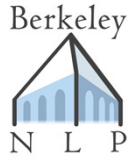


# Variational Inference for Structured NLP Models



ACL, August 4, 2013  
David Burkett and Dan Klein



# Tutorial Outline

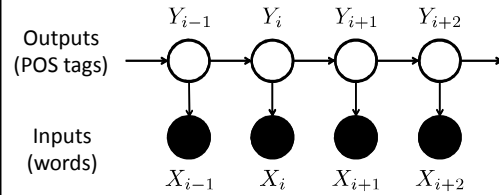
1. Structured Models and Factor Graphs
2. Mean Field
3. Structured Mean Field
4. Belief Propagation
5. Structured Belief Propagation
6. Wrap-Up

# Part 1: Structured Models and Factor Graphs



# Structured NLP Models

Example: Hidden Markov Model  
(Sample Application: Part of Speech Tagging)

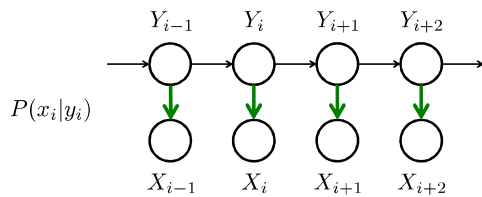


Goal: Queries from posterior  $P(Y = y|X = x)$  ( $P(y|x)$ )



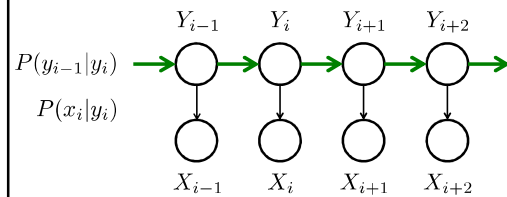
# Structured NLP Models

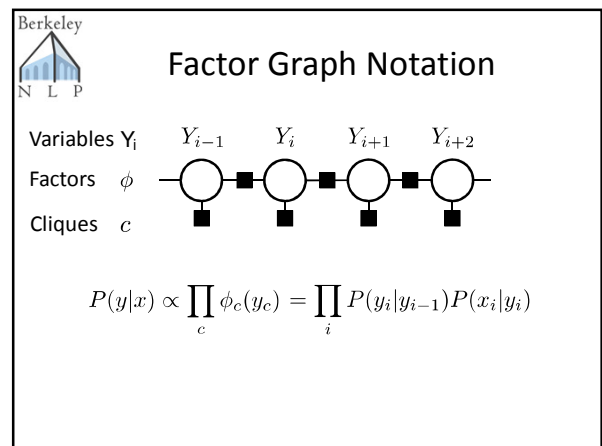
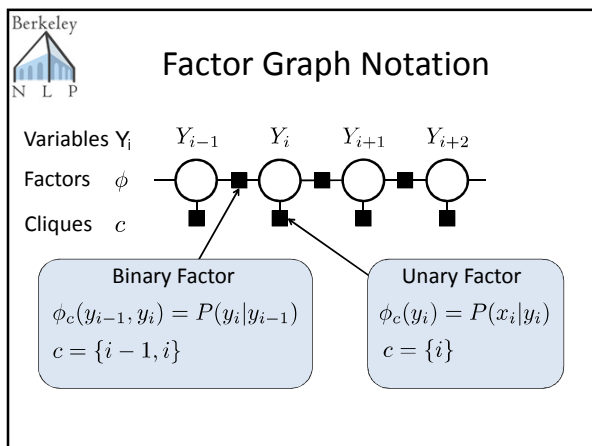
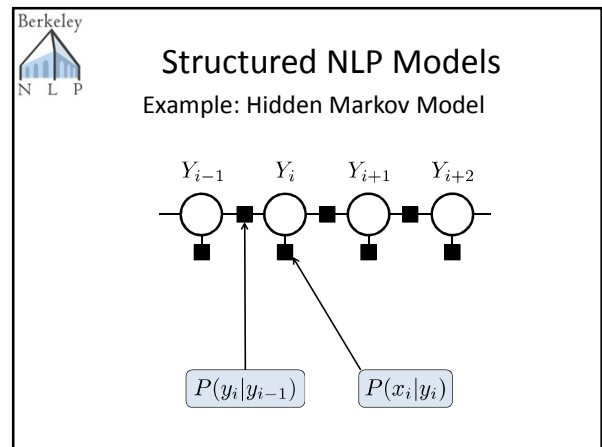
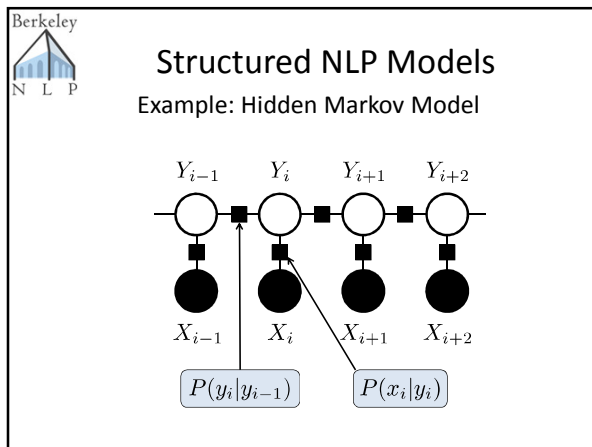
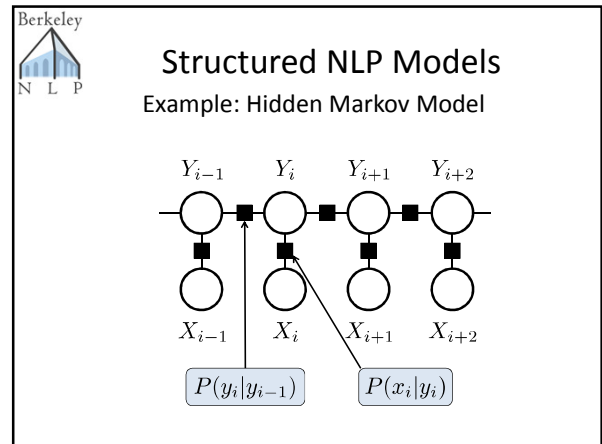
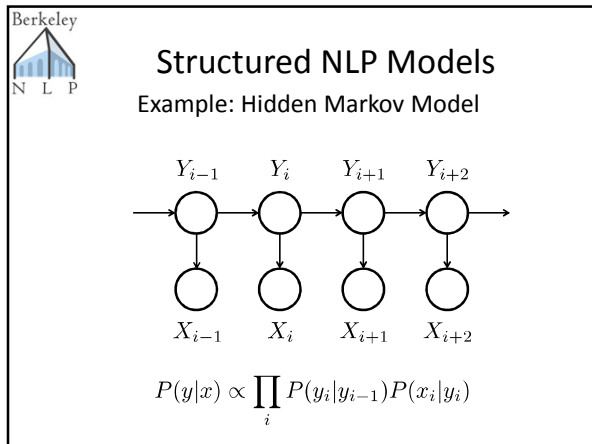
Example: Hidden Markov Model



# Structured NLP Models

Example: Hidden Markov Model





Berkeley

N L P

## Factor Graph Notation

Variables  $Y_i$

Factors  $\phi$

Cliques  $c$

Variables have factor (clique) neighbors:

$$\mathcal{N}(i) = \{c : i \in c\}$$

Factors have variable neighbors:

$$\mathcal{N}(\phi_c) = c$$

Berkeley

N L P

(Lafferty et al., 2001)

## Structured NLP Models

Example: Conditional Random Field  
(Sample Application: Named Entity Recognition)

$X = X_1, \dots, X_i, \dots, X_n$

$$P(y|x) \propto \exp \left( \sum_i w^\top f_i(y_i, x) + w^\top f_i(y_{i-1}, y_i, x) \right)$$

Berkeley

N L P

## Structured NLP Models

Example: Conditional Random Field

$\phi_i(y_i) = \exp(w^\top f_i(y_i, x))$

$\phi_{i-1,i}(y_{i-1}, y_i) = \exp(w^\top f_i(y_{i-1}, y_i, x))$

$$P(y|x) \propto \exp \left( \sum_i w^\top f_i(y_i, x) + w^\top f_i(y_{i-1}, y_i, x) \right)$$

Berkeley

N L P

## Structured NLP Models

Example: Conditional Random Field

$\phi(y_i) = \exp(w^\top f(y_i))$

$\phi(y_{i-1}, y_i) = \exp(w^\top f(y_{i-1}, y_i))$

$$P(y|x) \propto \exp \left( \sum_i w^\top f(y_i) + w^\top f(y_{i-1}, y_i) \right)$$

Berkeley

N L P

## Structured NLP Models

Example: Edge-Factored Dependency Parsing

$y_{ij} \in \{\text{left, right, off}\}$

(McDonald et al., 2005)

Berkeley

N L P

## Structured NLP Models

Example: Edge-Factored Dependency Parsing

$y_{ij} \in \{\text{left, right, off}\}$

Berkeley

**Structured NLP Models**  
Example: Edge-Factored Dependency Parsing

$y_{ij} \in \{\text{left, right, off}\}$

Berkeley

**Structured NLP Models**  
Example: Edge-Factored Dependency Parsing

$y_{ij} \in \{\text{left, right, off}\}$

Berkeley

**Structured NLP Models**  
Example: Edge-Factored Dependency Parsing

$y_{ij} \in \{\text{left, right, off}\}$

$\phi(y) = \begin{cases} 1 & y \text{ forms a tree} \\ 0 & \text{otherwise} \end{cases}$

$\phi(y_{ij}) = \begin{cases} \exp(w^\top f(i, j)) & y_{ij} = \text{left} \\ \exp(w^\top f(j, i)) & y_{ij} = \text{right} \\ 1 & y_{ij} = \text{off} \end{cases}$

Berkeley

**Inference**

▶ Input: Factor Graph

▶ Output: Marginals  $P(y_i|x)$

Berkeley

**Inference**

▶ Typical NLP Approach: Dynamic Programs!

▶ Examples:

- ▶ Sequence Models (Forward/Backward)
- ▶ Phrase Structure Parsing (CKY, Inside/Outside)
- ▶ Dependency Parsing (Eisner algorithm)
- ▶ ITG Parsing (Bitext Inside/Outside)

Berkeley

**Complex Structured Models**

POS Tagging

Joint

Named Entity Recognition

(Sutton et al., 2004)

Berkeley  
N L P

## Complex Structured Models

Dependency Parsing  
with Second Order Features

(McDonald & Pereira, 2006)  
(Carreras, 2007)

Berkeley  
N L P

## Complex Structured Models

### Word Alignment

$y_{ij} \in \{\text{on}, \text{off}\}$

(Taskar et al., 2005)

Berkeley  
N L P

## Complex Structured Models

### Word Alignment

$y_{ij} \in \{\text{on}, \text{off}\}$

$$\phi(y_{ij}) = \begin{cases} \exp(w^\top f(i, j)) & y_{ij} = \text{on} \\ 1 & y_{ij} = \text{off} \end{cases}$$

$$\phi(y_{i*}) = \begin{cases} 1 & |\{j : y_{ij} = \text{on}\}| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

$$\phi(y_{*j}) = \begin{cases} 1 & |\{i : y_{ij} = \text{on}\}| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

Berkeley  
N L P

## Variational Inference

- ▶ Approximate inference techniques that can be applied to any graphical model
- ▶ This tutorial:
  - ▶ Mean Field: Approximate the joint distribution with a product of marginals
  - ▶ Belief Propagation: Apply tree inference algorithms even if your graph isn't a tree
  - ▶ Structure: What changes when your factor graph has tractable substructures

## Part 2: Mean Field

Berkeley  
N L P

Berkeley  
N L P

## Mean Field Warmup

Wanted:  $\operatorname{argmax}_{a,b} P(a, b|x)$

Idea: coordinate ascent

Key object: assignments

Iterated Conditional Modes (Besag, 1986)

Berkeley  
N L P

### Mean Field Warmup

$\phi(a)$   
 $A = a^{(0)}$   $a^{(1)} = \operatorname{argmax}_a \phi(a)\phi(a,b)$   
 $\phi(a,b)$   
 $B = b^{(0)}$   
 $\phi(b)$

Wanted:  $\operatorname{argmax}_{a,b} P(a,b|x)$

Berkeley  
N L P

### Mean Field Warmup

$\phi(a)$   
 $A = a^{(1)}$   
 $\phi(a,b)$   
 $B$   
 $\phi(b)$   $b^{(1)} = \operatorname{argmax}_b \phi(b)\phi(a,b)$

Wanted:  $\operatorname{argmax}_{a,b} P(a,b|x)$

Berkeley  
N L P

### Mean Field Warmup

$\phi(a)$   
 $A$   $a^{(t)} = \operatorname{argmax}_a \phi(a)\phi(a,b)$   
 $\phi(a,b)$   
 $B = b^{(t-1)}$   
 $\phi(b)$

Wanted:  $\operatorname{argmax}_{a,b} P(a,b|x)$

Berkeley  
N L P

### Mean Field Warmup

$\phi(a)$   
 $A = a^{(t)}$   
 $\phi(a,b)$   
 $B$   
 $\phi(b)$   $b^{(t)} = \operatorname{argmax}_b \phi(b)\phi(a,b)$

Wanted:  $\operatorname{argmax}_{a,b} P(a,b|x)$

Berkeley  
N L P

### Mean Field Warmup

$\phi(a)$   
 $A$   $a^{(t)} = a^{(t-1)}$   
 $\phi(a,b)$   
 $B$   $b^{(t)} = b^{(t-1)}$   
 $\phi(b)$

Wanted:  $\operatorname{argmax}_{a,b} P(a,b|x)$

Approximate Result:  $(a^{(t)}, b^{(t)})$

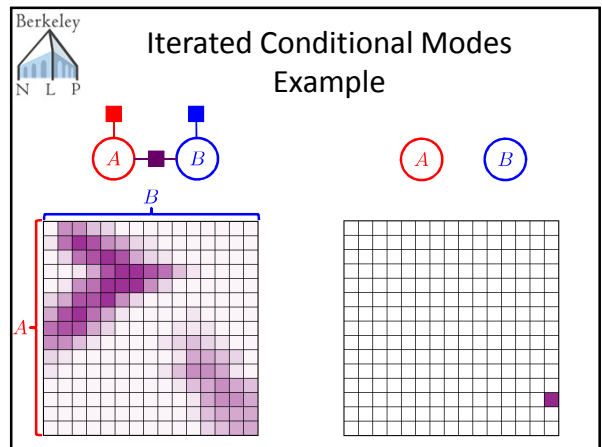
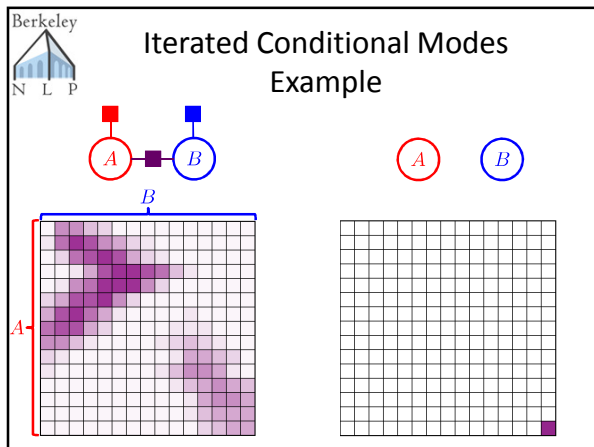
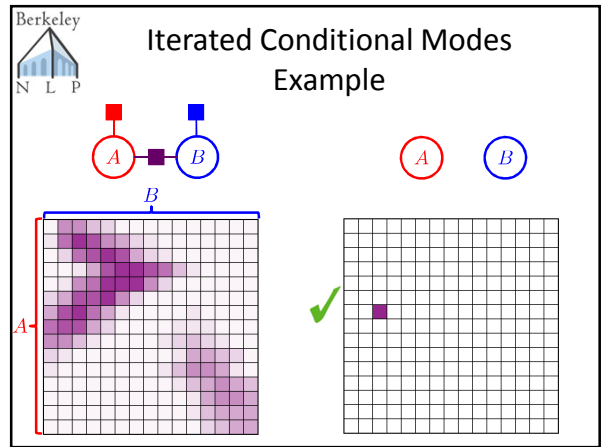
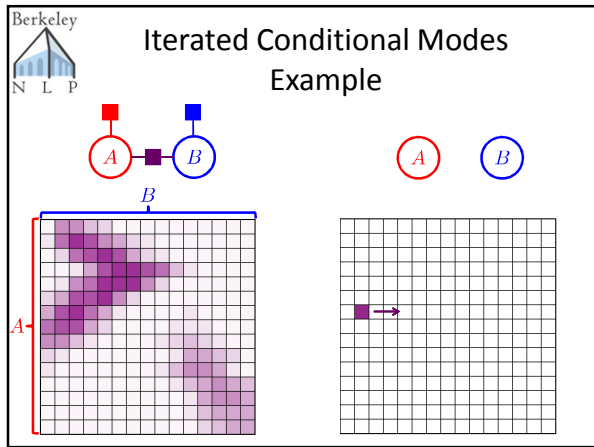
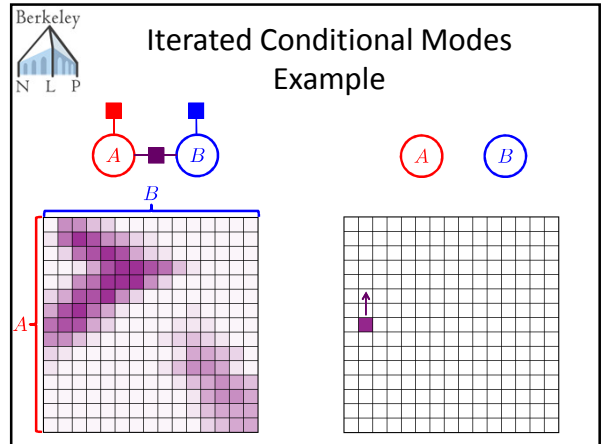
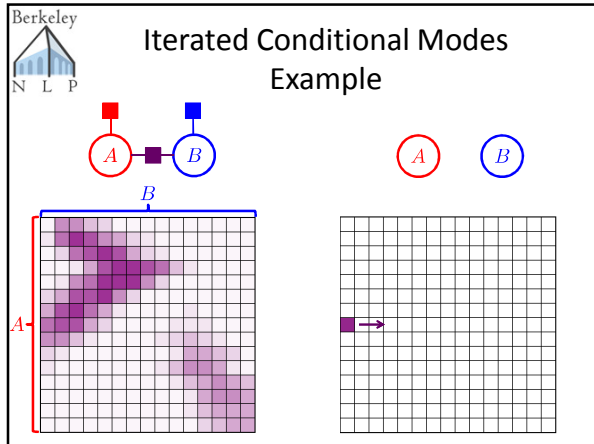
Berkeley  
N L P

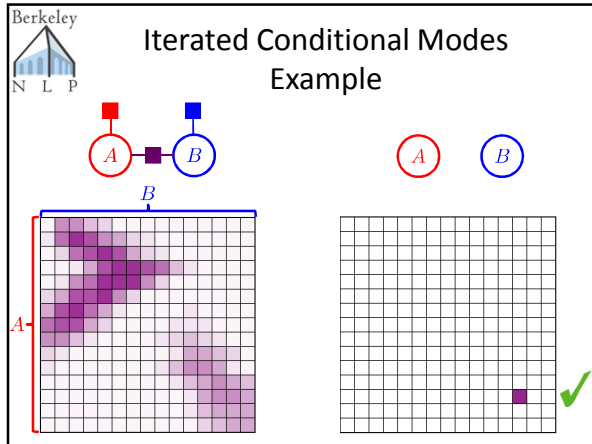
### Iterated Conditional Modes Example

$\phi(a)$   
 $A$   
 $\phi(a,b)$   
 $B$   
 $\phi(b)$

Wanted:  $\operatorname{argmax}_{a,b} P(a,b|x)$

Approximate Result:  $(a^{(t)}, b^{(t)})$





Berkeley  
N L P

### Mean Field Intro

Mean Field is coordinate ascent, just like Iterated Conditional Modes, but with soft assignments to each variable!

Berkeley  
N L P

### Mean Field Intro

$\phi(a)$    
  
 $\phi(a, b)$    
  
 $\phi(b)$

Wanted:  $P(a|x), P(b|x)$

Idea: coordinate ascent

Key object: (approx) marginals

Berkeley  
N L P

### Mean Field Intro

$\phi(a)$  =  $\exp(w^\top f(a))$    
  
 $\phi(a, b)$  =  $\exp(w^\top f(a, b))$    
  
 $\phi(b)$  =  $\exp(w^\top f(b))$

$P(a, b|x) \propto \phi(a)\phi(b)\phi(a, b)$

Berkeley  
N L P

### Mean Field Intro

$\phi(a)$  =  $\exp(w^\top f(a))$    
  
 $\phi(a, b)$  =  $\exp(w^\top f(a, b))$    
  
 $\phi(b)$  =  $\exp(w^\top f(b))$

$P(a, b|x) \propto \phi(a)\phi(b)\phi(a, b)$

$= \exp(w^\top f(a) + w^\top f(b) + w^\top f(a, b))$

Berkeley  
N L P

### Mean Field Intro

$w^\top f(a)$    
  
 $w^\top f(a, b)$    
 =  $b$    
 $w^\top f(b)$

$P(a|b, x) \propto \exp(w^\top f(a) + w^\top f(a, b))$

Wanted:  $P(a|x), P(b|x)$



Berkeley

### Mean Field Intro

$w^T f(a)$  ■  $P(a|b, x) \propto \exp(w^T f(a) + w^T f(a, b))$   
 $w^T f(a, b)$  ■  $q(b)$   
 $w^T f(b)$  ■

Wanted:  $P(a|x), P(b|x)$

Berkeley

### Mean Field Intro

$w^T f(a)$  ■  $q(a) \propto \exp(w^T f(a) + w^T \mathbb{E}_{q(b)} f(a, b))$   
 $w^T f(a, b)$  ■  $q(a)$   
 $w^T f(b)$  ■  $q(b)$

Wanted:  $P(a|x), P(b|x)$

Berkeley

### Mean Field Procedure

$w^T f(a)$  ■  $q^{(0)}(a)$   
 $w^T f(a, b)$  ■  $q^{(0)}(b)$   
 $w^T f(b)$  ■

Wanted:  $P(a|x), P(b|x)$

Berkeley

### Mean Field Procedure

$w^T f(a)$  ■  $q^{(t)}(a) \propto \exp(w^T f(a) + w^T \mathbb{E}_{q^{(t-1)}(b)} f(a, b))$   
 $w^T f(a, b)$  ■  $q^{(t-1)}(b)$   
 $w^T f(b)$  ■

Wanted:  $P(a|x), P(b|x)$

Berkeley

### Mean Field Procedure

$w^T f(a)$  ■  $q^{(t)}(a)$   
 $w^T f(a, b)$  ■  $q^{(t)}(b) \propto \exp(w^T f(b) + w^T \mathbb{E}_{q^{(t)}(a)} f(a, b))$   
 $w^T f(b)$  ■

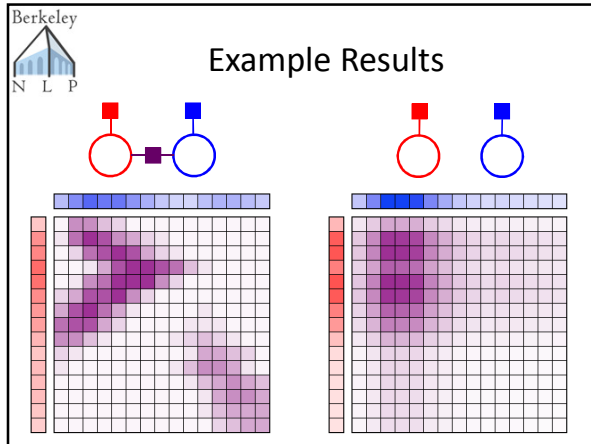
Wanted:  $P(a|x), P(b|x)$

Berkeley

### Mean Field Procedure

$w^T f(a)$  ■  $q^{(t)}(a)$   
 $w^T f(a, b)$  ■  $q^{(t)}(b)$   
 $w^T f(b)$  ■

Wanted:  $P(a|x), P(b|x)$



Berkeley

### Mean Field Derivation

- ▶ Goal:  $p(y) = P(y|x) \propto \exp\left(\sum_c w^\top f(y_c)\right)$
- ▶ Approximation:  $q(y) \approx p(y)$
- ▶ Constraint:  $q(y) = \prod_i q(y_i)$
- ▶ Objective:  $q(y) = \operatorname{argmin}_q KL(q||p)$
- ▶ Procedure: Coordinate ascent on each  $q(y_i)$
- ▶ What's the update?

Berkeley

### Mean Field Update

- 1)  $q(y_i) = \operatorname{argmin}_{q(y_i)} KL(q||p)$
- 2)  $\frac{\partial KL(q||p)}{\partial q(y_i)} = 0$
- 3-9) Lots of algebra
- 10)  $q(y_i) \propto \exp\left(\sum_{c \in \mathcal{N}(i)} w^\top \mathbb{E}_{q(y_{-i})} f_c(y_c)\right)$

Berkeley

### Approximate Expectations

$$\mathbb{E}_{q(y_{-i})} f(y_i, y_j, y_k) = \sum_{y_j} \sum_{y_k} q(y_j) q(y_k) f(y_i, y_j, y_k)$$

General:  $\mathbb{E}_{q(y_{-i})} f_c(y_c) = \sum_{y_c \in \{i\}} \left( \prod_{j \in \mathcal{N}(i)} q(y_j) \right) f_c(y_c)$

Berkeley

### General Update \*

Exponential Family:

$$q(y_i) \propto \exp\left(\sum_{c \in \mathcal{N}(i)} w^\top \mathbb{E}_{q(y_{-i})} f_c(y_c)\right)$$

Generic:

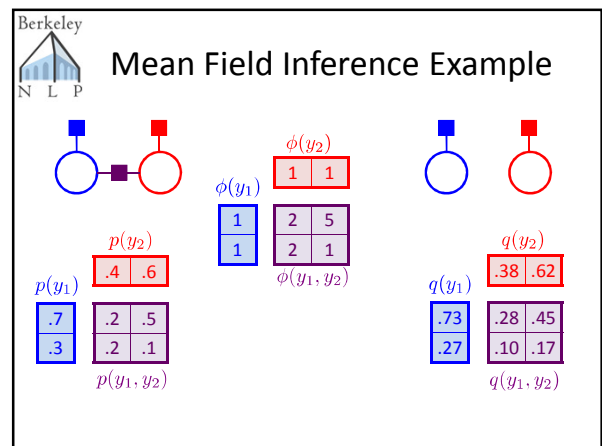
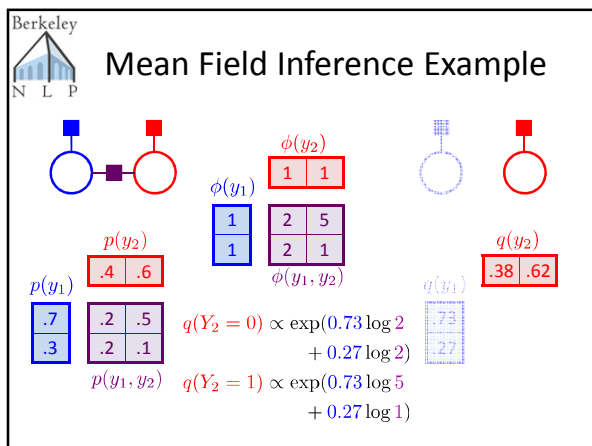
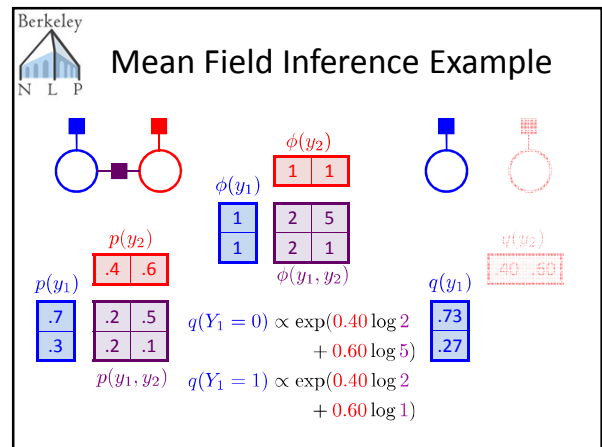
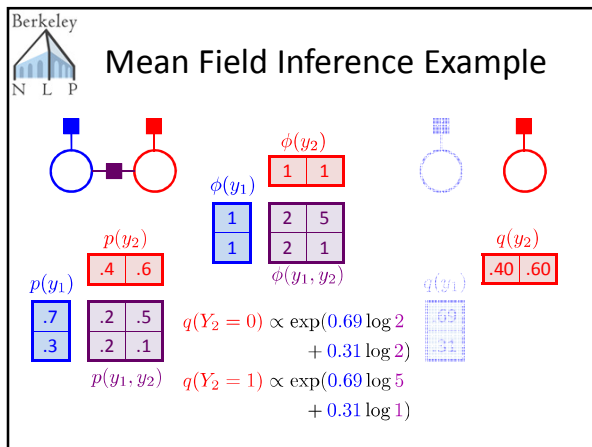
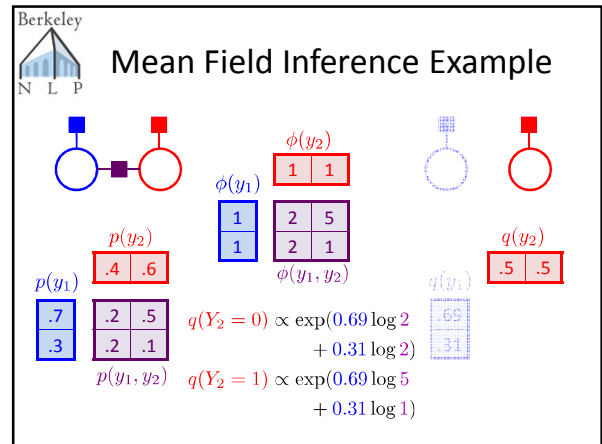
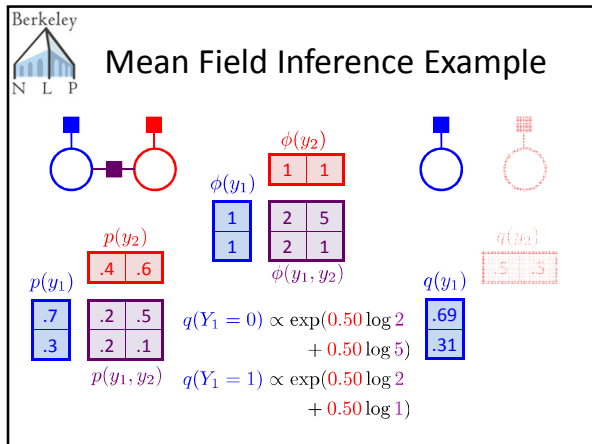
$$q(y_i) \propto \exp\left(\sum_{c \in \mathcal{N}(i)} \mathbb{E}_{q(y_{-i})} \log \phi_c(y_c)\right)$$

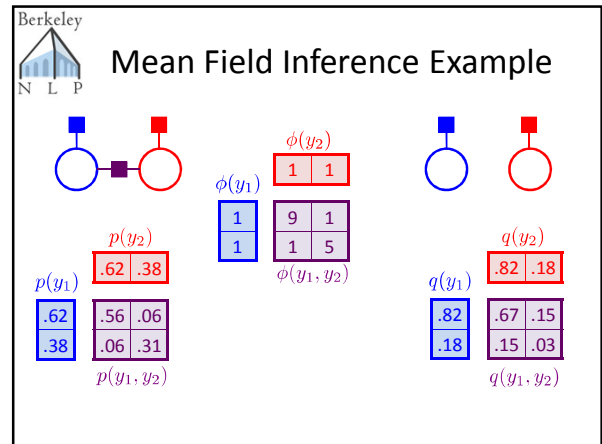
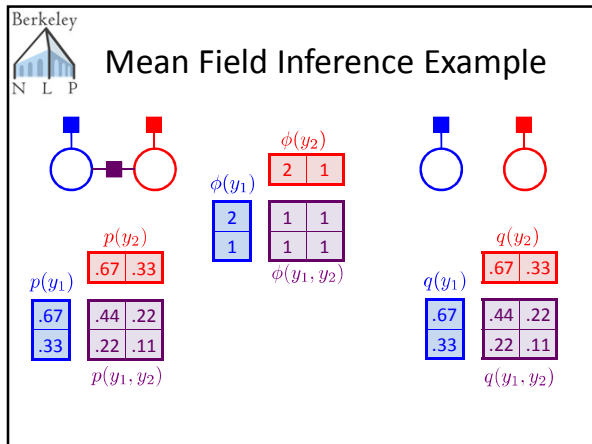
Berkeley

### Mean Field Inference Example

$p(y_1) = \begin{bmatrix} .7 & .3 \end{bmatrix}$ 
 $p(y_2) = \begin{bmatrix} .4 & .6 \end{bmatrix}$ 
 $\phi(y_1) = \begin{bmatrix} 1 & 2 \\ 1 & 2 \end{bmatrix}$ 
 $\phi(y_2) = \begin{bmatrix} 1 & 1 \\ 2 & 1 \end{bmatrix}$ 
 $q(y_1) = \begin{bmatrix} .5 & .5 \end{bmatrix}$ 
 $q(y_2) = \begin{bmatrix} .5 & .5 \end{bmatrix}$

$q(Y_1 = 0) \propto \exp(0.50 \log 2 + 0.50 \log 5)$   
 $q(Y_1 = 1) \propto \exp(0.50 \log 2 + 0.50 \log 1)$





Berkeley

### Mean Field Q&A

- Are the marginals guaranteed to converge to the right thing, like in sampling?  No
- Is the algorithm at least guaranteed to converge to something?  Yes
- So it's just like EM?  Yes

Berkeley

### Why Only Local Optima?!

Variables:  $Y_1, Y_2, \dots, Y_n$

Discrete distributions:  
e.g.  $P(0,1,0,\dots,0) = 1$

All distributions (all convex combos)

Mean field approximable (can represent all discrete ones, but not all)

### Part 3: Structured Mean Field

Berkeley

Berkeley

### Mean Field Approximation

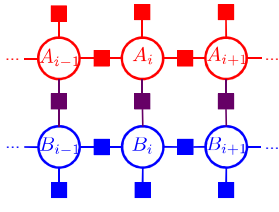
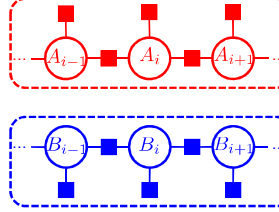
Model:

Approximate Graph:

$$p(a, b) \approx \prod_i q(a_i)q(b_i)$$

Berkeley

Structured Mean Field Approximation

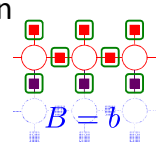
Model:  Approximate Graph: 

(Xing et al, 2003)

$$p(a, b) \approx q(a)q(b)$$

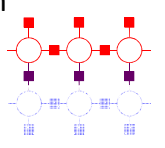
Berkeley

Structured Mean Field Approximation

$$P(a|b, x) \propto \exp \left( \sum_i w^\top f(a_i) + \sum_i w^\top f(a_{i-1}, a_i) + \sum_i w^\top f(a_i, b_i) \right)$$


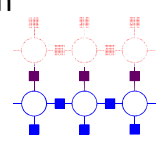
Berkeley

Structured Mean Field Approximation

$$q(a) \propto \exp \left( \sum_i w^\top f(a_i) + \sum_i w^\top f(a_{i-1}, a_i) + \sum_i w^\top \mathbb{E}_{q(b)} f(a_i, b_i) \right)$$


Berkeley

Structured Mean Field Approximation

$$q(b) \propto \exp \left( \sum_i w^\top f(b_i) + \sum_i w^\top f(b_{i-1}, b_i) + \sum_i w^\top \mathbb{E}_{q(a)} f(a_i, b_i) \right)$$


Berkeley

Computing Structured Updates

$$q(a) \propto \exp \left( \sum_i w^\top f(a_i) + \sum_i w^\top f(a_{i-1}, a_i) + \sum_i w^\top \mathbb{E}_{q(b)} f(a_i, b_i) \right)$$

$w$  ✓  
 $f(a_i)$  ✓  
 $f(a_{i-1}, a_i)$  ✓  
 $\mathbb{E}_{q(b)} f(a_i, b_i)$  ??

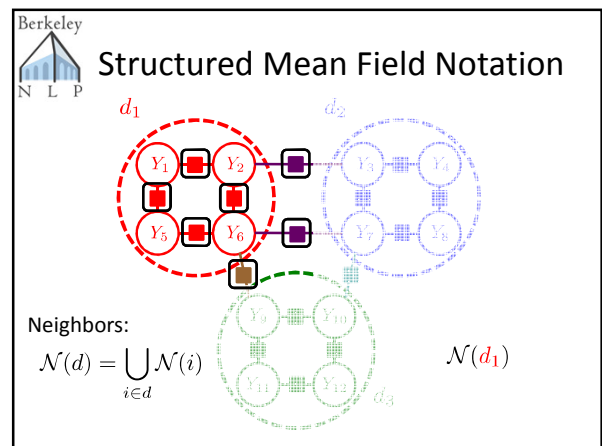
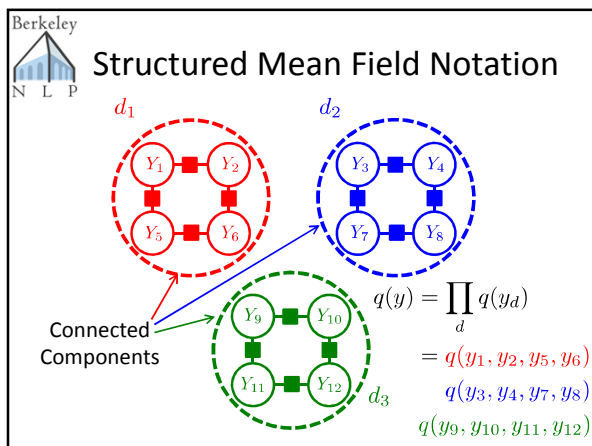
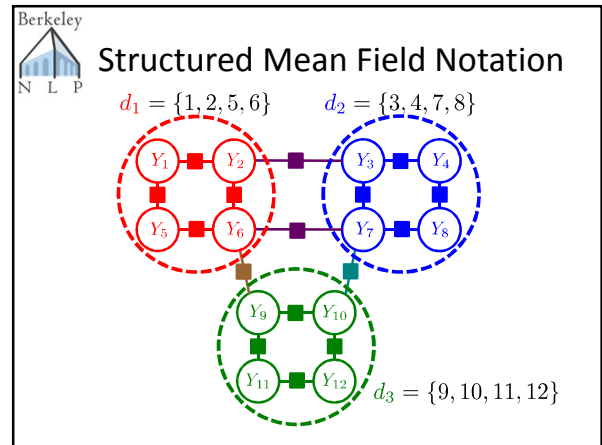
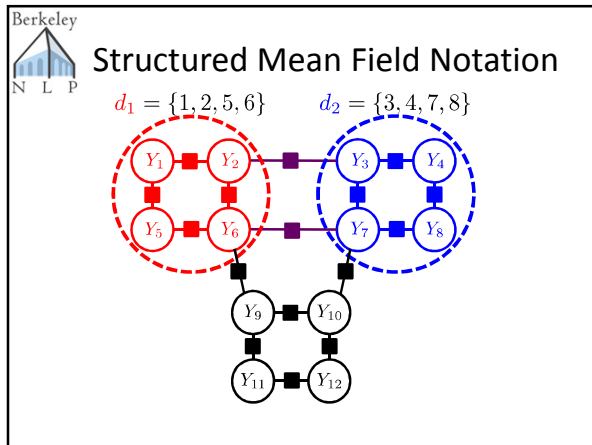
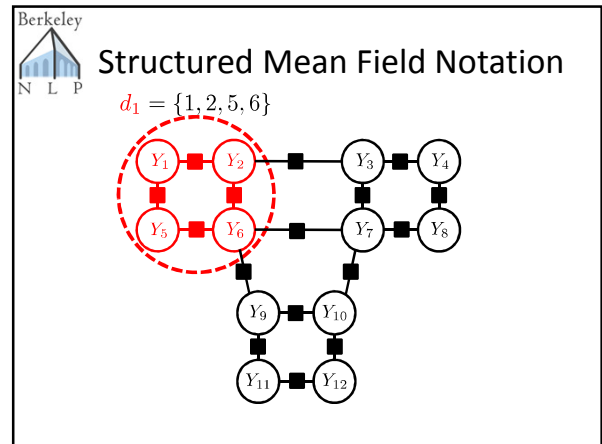
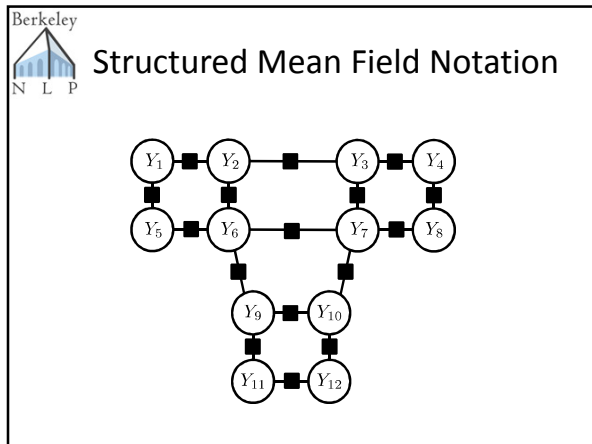
Berkeley

Computing Structured Updates

$$\mathbb{E}_{q(b)} f(a_i, b_i) = \sum_b q(b) f(a_i, b_i) = \sum_{b_i} q(b_i) f(a_i, b_i)$$

Updating  $q(a)$  consists of computing all marginals  $q(b_i)$

Marginal probability of  $b_i$  under  $q(b)$  Computed with forward-backward



Berkeley

## Structured Mean Field Updates

Naïve Mean Field:

$$q(y_i) \propto \exp \left( \sum_{c \in \mathcal{N}(i)} w^T \mathbb{E}_{q_{-i}} f(y_c) \right)$$

Structured Mean Field:

$$q(y_d) \propto \exp \left( \sum_{c \in \mathcal{N}(d)} w^T \mathbb{E}_{q_{-d}} f(y_c) \right)$$

Berkeley

## Expected Feature Counts

$q(d_1)$   
 $\mathcal{N}(d_1)$   
 $\mathbb{E}_{q_{-d_1}} f(y_6, y_7, y_9, y_{10})$   
 $= \sum_{y_7, y_9, y_{10}} f(y_6, y_7, y_9, y_{10})$

Berkeley

## Component Factorizability \*

Condition	Example Feature	Generic Condition
$f(a_i, b_i) = f(a_i)f(b_i)$	$f(a_i, b_i) = \begin{cases} 1 & a_i = \text{NNP} \ \& \ b_i = \text{B-PER} \\ 0 & \text{otherwise} \end{cases}$ $= \begin{cases} 1 & a_i = \text{NNP} \\ 0 & \text{otherwise} \end{cases} \cdot \begin{cases} 1 & b_i = \text{B-PER} \\ 0 & \text{otherwise} \end{cases}$ $= f(a_i)f(b_i)$	$f_c(y_c) = \prod_{d: c \cap d \neq \emptyset} f_{c \cap d}(y_{c \cap d})$ (pointwise product)

Berkeley

## Component Factorizability \*

(Abridged)

Use conjunctive indicator features

Berkeley

## Joint Parsing and Alignment

(Burkett et al, 2010)

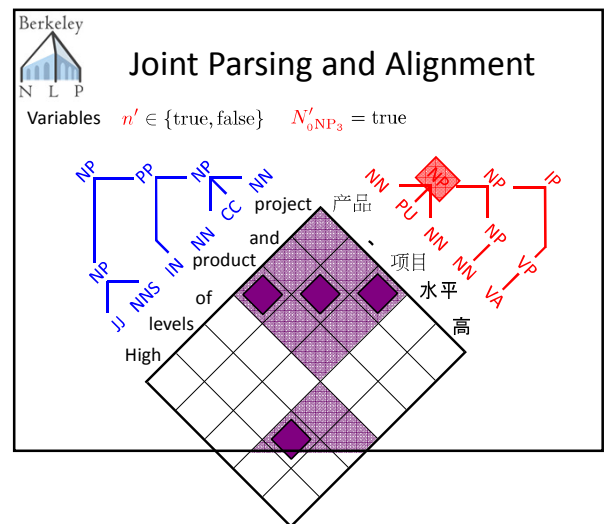
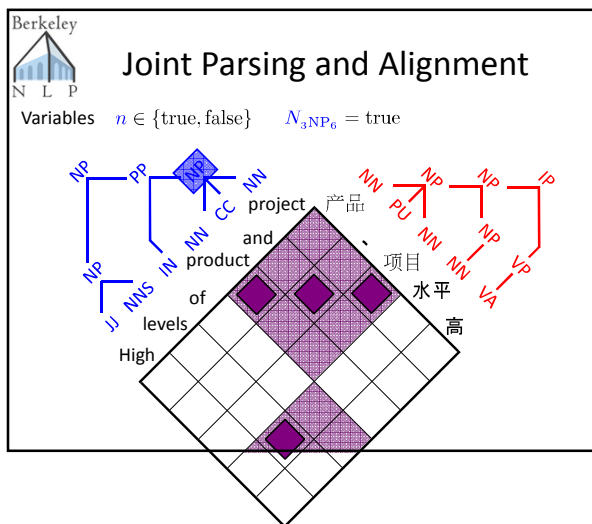
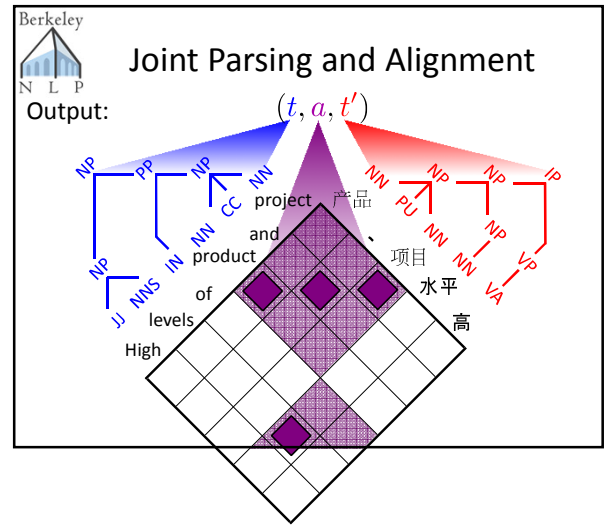
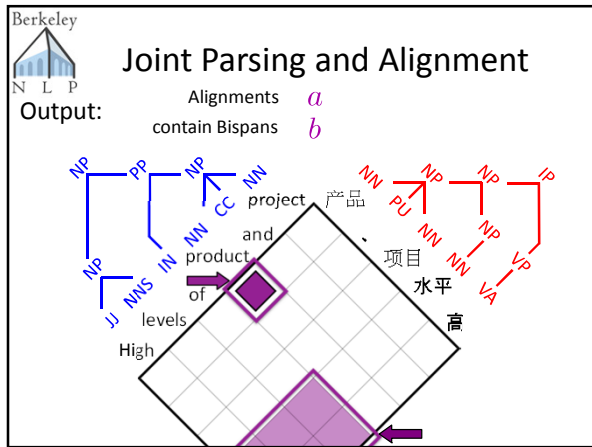
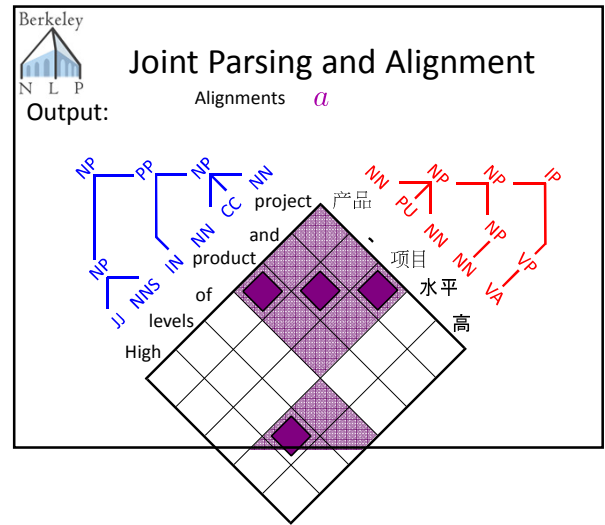
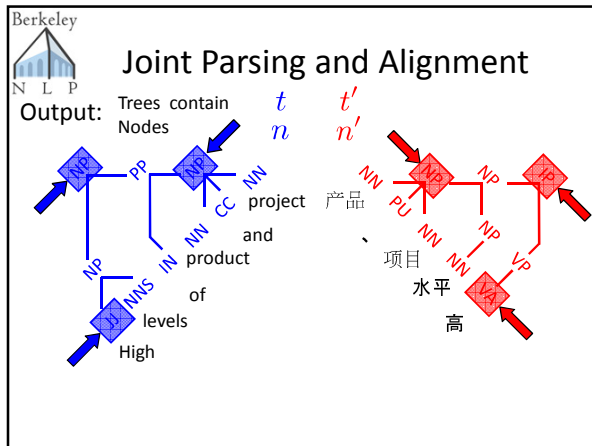
Berkeley

## Joint Parsing and Alignment

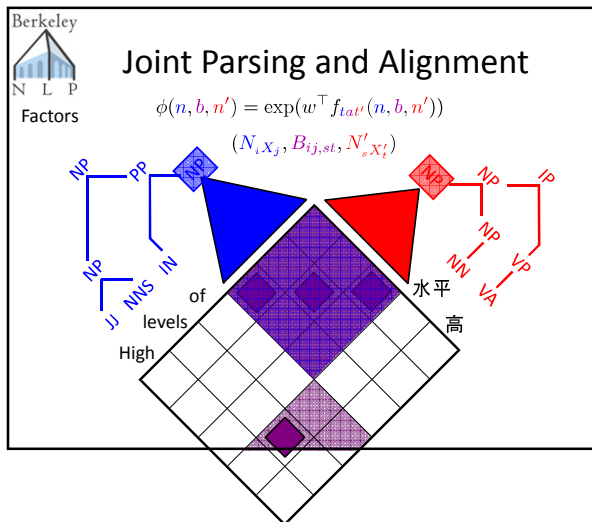
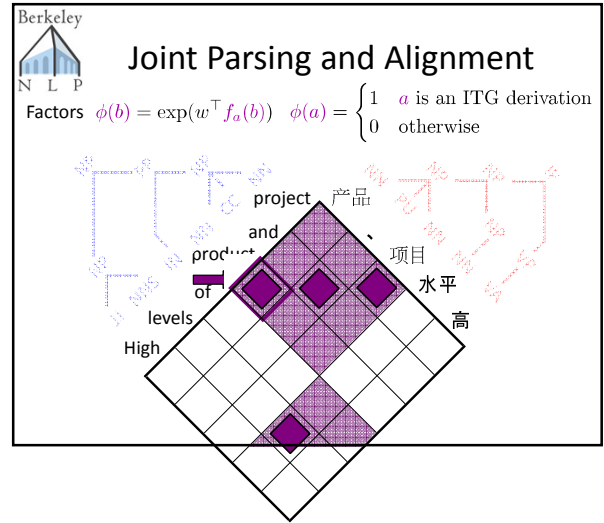
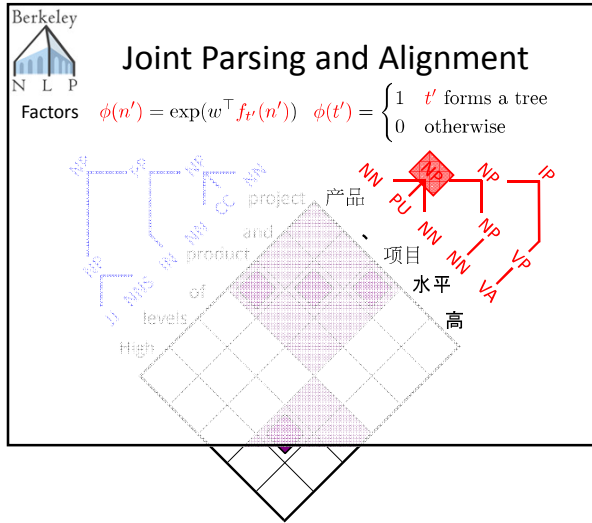
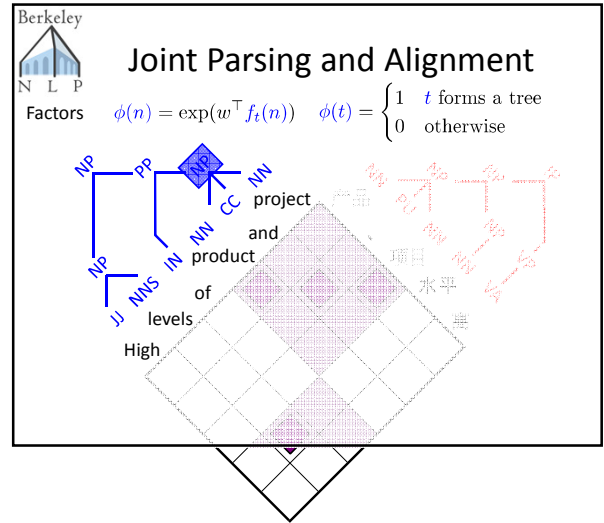
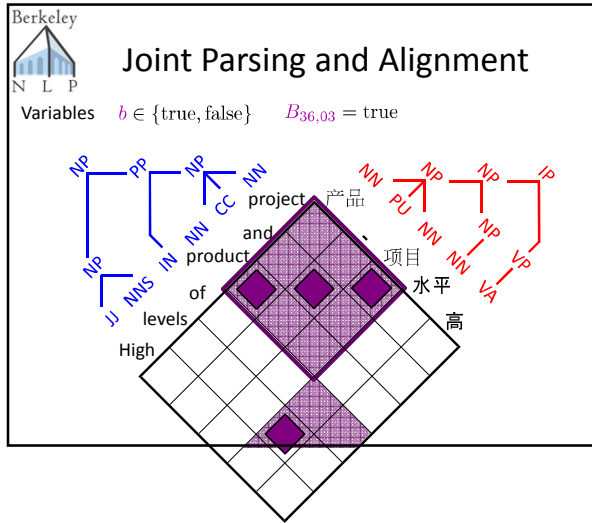
Input:

Sentences  $(s, s')$

project 产品  
and 项目  
product 项目  
of 水平  
levels 高  
High







Berkeley

N L P

### Notational Abuse


Subscript Omission:  
 $f_t(n) = f_t(n_i X_j)$

Shorthand:  
 $n \in t \Leftrightarrow N_i X_j = \text{true}$   
 $n \triangleright b \triangleleft n' \Leftrightarrow n \in t \ \& \ b \in a \ \& \ n' \in t' \ \& \ (N_i X_j, B_{ij, st}, N'_s X'_t) \text{ match up}$

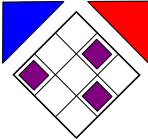
Skip Nonexistent Substructures:  
 $n \notin t \Rightarrow f_t(n) = 0$

Structural factors  $\phi(t), \phi(a), \phi(t')$  are implicit


Berkeley



## Model Form

$$P(t, a, t' | s, s') \propto \exp \left( \sum_{n \in t} w^\top f_t(n) + \sum_{b \in a} w^\top f_a(b) + \sum_{n' \in t'} w^\top f_{t'}(n') + \sum_{n \triangleright b \triangleleft n'} w^\top f_{ta t'}(n, b, n') \right)$$



Berkeley



## Training

Expected Feature Counts	Marginals
$\mathbb{E} f_t(n)$	$P(n \in t   s, s')$
$\mathbb{E} f_a(b)$	$P(b \in a   s, s')$
$\mathbb{E} f_{t'}(n')$	$P(n' \in t'   s, s')$
$\mathbb{E} f_{ta t'}(n, b, n')$	$P(n \triangleright b \triangleleft n'   s, s')$

Berkeley




## Structured Mean Field Approximation

$$P(t, a, t' | s, s') \propto \exp \left( \sum_{n \in t} w^\top f_t(n) + \sum_{b \in a} w^\top f_a(b) + \sum_{n' \in t'} w^\top f_{t'}(n') + \sum_{n \triangleright b \triangleleft n'} w^\top f_{ta t'}(n, b, n') \right)$$

$$\approx q(t)q(a)q(t')$$

Berkeley




## Approximate Component Scores

Monolingual parser:  
Score for  $n = w^\top f_t(n)$

If we knew  $(a, t')$ :  
Score for  $n = w^\top f_t(n) + w^\top f_{ta t'}(n, b, n')$

To compute  $q(t)$ :  
Score for  $n = w^\top f_t(n) + w^\top \mathbb{E}_{q(a, t')} f_{ta t'}(n, b, n')$

Berkeley



## Expected Feature Counts

For fixed  $n_i, X_j$ :

$$\begin{aligned} & \mathbb{E}_{q(a, t')} f_{ta t'}(n, b, n') \\ &= \sum_{s' X'_t} P_q(n_i, X_j \triangleright b_{ij, st} \triangleleft n'_s X'_t) f_{ta t'}(n, b, n') \\ &= \sum_{s' X'_t} q(b_{ij, st}) q(n'_s X'_t) f_{ta t'}(n, b, n') \end{aligned}$$


↑

Marginals computed with bitext inside-outside

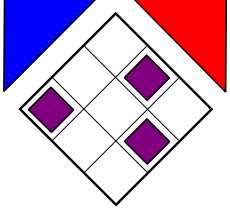
↑

Marginals computed with inside-outside

Berkeley



## Inference Procedure



Initialize:

$$q(t) \propto \exp \left( \sum_{n \in t} w^\top f_t(n) \right)$$

$$q(a) \propto \exp \left( \sum_{b \in a} w^\top f_a(b) \right)$$

$$q(t') \propto \exp \left( \sum_{n' \in t'} w^\top f_{t'}(n') \right)$$

Berkeley  
N L P

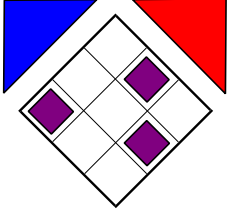
## Inference Procedure

Iterate marginal updates:

$q(n)$

$q(b)$

$q(n')$



...until convergence!

Berkeley  
N L P

## Approximate Marginals

$P(n \in t | s, s') \approx q(n)$

$P(b \in a | s, s') \approx q(b)$

$P(n' \in t' | s, s') \approx q(n')$

$P(n \triangleright b \triangleleft n' | s, s') \approx q(n)q(b)q(n')$

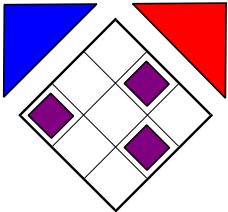
Berkeley  
N L P

## Decoding

$\hat{t} = \operatorname{argmax}_t q(t)$

$\hat{a} = \operatorname{argmax}_a q(a)$

$\hat{t}' = \operatorname{argmax}_{t'} q(t')$



(Minimum Risk)

Berkeley  
N L P

## Structured Mean Field Summary

- ▶ Split the model into pieces you have dynamic programs for
- ▶ Substitute expected feature counts for actual counts in cross-component factors
- ▶ Iterate computing marginals until convergence

Berkeley  
N L P

## Structured Mean Field Tips

- ▶ Try to make sure cross-component features are products of indicators
- ▶ You don't have to run all the way to convergence; marginals are usually pretty good after just a few rounds
- ▶ Recompute marginals for fast components more frequently than for slow ones
  - ▶ e.g. For joint parsing and alignment, the two monolingual tree marginals ( $O(n^3)$ ) were updated until convergence between each update of the ITG marginals ( $O(n^6)$ )

Berkeley  
N L P

## Break Time!



## Part 4: Belief Propagation

Berkeley

## Belief Propagation

Wanted:  $P(a|x), P(b|x)$

Idea: pretend graph is a tree

Key objects:  
Beliefs (marginals)  
Messages

## Belief Propagation Intro

Assume we have a tree

$$\begin{aligned}
 P(a|x) &\propto \sum_{y \setminus \{a\}} \text{score}(y) \\
 &= \left( \sum_{\mathcal{R}} \text{score}(\mathcal{R}, a) \right) \left( \sum_{\mathcal{G}} \text{score}(\mathcal{G}, a) \right) \left( \sum_{\mathcal{B}} \text{score}(\mathcal{B}, a) \right) \\
 &= m_{\mathcal{R} \rightarrow A}(a) \cdot m_{\mathcal{G} \rightarrow A}(a) \cdot m_{\mathcal{B} \rightarrow A}(a)
 \end{aligned}$$

## Belief Propagation Intro

$$\begin{aligned}
 m_{\phi \rightarrow A}(a) &= \sum_{\mathcal{R}} \text{score}(\mathcal{R}, a) \\
 &= \sum_r \phi(r, a) m_{\phi_1 \rightarrow R}(r) m_{\phi_2 \rightarrow R}(r) \\
 &= \sum_r \phi(r, a) m_{R \rightarrow \phi}(r) \\
 m_{R \rightarrow \phi}(r) &= m_{\phi_1 \rightarrow R}(r) m_{\phi_2 \rightarrow R}(r)
 \end{aligned}$$

## Messages

Variable to Factor

Factor to Variable

Both take form of "distribution" over  $Y_i$

## Messages General Form

▶ Messages from variables to factors:

$$m_{Y_i \rightarrow \phi_c}(y_i) \propto \prod_{c' \in \mathcal{N}(i) \setminus \{c\}} m_{\phi_{c'} \rightarrow Y_i}(y_i)$$

Berkeley  
N L P

## Messages General Form

▶ Messages from factors to variables:

$$m_{\phi_c \rightarrow Y_i}(y_i) \propto \sum_{y_{c \setminus \{i\}}} \phi_c(y_c) \prod_{i' \in c \setminus \{i\}} m_{Y_{i'} \rightarrow \phi_c}(y_{i'})$$

Berkeley  
N L P

## Marginal Beliefs

$$b_{Y_i}(y_i) \propto \prod_{c \in \mathcal{N}(i)} m_{\phi_c \rightarrow Y_i}(y_i)$$

Berkeley  
N L P

## Belief Propagation on Tree-Structured Graphs

▶ If the factor graph has no cycles, BP is exact

- ▶ Can always order message computations

▶ After one pass, marginal beliefs are correct

Berkeley  
N L P

## “Loopy” Belief Propagation

Problem: we no longer have a tree  
Solution: ignore problem

Berkeley  
N L P

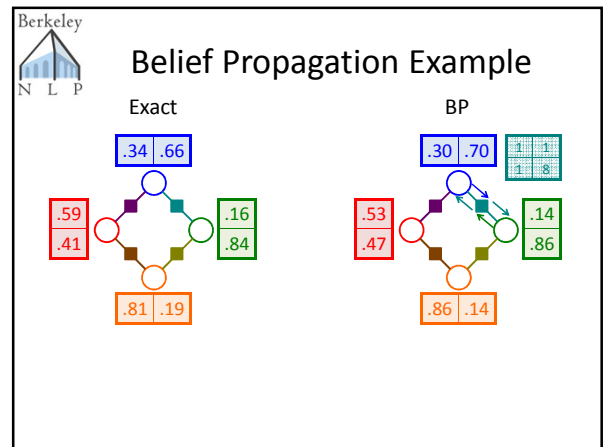
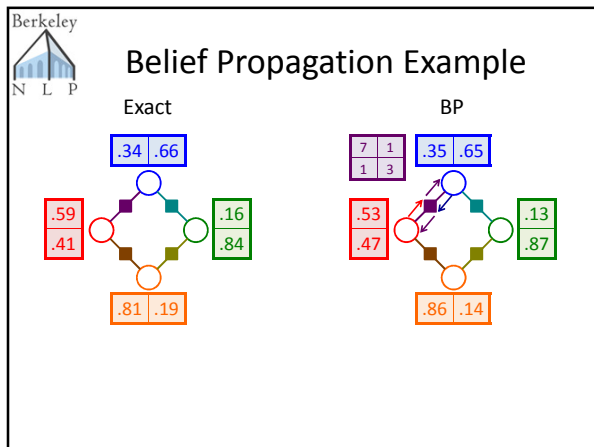
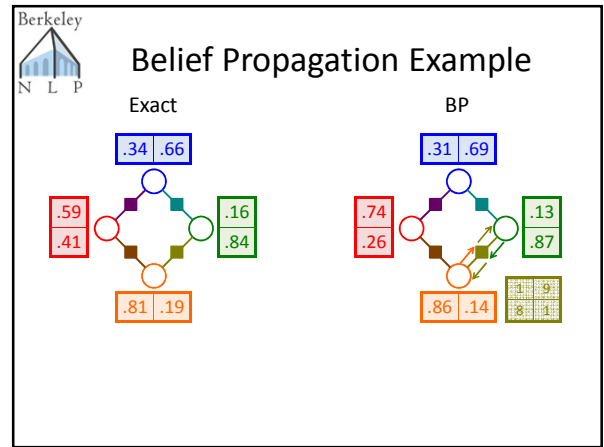
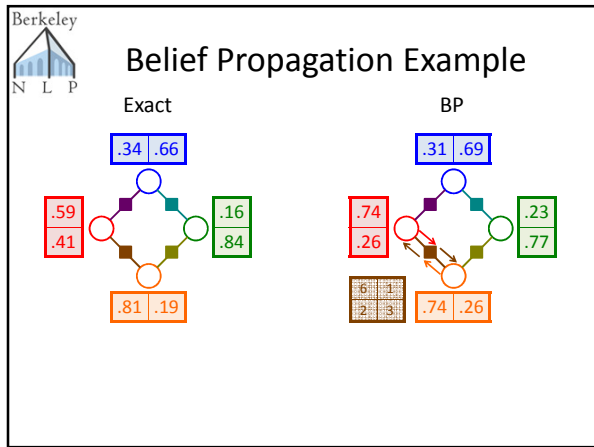
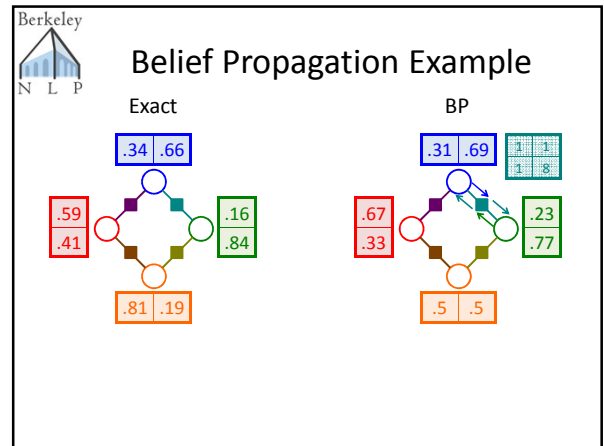
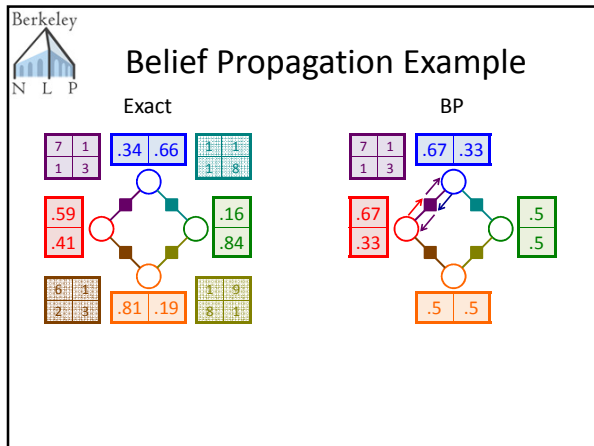
## “Loopy” Belief Propagation

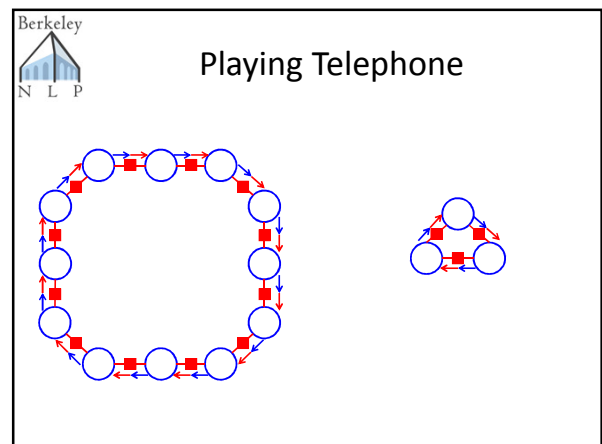
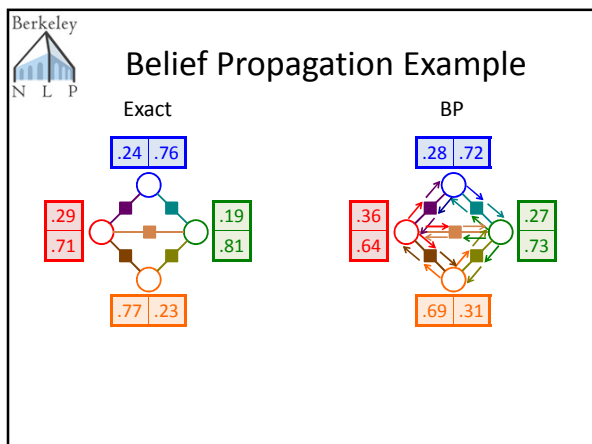
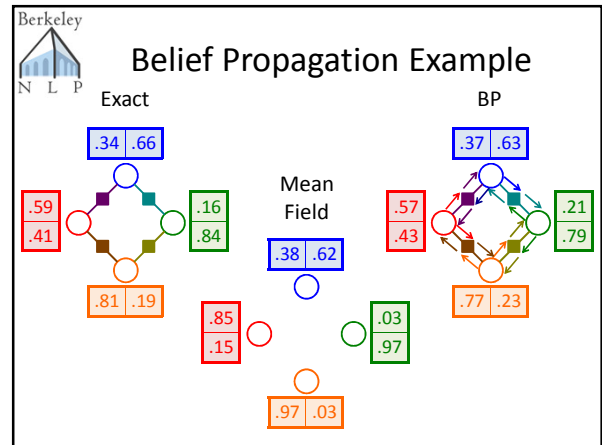
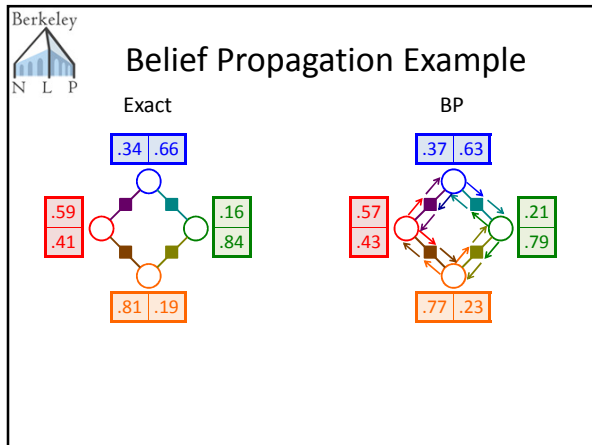
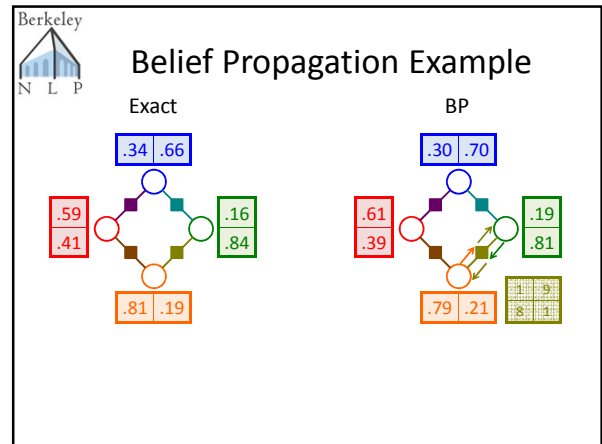
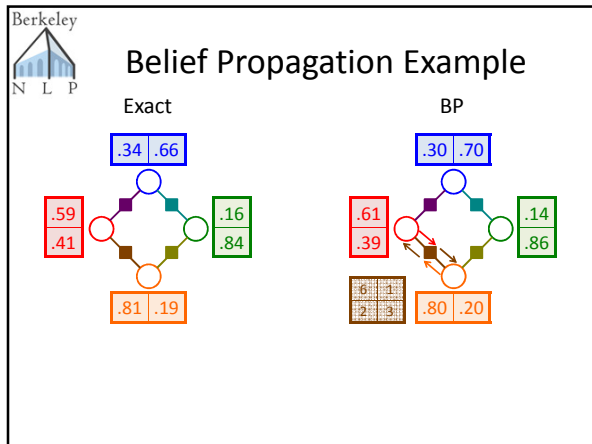
Just start passing messages anyway!

Berkeley  
N L P

## Belief Propagation Q&A

- ▶ Are the marginals guaranteed to converge to the right thing, like in sampling? No
- ▶ Well, is the algorithm at least guaranteed to converge to something, like mean field? No
- ▶ Will everything often work out more or less OK in practice? Maybe





## Part 5: Belief Propagation with Structured Factors



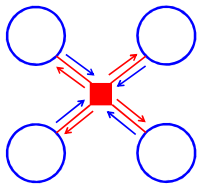
## Structured Factors

- ▶ **Problem:**
  - ▶ Computing factor messages is exponential in arity
  - ▶ Many models we care about have high-arity factors
- ▶ **Solution:**
  - ▶ Take advantage of NLP tricks for efficient sums
- ▶ **Examples:**
  - ▶ Word Alignment (at-most-one constraints)
  - ▶ Dependency Parsing (tree constraint)



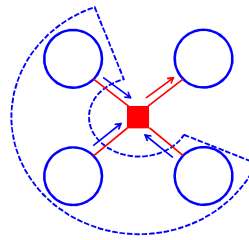
## Warm-up Exercise

$$m_{\phi_c \rightarrow y_i}(y_i) \propto \sum_{\mathbf{y}_{c \setminus \{i\}}} \phi_c(\mathbf{y}_c) \prod_{i' \in c \setminus \{i\}} m_{Y_{i'} \rightarrow \phi_c}(y_{i'})$$



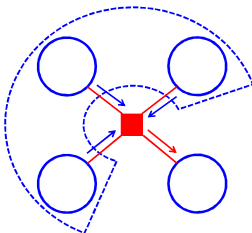
## Warm-up Exercise

$$m_{\phi_c \rightarrow y_i}(y_i) \propto \sum_{\mathbf{y}_{c \setminus \{i\}}} \phi_c(\mathbf{y}_c) \prod_{i' \in c \setminus \{i\}} m_{Y_{i'} \rightarrow \phi_c}(y_{i'})$$



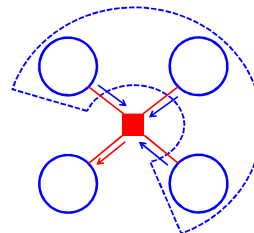
## Warm-up Exercise

$$m_{\phi_c \rightarrow y_i}(y_i) \propto \sum_{\mathbf{y}_{c \setminus \{i\}}} \phi_c(\mathbf{y}_c) \prod_{i' \in c \setminus \{i\}} m_{Y_{i'} \rightarrow \phi_c}(y_{i'})$$



## Warm-up Exercise

$$m_{\phi_c \rightarrow y_i}(y_i) \propto \sum_{\mathbf{y}_{c \setminus \{i\}}} \phi_c(\mathbf{y}_c) \prod_{i' \in c \setminus \{i\}} m_{Y_{i'} \rightarrow \phi_c}(y_{i'})$$





Berkeley  
N L P

### Warm-up Exercise

$$m_{\phi_c \rightarrow y_c}(y_c) \propto \sum_{\{y_{i'}\}} \phi_c(y_c) \prod_{i' \in c \setminus \{i\}} m_{Y_{i'} \rightarrow \phi_c}(y_{i'})$$

Berkeley  
N L P

### Warm-up Exercise

$$m_{\phi_c \rightarrow y_c}(y_c) \propto \sum_{\{y_{i'}\}} \phi_c(y_c) \prod_{i' \in c \setminus \{i\}} m_{Y_{i'} \rightarrow \phi_c}(y_{i'})$$

$$s(y_c) = \prod_{i \in c} m_{Y_i \rightarrow \phi_c}(y_i)$$

Berkeley  
N L P

### Warm-up Exercise

$$m_{\phi_c \rightarrow y_c}(y_c) \propto \sum_{\{y_{i'}\}} \phi_c(y_c) \prod_{i' \in c \setminus \{i\}} m_{Y_{i'} \rightarrow \phi_c}(y_{i'})$$

$$= \frac{s(y_c)}{m_{Y_i \rightarrow \phi_c}(y_i)}$$

$$s(y_c) = \prod_{i \in c} m_{Y_i \rightarrow \phi_c}(y_i)$$

Berkeley  
N L P

### Warm-up Exercise

$$m_{\phi_c \rightarrow y_c}(y_c) \propto \sum_{\{y_{i'}\}} \phi_c(y_c) \prod_{i' \in c \setminus \{i\}} m_{Y_{i'} \rightarrow \phi_c}(y_{i'})$$

$$= \frac{s(y_c)}{m_{Y_i \rightarrow \phi_c}(y_i)}$$

► **Benefits:**

- Cleans up notation
- Saves time multiplying
- Enables efficient summing

$$s(y_c) = \prod_{i \in c} m_{Y_i \rightarrow \phi_c}(y_i)$$

Berkeley  
N L P

### The Shape of Structured BP

- Isolate the combinatorial factors
- Figure out how to compute efficient sums
  - Directly exploiting sparsity
  - Dynamic programming
- Work out the bookkeeping
  - Or, use a reference!

Berkeley  
N L P

### Word Alignment with BP


$$\phi(y_{ij}) = \begin{cases} \exp(w^T f(i, j)) & y_{ij} = \text{on} \\ 1 & y_{ij} = \text{off} \end{cases} \quad y_{ij} \in \{\text{on}, \text{off}\}$$

$$\phi(y_{i*}) = \begin{cases} 1 & |\{j : y_{ij} = \text{on}\}| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

$$\phi(y_{*j}) = \begin{cases} 1 & |\{i : y_{ij} = \text{on}\}| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$



(Cromières & Kurohashi, 2009)  
(Burkett & Klein, 2012)

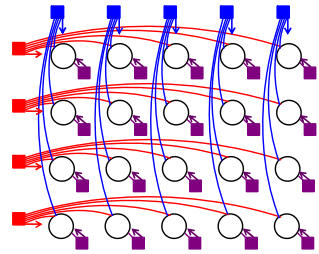
Berkeley




### Computing Messages from Factors

Exponential in arity of factor  
(have to sum over all assignments)

$\phi_{ij}(y_{ij})$  Arity 1   
 $\phi_i(y_{i*})$  Arity  $O(n)$   
 $\phi_j(y_{*j})$  Arity  $O(n)$  

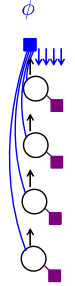


Berkeley




### Computing Constraint Factor Messages

▶ Input:  $m_{Y_j \rightarrow \phi}(y_j) \forall j$   
 ▶ Goal:  $m_{\phi \rightarrow Y_j}(y_j) \forall j$

$$m_{\phi \rightarrow Y_j}(y_j) \propto \frac{\sum_{y: Y_j=y_j} \phi(y) s(y)}{m_{Y_j \rightarrow \phi}(y_j)}$$



Berkeley




### Computing Constraint Factor Messages

$y(j)$ : Assignment to variables where  $Y_j = \text{on}$

$y(2) = \{Y_1 = \text{off},$   
 $Y_2 = \text{on},$   
 $Y_3 = \text{off},$   
 $Y_4 = \text{off}\}$

$$\phi(y) = \begin{cases} 1 & |\{j : y_j = \text{on}\}| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$


Berkeley

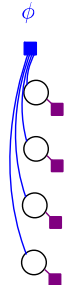


### Computing Constraint Factor Messages


$y(j)$ : Assignment to variables where  $Y_j = \text{on}$

$y(0)$ : Special case for all off

$y(0) = \{Y_1 = \text{off},$   
 $Y_2 = \text{off},$   
 $Y_3 = \text{off},$   
 $Y_4 = \text{off}\}$

$$\phi(y) = \begin{cases} 1 & |\{j : y_j = \text{on}\}| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$


Berkeley




### Computing Constraint Factor Messages


▶ Input:  $m_{Y_j \rightarrow \phi}(y_j) \forall j$   
 ▶ Goal:  $m_{\phi \rightarrow Y_j}(y_j) \forall j$

$$m_{\phi \rightarrow Y_j}(y_j) \propto \frac{\sum_{y: Y_j=y_j} \phi(y) s(y)}{m_{Y_j \rightarrow \phi}(y_j)}$$

Only need to consider  $y(j')$  for  $0 \leq j' \leq n$



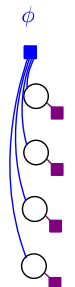
Berkeley



### Computing Constraint Factor Messages

$$s(y) = \prod_j m_{Y_j \rightarrow \phi}(y_j)$$

$$s(y(0)) = m_{Y_1 \rightarrow \phi}(\text{off}) \cdot m_{Y_2 \rightarrow \phi}(\text{off}) \cdot m_{Y_3 \rightarrow \phi}(\text{off}) \cdot m_{Y_4 \rightarrow \phi}(\text{off})$$

$$s(y(0)) = \prod_{1 \leq j \leq n} m_{Y_j \rightarrow \phi}(\text{off})$$


Berkeley  
N L P

### Computing Constraint Factor Messages

$$s(y) = \prod_j m_{Y_i \rightarrow \phi}(y_j)$$

$$s(y(1)) = m_{Y_1 \rightarrow \phi}(\text{on}) \cdot m_{Y_2 \rightarrow \phi}(\text{off}) \cdot m_{Y_3 \rightarrow \phi}(\text{off}) \cdot m_{Y_4 \rightarrow \phi}(\text{off})$$

Berkeley  
N L P

### Computing Constraint Factor Messages

$$s(y) = \prod_j m_{Y_i \rightarrow \phi}(y_j)$$

$$s(y(2)) = m_{Y_1 \rightarrow \phi}(\text{off}) \cdot m_{Y_2 \rightarrow \phi}(\text{on}) \cdot m_{Y_3 \rightarrow \phi}(\text{off}) \cdot m_{Y_4 \rightarrow \phi}(\text{off})$$

Berkeley  
N L P

### Computing Constraint Factor Messages

$$s(y) = \prod_j m_{Y_i \rightarrow \phi}(y_j)$$

$$s(y(3)) = m_{Y_1 \rightarrow \phi}(\text{off}) \cdot m_{Y_2 \rightarrow \phi}(\text{off}) \cdot m_{Y_3 \rightarrow \phi}(\text{on}) \cdot m_{Y_4 \rightarrow \phi}(\text{off})$$

Berkeley  
N L P

### Computing Constraint Factor Messages

$$s(y) = \prod_j m_{Y_i \rightarrow \phi}(y_j)$$

$$s(y(4)) = m_{Y_1 \rightarrow \phi}(\text{off}) \cdot m_{Y_2 \rightarrow \phi}(\text{off}) \cdot m_{Y_3 \rightarrow \phi}(\text{off}) \cdot m_{Y_4 \rightarrow \phi}(\text{on})$$

Berkeley  
N L P

### Computing Constraint Factor Messages

$$s(y) = \prod_j m_{Y_i \rightarrow \phi}(y_j)$$

$$s(y(0)) = \prod_{1 \leq j \leq n} m_{Y_j \rightarrow \phi}(\text{off})$$

$\forall j > 0 :$

$$s(y(j)) = s(y(0)) \frac{m_{Y_j \rightarrow \phi}(\text{on})}{m_{Y_j \rightarrow \phi}(\text{off})}$$

Berkeley  
N L P


### Computing Constraint Factor Messages

$$m_{\phi \rightarrow Y_1}(\text{on}) \propto \frac{s(y(1))}{m_{Y_1 \rightarrow \phi}(\text{on})}$$

$$m_{\phi \rightarrow Y_1}(\text{off}) \propto \frac{s(*) - s(y(1))}{m_{Y_1 \rightarrow \phi}(\text{off})}$$


$$s(*) = \sum_{0 \leq j \leq n} s(y(j))$$

Berkeley




## Computing Constraint Factor Messages

1. Precompute:  $s(y(0)) = \prod_{1 \leq j \leq n} m_{Y_j \rightarrow \phi}(\text{off})$   
 $O(n)$
2.  $\forall j > 0: s(y(j)) = s(y(0)) \frac{m_{Y_j \rightarrow \phi}(\text{on})}{m_{Y_j \rightarrow \phi}(\text{off})}$   
 $O(n)$
3. Partition:  $s(*) = \sum_{0 \leq j \leq n} s(y(j))$
4. Messages:  $m_{\phi \rightarrow Y_j}(\text{on}) \propto \frac{s(y(j))}{m_{Y_j \rightarrow \phi}(\text{on})}$   
 $m_{\phi \rightarrow Y_j}(\text{off}) \propto \frac{s(*) - s(y(j))}{m_{Y_j \rightarrow \phi}(\text{off})}$



Berkeley

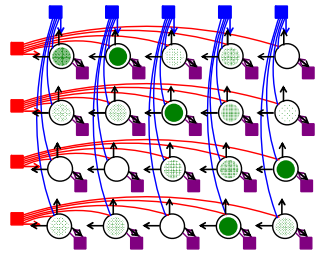


## Using BP Marginals


$P(y_{ij}|x) \approx b_{Y_{ij}}(y_{ij})$

Expected Feature Counts:  
 $\mathbb{E}f(i, j) \approx b_{Y_{ij}}(\text{on})f(i, j)$

Marginal Decoding:  
 $\hat{y}_{ij} = \begin{cases} \text{on} & b_{Y_{ij}}(\text{on}) \geq \tau \\ \text{off} & \text{otherwise} \end{cases}$



Berkeley



## Dependency Parsing with BP

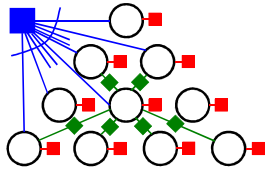
$y_{ij} \in \{\text{left}, \text{right}, \text{off}\}$

$\phi(y) = \begin{cases} 1 & y \text{ forms a tree} \\ 0 & \text{otherwise} \end{cases}$


$\phi(y_{ij}) = \begin{cases} \exp(w^\top f(i, j)) & y_{ij} = \text{left} \\ \exp(w^\top f(j, i)) & y_{ij} = \text{right} \\ 1 & y_{ij} = \text{off} \end{cases}$

$\phi(y_{ij}, y_{jk})$

(Smith & Eisner, 2008)  
(Martins et al., 2010)



Berkeley



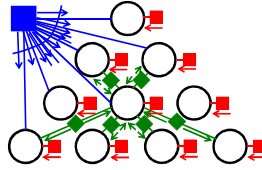
## Dependency Parsing with BP

$\phi_{ij}(y_{ij})$  Arity 1 ✓


$\phi_{ijk}(y_{ij}, y_{jk})$  Arity 2 ✓

$\phi_{\text{TREE}}(y)$  Arity  $O(n^2)$  ✗

Exponential in arity of factor



Berkeley



## Messages from the Tree Factor

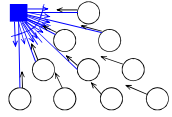
▶ Input:  $m_{Y_{ij} \rightarrow \phi_{\text{TREE}}}(y_{ij})$  for all variables

▶ Goal:  $m_{\phi_{\text{TREE}} \rightarrow Y_{ij}}(y_{ij})$  for all variables


$m_{\phi_{\text{TREE}} \rightarrow Y_{ij}}(y_{ij}) \propto \sum \phi_{\text{TREE}}(y) s(y) \frac{1}{m_{Y_{ij} \rightarrow \phi_{\text{TREE}}}(y_{ij})}$

$\phi_{\text{TREE}}(y) = \begin{cases} 1 & y \text{ forms a tree} \\ 0 & \text{otherwise} \end{cases}$

$T = \{y : y \text{ forms a tree}\}$



Berkeley



## What Do Parsers Do?


▶ Initial state:

- ▶ Value of an edge ( $i$  has parent  $j$ ):  $v(i, j)$
- ▶ Value of a tree:  $v(t) = \prod_{(i,j) \in t} v(i, j)$

▶ Run inside-outside to compute:

- ▶ Total score for all trees:  $Z = \sum_t v(t)$
- ▶ Total score for an edge:  $Z(i, j) = \sum_{t: (i,j) \in t} v(t)$

Berkeley (Klein & Manning, 2002)



## Initializing the Parser

**Problem:**

$$v(t) = \prod_{(i,j) \in t} v(i,j) \qquad s(y) = \prod_{ij} m_{Y_{ij} \rightarrow \phi_{\text{TREE}}}(y_{ij})$$

Product over edges in  $t$ :  $y_{ij} = \text{left or } y_{ji} = \text{right}$       Product over ALL edges, including  $y_{ij} = \text{off}$


**Solution: Use odds ratios**

$$v(y_{ij}) = \begin{cases} \frac{m_{Y_{ij} \rightarrow \phi_{\text{TREE}}(\text{left})}}{m_{Y_{ij} \rightarrow \phi_{\text{TREE}}(\text{off})}} & y_{ij} \neq \text{off} \\ 1 & y_{ij} = \text{off} \end{cases}$$

$$\pi = \prod_{ij} m_{Y_{ij} \rightarrow \phi_{\text{TREE}}}(\text{off})$$

$$\pi v(t) = s(y)$$

Berkeley



## Running the Parser


$$Z = \sum_t v(t) \qquad \pi v(t) = s(y) \qquad \pi Z = \sum_{y \in T} s(y)$$

**Sums we want:**

$$\pi Z(i,j) = \sum_{\substack{y \in T \\ y_{ij} = \text{left}}} s(y) \qquad \pi Z(j,i) = \sum_{\substack{y \in T \\ y_{ji} = \text{right}}} s(y)$$

$$\pi(Z - Z(i,j) - Z(j,i)) = \sum_{\substack{y \in T \\ y_{ij} = \text{off}}} s(y)$$


Berkeley



## Computing Tree Factor Messages

- Precompute:  $\pi = \prod_{ij} m_{Y_{ij} \rightarrow \phi_{\text{TREE}}}(\text{off})$
- Initialize:  $v(i,j) = \begin{cases} \frac{m_{Y_{ij} \rightarrow \phi_{\text{TREE}}(\text{left})}}{m_{Y_{ij} \rightarrow \phi_{\text{TREE}}(\text{off})}} & i < j \\ \frac{m_{Y_{ji} \rightarrow \phi_{\text{TREE}}(\text{right})}}{m_{Y_{ji} \rightarrow \phi_{\text{TREE}}(\text{off})}} & j < i \end{cases}$
- Run inside-outside
- Messages:  $m_{\phi_{\text{TREE}} \rightarrow Y_{ij}}(y_{ij}) \propto \begin{cases} \frac{\pi Z(i,j)}{m_{Y_{ij} \rightarrow \phi_{\text{TREE}}}(y_{ij})} & y_{ij} = \text{left} \\ \frac{\pi Z(j,i)}{m_{Y_{ij} \rightarrow \phi_{\text{TREE}}}(y_{ij})} & y_{ij} = \text{right} \\ \frac{\pi(Z - Z(i,j) - Z(j,i))}{m_{Y_{ij} \rightarrow \phi_{\text{TREE}}}(y_{ij})} & y_{ij} = \text{off} \end{cases}$

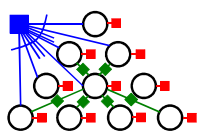
Berkeley




## Using BP Marginals

$$P(y_{ij}|x) \approx b_{Y_{ij}}(y_{ij})$$

- Expected Feature Counts:  $\mathbb{E}f(i,j) \approx \begin{cases} b_{Y_{ij}}(\text{left})f(i,j) & i < j \\ b_{Y_{ji}}(\text{right})f(i,j) & j < i \end{cases}$
- Minimum Risk Decoding:
  - Initialize:  $v(i,j) = \begin{cases} \frac{b_{Y_{ij}}(\text{left})}{b_{Y_{ij}}(\text{off})} & i < j \\ \frac{b_{Y_{ji}}(\text{right})}{b_{Y_{ji}}(\text{off})} & j < i \end{cases}$
  - Run parser:  $\hat{t} = \text{argmax}_t s(t)$




Berkeley



## Structured BP Summary

- Tricky part is factors whose arity grows with input size
- Simplify the problem by focusing on sums of total scores
- Exploit problem-specific structure to compute sums efficiently
- Use odds ratios to eliminate “default” values that don’t appear in dynamic program sums

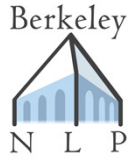
Berkeley



## Belief Propagation Tips

- Don’t compute unary messages multiple times
- Store variable beliefs to save time computing variable to factor messages (divide one out)
- Update the slowest messages less frequently
- You don’t usually need to run to convergence; measure the speed/performance tradeoff

## Part 6: Wrap-Up



## Mean Field vs Belief Propagation

- ▶ **When to use Mean Field:**
  - ▶ Models made up of weakly interacting structures that are individually tractable
  - ▶ Joint models often have this flavor
- ▶ **When to use Belief Propagation:**
  - ▶ Models with intersecting factors that are tractable in isolation but interact badly
  - ▶ You often get models like this when adding non-local features to an existing tractable model



## Mean Field vs Belief Propagation

- ▶ **Mean Field Advantages**
  - ▶ For models where it applies, the coordinate ascent procedure converges quite quickly
- ▶ **Belief Propagation Advantages**
  - ▶ More broadly applicable
  - ▶ More freedom to focus on factor graph design when modeling
- ▶ **Advantages of Both**
  - ▶ Work pretty well when the real posterior is peaked (like in NLP models!)



## Other Variational Techniques

- ▶ **Variational Bayes**
  - ▶ Mean Field for models with parametric forms (e.g. Liang et al., 2007; Cohen et al., 2010)
- ▶ **Expectation Propagation**
  - ▶ Theoretical generalization of BP
  - ▶ Works kind of like Mean Field in practice; good for product models (e.g. Hall and Klein, 2012)
- ▶ **Convex Relaxation**
  - ▶ Optimize a convex approximate objective

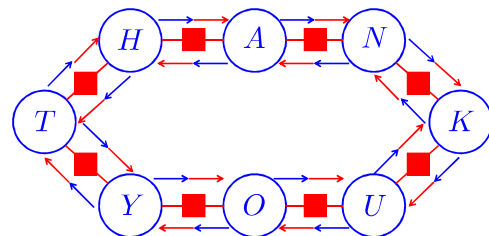


## Related Techniques

- ▶ **Dual Decomposition**
  - ▶ Not probabilistic, but good for finding maxes in similar models (e.g. Koo et al., 2010; DeNero & Machery, 2011)
- ▶ **Search approximations**
  - ▶ E.g. pruning, beam search, reranking
  - ▶ Orthogonal to approximate inference techniques (and often stackable!)



## Thank You



## Appendix A: Bibliography



## References

- ▶ **Conditional Random Fields**
  - ▶ John D. Lafferty, Andrew McCallum, and Fernando C. N. Pereira (2001). Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data. In ICML.
- ▶ **Edge-Factored Dependency Parsing**
  - ▶ Ryan McDonald, Koby Crammer, and Fernando Pereira (2005). Online Large-Margin Training of Dependency Parsers. In ACL.
  - ▶ Ryan McDonald, Fernando Pereira, Kiril Ribarov, and Jan Hajič (2005). Non-projective Dependency Parsing using Spanning Tree Algorithms. In HLT/EMNLP.



## References

- ▶ **Factorial Chain CRF**
  - ▶ Charles Sutton, Khashayar Rohanimanesh, and Andrew McCallum (2004). Dynamic Conditional Random Fields: Factorized Probabilistic Models for Labeling and Segmenting Sequence Data. In ICML.
- ▶ **Second-Order Dependency Parsing**
  - ▶ Ryan McDonald and Fernando Pereira (2006). Online Learning of Approximate Dependency Parsing Algorithms. In EACL.
  - ▶ Xavier Carreras (2007). Experiments with a Higher-Order Projective Dependency Parser. In CoNLL Shared Task Session.



## References

- ▶ **Max Matching Word Alignment**
  - ▶ Ben Taskar, Simon, Lacoste-Julien, and Dan Klein (2005). A discriminative matching approach to word alignment. In HLT/EMNLP.
- ▶ **Iterated Conditional Modes**
  - ▶ Julian Besag (1986). On the Statistical Analysis of Dirty Pictures. *Journal of the Royal Statistical Society, Series B*. Vol. 48, No. 3, pp. 259-302.
- ▶ **Structured Mean Field**
  - ▶ Eric P. Xing, Michael I. Jordan, and Stuart Russell (2003). A Generalized Mean Field Algorithm for Variational Inference in Exponential Families. In UAI.



## References

- ▶ **Joint Parsing and Alignment**
  - ▶ David Burkett, John Blitzer, and Dan Klein (2010). Joint Parsing and Alignment with Weakly Synchronized Grammars. In NAACL.
- ▶ **Word Alignment with Belief Propagation**
  - ▶ Jan Niehues and Stephan Vogel (2008). Discriminative Word Alignment via Alignment Matrix Modelling. In ACL:HLT.
  - ▶ Fabien Cromières and Sadao Kurohashi (2009). An Alignment Algorithm using Belief Propagation and a Structure-Based Distortion Model. In EACL.
  - ▶ David Burkett and Dan Klein (2012). Fast Inference in Phrase Extraction Models with Belief Propagation. In NAACL.



## References

- ▶ **Dependency Parsing with Belief Propagation**
  - ▶ David A. Smith and Jason Eisner (2008). Dependency Parsing by Belief Propagation. In EMNLP.
  - ▶ André F. T. Martins, Noah A. Smith, Eric P. Xing, Pedro M. Q. Aguiar, and Mário A. T. Figueiredo (2010). Turbo Parsers: Dependency Parsing by Approximate Variational Inference. In EMNLP.
- ▶ **Odds Ratios**
  - ▶ Dan Klein and Chris Manning (2002). A Generative Constituent-Context Model for Improved Grammar Induction. In ACL.
- ▶ **Variational Bayes**
  - ▶ Percy Liang, Slav Petrov, Michael I. Jordan, and Dan Klein (2007). The Infinite PCFG using Hierarchical Dirichlet Processes. In EMNLP/CoNLL.
  - ▶ Shay B. Cohen, David M. Blei, and Noah A. Smith (2010). Variational Inference for Adaptor Grammars. In NAACL.

Berkeley  
N L P

## References

- ▶ Expectation Propagation
  - ▶ David Hall and Dan Klein (2012). Training Factored PCFGs with Expectation Propagation. In EMNLP-CoNLL.
- ▶ Dual Decomposition
  - ▶ Terry Koo, Alexander M. Rush, Michael Collins, Tommi Jaakkola, and David Sontag (2010). Dual Decomposition for Parsing with Non-Projective Head Automata. In EMNLP.
  - ▶ Alexander M. Rush, David Sontag, Michael Collins, and Tommi Jaakkola (2010). On Dual Decomposition and Linear Programming Relaxations for Natural Language Processing. In EMNLP.
  - ▶ John DeNero and Klaus Macherey (2011). Model-Based Aligner Combination Using Dual Decomposition. In ACL.

Berkeley  
N L P

## Further Reading

- ▶ Theoretical Background
  - ▶ Martin J. Wainwright and Michael I. Jordan (2008). Graphical Models, Exponential Families, and Variational Inference. *Foundations and Trends in Machine Learning*, Vol. 1, No. 1-2, pp. 1-305.
- ▶ Gentle Introductions
  - ▶ Christopher M. Bishop (2006). *Pattern Recognition and Machine Learning*. Springer.
  - ▶ David J.C. MacKay (2003). *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press.

Berkeley  
N L P

## Further Reading

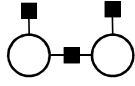
- ▶ More Variational Inference for Structured NLP
  - ▶ Zhifei Li, Jason Eisner, and Sanjeev Khudanpur (2009). Variational Decoding for Statistical Machine Translation. In ACL.
  - ▶ Michael Auli and Adam Lopez (2011). A Comparison of Loopy Belief Propagation and Dual Decomposition for Integrated CCG Supertagging and Parsing. In ACL.
  - ▶ Veselin Stoyanov and Jason Eisner (2012). Minimum-Risk Training of Approximate CRF-Based NLP Systems. In NAACL.
  - ▶ Jason Naradowsky, Sebastian Riedel, and David A. Smith (2012). Improving NLP through Marginalization of Hidden Syntactic Structure. In EMNLP-CoNLL.
  - ▶ Greg Durrett, David Hall, and Dan Klein (2013). Decentralized Entity-Level Modeling for Coreference Resolution. In ACL.

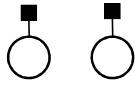
## Appendix B: Mean Field Update Derivation

Berkeley  
N L P

Berkeley  
N L P

## Mean Field Update Derivation

**Model:**   $p(y) \propto \phi(y_1)\phi(y_2)\phi(y_1, y_2)$

**Approximate Graph:**   $q(y) = q(y_1)q(y_2)$

**Goal:**  $q(y_1) = \operatorname{argmin}_{q(y_1)} KL(q||p)$


Berkeley  
N L P

## Mean Field Update Derivation

$$KL(q||p) = \sum_{y_1, y_2} q(y_1)q(y_2) \frac{\log q(y_1)q(y_2)}{\log p(y_1, y_2)}$$



Berkeley




## Mean Field Update Derivation

$$KL(q||p) = \sum_{y_1, y_2} q(y_1)q(y_2) \frac{\log q(y_1)q(y_2)}{\log p(y_1, y_2)}$$

$$= \sum_{y_1, y_2} q(y_1)q(y_2)$$

Berkeley




## Mean Field Update Derivation

$$KL(q||p) = \sum_{y_1, y_2} q(y_1)q(y_2) \frac{\log q(y_1)q(y_2)}{\log p(y_1, y_2)}$$

$$= \sum_{y_1, y_2} q(y_1)q(y_2) (\log q(y_1) + \log q(y_2))$$

Berkeley




## Mean Field Update Derivation

$$KL(q||p) = \sum_{y_1, y_2} q(y_1)q(y_2) \frac{\log q(y_1)q(y_2)}{\log p(y_1, y_2)}$$

$$= \sum_{y_1, y_2} q(y_1)q(y_2) (\log q(y_1) + \log q(y_2) - \log \phi(y_1) - \log \phi(y_2) - \log \phi(y_1, y_2) + \log Z_x)$$

Berkeley




## Mean Field Update Derivation

$$KL(q||p) = \sum_{y_1, y_2} q(y_1)q(y_2) \frac{\log q(y_1)q(y_2)}{\log p(y_1, y_2)}$$

$$= \sum_{y_1, y_2} q(y_1)q(y_2) (\log q(y_1) + \log q(y_2) - \log \phi(y_1) - \log \phi(y_2) - \log \phi(y_1, y_2) + \log Z_x)$$

$$= \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log q(y_1) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_1) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_2) \right) + \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log q(y_2) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_2) \right) + \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log Z_x \right)$$

Berkeley



## Mean Field Update Derivation


$$KL(q||p) = \sum_{y_1, y_2} q(y_1)q(y_2) \frac{\log q(y_1)q(y_2)}{\log p(y_1, y_2)}$$

$$= \sum_{y_1, y_2} q(y_1)q(y_2) (\log q(y_1) + \log q(y_2) - \log \phi(y_1) - \log \phi(y_2) - \log \phi(y_1, y_2) + \log Z_x)$$

$$= \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log q(y_1) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_1) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_2) \right) + \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log q(y_2) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_2) \right) + \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log Z_x \right)$$

$$= \left( \sum_{y_1} q(y_1) \log q(y_1) \right) - \left( \sum_{y_1} q(y_1) \log \phi(y_1) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_1, y_2) \right) + \left( \sum_{y_2} q(y_2) \log q(y_2) \right) - \left( \sum_{y_2} q(y_2) \log \phi(y_2) \right) + \log Z_x$$

Berkeley



## Mean Field Update Derivation

$$KL(q||p) = \sum_{y_1, y_2} q(y_1)q(y_2) \frac{\log q(y_1)q(y_2)}{\log p(y_1, y_2)}$$


$$= \sum_{y_1, y_2} q(y_1)q(y_2) (\log q(y_1) + \log q(y_2) - \log \phi(y_1) - \log \phi(y_2) - \log \phi(y_1, y_2) + \log Z_x)$$

$$= \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log q(y_1) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_1) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_2) \right) + \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log q(y_2) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_2) \right) + \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log Z_x \right)$$

$$= \left( \sum_{y_1} q(y_1) \log q(y_1) \right) - \left( \sum_{y_1} q(y_1) \log \phi(y_1) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_1, y_2) \right) + \left( \sum_{y_2} q(y_2) \log q(y_2) \right) - \left( \sum_{y_2} q(y_2) \log \phi(y_2) \right) + \log Z_x$$

$$\frac{\partial KL(q||p)}{\partial q(y_1)} =$$

Berkeley




## Mean Field Update Derivation

$$\begin{aligned}
 KL(q||p) &= \sum_{y_1, y_2} q(y_1)q(y_2) \frac{\log q(y_1)q(y_2)}{\log p(y_1, y_2)} \\
 &= \sum_{y_1, y_2} q(y_1)q(y_2) (\log q(y_1) + \log q(y_2) - \log \phi(y_1) - \log \phi(y_2) - \log \phi(y_1, y_2) + \log Z_x) \\
 &= \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log q(y_1) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_1) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_2) \right) + \\
 &\quad \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log q(y_2) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_2) \right) + \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log Z_x \right) \\
 &= \left( \sum_{y_1} q(y_1) \log q(y_1) \right) - \left( \sum_{y_1} q(y_1) \log \phi(y_1) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_1, y_2) \right) + \\
 &\quad \left( \sum_{y_2} q(y_2) \log q(y_2) \right) - \left( \sum_{y_2} q(y_2) \log \phi(y_2) \right) + \log Z_x
 \end{aligned}$$

$$\frac{\partial KL(q||p)}{\partial q(y_1)} = (\log q(y_1) + 1)$$

Berkeley




## Mean Field Update Derivation

$$\begin{aligned}
 KL(q||p) &= \sum_{y_1, y_2} q(y_1)q(y_2) \frac{\log q(y_1)q(y_2)}{\log p(y_1, y_2)} \\
 &= \sum_{y_1, y_2} q(y_1)q(y_2) (\log q(y_1) + \log q(y_2) - \log \phi(y_1) - \log \phi(y_2) - \log \phi(y_1, y_2) + \log Z_x) \\
 &= \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log q(y_1) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_1) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_2) \right) + \\
 &\quad \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log q(y_2) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_2) \right) + \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log Z_x \right) \\
 &= \left( \sum_{y_1} q(y_1) \log q(y_1) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_1) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_2) \right) + \\
 &\quad \left( \sum_{y_2} q(y_2) \log q(y_2) \right) - \left( \sum_{y_2} q(y_2) \log \phi(y_2) \right) + \log Z_x
 \end{aligned}$$

$$\frac{\partial KL(q||p)}{\partial q(y_1)} = (\log q(y_1) + 1) - \log \phi(y_1)$$

Berkeley




## Mean Field Update Derivation

$$\begin{aligned}
 KL(q||p) &= \sum_{y_1, y_2} q(y_1)q(y_2) \frac{\log q(y_1)q(y_2)}{\log p(y_1, y_2)} \\
 &= \sum_{y_1, y_2} q(y_1)q(y_2) (\log q(y_1) + \log q(y_2) - \log \phi(y_1) - \log \phi(y_2) - \log \phi(y_1, y_2) + \log Z_x) \\
 &= \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log q(y_1) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_1) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_2) \right) + \\
 &\quad \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log q(y_2) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_2) \right) + \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log Z_x \right) \\
 &= \left( \sum_{y_1} q(y_1) \log q(y_1) \right) - \left( \sum_{y_1} q(y_1) \log \phi(y_1) \right) - \left( \sum_{y_1, y_2} q(y_1)q(y_2) \log \phi(y_1, y_2) \right) + \\
 &\quad \left( \sum_{y_2} q(y_2) \log q(y_2) \right) - \left( \sum_{y_2} q(y_2) \log \phi(y_2) \right) + \log Z_x
 \end{aligned}$$

$$\frac{\partial KL(q||p)}{\partial q(y_1)} = (\log q(y_1) + 1) - \log \phi(y_1) - \sum_{y_2} q(y_2) \log \phi(y_1, y_2)$$


Berkeley



## Mean Field Update Derivation

$$0 = (\log q(y_1) + 1) - \log \phi(y_1) - \sum_{y_2} q(y_2) \log \phi(y_1, y_2)$$

Berkeley




## Mean Field Update Derivation

$$0 = (\log q(y_1) + 1) - \log \phi(y_1) - \sum_{y_2} q(y_2) \log \phi(y_1, y_2)$$

$$\log q(y_1) = \log \phi(y_1) + \sum_{y_2} q(y_2) \log \phi(y_1, y_2) - 1$$

Berkeley




## Mean Field Update Derivation

$$0 = (\log q(y_1) + 1) - \log \phi(y_1) - \sum_{y_2} q(y_2) \log \phi(y_1, y_2)$$

$$\log q(y_1) = \log \phi(y_1) + \sum_{y_2} q(y_2) \log \phi(y_1, y_2) - 1$$

$$q(y_1) = \exp \left( \log \phi(y_1) + \sum_{y_2} q(y_2) \log \phi(y_1, y_2) - 1 \right)$$

Berkeley




### Mean Field Update Derivation

$$0 = (\log q(y_1) + 1) - \log \phi(y_1) - \sum_{y_2} q(y_2) \log \phi(y_1, y_2)$$

$$\log q(y_1) = \log \phi(y_1) + \sum_{y_2} q(y_2) \log \phi(y_1, y_2) - 1$$

$$q(y_1) \propto \exp \left( \log \phi(y_1) + \sum_{y_2} q(y_2) \log \phi(y_1, y_2) - 1 \right)$$

Berkeley



### Mean Field Update Derivation

$$0 = (\log q(y_1) + 1) - \log \phi(y_1) - \sum_{y_2} q(y_2) \log \phi(y_1, y_2)$$


$$\log q(y_1) = \log \phi(y_1) + \sum_{y_2} q(y_2) \log \phi(y_1, y_2) - 1$$

$$q(y_1) \propto \exp \left( \log \phi(y_1) + \sum_{y_2} q(y_2) \log \phi(y_1, y_2) \right)$$


$$q(y_i) \propto \exp \left( \sum_{c: i \in c} \mathbb{E}_{q_{-y_i}} \log \phi_c(y_c) \right)$$

### Appendix C: Joint Parsing and Alignment Component Distributions

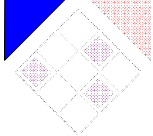
Berkeley



Berkeley




### Joint Parsing and Alignment Component Distributions

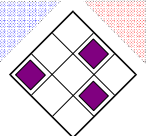


$$q(t) \propto \exp \left( \sum_{n_i X_j \in t} w^\top f_t(n) + \sum_{n_i X_j \in t} \sum_{n'_s X'_t} q(b_{ij, st}) q(n'_s X'_t) w^\top f_{t a t'}(n, b, n') \right)$$

Berkeley




### Joint Parsing and Alignment Component Distributions

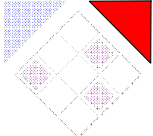


$$q(a) \propto \exp \left( \sum_{b_{ij, st} \in a} w^\top f_a(b) + \sum_{b_{ij, st} \in a} \sum_{n_i X_j \in X, X'} q(n_i X_j) q(n'_s X'_t) w^\top f_{t a t'}(n, b, n') \right)$$

Berkeley

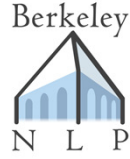


### Joint Parsing and Alignment Component Distributions




$$q(t') \propto \exp \left( \sum_{n'_s X'_t \in t'} w^\top f_{t'}(n') + \sum_{n'_s X'_t \in t'} \sum_{n_i X_j \in X_j} q(n_i X_j) q(b_{ij, st}) w^\top f_{t a t'}(n, b, n') \right)$$

## Appendix D: Forward-Backward as Belief Propagation

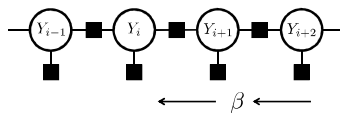


Berkeley



### Forward-Backward as Belief Propagation

→  $\alpha$  →




←  $\beta$  ←

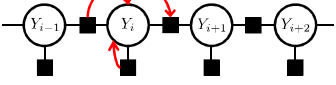
$$\alpha_i(y_i) = \phi_i(y_i) \sum_{y_{i-1}} \alpha_{i-1}(y_{i-1}) \phi_{i-1,i}(y_{i-1}, y_i)$$

$$\beta_i(y_i) = \sum_{y_{i+1}} \beta_{i+1}(y_{i+1}) \phi_{i,i+1}(y_i, y_{i+1}) \phi_{i+1}(y_{i+1})$$

Berkeley



### Forward-Backward as Belief Propagation




$$\alpha_i(y_i) = \phi_i(y_i) \sum_{y_{i-1}} \alpha_{i-1}(y_{i-1}) \phi_{i-1,i}(y_{i-1}, y_i)$$

$$= m_{\phi_i \rightarrow Y_i}(y_i) m_{\phi_{i-1,i} \rightarrow Y_i}(y_i)$$

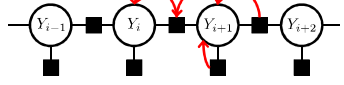
$$m_{\phi_i \rightarrow Y_i}(y_i) = \phi_i(y_i) \quad m_{Y_i \rightarrow \phi_{i,i+1}}(y_i) = \alpha_i(y_i)$$

$$m_{\phi_{i-1,i} \rightarrow Y_i}(y_i) = \sum_{y_{i-1}} \alpha_{i-1}(y_{i-1}) \phi_{i-1,i}(y_{i-1}, y_i)$$

Berkeley



### Forward-Backward as Belief Propagation




$$\beta_i(y_i) = \sum_{y_{i+1}} \beta_{i+1}(y_{i+1}) \phi_{i,i+1}(y_i, y_{i+1}) \phi_{i+1}(y_{i+1})$$

$$= m_{\phi_{i,i+1} \rightarrow Y_i}(y_i)$$

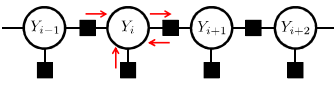
$$= \sum_{y_{i+1}} m_{Y_{i+1} \rightarrow \phi_{i,i+1}}(y_{i+1}) \phi_{i,i+1}(y_i, y_{i+1})$$

$$m_{Y_{i+1} \rightarrow \phi_{i,i+1}}(y_{i+1}) = m_{\phi_{i+1} \rightarrow Y_{i+1}}(y_{i+1}) * m_{\phi_{i+1,i+2} \rightarrow Y_{i+1}}(y_{i+1})$$

Berkeley



### Forward-Backward Marginal Beliefs



$$P(y_i|x) \propto \alpha_i(y_i) \beta_i(y_i)$$

$$= m_{Y_i \rightarrow \phi_{i,i+1}}(y_i) m_{\phi_{i,i+1} \rightarrow Y_i}(y_i)$$

$$= m_{\phi_{i-1,i} \rightarrow Y_i}(y_i) m_{\phi_i \rightarrow Y_i}(y_i) m_{\phi_{i,i+1} \rightarrow Y_i}(y_i)$$