

# Natural Language Processing (CSE 517): Graphical Models

Noah Smith

© 2018

University of Washington  
nasmith@cs.washington.edu

May 2–4, 2018

## Notation

Let  $\mathbf{V} = \langle V_1, V_2, \dots, V_\ell \rangle$  be a collection of random variables (not necessarily a sequence).

$\text{Val}(V)$  will denote the values of a r.v.  $V$ .

$\mathbf{V}_I$  denotes a subset of the r.v.s  $\mathbf{V}$  with indices  $i \in I$ .

$$\mathbf{V}_{\neg I} = \mathbf{V} \setminus \mathbf{V}_I$$

Recall:

- ▶  $p(\mathbf{V}) = \prod_{i=1}^{\ell} p(V_i | V_1, \dots, V_{i-1})$  (always true, for any ordering)
- ▶  $p(\mathbf{V}_I, \mathbf{V}_J | \mathbf{V}_K) = p(\mathbf{V}_I | \mathbf{V}_K) \cdot p(\mathbf{V}_J | \mathbf{V}_K)$  if and only if  $\mathbf{V}_I \perp \mathbf{V}_J | \mathbf{V}_K$  (conditional independence)
- ▶  $p(\mathbf{V}_I = \mathbf{v}_I) = \sum_{\mathbf{v}_{\neg I} \in \text{Val}(\mathbf{V}_{\neg I})} p(\mathbf{V}_I = \mathbf{v}_I, \mathbf{V}_{\neg I} = \mathbf{v}_{\neg I})$  (marginalization)

# Factor Graphs

Two kinds of vertices:

- ▶ Random variables (denoted by circles, “ $V_i$ ”)
- ▶ Factors (denoted by squares, “ $f_j$ ”)

The graph is *bipartite*; every edge connects some variable to some factor. Let  $I_j \subseteq \{1, \dots, \ell\}$  be the set of variables  $f_j$  is connected to.

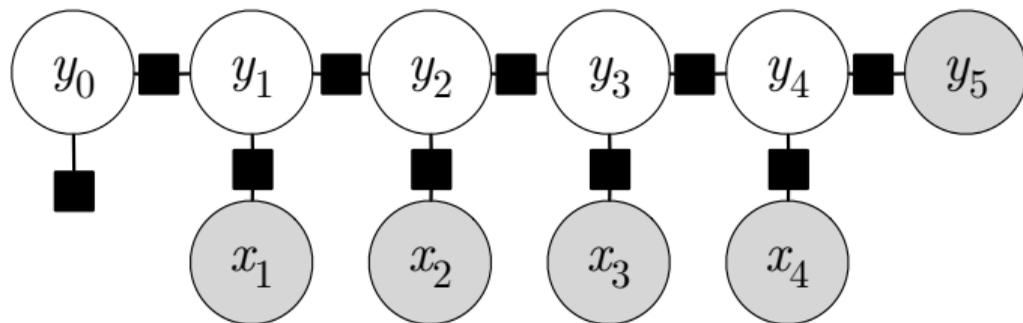
Factor  $f_j$  defines a map  $\text{Val}(\mathbf{V}_{I_j}) \rightarrow \mathbb{R}_{\geq 0}$ .

The graph and factors define a probability distribution:

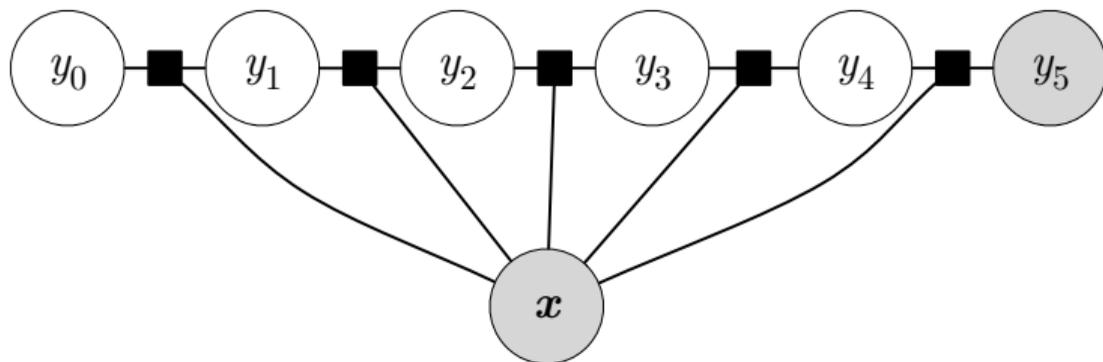
$$p(\mathbf{V} = \mathbf{v}) \propto \prod_j f_j(\mathbf{v}_{I_j})$$

## Factor Graphs We've Seen Before

Hidden Markov model:



General first-order sequence model:



## Two Kinds of Factors

Conditional probability tables. E.g., if  $I_j = \{1, 2, 3\}$ :

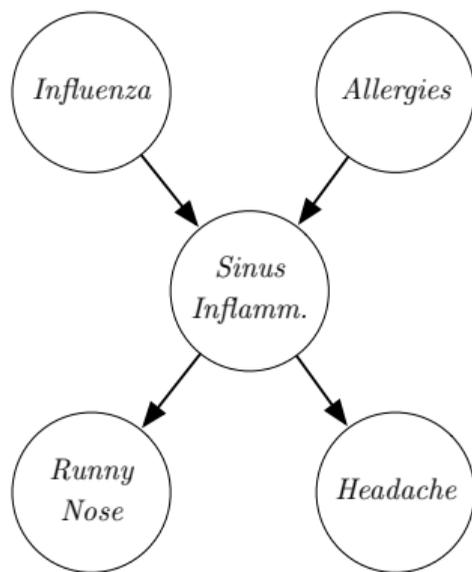
$$f_j(v_1, v_2, v_3) = p(V_3 = v_3 \mid V_1 = v_1, V_2 = v_2)$$

Lead to **Bayesian networks** (with some constraints).

Potential functions (arbitrary nonnegative values).

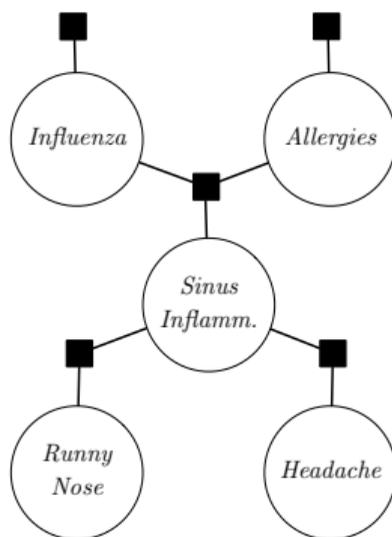
Lead to **Markov random fields** (a.k.a. Markov networks).

## Yucky Bayesian Network



Sinus inflammation is caused by flu, but also by allergies.  
Runny nose and headache are both caused by sinus inflammation.

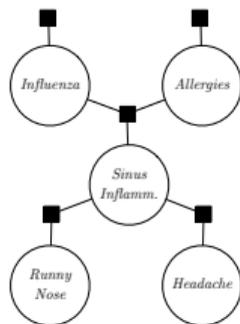
## Yucky Factor Graph



~~Sinus inflammation is caused by flu, but also by allergies.~~

~~Runny nose and headache are both caused by sinus inflammation.~~

# Yucky Factor Graph



$I$	$f_I$
0	
1	

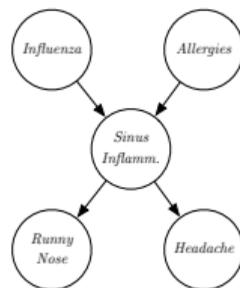
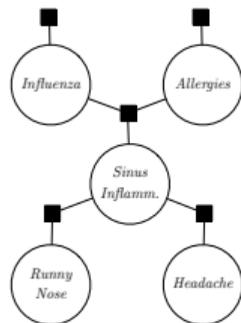
$I$	$f_A$
0	
1	

$S$	$I$	$A$	$f_{S,I,A}$
0	0	0	
0	0	1	
0	1	0	
0	1	1	
1	0	0	
1	0	1	
1	1	0	
1	1	1	

$R$	$S$	$f_{R,S}$
0	0	
0	1	
1	0	
1	1	

$H$	$S$	$f_{H,S}$
0	0	
0	1	
1	0	
1	1	

# Yucky Factor Graph



$I$	$f_I$
0	
1	

$I$	$f_A$
0	
1	

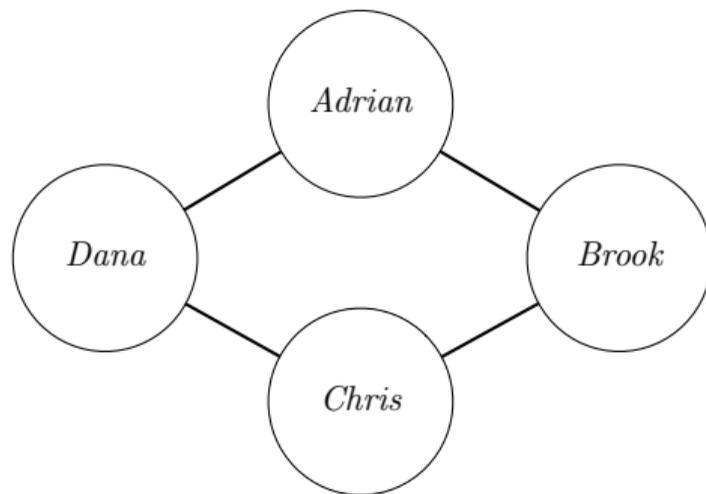
$S$	$I$	$A$	$f_{S,I,A}$
0	0	0	
0	0	1	
0	1	0	
0	1	1	
1	0	0	
1	0	1	
1	1	0	
1	1	1	

$R$	$S$	$f_{R,S}$
0	0	
0	1	
1	0	
1	1	

$H$	$S$	$f_{H,S}$
0	0	
0	1	
1	0	
1	1	

$$\begin{aligned}
 p(i, a, s, r, h) &= f_I(i) \cdot f_A(a) \cdot f_{S,I,A}(s, i, a) \cdot f_{R,S}(r, s) \cdot f_{H,S}(h, s) \\
 &= p(i) \cdot p(a) \cdot p(s \mid i, a) \cdot p(r \mid s) \cdot p(h \mid s)
 \end{aligned}$$

# Naughty Markov Random Field



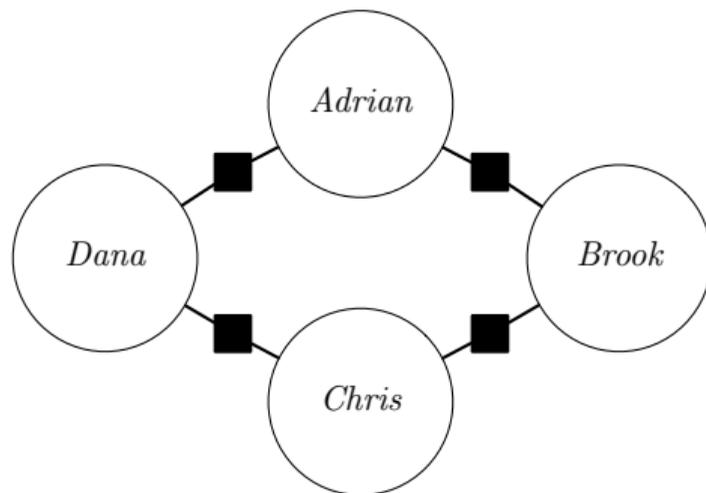
Independencies:  $A \perp C \mid B, D$ ;

$B \perp D \mid A, C$ ;

$\neg A \perp C$ ;

$\neg B \perp D$

# Naughty Factor Graph



A	B	$f_{A,B}$
0	0	
0	1	
1	0	
1	1	

B	C	$f_{B,C}$
0	0	
0	1	
1	0	
1	1	

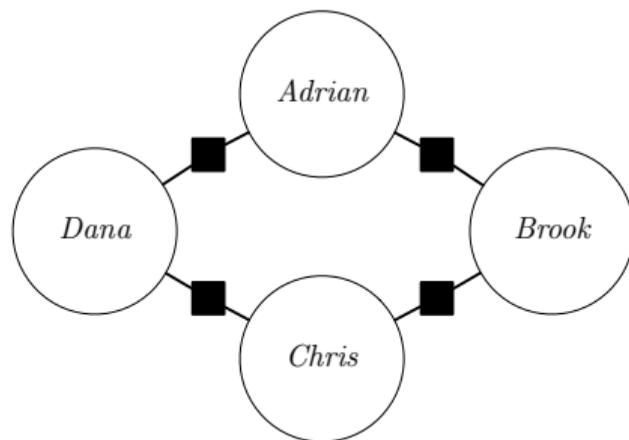
C	D	$f_{C,D}$
0	0	
0	1	
1	0	
1	1	

D	A	$f_{D,A}$
0	0	
0	1	
1	0	
1	1	

$$p(a, b, c, d) =$$

$$\frac{f_{A,B}(a, b) \cdot f_{B,C}(b, c) \cdot f_{C,D}(c, d) \cdot f_{D,A}(d, a)}{\sum_{a'} \sum_{b'} \sum_{c'} \sum_{d'} f_{A,B}(a', b') \cdot f_{B,C}(b', c') \cdot f_{C,D}(c', d') \cdot f_{D,A}(d', a')}$$

# Assignment Probabilities: Examples



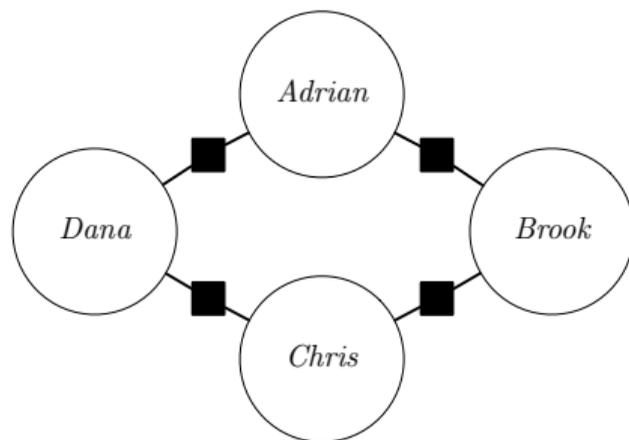
<i>A</i>	<i>B</i>	$f_{A,B}$
0	0	30
0	1	5
1	0	1
1	1	10

<i>B</i>	<i>C</i>	$f_{B,C}$
0	0	100
0	1	1
1	0	1
1	1	100

<i>C</i>	<i>D</i>	$f_{C,D}$
0	0	1
0	1	100
1	0	100
1	1	1

<i>D</i>	<i>A</i>	$f_{D,A}$
0	0	100
0	1	1
1	0	1
1	1	100

# Assignment Probabilities: Examples



A	B	$f_{A,B}$
0	0	30
0	1	5
1	0	1
1	1	10

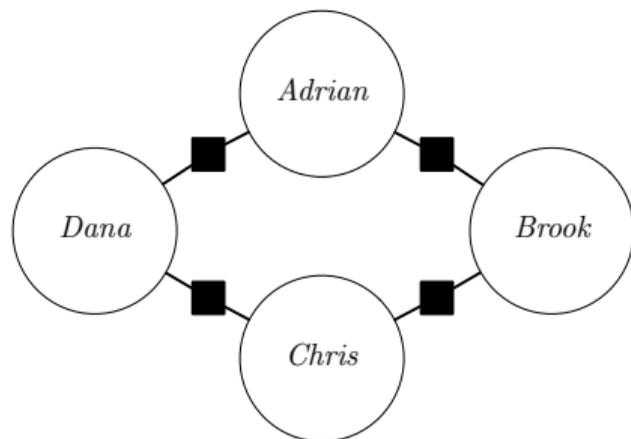
B	C	$f_{B,C}$
0	0	100
0	1	1
1	0	1
1	1	100

C	D	$f_{C,D}$
0	0	1
0	1	100
1	0	100
1	1	1

D	A	$f_{D,A}$
0	0	100
0	1	1
1	0	1
1	1	100

$$\sum_{\substack{a' \in \\ \text{Val}(A)}} \sum_{\substack{b' \in \\ \text{Val}(B)}} \sum_{\substack{c' \in \\ \text{Val}(C)}} \sum_{\substack{d' \in \\ \text{Val}(D)}} f_{A,B}(a', b') \cdot f_{B,C}(b', c') \cdot f_{C,D}(c', d') \cdot f_{D,A}(d', a')$$
$$= 7,201,840$$

# Assignment Probabilities: Examples



<i>A</i>	<i>B</i>	$f_{A,B}$
0	0	30
0	1	5
1	0	1
1	1	10

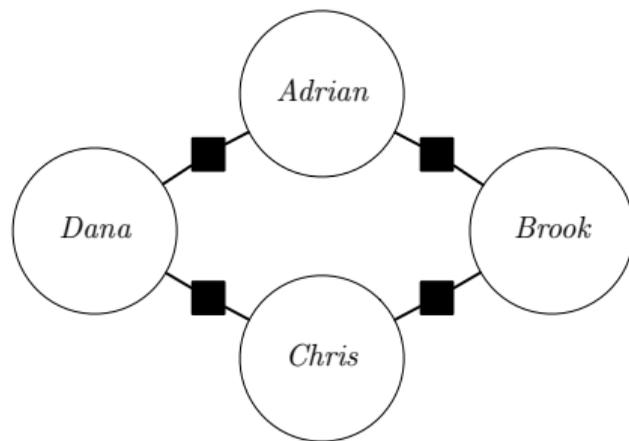
<i>B</i>	<i>C</i>	$f_{B,C}$
0	0	100
0	1	1
1	0	1
1	1	100

<i>C</i>	<i>D</i>	$f_{C,D}$
0	0	1
0	1	100
1	0	100
1	1	1

<i>D</i>	<i>A</i>	$f_{D,A}$
0	0	100
0	1	1
1	0	1
1	1	100

$$p(A = 0, B = 1, C = 1, D = 0) = \frac{5,000,000}{7,201,840} \approx 0.69$$

# Assignment Probabilities: Examples



A	B	$f_{A,B}$
0	0	30
0	1	5
1	0	1
1	1	10

B	C	$f_{B,C}$
0	0	100
0	1	1
1	0	1
1	1	100

C	D	$f_{C,D}$
0	0	1
0	1	100
1	0	100
1	1	1

D	A	$f_{D,A}$
0	0	100
0	1	1
1	0	1
1	1	100

$$p(A = 1, B = 1, C = 0, D = 0) = \frac{10}{7,201,840} \approx 0.0000014$$

# Structure and Independence

Bayesian networks:

- ▶ A variable is conditionally independent of its non-descendants given its parents.

Markov networks:

- ▶ Conditional independence derived from “Markov blanket” and separation properties.

Local configurations can be used to check *all* conditional independence questions; almost no need to look at the values in the factors!

# Independence “Spectrum”

$$\prod_{i=1}^{\ell} f_{V_i}(V_i)$$

everything is independent

minimal expressive power

fewer parameters

$$f_{\mathbf{V}}(\mathbf{V})$$

everything can be interdependent

arbitrary expressive power

more parameters

## Operations on Factors: Multiplication

Given two factors  $f_U$  and  $f_V$ , we can create a new “product” factor such that:

$$f_{U \cup V}(\mathbf{u} \cup \mathbf{v}) = f_U(\mathbf{u}) \cdot f_V(\mathbf{v})$$

for all  $\mathbf{u} \in \text{Val}(U)$  and all  $\mathbf{v} \in \text{Val}(V)$ .

A	B	$f_{A,B}$
0	0	30
0	1	5
1	0	1
1	1	10

.

B	C	$f_{B,C}$
0	0	100
0	1	1
1	0	1
1	1	100

=

A	B	C	$f_{A,B,C}$
0	0	0	3,000
0	0	1	30
0	1	0	5
0	1	1	500
1	0	0	100
1	0	1	1
1	1	0	10
1	1	1	1,000

## Operations on Factors: Multiplication

Given two factors  $f_U$  and  $f_V$ , we can create a new “product” factor such that:

$$f_{U \cup V}(u \cup v) = f_U(u) \cdot f_V(v)$$

for all  $u \in \text{Val}(U)$  and all  $v \in \text{Val}(V)$ .

$A$	$B$	$f_{A,B}$
0	0	30
0	1	5
1	0	1
1	1	10

.

$B$	$C$	$f_{B,C}$
0	0	100
0	1	1
1	0	1
1	1	100

=

$A$	$B$	$C$	$f_{A,B,C}$
0	0	0	3,000
0	0	1	30
0	1	0	5
0	1	1	500
1	0	0	100
1	0	1	1
1	1	0	10
1	1	1	1,000

This might remind you of a **join** operation on a database.

## Operations on Factors: Multiplication

Given two factors  $f_U$  and  $f_V$ , we can create a new “product” factor such that:

$$f_{U \cup V}(u \cup v) = f_U(u) \cdot f_V(v)$$

for all  $u \in \text{Val}(U)$  and all  $v \in \text{Val}(V)$ .

$A$	$B$	$f_{A,B}$
0	0	30
0	1	5
1	0	1
1	1	10

.

$B$	$C$	$f_{B,C}$
0	0	100
0	1	1
1	0	1
1	1	100

=

$A$	$B$	$C$	$f_{A,B,C}$
0	0	0	3,000
0	0	1	30
0	1	0	5
0	1	1	500
1	0	0	100
1	0	1	1
1	1	0	10
1	1	1	1,000

What happens if you multiply out all the factors in a factor graph?

## Operations on Factors: Maximization

Given a factor  $f_U$  and a variable  $V \notin U$ , we can transform  $f_{U,V}$  into  $f_U$  by:

$$f_U(\mathbf{u}) = \max_{v \in \text{Val}(V)} f_{U,V}(\mathbf{u}, v)$$

for all  $\mathbf{u} \in \text{Val}(U)$ .

A	C	$f_{A,C}$	
0	0	3,000	$B = 0$
0	1	500	$B = 1$
1	0	100	$B = 0$
1	1	1,000	$B = 1$

=

$\max_B$

A	B	C	$f_{A,B,C}$
0	0	0	3,000
0	0	1	30
0	1	0	5
0	1	1	500
1	0	0	100
1	0	1	1
1	1	0	10
1	1	1	1,000

## Operations on Factors: Marginalization

Given a factor  $f_U$  and a variable  $V \notin U$ , we can transform  $f_{U,V}$  into  $f_U$  by:

$$f_U(\mathbf{u}) = \sum_{v \in \text{Val}(V)} f_{U,V}(\mathbf{u}, v)$$

for all  $\mathbf{u} \in \text{Val}(U)$ .

$A$	$C$	$f_{A,C}$
0	0	3,000 + 5
0	1	30 + 500
1	0	100 + 10
1	1	1 + 1,000

=

$\sum_B$

$A$	$B$	$C$	$f_{A,B,C}$
0	0	0	3,000
0	0	1	30
0	1	0	5
0	1	1	500
1	0	0	100
1	0	1	1
1	1	0	10
1	1	1	1,000

## Operations on Factors: Marginalization

Given a factor  $f_U$  and a variable  $V \notin U$ , we can transform  $f_{U,V}$  into  $f_U$  by:

$$f_U(\mathbf{u}) = \sum_{v \in \text{Val}(V)} f_{U,V}(\mathbf{u}, v)$$

for all  $\mathbf{u} \in \text{Val}(U)$ .

$A$	$C$	$f_{A,C}$
0	0	3,000 + 5
0	1	30 + 500
1	0	100 + 10
1	1	1 + 1,000

=

$\sum_B$

$A$	$B$	$C$	$f_{A,B,C}$
0	0	0	3,000
0	0	1	30
0	1	0	5
0	1	1	500
1	0	0	100
1	0	1	1
1	1	0	10
1	1	1	1,000

If you multiply out all the factors in a factor graph, then sum out each variable, one by one, until none are left, what do you get?

## Factors are like numbers.

- ▶ Products are commutative:  $f_1 \cdot f_2 = f_2 \cdot f_1$

## Factors are like numbers.

- ▶ Products are commutative:  $f_1 \cdot f_2 = f_2 \cdot f_1$
- ▶ Products are associative:  $(f_1 \cdot f_2) \cdot f_3 = f_1 \cdot (f_2 \cdot f_3)$

## Factors are like numbers.

- ▶ Products are commutative:  $f_1 \cdot f_2 = f_2 \cdot f_1$
- ▶ Products are associative:  $(f_1 \cdot f_2) \cdot f_3 = f_1 \cdot (f_2 \cdot f_3)$
- ▶ Sums are commutative:  $\sum_X \sum_Y f = \sum_Y \sum_X f$

## Factors are like numbers.

- ▶ Products are commutative:  $f_1 \cdot f_2 = f_2 \cdot f_1$
- ▶ Products are associative:  $(f_1 \cdot f_2) \cdot f_3 = f_1 \cdot (f_2 \cdot f_3)$
- ▶ Sums are commutative:  $\sum_X \sum_Y f = \sum_Y \sum_X f$
- ▶ Maximizations are commutative:  $\max_X \max_Y f = \max_Y \max_X f$

## Factors are like numbers.

- ▶ Products are commutative:  $f_1 \cdot f_2 = f_2 \cdot f_1$
- ▶ Products are associative:  $(f_1 \cdot f_2) \cdot f_3 = f_1 \cdot (f_2 \cdot f_3)$
- ▶ Sums are commutative:  $\sum_X \sum_Y f = \sum_Y \sum_X f$
- ▶ Maximizations are commutative:  $\max_X \max_Y f = \max_Y \max_X f$
- ▶ Multiplication distributes over marginalization and maximization:

$$\sum_X (f_1 \cdot f_2) = f_1 \cdot \sum_X f_2$$
$$\max_X (f_1 \cdot f_2) = f_1 \cdot \max_X f_2$$

(assuming  $X$  is not in the scope of  $f_1$ ).

# Inference

Most general definition: “reason about some variables, optionally given values of some others.” Let  $\mathbf{O}$  be the observed variables and  $\mathbf{U}$  be the unobserved ones;  $\mathbf{V} = \mathbf{O} \cup \mathbf{U}$ .

Three inference problems, all given  $\mathbf{O} = \mathbf{o} \dots$

# Inference

Most general definition: “reason about some variables, optionally given values of some others.” Let  $\mathbf{O}$  be the observed variables and  $\mathbf{U}$  be the unobserved ones;  $\mathbf{V} = \mathbf{O} \cup \mathbf{U}$ .

Three inference problems, all given  $\mathbf{O} = \mathbf{o} \dots$

- ▶ **Marginal inference:** what is the marginal distribution over  $\mathbf{Q} \subset \mathbf{U}$ ? ( $p(\mathbf{Q} \mid \mathbf{o})$ , marginalizing out the rest.)

# Inference

Most general definition: “reason about some variables, optionally given values of some others.” Let  $\mathbf{O}$  be the observed variables and  $\mathbf{U}$  be the unobserved ones;  $\mathbf{V} = \mathbf{O} \cup \mathbf{U}$ .

Three inference problems, all given  $\mathbf{O} = \mathbf{o} \dots$

- ▶ **Marginal inference:** what is the marginal distribution over  $\mathbf{Q} \subset \mathbf{U}$ ? ( $p(\mathbf{Q} \mid \mathbf{o})$ , marginalizing out the rest.)
  - ▶ Related: draw samples from that distribution.

# Inference

Most general definition: “reason about some variables, optionally given values of some others.” Let  $\mathbf{O}$  be the observed variables and  $\mathbf{U}$  be the unobserved ones;  $\mathbf{V} = \mathbf{O} \cup \mathbf{U}$ .

Three inference problems, all given  $\mathbf{O} = \mathbf{o} \dots$

- ▶ **Marginal inference:** what is the marginal distribution over  $\mathbf{Q} \subset \mathbf{U}$ ? ( $p(\mathbf{Q} \mid \mathbf{o})$ , marginalizing out the rest.)
  - ▶ Related: draw samples from that distribution.
- ▶ **Most probable explanation (MPE):** what is the most probable assignment to  $\mathbf{U}$ ? ( $\operatorname{argmax}_{\mathbf{u}} p(\mathbf{u} \mid \mathbf{o})$ )

# Inference

Most general definition: “reason about some variables, optionally given values of some others.” Let  $\mathbf{O}$  be the observed variables and  $\mathbf{U}$  be the unobserved ones;  $\mathbf{V} = \mathbf{O} \cup \mathbf{U}$ .

Three inference problems, all given  $\mathbf{O} = \mathbf{o} \dots$

- ▶ **Marginal inference:** what is the marginal distribution over  $\mathbf{Q} \subset \mathbf{U}$ ? ( $p(\mathbf{Q} \mid \mathbf{o})$ , marginalizing out the rest.)
  - ▶ Related: draw samples from that distribution.
- ▶ **Most probable explanation (MPE):** what is the most probable assignment to  $\mathbf{U}$ ? ( $\operatorname{argmax}_{\mathbf{u}} p(\mathbf{u} \mid \mathbf{o})$ )
  - ▶ Related: what is the most *dangerous* assignment to  $\mathbf{U}$ ?

# Inference

Most general definition: “reason about some variables, optionally given values of some others.” Let  $\mathbf{O}$  be the observed variables and  $\mathbf{U}$  be the unobserved ones;  $\mathbf{V} = \mathbf{O} \cup \mathbf{U}$ .

Three inference problems, all given  $\mathbf{O} = \mathbf{o} \dots$

- ▶ **Marginal inference:** what is the marginal distribution over  $\mathbf{Q} \subset \mathbf{U}$ ? ( $p(\mathbf{Q} \mid \mathbf{o})$ , marginalizing out the rest.)
  - ▶ Related: draw samples from that distribution.
- ▶ **Most probable explanation (MPE):** what is the most probable assignment to  $\mathbf{U}$ ? ( $\operatorname{argmax}_{\mathbf{u}} p(\mathbf{u} \mid \mathbf{o})$ )
  - ▶ Related: what is the most *dangerous* assignment to  $\mathbf{U}$ ?
- ▶ **Maximum a posteriori (MAP):** what is the most probable assignment to  $\mathbf{Q} \subset \mathbf{U}$ ? ( $\operatorname{argmax}_{\mathbf{q}} p(\mathbf{q} \mid \mathbf{o})$ )

# Inference

Most general definition: “reason about some variables, optionally given values of some others.” Let  $O$  be the observed variables and  $U$  be the unobserved ones;  $V = O \cup U$ .

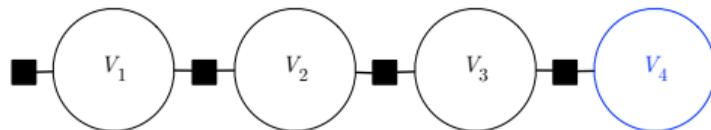
Three inference problems, all given  $O = o \dots$

- ▶ **Marginal inference:** what is the marginal distribution over  $Q \subset U$ ? ( $p(Q | o)$ , marginalizing out the rest.)
  - ▶ Related: draw samples from that distribution.
- ▶ **Most probable explanation (MPE):** what is the most probable assignment to  $U$ ? ( $\operatorname{argmax}_{\mathbf{u}} p(\mathbf{u} | o)$ )
  - ▶ Related: what is the most *dangerous* assignment to  $U$ ?
- ▶ **Maximum a posteriori (MAP):** what is the most probable assignment to  $Q \subset U$ ? ( $\operatorname{argmax}_{\mathbf{q}} p(\mathbf{q} | o)$ )
  - ▶ Related: what values of  $Q$  have the lowest expected cost?

# Marginal Inference

Given a factor graph with variables  $\mathbf{V}$ , find the marginal distribution over some  $V_i \in \mathbf{V}$ ,  $p(V_i)$ .

Simple chain example, focusing on  $i = 4$ :



$V_1$	$f_{V_1}$
0	
1	

$V_1$	$V_2$	$f_{V_1, V_2}$
0	0	
0	1	
1	0	
1	1	

$V_2$	$V_3$	$f_{V_2, V_3}$
0	0	
0	1	
1	0	
1	1	

$V_3$	$V_4$	$f_{V_3, V_4}$
0	0	
0	1	
1	0	
1	1	

## Observations

- ▶ If we had a single  $f_{V_4}$ , we could easily renormalize it to get  $p(V_4)$ .

## Observations

- ▶ If we had a single  $f_{V_4}$ , we could easily renormalize it to get  $p(V_4)$ .
- ▶ Correct:  $f_{V_4} = \sum_{V_1} \sum_{V_2} \sum_{V_3} f_{V_1} \cdot f_{V_1, V_2} \cdot f_{V_2, V_3} \cdot f_{V_3, V_4}$

## Observations

- ▶ If we had a single  $f_{V_4}$ , we could easily renormalize it to get  $p(V_4)$ .
- ▶ Correct:  $f_{V_4} = \sum_{V_1} \sum_{V_2} \sum_{V_3} f_{V_1} \cdot f_{V_1, V_2} \cdot f_{V_2, V_3} \cdot f_{V_3, V_4}$ 
  - ▶ But that multiplied-out factor would have  $\prod_i |\text{Val}(V_i)|$  values!

## Observations

► If we had a single  $f_{V_4}$ , we could easily renormalize it to get  $p(V_4)$ .

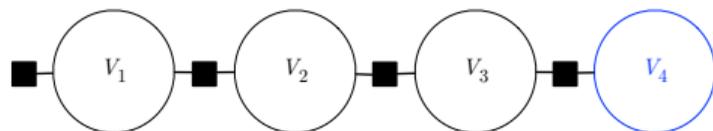
► Correct:  $f_{V_4} = \sum_{V_1} \sum_{V_2} \sum_{V_3} f_{V_1} \cdot f_{V_1, V_2} \cdot f_{V_2, V_3} \cdot f_{V_3, V_4}$

► But that multiplied-out factor would have  $\prod_i |\text{Val}(V_i)|$  values!

► Reorganize calculations:

$$\begin{aligned} & \sum_{V_1} \sum_{V_2} \sum_{V_3} f_{V_1} \cdot f_{V_1, V_2} \cdot f_{V_2, V_3} \cdot f_{V_3, V_4} \\ &= \sum_{V_3} f_{V_3, V_4} \cdot \left( \sum_{V_2} f_{V_2, V_3} \cdot \left( \sum_{V_1} f_{V_1, V_2} \cdot f_{V_1} \right) \right) \end{aligned}$$

# Marginal Inference



$V_1$	$f_{V_1}$
0	
1	

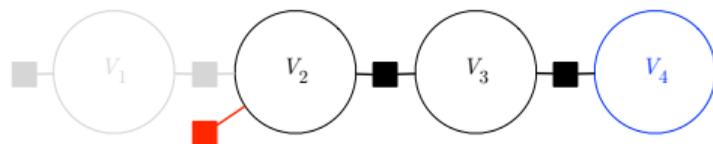
$V_1$	$V_2$	$f_{V_1, V_2}$
0	0	
0	1	
1	0	
1	1	

$V_2$	$V_3$	$f_{V_2, V_3}$
0	0	
0	1	
1	0	
1	1	

$V_3$	$V_4$	$f_{V_3, V_4}$
0	0	
0	1	
1	0	
1	1	

$$\begin{aligned} & \sum_{V_1} \sum_{V_2} \sum_{V_3} f_{V_1} \cdot f_{V_1, V_2} \cdot f_{V_2, V_3} \cdot f_{V_3, V_4} \\ &= \sum_{V_3} f_{V_3, V_4} \cdot \left( \sum_{V_2} f_{V_2, V_3} \cdot \left( \sum_{V_1} f_{V_1, V_2} \cdot f_{V_1} \right) \right) \end{aligned}$$

# Marginal Inference



$V_2$	$f_{V_2}$
0	
1	

$V_2$	$V_3$	$f_{V_2, V_3}$
0	0	
0	1	
1	0	
1	1	

$V_3$	$V_4$	$f_{V_3, V_4}$
0	0	
0	1	
1	0	
1	1	

$$\begin{aligned} & \sum_{V_1} \sum_{V_2} \sum_{V_3} f_{V_1} \cdot f_{V_1, V_2} \cdot f_{V_2, V_3} \cdot f_{V_3, V_4} \\ &= \sum_{V_3} f_{V_3, V_4} \cdot \left( \sum_{V_2} f_{V_2, V_3} \cdot f_{V_2} \right) \end{aligned}$$

# Marginal Inference



$V_3$	$f_{V_3}$
0	
1	

$V_3$	$V_4$	$f_{V_3, V_4}$
0	0	
0	1	
1	0	
1	1	

$$\begin{aligned} & \sum_{V_1} \sum_{V_2} \sum_{V_3} f_{V_1} \cdot f_{V_1, V_2} \cdot f_{V_2, V_3} \cdot f_{V_3, V_4} \\ &= \sum_{V_3} f_{V_3, V_4} \cdot f_{V_3} \end{aligned}$$

# Marginal Inference



$V_4$	$f_{V_4}$
0	
1	

$$\sum_{V_1} \sum_{V_2} \sum_{V_3} f_{V_1} \cdot f_{V_1, V_2} \cdot f_{V_2, V_3} \cdot f_{V_3, V_4}$$
$$= f_{V_4}$$

# Variable Elimination

Given a factor graph with factors  $\mathbf{f}$ , eliminate variable  $V$ .

1. Let  $\mathbf{f}_{elim} \subset \mathbf{f}$  be the factors connected to  $V$
2. Let  $\mathbf{f}_{keep} = \mathbf{f} \setminus \mathbf{f}_{elim}$  be the rest
3. Let  $f_{new} = \sum_V \prod_{f \in \mathbf{f}_{elim}} f$
4. Return  $\mathbf{f}_{keep} \cup \{f_{new}\}$

Uses the graph structure to avoid exponential blowup; this is an example of dynamic programming.

## Marginal Inference by Variable Elimination (No Evidence)

Given a factor graph with variables  $\mathbf{V}$  and factors  $\mathbf{f}$ , find the marginal distribution over some  $\mathbf{V}_{keep} \subset \mathbf{V}$ .

1. Order the variables in  $\mathbf{V} \setminus \mathbf{V}_{keep}$ .
2. For each  $V \in \mathbf{V} \setminus \mathbf{V}_{keep}$ :
  - ▶ Eliminate  $V$ ; i.e., remove factors connected to  $V$  and replace with the derived  $f_{new}$ .

The resulting factor graph is proportional to  $p(\mathbf{V}_{keep})$ .

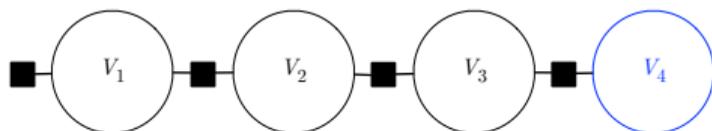
## Marginal Inference by Variable Elimination (No Evidence)

Given a factor graph with variables  $\mathbf{V}$  and factors  $\mathbf{f}$ , find the marginal distribution over some  $\mathbf{V}_{keep} \subset \mathbf{V}$ .

1. Order the variables in  $\mathbf{V} \setminus \mathbf{V}_{keep}$ .  
The ordering can make a huge difference!
2. For each  $V \in \mathbf{V} \setminus \mathbf{V}_{keep}$ :
  - ▶ Eliminate  $V$ ; i.e., remove factors connected to  $V$  and replace with the derived  $f_{new}$ .

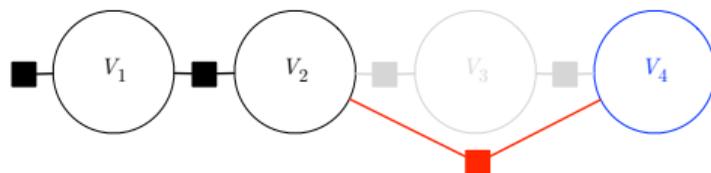
The resulting factor graph is proportional to  $p(\mathbf{V}_{keep})$ .

## A Less Good Ordering



$$\begin{aligned} & \sum_{V_1} \sum_{V_2} \sum_{V_3} f_{V_1} \cdot f_{V_1, V_2} \cdot f_{V_2, V_3} \cdot f_{V_3, V_4} \\ &= \sum_{V_1} f_{V_1} \cdot \left( \sum_{V_2} f_{V_1, V_2} \cdot \left( \sum_{V_3} f_{V_2, V_3} \cdot f_{V_3, V_4} \right) \right) \end{aligned}$$

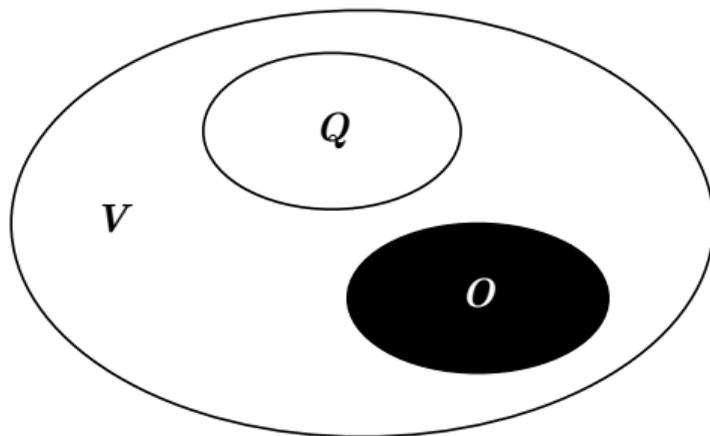
## A Less Good Ordering



$$\begin{aligned} & \sum_{V_1} \sum_{V_2} \sum_{V_3} f_{V_1} \cdot f_{V_1, V_2} \cdot f_{V_2, V_3} \cdot f_{V_3, V_4} \\ &= \sum_{V_1} f_{V_1} \cdot \left( \sum_{V_2} f_{V_1, V_2} \cdot \left( \sum_{V_3} f_{V_2, V_3} \cdot f_{V_3, V_4} \right) \right) \\ &= \sum_{V_1} f_{V_1} \cdot \left( \sum_{V_2} f_{V_1, V_2} \cdot f_{V_2, V_4} \right) \end{aligned}$$

## What About Evidence?

Original problem: given  $O = o$ , what is the marginal distribution over  $Q \subset U$ ? (i.e.,  $p(Q | O = o)$ .)



## What About Evidence?

Original problem: given  $\mathbf{O} = \mathbf{o}$ , what is the marginal distribution over  $\mathbf{Q} \subset \mathbf{U}$ ? (i.e.,  $p(\mathbf{Q} \mid \mathbf{O} = \mathbf{o})$ .)

This adds a step at the beginning: **reduce** factors to “respect the evidence.”

## What About Evidence?

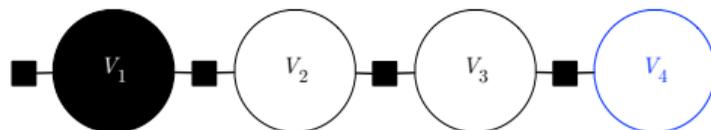
Original problem: given  $O = o$ , what is the marginal distribution over  $Q \subset U$ ? (i.e.,  $p(Q | O = o)$ .)

This adds a step at the beginning: **reduce** factors to “respect the evidence.”

This will remind you of a **select ... where** operation in a database.

# Marginal Inference

Suppose  $V_1$  is observed to take value 1.



$V_1$	$f_{V_1}$
0	
1	

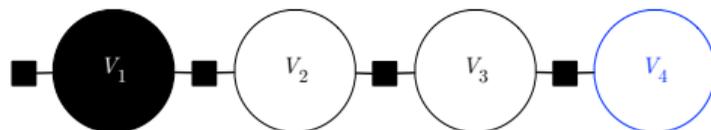
$V_1$	$V_2$	$f_{V_1, V_2}$
0	0	
0	1	
1	0	
1	1	

$V_2$	$V_3$	$f_{V_2, V_3}$
0	0	
0	1	
1	0	
1	1	

$V_3$	$V_4$	$f_{V_3, V_4}$
0	0	
0	1	
1	0	
1	1	

# Marginal Inference

Suppose  $V_1$  is observed to take value 1.



$V_1$	$f_{V_1}$
1	

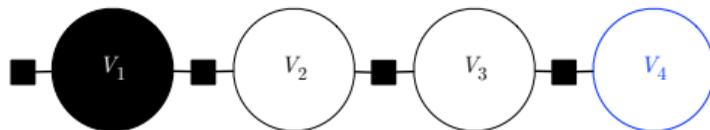
$V_1$	$V_2$	$f_{V_1, V_2}$
1	0	
1	1	

$V_2$	$V_3$	$f_{V_2, V_3}$
0	0	
0	1	
1	0	
1	1	

$V_3$	$V_4$	$f_{V_3, V_4}$
0	0	
0	1	
1	0	
1	1	

# Marginal Inference

Suppose  $V_1$  is observed to take value 1.



$V_1$	$f_{V_1}$
1	

$V_1$	$V_2$	$f_{V_1, V_2}$
1	0	
1	1	

$V_2$	$V_3$	$f_{V_2, V_3}$
0	0	
0	1	
1	0	
1	1	

$V_3$	$V_4$	$f_{V_3, V_4}$
0	0	
0	1	
1	0	
1	1	

Note that  $f_{V_1}$  is now a constant; since we renormalize at the end, we can ignore it. Observed nodes may create a “separation” between variables of interest and some factors.

# Marginal Inference by Variable Elimination with Evidence

Given a factor graph with variables  $V$  and factors  $f$ , and given  $O = o$  (where  $O \subset V$ ), find the marginal distribution over  $Q \subseteq U = V \setminus O$ .

1. Reduce factors connected to  $O$  to respect the evidence.
2. Order the variables in  $U \setminus Q$ .
3. For each  $V \in U \setminus Q$ :
  - ▶ Eliminate  $V$ ; i.e., remove factors connected to  $V$  and replace with the derived  $f_{new}$ .

The resulting factor graph is proportional to  $p(Q \mid O = o)$ .

## Remarks on Computational Complexity

In general, denser graphs are more expensive.

Runtime and space depend on the size of the original and intermediate factors. (This is why ordering matters so much.)

Finding the best ordering is NP-hard.

Certain graphical structures allow inference in linear time with respect to the size of the *original* factors.

- ▶ Bayesian networks: polytrees
- ▶ Markov networks: chordal graphs

## Return to Hidden Markov Models

- ▶ Hidden Markov models are not (quite) Bayesian networks.

## Return to Hidden Markov Models

- ▶ Hidden Markov models are not (quite) Bayesian networks.
  - ▶ Given an observed sequence  $x$ , however, an HMM provides a pattern to construct a Bayesian network.

## Return to Hidden Markov Models

- ▶ Hidden Markov models are not (quite) Bayesian networks.
  - ▶ Given an observed sequence  $x$ , however, an HMM provides a pattern to construct a Bayesian network.
  - ▶ Sometimes called “dynamic graphical models.”

## Return to Hidden Markov Models

- ▶ Hidden Markov models are not (quite) Bayesian networks.
  - ▶ Given an observed sequence  $x$ , however, an HMM provides a pattern to construct a Bayesian network.
  - ▶ Sometimes called “dynamic graphical models.”
- ▶ Marginal inference for every  $Y_i$  in an HMM can be accomplished by variable elimination.

## Return to Hidden Markov Models

- ▶ Hidden Markov models are not (quite) Bayesian networks.
  - ▶ Given an observed sequence  $x$ , however, an HMM provides a pattern to construct a Bayesian network.
  - ▶ Sometimes called “dynamic graphical models.”
- ▶ Marginal inference for every  $Y_i$  in an HMM can be accomplished by variable elimination.
  - ▶ All variables share some computation with those to their right and those to their left.

## Return to Hidden Markov Models

- ▶ Hidden Markov models are not (quite) Bayesian networks.
  - ▶ Given an observed sequence  $x$ , however, an HMM provides a pattern to construct a Bayesian network.
  - ▶ Sometimes called “dynamic graphical models.”
- ▶ Marginal inference for every  $Y_i$  in an HMM can be accomplished by variable elimination.
  - ▶ All variables share some computation with those to their right and those to their left.
  - ▶ This is called the **forward-backward** algorithm.

## Return to Hidden Markov Models

- ▶ Hidden Markov models are not (quite) Bayesian networks.
  - ▶ Given an observed sequence  $x$ , however, an HMM provides a pattern to construct a Bayesian network.
  - ▶ Sometimes called “dynamic graphical models.”
- ▶ Marginal inference for every  $Y_i$  in an HMM can be accomplished by variable elimination.
  - ▶ All variables share some computation with those to their right and those to their left.
  - ▶ This is called the **forward-backward** algorithm.
  - ▶ This is useful when we want to apply EM to HMMs (unsupervised sequence modeling).

## Return to Hidden Markov Models

- ▶ Hidden Markov models are not (quite) Bayesian networks.
  - ▶ Given an observed sequence  $x$ , however, an HMM provides a pattern to construct a Bayesian network.
  - ▶ Sometimes called “dynamic graphical models.”
- ▶ Marginal inference for every  $Y_i$  in an HMM can be accomplished by variable elimination.
  - ▶ All variables share some computation with those to their right and those to their left.
  - ▶ This is called the **forward-backward** algorithm.
  - ▶ This is useful when we want to apply EM to HMMs (unsupervised sequence modeling).
  - ▶ It is also useful in supervised learning.

## Related Topics

- ▶ Conditional random fields
- ▶ MPE inference
- ▶ MAP inference
- ▶ Inexact inference

# Conditional Random Fields (Sequence Version)

Lafferty et al. (2001)

A nice confluence:

- ▶ Probabilistic graphical model-style reasoning, as in HMMs.
- ▶ Discriminative training, as with structured perceptron.

Local factors:  $f_i(\mathbf{x}, y, y') = \exp(\mathbf{w} \cdot \phi(\mathbf{x}, i, y, y'))$

Log loss, where the graphical model parameterizes the probability distribution:

$$\sum_{i=1}^n \log \underbrace{\sum_{\mathbf{y} \in \mathcal{L}^{\ell_i+1}} \exp \left( \mathbf{w} \cdot \sum_{j=1}^{\ell_i+1} \phi(\mathbf{x}_i, j, y_j, y_{j-1}) \right)}_{\text{fear}}$$
$$- \underbrace{\mathbf{w} \cdot \sum_{j=1}^{\ell_i+1} \phi(\mathbf{x}_i, j, y_{i,j}, y_{i,j-1})}_{\text{hope}}$$

## Conditional Random Fields (General Version)

Factor graph consisting of “input” variables  $\mathbf{X}$  (always observed) and “output” variables  $\mathbf{Y}$ .

$$p(\mathbf{Y} = \mathbf{y} \mid \mathbf{X} = \mathbf{x}) = \frac{\prod_j f_j(\mathbf{x}, \mathbf{y}_{I_j})}{\sum_{\mathbf{y}' \in \text{Val}(\mathbf{Y})} \prod_j f_j(\mathbf{x}, \mathbf{y}'_{I_j})}$$

MLE:

$$\sum_{i=1}^n \log \underbrace{\sum_{\mathbf{y} \in \text{Val}(\mathbf{Y})} \prod_j f_j(\mathbf{x}_i, \mathbf{y}_{I_j})}_{\text{fear}} - \log \underbrace{\prod_j f_j(\mathbf{x}_i, \mathbf{y}_i_{I_j})}_{\text{hope}}$$

Marginal inference is required for calculating the left term and its gradient with respect to  $\mathbf{w}$ .

# MPE Inference

$$\operatorname{argmax}_{\mathbf{u} \in \text{Val}(\mathbf{U})} p(\mathbf{U} = \mathbf{u} \mid \mathbf{O} = \mathbf{o})$$

# MPE Inference

$$\operatorname{argmax}_{\mathbf{u} \in \text{Val}(\mathbf{U})} p(\mathbf{U} = \mathbf{u} \mid \mathbf{O} = \mathbf{o})$$

Variable elimination and exact inference are identical to the marginal case!

# MPE Inference

$$\operatorname{argmax}_{\mathbf{u} \in \text{Val}(\mathbf{U})} p(\mathbf{U} = \mathbf{u} \mid \mathbf{O} = \mathbf{o})$$

Variable elimination and exact inference are identical to the marginal case!

Just replace each sum operation with a max operation, and add bookkeeping to recover the most probable assignment.

# MPE Inference

$$\operatorname{argmax}_{\mathbf{u} \in \text{Val}(\mathbf{U})} p(\mathbf{U} = \mathbf{u} \mid \mathbf{O} = \mathbf{o})$$

Variable elimination and exact inference are identical to the marginal case!

Just replace each sum operation with a max operation, and add bookkeeping to recover the most probable assignment.

The Viterbi algorithm is, of course, an instance of this. Each “ $s_i(*)$ ” is an intermediate factor.

# MPE Inference

$$\operatorname{argmax}_{\mathbf{u} \in \text{Val}(U)} p(U = \mathbf{u} \mid O = \mathbf{o})$$

Variable elimination and exact inference are identical to the marginal case!

Just replace each sum operation with a max operation, and add bookkeeping to recover the most probable assignment.

The Viterbi algorithm is, of course, an instance of this. Each “ $s_i(*)$ ” is an intermediate factor.

Specifically for sequence models, it should be clear how factors/features that depend on the observed sequence  $\mathbf{X}$  don't affect the asymptotics of exact inference.

## Rocket Science: True MAP

Given a factor graph with variables  $V$  and factors  $f$ , and given  $O = o$  (where  $O \subset V$ ), find the most probable assignment of  $Q \subset U = V \setminus O$ .

Let  $R = U \setminus Q$ .

$$\begin{aligned} & \operatorname{argmax}_{q \in \text{Val}(Q)} p(Q = q \mid O = o) \\ &= \operatorname{argmax}_{q \in \text{Val}(Q)} \sum_{r \in \text{Val}(R)} p(Q = q, R = r \mid O = o) \end{aligned}$$

Solution: first use marginal inference to eliminate  $R$ , then use max inference to solve for  $Q$ .

# Alternative Inference Methods

Huge range of techniques!

Exact:

- ▶ Integer linear programming

Inexact:

- ▶ randomized (e.g., Gibbs sampling, importance sampling, simulated annealing, stochastic variational)
- ▶ deterministic (e.g., mean field variational, loopy belief propagation, linear programming relaxations, dual decomposition, beam search)

# References I

John D. Lafferty, Andrew McCallum, and Fernando C. N. Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proc. of ICML*, 2001.