

# Sequence Labeling, Contd

Prof. Sameer Singh

---

CS 295: STATISTICAL NLP

WINTER 2017

February 2, 2017

# Outline

---

Marginal Inference in HMMs: F/B algorithm

Maximum Entropy Markov Models

Conditional Random Fields

Neural Sequence Tagging

# Outline

---

Marginal Inference in HMMs: F/B algorithm

Maximum Entropy Markov Models

Conditional Random Fields

Neural Sequence Tagging

# Expectation Maximization

---

Initialization

$$e(x_i | y_i) \sim \text{uniform}$$
$$t(y_i | y_{i-1}) \sim \text{uniform}$$

K-Means

Pick K random centroids

Label Data from the Model

Fix  $e, t$

compute  $p(\vec{y} | \vec{x})$

Cluster all the points

Update the Model from Data

Fix  $p(y|x)$

Update  $e, t$

Update centroids

# Label Data from the Model

Hard-EM

$$y^* = \operatorname{argmax}_y p(y|x) = \operatorname{argmax}_y p(y, x)$$

$$\hat{e}(x_i | y_i) = \frac{\# x_i \wedge y_i \text{ in } y^*}{\# y_i \text{ in } y^*} \quad \vec{t} = \dots$$

Soft-EM

$$p(y_i | \vec{x}) = \sum_{\substack{\vec{y}_j \\ y_{ij} = y_i}} p(\vec{y}_j | \vec{x})$$

$$p(x_i, y_i) = \sum_{y_j} p(x_j = x_i, y_{ij} = y_i)$$

# Dynamic Programming

---

$$\begin{aligned} P(x_1 \dots x_n, y_i) &= \sum_{y_1 \dots y_{i-1}} \sum_{y_{i+1} \dots y_n} P(x_1 \dots x_n, y_1 \dots y_n) \\ &= P(x_1 \dots x_i, y_i) P(x_{i+1} \dots x_n | y_i) \end{aligned}$$

$$\begin{aligned} \alpha(i, y_i) &= P(x_1 \dots x_i, y_i) = \sum_{y_1 \dots y_{i-1}} P(x_1 \dots x_i, y_1 \dots y_i) \\ &= \sum_{y_{i-1}} e(x_i | y_i) t(y_i | y_{i-1}) \alpha(i-1, y_{i-1}) \end{aligned}$$

$$\begin{aligned} \beta(i, y_i) &= P(x_{i+1} \dots x_n | y_i) = \sum_{y_{i+1} \dots y_n} P(x_{i+1} \dots x_n, y_{i+1} \dots y_n) \\ &= \sum_{y_{i+1}} e(x_{i+1} | y_{i+1}) t(y_{i+1} | y_i) \beta(i+1, y_{i+1}) \end{aligned}$$

# Forward Backward Algorithm

---

Forward

$$\alpha(0, s) = 1$$

$$\alpha(i, y_i) = \sum_{y_{i-1}} e(x_i | y_i) t(y_i | y_{i-1}) \alpha(i-1, y_{i-1})$$

Backward

$$\beta(n, e) = 1$$

$$\beta(i, y_i) = \sum_{y_{i+1}} e(x_{i+1} | y_{i+1}) t(y_{i+1} | y_i) \beta(i+1, y_{i+1})$$

# Updating the Model from Data

Hard-EM

$$e(x_i | y_i) = \frac{\# x_i \text{ tagged as } y_i}{\# y_i} \quad t(y_i | y_{i-1}) = \frac{\# y_{i-1} \rightarrow y_i}{\# y_{i-1}}$$

Soft-EM

$$e(x_i | y_i) = \sum_j p(x_j = x_i, y_j = y_i | x_1 \dots x_n)$$
$$p(x_1 \dots x_n, y_i) = \alpha(i, y_i) \beta(i, y_i) = \sum_{\substack{j \\ x_j = x_i}} p(y_j = y_i | x_1 \dots x_n)$$

$$p(x_1 \dots x_n, y_i, y_{i+1}) = \alpha(i, y_i) t(y_{i+1} | y_i) e(x_{i+1} | y_{i+1}) \beta(i+1, y_{i+1})$$



# Outline

---

Marginal Inference in HMMs: F/B algorithm

Maximum Entropy Markov Models

Conditional Random Fields

Neural Sequence Tagging

# Named Entity Recognition

---

George W. Bush spoke from the White House today .

PER PER PER O O O LOC LOC O O

B-PER I-PER I-PER O O O B-LOC I-LOC O O

B-PER I-PER E-PER O O O B-LOC E-LOC O O

# Max. Entropy Markov Models

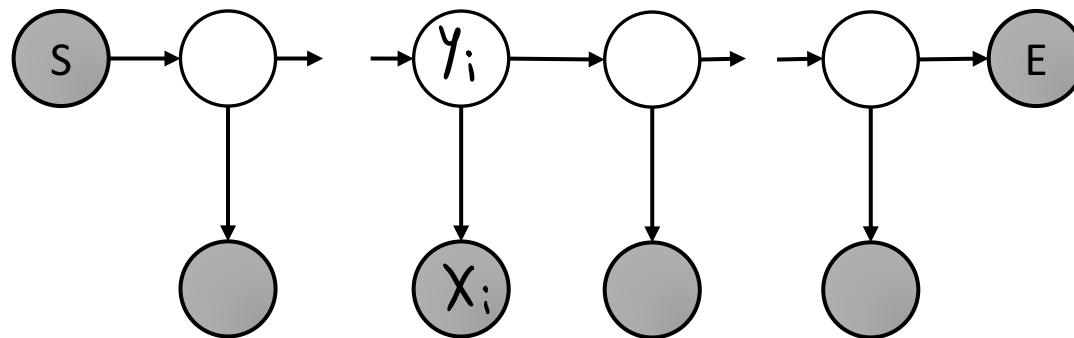
---

$$p(\vec{y}|\vec{x}) = \prod_i p(y_i|x_i, y_{i-1})$$

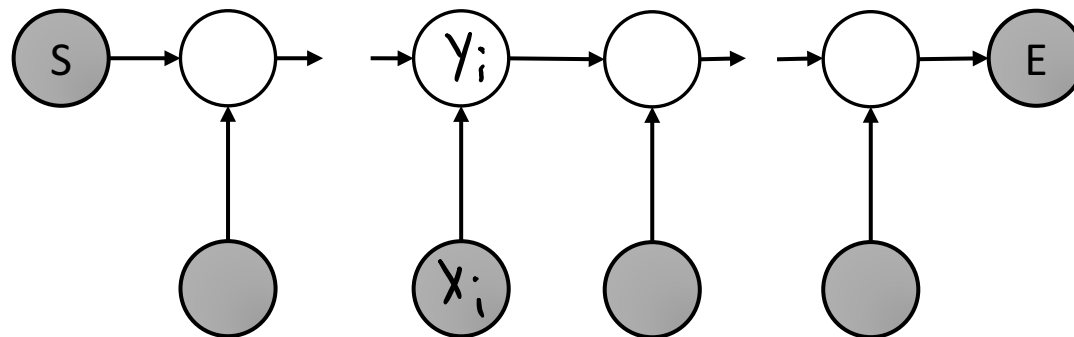
$$p(y_i|x_i, y_{i-1}) = \frac{e^{\theta \cdot \phi(x_i, y_i, y_{i-1})}}{\sum_y e^{\theta \cdot \phi(x_i, y, y_{i-1})}}$$

# Graphical Model Notation

HMMs



MEMMs



# Adding Features (for POS)

---

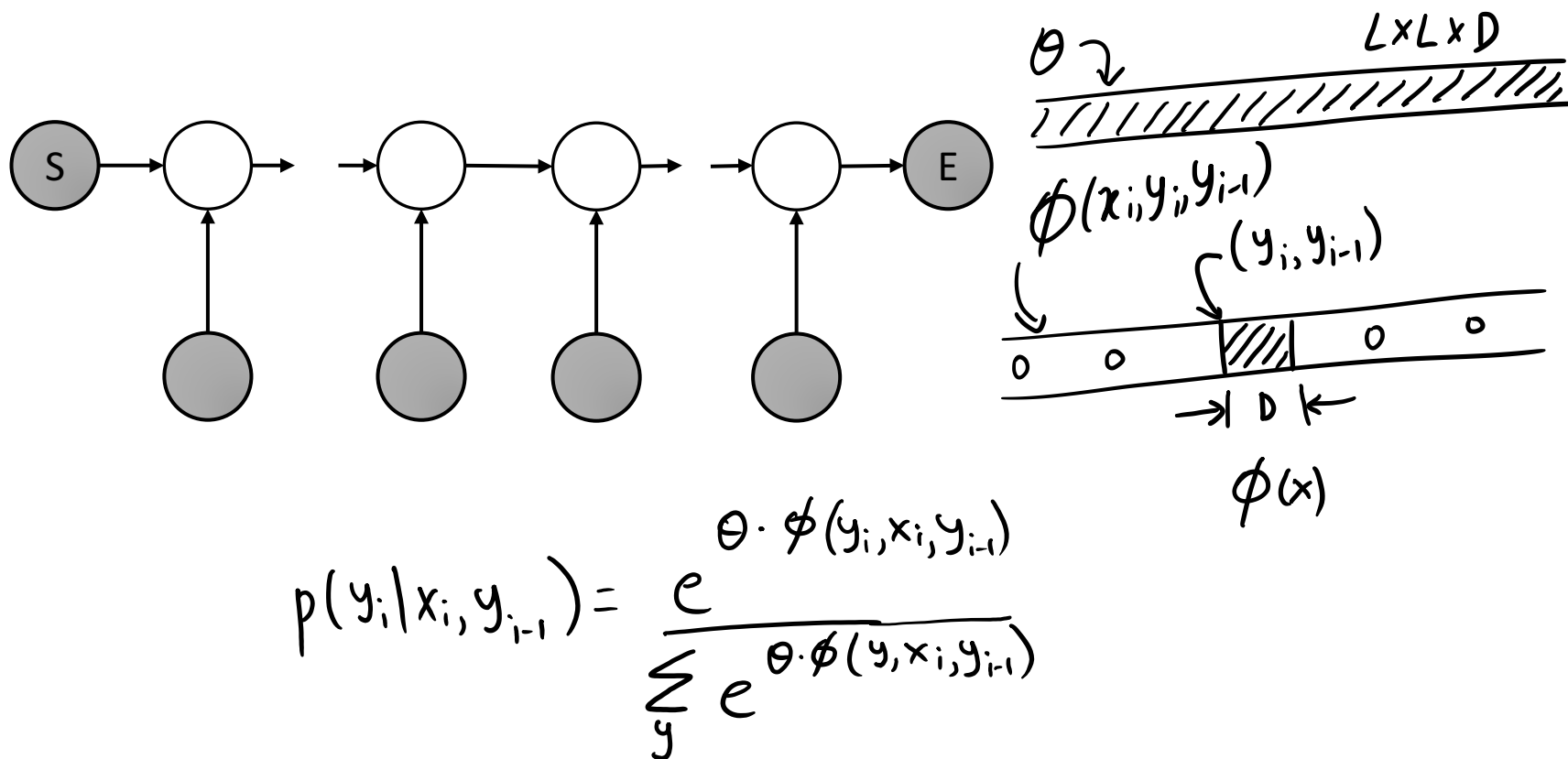
## Current Word

- Word                      the: the → DT
- Lowercased word        Importantly: importantly → RB
- Prefixes                unfathomable: un- → JJ
- Suffixes                Surprisingly: -ly → RB
- Capitalization         Meridian: CAP → NNP
- Word shapes            35-year: d-x → JJ

## Window Words

- Add in previous / next word the \_\_\_\_
- Previous / next word shapes        X \_\_\_\_ X
- Occurrence pattern features        [X: x X occurs]
- Crude entity detection                \_\_\_\_ ..... (Inc. | Co.)
- Phrasal verb in sentence?    put ..... \_\_\_\_
- Conjunctions of these things

# Adding Features



# Predictions Using MEMMs

Greedy

$$y_1 = \operatorname{argmax}_y p(y|x_1, s) \quad y_2 = \operatorname{argmax}_y p(y|x_2, y_1) \quad \dots$$

Viterbi Decoding

$$\vec{y} = \operatorname{argmax}_y \prod_{i=1}^n p(y_i | x_i, y_{i-1})$$

$$\pi(i, y_i) = \operatorname{argmax}_{y_1 \dots y_{i-1}} \prod_{j=1}^{i-1} p(y_j | x_j, y_{j-1})$$

$$= \operatorname{argmax}_{y_{i-1}} p(y_i | x_i, y_{i-1}) \pi(i-1, y_{i-1})$$

# Training MEMMs

---

$$\begin{aligned} \mathcal{L}(\theta, D) &= \sum_{d \in D} \log P(\vec{y} | \vec{x}) \\ &= \sum_d \sum_i \log \underbrace{P(y_i | x_i, y_{i-1})}_{\text{Independent}} \end{aligned}$$

Train using off the shelf classifiers!



# Outline

---

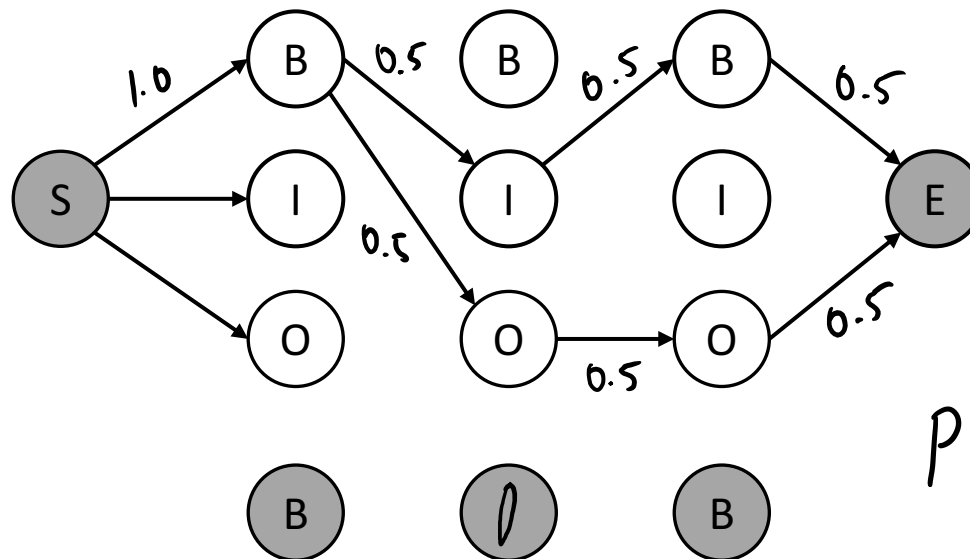
Marginal Inference in HMMs: F/B algorithm

Maximum Entropy Markov Models

Conditional Random Fields

Neural Sequence Tagging

# Label Bias Problem



$$P(B|B) = 1$$

$$P(I|O, B) = 0.5$$

$$P(O|O, B) = 0.5$$

$$P(B|B) = 0.5$$

$$P(B|O) = 0.5$$

$$P(B|B, I) = 0.5$$

$$P(O|B, O) = 0.5$$

XX → ignores "B"

# Conditional Random Fields

$$p(\vec{y}|\vec{x}) = \frac{e^{\theta \cdot \Phi(\vec{x}, \vec{y})}}{\sum_{\vec{y}'} e^{\theta \cdot \Phi(\vec{x}, \vec{y}')}}}$$

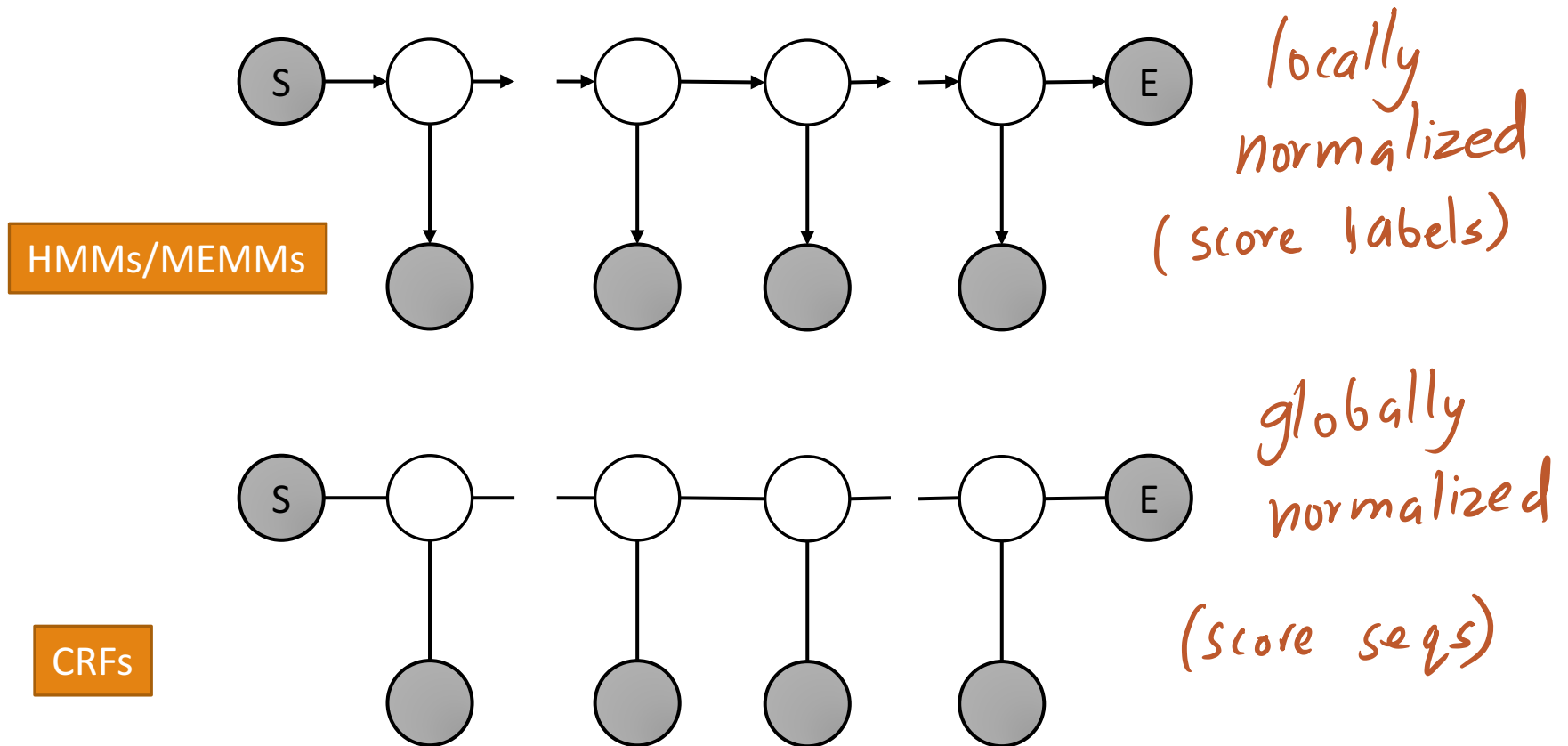
sometimes  
combined  
 $\phi(x_i, y_i, y_{i-1})$

$$\Phi(\vec{x}, \vec{y}) = \underbrace{\sum_{i=1}^n \phi_x(x_i, y_i)}_{\text{arbitrary features}} + \underbrace{\sum_{i=2}^n \phi_t(y_i, y_{i-1})}_{\text{transition table}}$$

arbitrary  
features

transition  
table

# Graphical Model Notation



# Predictions Using CRFs

---

$$\vec{y} = \operatorname{argmax}_y p(y|x)$$

$$= \operatorname{argmax}_y e^{\theta \cdot \Phi(x, y)} = \operatorname{argmax}_y \theta \cdot \Phi(x, y)$$

$$= \operatorname{argmax}_y \theta \cdot \sum_i \phi(x_i, y_i, y_{i-1}) = \operatorname{argmax}_y \sum_i \theta \cdot \phi(x_i, y_i, y_{i-1})$$

Viterbi

$$\Pi(i, y_i) = \max_{y_1, \dots, y_i} \sum_{j=1}^{i-1} \theta \cdot \phi(x_j, y_j, y_{j-1})$$

$$= \max_{y_i} \theta \cdot \phi(x_i, y_i, y_{i-1}) + \Pi(i-1, y_{i-1})$$

# Likelihood Training of CRFs

$$\mathcal{L}(\theta, D) = \log \prod_j P(\vec{y}_j, \vec{x}_j) = \sum_j \log \frac{e^{\theta \cdot \Phi(\vec{x}_j, \vec{y}_j)}}{\sum_y e^{\theta \cdot \Phi(\vec{x}_j, y)}}$$

$$= \sum_j \theta \cdot \Phi(\vec{x}_j, \vec{y}_j) - \log \sum_y e^{\theta \cdot \Phi(\vec{x}_j, y)}$$

$$\frac{\partial \mathcal{L}(\theta, D)}{\partial \theta_k} = \sum_j \Phi_k(\vec{x}_j, \vec{y}_j) - \sum_y p(y | \vec{x}_j) \Phi_k(\vec{x}_j, y)$$

# Likelihood Training of CRFs

$$\begin{aligned}
 \sum_{\vec{y}} p(\vec{y} | \vec{x}) \Phi_k(\vec{x}, \vec{y}) &= \sum_{\vec{y}} p(\vec{y} | \vec{x}) \sum_i \phi_k(x_i, y_i, y_{i-1}) \\
 &= \sum_i \sum_{\vec{y}} p(y_i | x) \phi_k(x_i, y_i, y_{i-1}) = \sum_i \sum_{\substack{y_i \\ y_{i-1}}} \sum_{y/y_{i-1}} p(y_i | x) \phi_k(x_i, y_i, y_{i-1}) \\
 &= \sum_i \sum_{\substack{y_i \\ y_{i-1}}} \phi_k(x_i, y_i, y_{i-1}) \underbrace{q_i(y_i, y_{i-1})}_{e^{\theta \Phi_k(\vec{x}, \vec{y})}} \quad q_i(a, b) = \sum_{\substack{y \\ y_i=a \\ y_{i+1}=b}} p(y | \vec{x}) \\
 \mu_i(a, b) &= \sum_{\substack{y \\ y_i=a \\ y_{i+1}=b}} \Psi(\vec{y}) \quad q_i(a, b) = \frac{\mu_i(a, b)}{\sum_{a, b} \mu_i(a, b)}
 \end{aligned}$$

# Forward-Backward Algorithm

$$\alpha(i, y_i) = P(y_i | x_1 \dots x_i) = \sum_{y_1 \dots y_{i-1}} \psi(x_i, y_i, y_{i-1}) \prod_{j=1}^{i-1} \psi(x_j, y_j, y_{j-1})$$

$\underbrace{e^{\theta \cdot \phi(x_j, y_j, y_{j-1})}}$

$$= \sum_{y_{i-1}} \psi(x_i, y_i, y_{i-1}) \alpha(i-1, y_{i-1})$$

$$\beta(i, y_i) = P(y_i | x_{i+1} \dots x_n) = \sum_{y_{i+1} \dots y_n} \prod_{j=i+1}^n \psi(x_j, y_j, y_{j-1})$$

$Z = \sum_y \alpha(n, y)$

$$= \sum_{y_{i+1}} \psi(x_{i+1}, y_{i+1}, y_i) \beta(i+1, y_{i+1})$$

$\mu_j(a, b) = \alpha(j, a) \psi(x_j, b, a) \beta(j+1, b)$



# Outline

---

Marginal Inference in HMMs: F/B algorithm

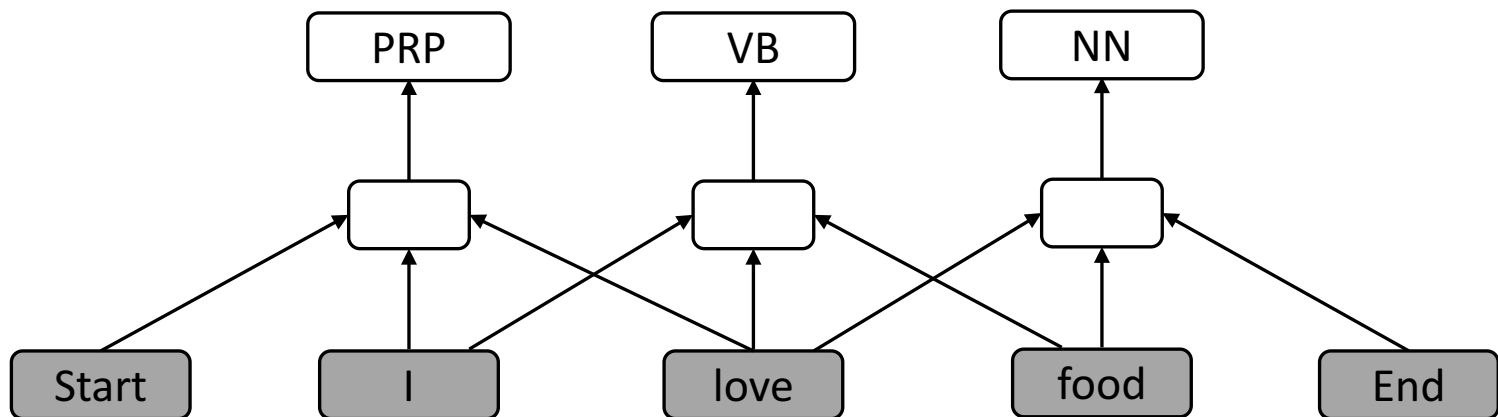
Maximum Entropy Markov Models

Conditional Random Fields

Neural Sequence Tagging

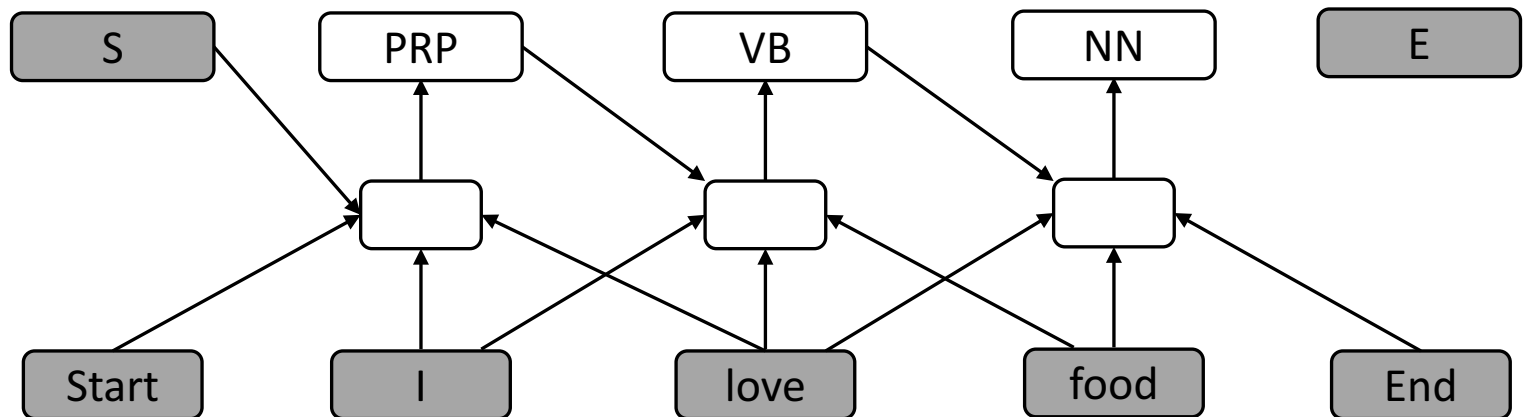
# Simple Neural Tagger

---



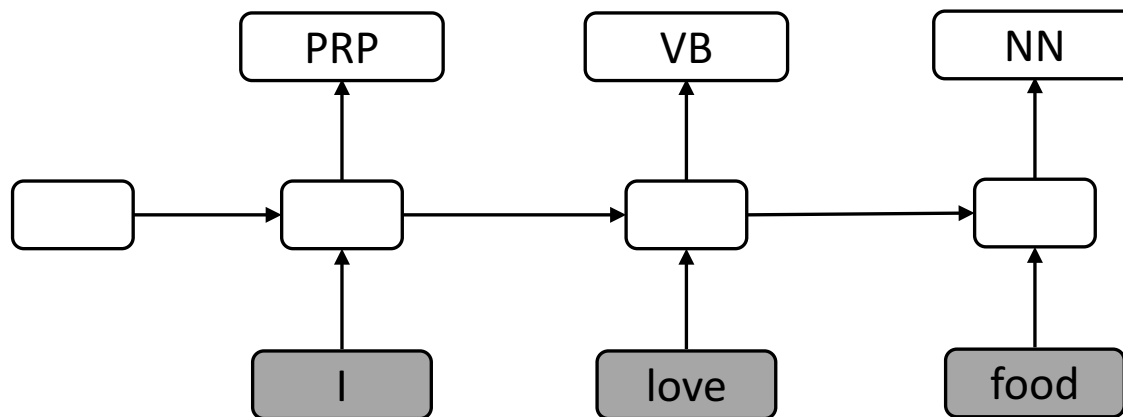
# MEMM-*ish* Neural Tagger

---



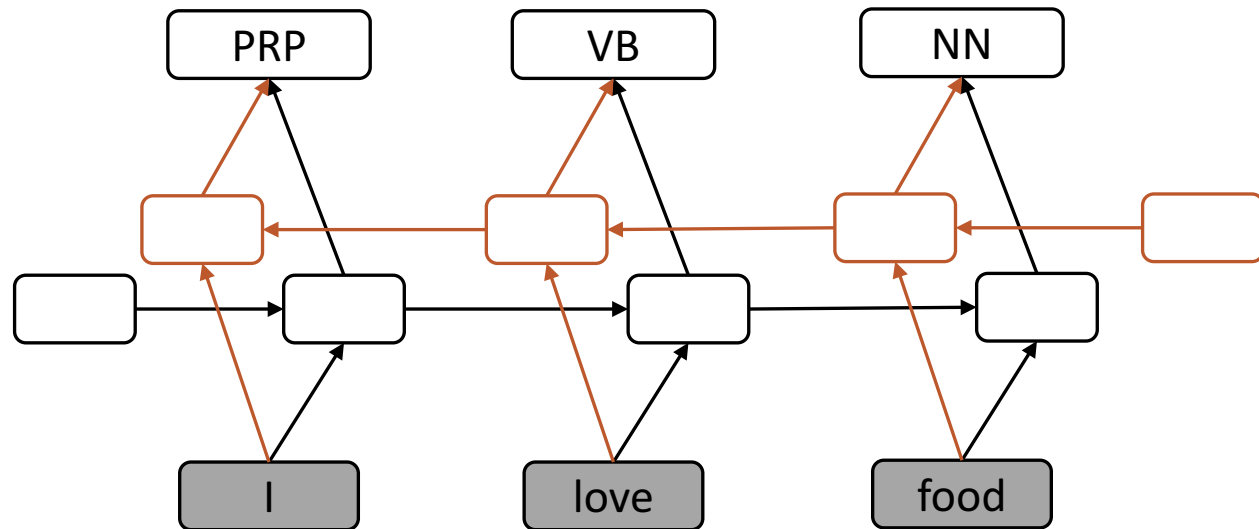
# Recurrent Neural Tagger

---



# Bidirectional RNN Tagger

---



# Upcoming...

---

## Homework

- Homework 2 is due (~10 days): **February 13, 2017**
- Write-up, data, and code for Homework 2 is up
- Ask questions early!

## Project

- Proposal is due on Tuesday: **February 7, 2017**
- Only **2 pages**

## Summaries

- Paper summaries: **February 17, February 28, March 14**
- Only **1 page** each