

# CS 6501 Natural Language Processing

## Sequence Labeling (II)

---

Yangfeng Ji

September 17, 2018

Department of Computer Science  
University of Virginia



ENGINEERING

# Important Dates

- ▶ Project 1 due: Sept. 23, 11:59PM
- ▶ Group project proposal due: Oct. 7, 11:59PM
- ▶ Group project proposal presentation: Oct. 10 & Oct. 15

## Other things

- ▶ Project 1 submission is open on Collab

# Overview

1. Conditional Random Fields
2. Inference
3. Parameter Estimation

# POS Tagging

## Example

Rain and wind from Florence will pop up on Tuesday

## POS tagging

Rain<sub>NN</sub> and<sub>CC</sub> wind<sub>NN</sub> from<sub>IN</sub> Florence<sub>NNP</sub> will<sub>MD</sub> pop<sub>VB</sub> up<sub>RP</sub>  
on<sub>IN</sub> Tuesday<sub>NNP</sub>

- ▶ NN: Noun, singular or mass
- ▶ NNP: Proper noun, singular
- ▶ IN: Preposition or subordinating conjunction
- ▶ CC: Coordinating conjunction
- ▶ MD: Modal
- ▶ VB: Verb, base form
- ▶ RP: Particle

# Sequence Labeling

## Example

[Atlantis]<sub>MSIC</sub> touched down at [Kennedy Space  
Center]<sub>LOC</sub>

Atlantis	touched	down	at	Kennedy	Space	Center	.
B <sub>MSIC</sub>	O	O	O	B <sub>LOC</sub>	I <sub>LOC</sub>	I <sub>LOC</sub>	O

Tag set

- ▶ B: beginning
- ▶ I: inside
- ▶ O: outside

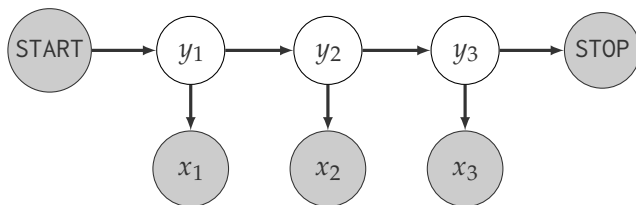
Category

- ▶ Person
- ▶ Location
- ▶ Organization
- ▶ Msic

# Hidden Markov Models

$$P(\mathbf{x}, \mathbf{y}) = \prod_{i=1} \left\{ P(y_i | y_{i-1}) P(x_i | y_i) \right\} \quad (1)$$

Graphical model



- ▶  $\mathbf{x}$ : observation (e.g., sentences)
- ▶  $\mathbf{y}$ : **hidden** variables (e.g., POS sequences)

# Generative Models

$$P(x, y) = P(x|y) \cdot P(y) \quad (2)$$

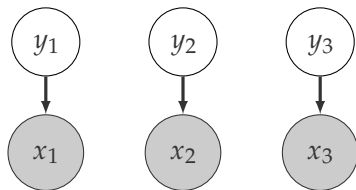
# Generative Models

$$P(x, y) = P(x|y) \cdot P(y) \quad (2)$$

Factorization

$$P(x|y) = \prod_{i=1} \underbrace{P(x_i|y_i)}_{\text{Emission probability}} \quad (3)$$

Graphical model





# Discriminative Models: Logistic Regression

$$P(y|x) = \frac{\exp(\boldsymbol{\theta}^\top \mathbf{f}(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(\boldsymbol{\theta}^\top \mathbf{f}(x, y'))} \quad (4)$$

# Question

How to build a discriminative sequential model  $p(\mathbf{y}|\mathbf{x})$ ?

- ▶  $\mathbf{y} \in \mathcal{Y}^T$  is a sequence

# Conditional Random Fields

---

# Logistic Regression

$$P(\mathbf{y}|\mathbf{x}) = \frac{\exp(\boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{Y}^T} \exp(\boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y}'))} \quad (5)$$

Huge  $\mathcal{Y}^T$  causes the problems on

- ▶ decoding  $\arg \max_{\mathbf{y}' \in \mathcal{Y}^T} P(\mathbf{y}'|\mathbf{x})$

# Logistic Regression

$$P(\mathbf{y}|\mathbf{x}) = \frac{\exp(\boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y}))}{\underbrace{\sum_{\mathbf{y}' \in \mathcal{Y}^T} \exp(\boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y}'))}_{\text{partition function } Z}} \quad (5)$$

Huge  $\mathcal{Y}^T$  causes the problems on

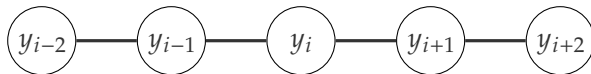
- ▶ decoding  $\arg \max_{\mathbf{y}' \in \mathcal{Y}^T} P(\mathbf{y}'|\mathbf{x})$
- ▶ computing the partition function with  $|\mathcal{Y}^T| = K^T$  possible values

# Markov Property

Global feature function:

$$f(x, y) \quad (6)$$

Markov assumption:



- ▶ Conditional independence
- ▶ Factorization over cliques

# Decomposition of $f(x, y)$

$$f(x, y) = \sum_{i=1}^T \underbrace{f_i(x, y_{i-1}, y_i)}_{\text{local feature function}} \quad (7)$$

- ▶  $f_i(x, y_{i-1}, y_i)$  captures the transition from  $y_i$  to  $y_i$
- ▶  $i$ : the position to be tagged
- ▶  $y_i \in \mathcal{Y}$ : POS tag at position  $i$
- ▶  $y_{i-1} \in \mathcal{Y}$ : POS tag at position  $i - 1$
- ▶  $x$ : the entire sentence

# Local Feature Function: Example

- ▶ standard features

[Lafferty et al., 2001]



# Local Feature Function: Example

- ▶ standard features
- ▶ whether a spelling begins with upper case letter,
  - ▶ IBM, Virginia: PROPER NOUN

[Lafferty et al., 2001]

# Local Feature Function: Example

- ▶ standard features
- ▶ whether a spelling begins with upper case letter,
  - ▶ IBM, Virginia: PROPER NOUN
- ▶ whether it ends in one of the following suffixes:
  - ▶ -ies e.g., parties: PROPER NOUN, PLURAL
  - ▶ -ly e.g., extremely, loudly: ADVERB
  - ▶ -ing e.g., : VERB, GERUND OR PRESENT PARTICIPLE
  - ▶ ...

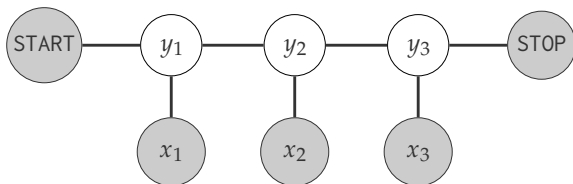
[Lafferty et al., 2001]

# Example

Rain<sub>NN</sub> and<sub>CC</sub> wind<sub>NN</sub> from<sub>IN</sub> Florence<sub>NNP</sub> will<sub>MD</sub> pop<sub>VB</sub> up<sub>RP</sub>  
on<sub>IN</sub> Tuesday<sub>NNP</sub>

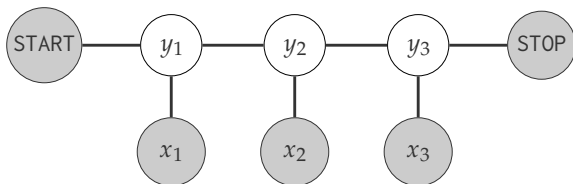
# Graphical Model Representation

Conditional Random Fields:

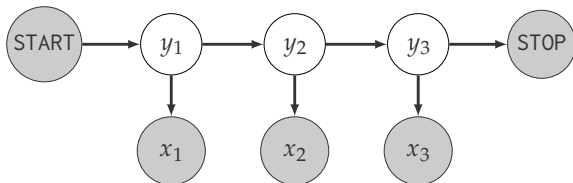


# Graphical Model Representation

Conditional Random Fields:



Hidden Markov Models:



# CRF vs HMM

HMM factorization

$$P(\mathbf{x}, \mathbf{y}) = \prod_{i=1} \left\{ P(y_i | y_{i-1}) P(x_i | y_i) \right\} \quad (8)$$

Let

- ▶  $\psi(y_1, y_2) = P(y_1)P(y_2|y_1)P(x_1|y_1)P(x_2|y_2)$
- ▶  $\psi(y_{i-1}, y_i) = P(y_i|y_{i-1})P(x_i|y_i), \forall i > 2$

# CRF vs HMM

HMM factorization

$$P(\mathbf{x}, \mathbf{y}) = \prod_{i=1} \left\{ P(y_i | y_{i-1}) P(x_i | y_i) \right\} \quad (8)$$

Let

- ▶  $\psi(y_1, y_2) = P(y_1)P(y_2|y_1)P(x_1|y_1)P(x_2|y_2)$
- ▶  $\psi(y_{i-1}, y_i) = P(y_i|y_{i-1})P(x_i|y_i), \forall i > 2$

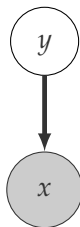
$$P(\mathbf{x}, \mathbf{y}) = \prod_{i=1}^T \psi(y_{i-1}, y_i) \quad (9)$$

with  $Z = 1$

# Directed vs. Undirected Graphical Models

- ▶ Undirected GMs allow more flexible definitions
- ▶ Directed GMs can always be converted into an undirected GMs

$$P(x, y) = P(x|y)P(y)$$



$$P(x, y) = \psi(x, y)/Z$$



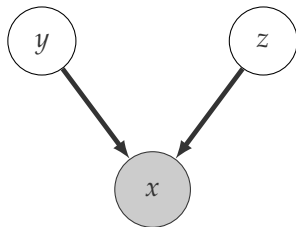
[Bishop, 2006, Chap. 8]



# Directed vs. Undirected Graphical Models (II)

- ▶ Directed GMs can give more concise definitions

$$P(x, y, z) = P(x|y, z)P(y)P(z)$$

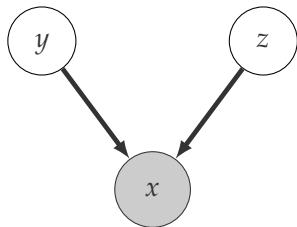


[Bishop, 2006, Chap. 8]

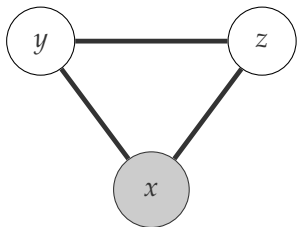
## Directed vs. Undirected Graphical Models (II)

- ▶ Directed GMs can give more concise definitions

$$P(x, y, z) = P(x|y, z)P(y)P(z)$$



$$P(x, y, z) = \psi(x, y, z)/Z$$



[Bishop, 2006, Chap. 8]

# Inference

---

# Decode $P(\boldsymbol{y}|\boldsymbol{x})$

$$f(\boldsymbol{x}, \boldsymbol{y}) = \sum_{i=1}^T f_i(\boldsymbol{x}, y_{i-1}, y_i) \quad (10)$$

# Decode $P(\mathbf{y}|\mathbf{x})$

$$f(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^T f_i(\mathbf{x}, y_{i-1}, y_i) \quad (10)$$

$$\begin{aligned} \arg \max_{\mathbf{y} \in \mathcal{Y}^T} P(\mathbf{y}|\mathbf{x}) &= \arg \max_{\mathbf{y} \in \mathcal{Y}^T} \frac{\exp(\boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{Y}^T} \exp(\boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y}'))} \\ &= \arg \max_{\mathbf{y} \in \mathcal{Y}^T} \exp(\boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y})) \end{aligned}$$

# Decode $P(\mathbf{y}|\mathbf{x})$

$$f(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^T f_i(\mathbf{x}, y_{i-1}, y_i) \quad (10)$$

$$\begin{aligned} \arg \max_{\mathbf{y} \in \mathcal{Y}^T} P(\mathbf{y}|\mathbf{x}) &= \arg \max_{\mathbf{y} \in \mathcal{Y}^T} \frac{\exp(\boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{Y}^T} \exp(\boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y}'))} \\ &= \arg \max_{\mathbf{y} \in \mathcal{Y}^T} \exp(\boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y})) \\ &= \arg \max_{\mathbf{y} \in \mathcal{Y}^T} \boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y}) \end{aligned}$$

# Decode $P(\mathbf{y}|\mathbf{x})$

$$f(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^T f_i(\mathbf{x}, y_{i-1}, y_i) \quad (10)$$

$$\begin{aligned} \arg \max_{\mathbf{y} \in \mathcal{Y}^T} P(\mathbf{y}|\mathbf{x}) &= \arg \max_{\mathbf{y} \in \mathcal{Y}^T} \frac{\exp(\boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{Y}^T} \exp(\boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y}'))} \\ &= \arg \max_{\mathbf{y} \in \mathcal{Y}^T} \exp(\boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y})) \\ &= \arg \max_{\mathbf{y} \in \mathcal{Y}^T} \boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y}) \\ &= \arg \max_{\mathbf{y} \in \mathcal{Y}^T} \boldsymbol{\theta}^\top \sum_{i=1}^T f_i(\mathbf{x}, y_{i-1}, y_i) \end{aligned}$$

# Decode $P(\mathbf{y}|\mathbf{x})$

$$f(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^T f_i(\mathbf{x}, y_{i-1}, y_i) \quad (10)$$

$$\begin{aligned} \arg \max_{\mathbf{y} \in \mathcal{Y}^T} P(\mathbf{y}|\mathbf{x}) &= \arg \max_{\mathbf{y} \in \mathcal{Y}^T} \frac{\exp(\boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{Y}^T} \exp(\boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y}'))} \\ &= \arg \max_{\mathbf{y} \in \mathcal{Y}^T} \exp(\boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y})) \\ &= \arg \max_{\mathbf{y} \in \mathcal{Y}^T} \boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y}) \\ &= \arg \max_{\mathbf{y} \in \mathcal{Y}^T} \boldsymbol{\theta}^\top \sum_{i=1}^T f_i(\mathbf{x}, y_{i-1}, y_i) \\ &= \arg \max_{\mathbf{y} \in \mathcal{Y}^T} \sum_{i=1}^T \boldsymbol{\theta}^\top f_i(\mathbf{x}, y_{i-1}, y_i) \end{aligned}$$



# Factorization

Factorize  $\theta^\top f(x, y)$  with respect to timestep  $i$

$$\begin{aligned} \sum_{i=1}^T \theta^\top f_i(x, y_{i-1}, y_i) &= \underbrace{\sum_{j \leq i-1} \theta^\top f_j(x, y_{j-1}, y_j)}_{\text{past}} \\ &+ \underbrace{\theta^\top f_i(x, y_{i-1}, y_i)}_{\text{present}} \\ &+ \underbrace{\sum_{k \geq i+1} \theta^\top f_k(x, y_{k-1}, y_k)}_{\text{future}} \end{aligned} \quad (11)$$

# Viterbi Algorithm

$$s_i(k, k') = \theta^\top f_i(x, y_{i-1} = k', y_i = k)$$

---

**Algorithm 11** The Viterbi algorithm. Each  $s_m(k, k')$  is a local score for tag  $y_m = k$  and  $y_{m-1} = k'$ .

---

```
for  $k \in \{0, \dots, K\}$  do
     $v_1(k) = s_1(k, \diamond)$ 
for  $m \in \{2, \dots, M\}$  do
    for  $k \in \{0, \dots, K\}$  do
         $v_m(k) = \max_{k'} s_m(k, k') + v_{m-1}(k')$ 
         $b_m(k) = \operatorname{argmax}_{k'} s_m(k, k') + v_{m-1}(k')$ 
 $y_M = \operatorname{argmax}_k s_{M+1}(\diamond, k) + v_M(k)$ 
for  $m \in \{M-1, \dots, 1\}$  do
     $y_m = b_m(y_{m+1})$ 
return  $y_{1:M}$ 
```

---

[Eisenstein, 2018]

# Parameter Estimation

---

# Parameter Estimation: Logistic regression

When label  $y$  is still a single component

$$\frac{\partial \log P(y|x; \theta)}{\partial \theta} = f(x, y) - \mathbb{E}_{Y|X}[f(x, y)] \quad (12)$$

where

$$\mathbb{E}_{Y|X}[f(x, y)] = \sum_{y \in \mathcal{Y}} \{P(y|x)f(x, y)\} \quad (13)$$

# Parameter Estimation: CRFs

When label  $\mathbf{y}$  is a sequence

$$\frac{\partial \log P(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbf{f}(\mathbf{x}, \mathbf{y}) - \mathbb{E}_{Y|X}[\mathbf{f}(\mathbf{x}, \mathbf{y})] \quad (14)$$

where

$$\mathbf{f}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^T f_i(\mathbf{x}, y_{i-1}, y_i) \quad (15)$$

and

$$\mathbb{E}_{Y|X}[\mathbf{f}(\mathbf{x}, \mathbf{y})] = \sum_{\mathbf{y} \in \mathcal{Y}^T} \left\{ P(\mathbf{y}|\mathbf{x}) \mathbf{f}(\mathbf{x}, \mathbf{y}) \right\} \quad (16)$$

# Expectation

$$\begin{aligned}\mathbb{E}_{Y|X}[f(x, y)] &= \sum_{y \in \mathcal{Y}^T} \left\{ P(y|x) f(x, y) \right\} \\&= \sum_{y \in \mathcal{Y}^T} \left\{ P(y|x) \sum_{i=1}^T f_i(x, y_{i-1}, y_i) \right\} \\&= \sum_{y \in \mathcal{Y}^T} \sum_{i=1}^T \left\{ P(y|x) f_i(x, y_{i-1}, y_i) \right\} \\&= \sum_{i=1}^T \sum_{y \in \mathcal{Y}^T} \left\{ P(y|x) f_i(x, y_{i-1}, y_i) \right\} \\&= \sum_{i=1}^T \sum_{y_{i-1} \in \mathcal{Y}; y_i \in \mathcal{Y}} \left\{ P(y_{i-1}, y_i|x) f_i(x, y_{i-1}, y_i) \right\}\end{aligned}$$

# Expectation

$$\begin{aligned}\mathbb{E}_{Y|X}[f(x, y)] &= \sum_{y \in \mathcal{Y}^T} \left\{ P(y|x) f(x, y) \right\} \\&= \sum_{y \in \mathcal{Y}^T} \left\{ P(y|x) \sum_{i=1}^T f_i(x, y_{i-1}, y_i) \right\} \\&= \sum_{y \in \mathcal{Y}^T} \sum_{i=1}^T \left\{ P(y|x) f_i(x, y_{i-1}, y_i) \right\} \\&= \sum_{i=1}^T \sum_{y \in \mathcal{Y}^T} \left\{ P(y|x) f_i(x, y_{i-1}, y_i) \right\} \\&= \sum_{i=1}^T \sum_{y_{i-1} \in \mathcal{Y}; y_i \in \mathcal{Y}} \left\{ P(y_{i-1}, y_i|x) f_i(x, y_{i-1}, y_i) \right\}\end{aligned}$$

# Basic Operation

$$P(y_{i-1}, y_i | x) = \frac{\sum_{y \setminus \{y_{i-1}, y_i\}} \exp(\theta^\top f(x, y))}{\sum_y \exp(\theta^\top f(s, y))} \quad (17)$$

Basic operation

$$\sum_{\tilde{y}} \exp(\theta^\top f(s, \tilde{y})) \quad (18)$$

where  $\tilde{y}$  could be

- ▶  $\tilde{y} = y \setminus \{y_{i-1}, y_i\}$
- ▶  $\tilde{y} = y$



$$\psi_i(y_{i-1}, y_i)$$

$$\begin{aligned}\psi(\mathbf{y}) &= \exp(\boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y})) \\ &= \exp\left(\sum_{i=1}^T \boldsymbol{\theta}^\top f_i(\mathbf{x}, y_{i-1}, y_i)\right) \\ &= \prod_{i=1}^T \psi_i(y_{i-1}, y_i)\end{aligned}$$

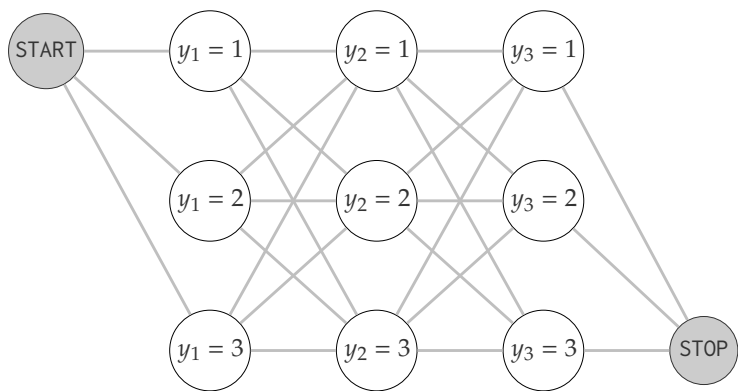
where

$$\psi_i(y_{i-1}, y_i) = \exp(\boldsymbol{\theta}^\top f_i(\mathbf{x}, y_{i-1}, y_i)) \quad (19)$$

# What We Need?

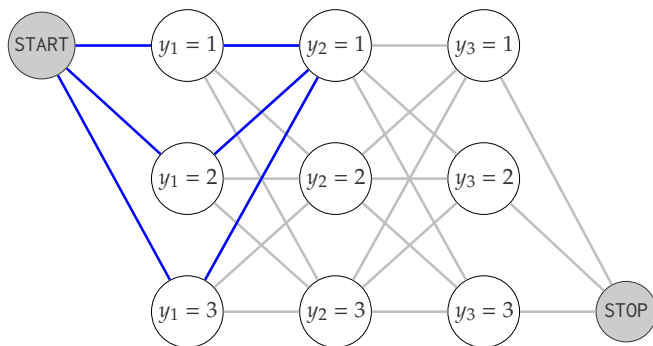
$$\begin{aligned} P(y_{i-1}, y_i | \mathbf{x}) &= \frac{\sum_{\mathbf{y} \setminus \{y_{i-1}, y_i\}} \exp(\boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y}))}{\sum_{\mathbf{y}} \exp(\boldsymbol{\theta}^\top f(\mathbf{x}, \mathbf{y}))} \\ &= \frac{\sum_{\mathbf{y} \setminus \{y_{i-1}, y_i\}} \prod_{i=1}^T \psi_i(y_{i-1}, y_i)}{\sum_{\mathbf{y}} \prod_{i=1}^T \psi_i(y_{i-1}, y_i)} \end{aligned}$$

# Forward-Backward Algorithm: Basic idea



# Forward-Backward Algorithm: Forward term

$$P(y_2 = 1, y_3 = 1) \propto \sum_{y_1} \psi_1(\text{START}, y_1) \psi_2(y_1, y_2 = 1) \cdot \psi_3(y_2 = 1, y_3 = 1) \quad (20)$$



# Forward-Backward Algorithm: Forward term

Forward term

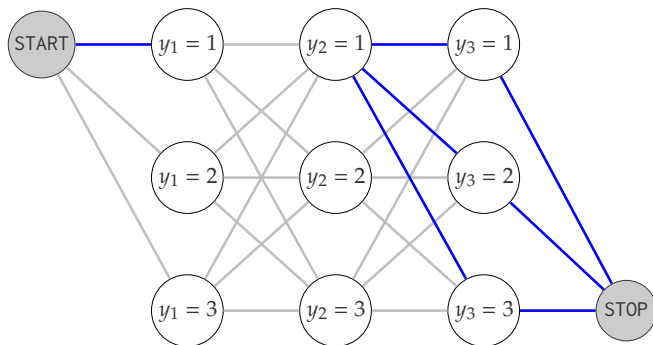
$$\alpha_i(y) = \sum_{y' \in \mathcal{Y}} \alpha_{i-1}(y') \psi_i(y', y) \quad (21)$$

Base case

$$\alpha_1(y) = \psi(\text{START}, y) \quad (22)$$

# Forward-Backward Algorithm: Back term

$$P(y_1 = 1, y_2 = 1) \propto \sum_{y_3} \psi_1(\text{START}, y_1) \psi_2(y_1 = 1, y_2 = 1) \cdot \psi_3(y_2 = 1, y_3) \quad (23)$$



# Forward-Backward Algorithm: Backward term

Backward term

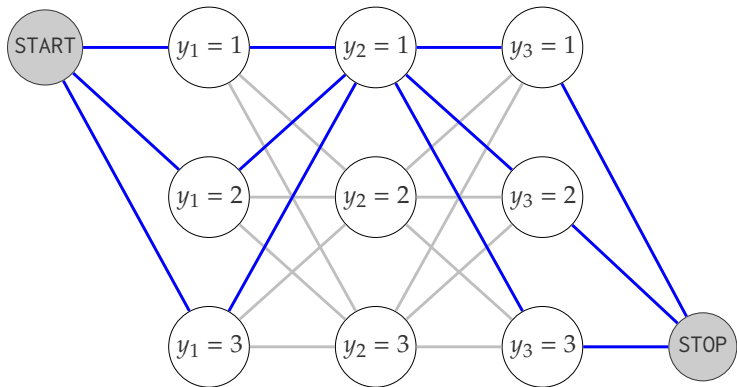
$$\beta_i(y) = \sum_{y' \in \mathcal{Y}} \psi_{i+1}(y, y') \beta_{i+1}(y') \quad (24)$$

Base case

$$\beta_T(y) = 1 \quad (25)$$

# Test

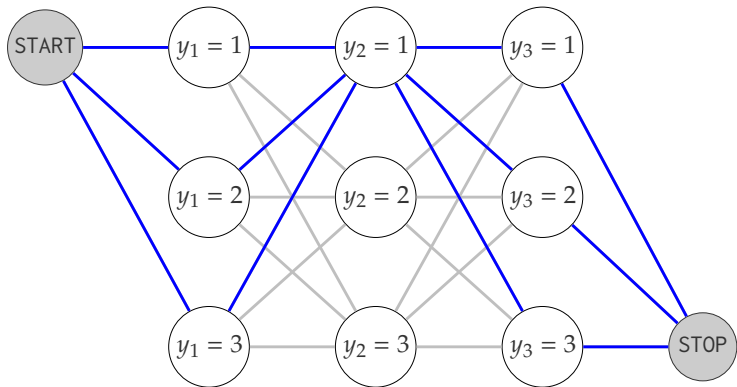
What we compute here?





# Test

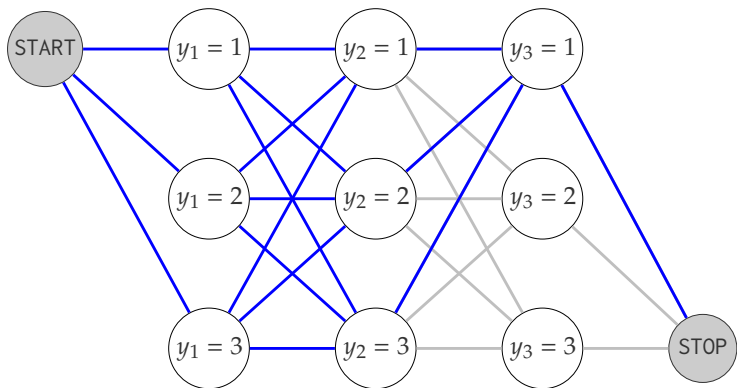
What we compute here?



Answer:  $P(y_2 = 1)$

# Normalization Term $Z$

$$Z = \sum_{\mathbf{y} \in \mathcal{Y}^T} \psi(\mathbf{y}) = \sum_{\mathbf{y} \in \mathcal{Y}} \alpha_T(\mathbf{y}) \beta_T(\mathbf{y}) = \sum_{\mathbf{y} \in \mathcal{Y}} \alpha_T(\mathbf{y}) \quad (26)$$



# Expectation

$$\begin{aligned}\mathbb{E}_{Y|X}[f(x, y)] &= \sum_{y \in \mathcal{Y}^T} \left\{ P(y|x) f(x, y) \right\} \\&= \sum_{y \in \mathcal{Y}^T} \left\{ P(y|x) \sum_{i=1}^T f_i(x, y_{i-1}, y_i) \right\} \\&= \sum_{y \in \mathcal{Y}^T} \sum_{i=1}^T \left\{ P(y|x) f_i(x, y_{i-1}, y_i) \right\} \\&= \sum_{i=1}^T \sum_{y \in \mathcal{Y}^T} \left\{ P(y|x) f_i(x, y_{i-1}, y_i) \right\} \\&= \sum_{i=1}^T \sum_{y_{i-1} \in \mathcal{Y}; y_i \in \mathcal{Y}} \left\{ P(y_{i-1}, y_i|x) f_i(x, y_{i-1}, y_i) \right\}\end{aligned}$$

# Parameter Estimation: CRFs

$$\frac{\partial \log P(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbf{f}(\mathbf{x}, \mathbf{y}) - \mathbb{E}_{Y|\mathbf{X}}[\mathbf{f}(\mathbf{x}, \mathbf{y})] \quad (27)$$

where

$$\mathbf{f}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^T f_i(\mathbf{x}, y_{i-1}, y_i) \quad (28)$$

and

$$\mathbb{E}_{Y|\mathbf{X}}[\mathbf{f}(\mathbf{x}, \mathbf{y})] = \sum_{\mathbf{y} \in \mathcal{Y}^T} \left\{ P(\mathbf{y}|\mathbf{x}) \mathbf{f}(\mathbf{x}, \mathbf{y}) \right\} \quad (29)$$

# Summary

1. Conditional Random Fields
  - 1.1 Logistic Regression
  - 1.2 Decomposition and CRF Formulation
  - 1.3 Directed vs. Undirected Graphical Models
2. Inference
  - 2.1 Viterbi Algorithm
3. Parameter Estimation
  - 3.1 Gradient based Learning
  - 3.2 Forward-Backward Algorithm

# Reference



Bishop, C. M. (2006).  
*Pattern Recognition and Machine Learning*.  
Springer-Verlag.



Eisenstein, J. (2018).  
*Natural Language Processing*.  
MIT Press.



Lafferty, J., McCallum, A., and Pereira, F. (2001).  
Conditional random fields: Probabilistic models for segmenting and labeling sequence data.  
In *ICML*.