}$ and $y_i = I$	0	0
	count of i s.t. $x_i = \textit{Britain}$ and $y_i = O$	0	1
<i>downcased:</i>	count of i s.t. $lc(x_i) = \textit{britain}$ and $y_i = B$	1	0
	count of i s.t. $lc(x_i) = \textit{britain}$ and $y_i = I$	0	0
	count of i s.t. $lc(x_i) = \textit{britain}$ and $y_i = O$	0	1
	count of i s.t. $lc(x_i) = \textit{sent}$ and $y_i = O$	1	1
	count of i s.t. $lc(x_i) = \textit{warships}$ and $y_i = O$	1	1
<i>shape:</i>	count of i s.t. $shape(x_i) = Aaaaaaa$ and $y_i = B$	3	2
	count of i s.t. $shape(x_i) = Aaaaaaa$ and $y_i = I$	1	1
	count of i s.t. $shape(x_i) = Aaaaaaa$ and $y_i = O$	0	1
<i>prefix:</i>	count of i s.t. $pre_1(x_i) = B$ and $y_i = B$	2	1
	count of i s.t. $pre_1(x_i) = B$ and $y_i = I$	0	0
	count of i s.t. $pre_1(x_i) = B$ and $y_i = O$	0	1
	count of i s.t. $pre_1(x_i) = s$ and $y_i = O$	2	2
	count of i s.t. $shape(pre_1(x_i)) = A$ and $y_i = B$	5	4
	count of i s.t. $shape(pre_1(x_i)) = A$ and $y_i = I$	1	1
	count of i s.t. $shape(pre_1(x_i)) = A$ and $y_i = O$	0	1
	$\llbracket shape(pre_1(x_1)) = A \wedge y_1 = B \rrbracket$	1	0
	$\llbracket shape(pre_1(x_1)) = A \wedge y_1 = O \rrbracket$	0	1
<i>gazetteer:</i>	count of i s.t. x_i is in the gazetteer and $y_i = B$	2	1
	count of i s.t. x_i is in the gazetteer and $y_i = I$	0	0
	count of i s.t. x_i is in the gazetteer and $y_i = O$	0	1
	count of i s.t. $x_i = \textit{sent}$ and $y_i = O$	1	1

(What is *Not* A Linear Model?)

- Models with hidden variables

$$\arg \max_{\mathbf{y}} p(\mathbf{y} \mid \mathbf{x}) = \arg \max_{\mathbf{y}} \sum_{\mathbf{z}} p(\mathbf{y}, \mathbf{z} \mid \mathbf{x})$$

- Models based on non-linear kernels

$$\arg \max_{\mathbf{y}} \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}) = \arg \max_{\mathbf{y}} \sum_{i=1}^N \alpha_i K(\langle \mathbf{x}_i, \mathbf{y}_i \rangle, \langle \mathbf{x}, \mathbf{y} \rangle)$$

Decoding

- For HMMs, the decoding algorithm we usually think of first is the Viterbi algorithm.
 - This is just one example.
- We will view decoding in five different ways.
 - Sequence models as a running example.
 - These views are not just for HMMs.
 - Sometimes they will lead us back to Viterbi!

Five Views of Decoding

1. Probabilistic Graphical Models

- View the linguistic structure as a collection of random variables that are interdependent.
- Represent interdependencies as a directed or undirected graphical model.
- Conditional probability tables (BNs) or factors (MNs) encode the probability distribution.

Inference in Graphical Models

- General algorithm for exact MAP inference: **variable elimination**.
 - Iteratively solve for the best values of each variable conditioned on values of “preceding” neighbors.
 - Then trace back.

The Viterbi algorithm is an instance of max-product variable elimination!

MAP is Linear Decoding

- Bayesian network:

$$\sum_i \log p(x_i \mid \text{parents}(X_i)) \\ + \sum_j \log p(y_j \mid \text{parents}(Y_j))$$

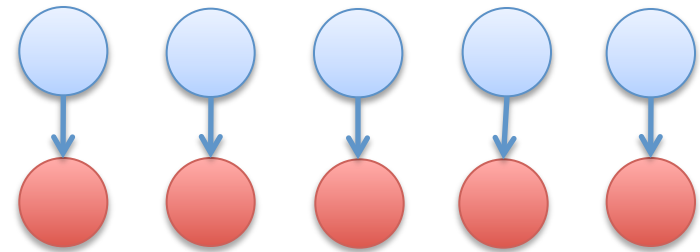
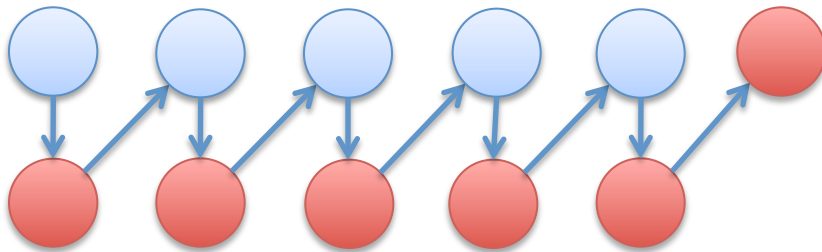
- Markov network:

$$\sum_C \log \phi_C (\{x_i\}_{i \in C}, \{y_j\}_{j \in C})$$

- This only works if every variable is in **X** or **Y**.

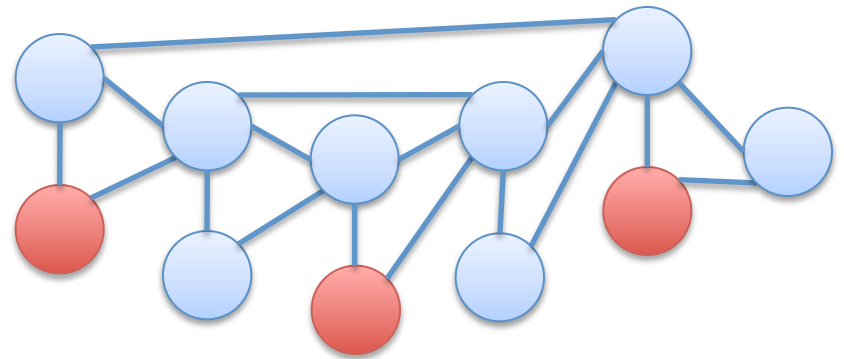
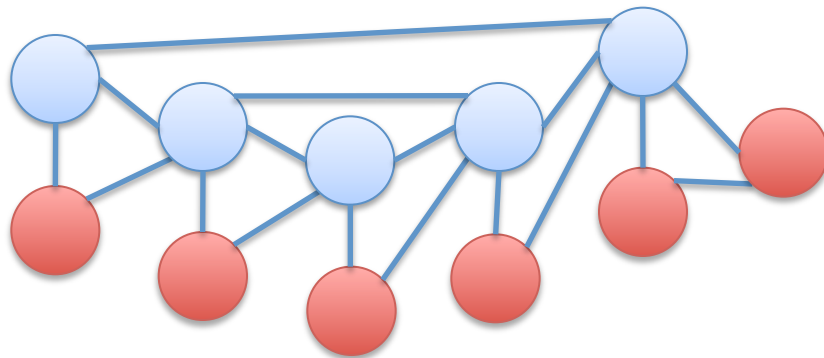
Inference in Graphical Models

- Remember: more edges make inference more expensive.
 - Fewer edges means stronger independence.
- Really pleasant:



Inference in Graphical Models

- Remember: more edges make inference more expensive.
 - Fewer edges means stronger independence.
- Really unpleasant:



2. Polytopes











“Parts”

- Assume that feature function \mathbf{g} breaks down into local parts.

$$\mathbf{g}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{\#parts(\mathbf{x})} \mathbf{f}(\Pi_i(\mathbf{x}, \mathbf{y}))$$

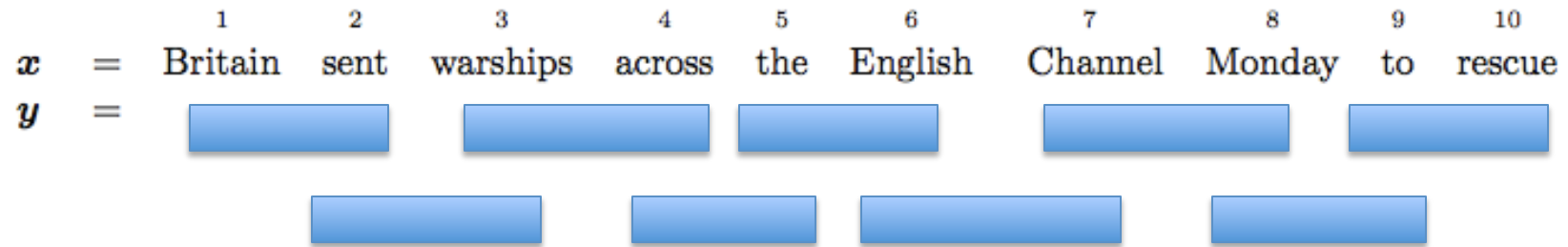
- Each part has an alphabet of possible values.
 - Decoding is choosing values for all parts, with **consistency** constraints.
 - (In the graphical models view, a part is a clique.)

Example

		1	2	3	4	5	6	7	8	9	10
x	=	Britain	sent	warships	across	the	English	Channel	Monday	to	rescue
y	=										

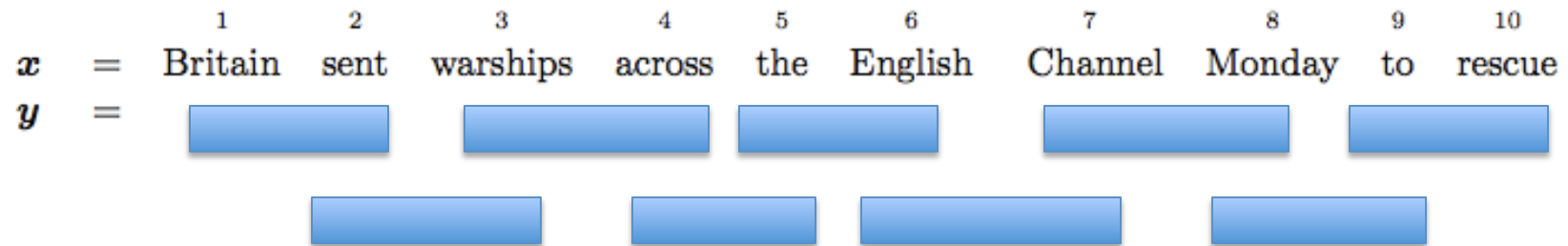
- One part per word, each is in $\{B, I, O\}$
- No features look at multiple parts
 - Fast inference
 - Not very expressive

Example



- One part per bigram, each is in $\{BB, BI, BO, IB, II, IO, OB, OO\}$
- Features and constraints can look at pairs
 - Slower inference
 - A bit more expressive

Geometric View



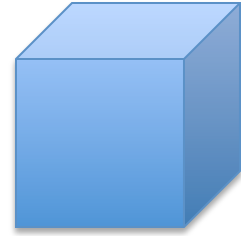
- Let $z_{i,\pi}$ be 1 if part i takes value π and 0 otherwise.
- \mathbf{z} is a vector in $\{0, 1\}^N$
 - N = total number of localized part values
 - Each \mathbf{z} is a vertex of the unit cube

Score is Linear in \mathbf{z}

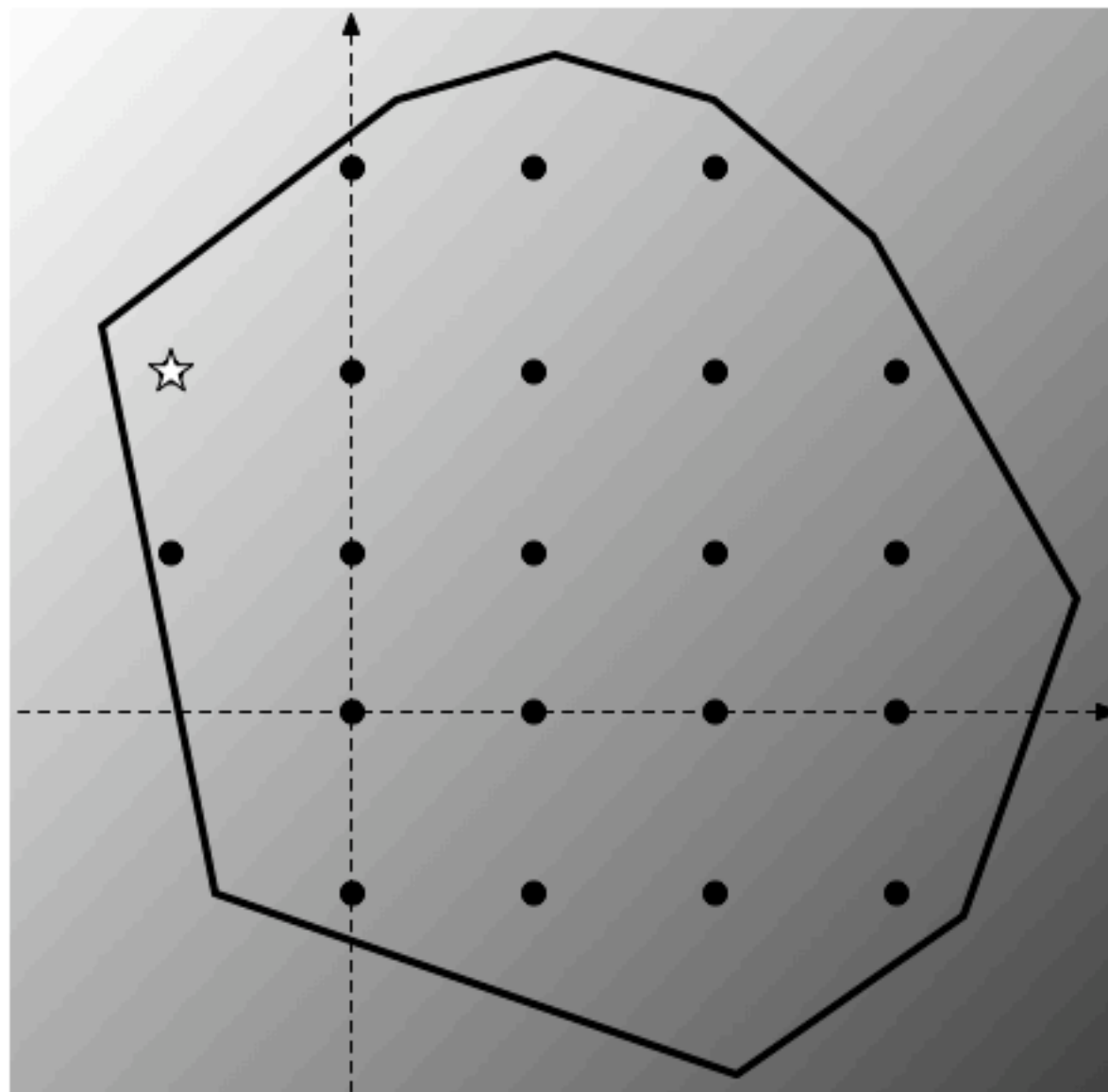
$$\begin{aligned}
 \arg \max_{\mathbf{y}} \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}) &= \arg \max_{\mathbf{y}} \mathbf{w}^\top \sum_{i=1}^{\#parts(\mathbf{x})} \mathbf{f}(\Pi_i(\mathbf{x}, \mathbf{y})) \\
 &= \arg \max_{\mathbf{y}} \mathbf{w}^\top \sum_{i=1}^{\#parts(\mathbf{x})} \sum_{\boldsymbol{\pi} \in \text{Values}(\Pi_i)} \mathbf{f}(\boldsymbol{\pi}) \mathbf{1}\{\Pi_i(\mathbf{x}, \mathbf{y}) = \boldsymbol{\pi}\} \\
 &= \arg \max_{\mathbf{z} \in \mathcal{Z}_{\mathbf{x}}} \mathbf{w}^\top \sum_{i=1}^{\#parts(\mathbf{x})} \sum_{\boldsymbol{\pi} \in \text{Values}(\Pi_i)} \mathbf{f}(\boldsymbol{\pi}) z_{i,\boldsymbol{\pi}} \\
 &= \arg \max_{\mathbf{z} \in \mathcal{Z}_{\mathbf{x}}} \mathbf{w}^\top \mathbf{F}_{\mathbf{x}} \mathbf{z} \\
 &= \arg \max_{\mathbf{z} \in \mathcal{Z}_{\mathbf{x}}} (\mathbf{w}^\top \mathbf{F}_{\mathbf{x}}) \mathbf{z}
 \end{aligned}$$

not really
equal; need
to transform
back to get \mathbf{y}

Polyhedra

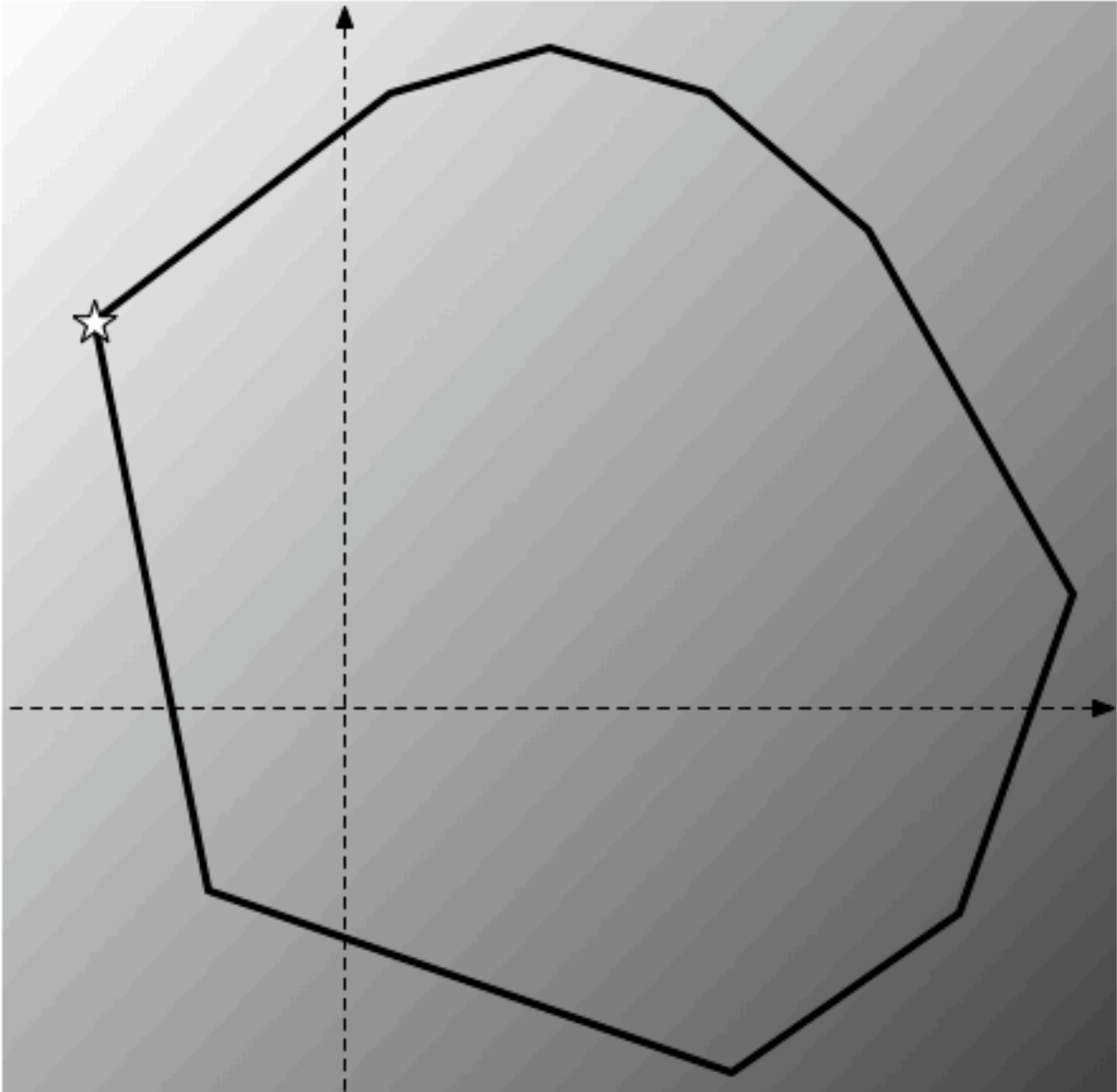


- Not all vertices of the N -dimensional unit cube satisfy the constraints.
 - E.g., can't have $z_{1,BI} = 1$ and $z_{2,BI} = 1$
- Sometimes we can write down a small (polynomial number) of linear constraints on \mathbf{z} .
- Result: linear objective, linear constraints, integer constraints ...



Integer Linear Programming

- Very easy to add new constraints and non-local features.
- Many decoding problems have been mapped to ILP (sequence labeling, parsing, ...), but it's *not* always trivial.
- NP-hard in general.
 - But there are packages that often work well in practice (e.g., CPLEX)
 - Specialized algorithms in some cases
 - LP relaxation for approximate solutions



Remark

- Graphical models assumed a probabilistic interpretation
 - Though they are not always learned using a probabilistic interpretation!
- The polytope view is agnostic about how you interpret the weights.
 - It only says that the decoding problem is an ILP.