

# CS 6501 Natural Language Processing

## Sequence Labeling (I)

---

Yangfeng Ji

September 19, 2019

Department of Computer Science  
University of Virginia



ENGINEERING

# Overview

1. Final Project Suggestion
2. Problem Definition
3. Hidden Markov Models
4. Parameter Estimation
5. Viterbi Decoding

# Final Project Suggestion

---

# Text Classification and Sentiment Analysis

Projects from last year:

- ▶ A Deeper Look into Social Media-focused Sentiment Analysis
- ▶ Citation Recommendation by Abstract

# Convation Modeling

Projects from last year:

- ▶ Neural Dialog System with Personality
- ▶ Chit-chat Bot Modeling

Projects from last year:

- ▶ A Deep Learning Approach for Meme Generation
- ▶ Variational Image Captioning Using Deterministic Attention
- ▶ Neural Style Transfer for Natural Language
- ▶ Generating Subject Line from Email Text

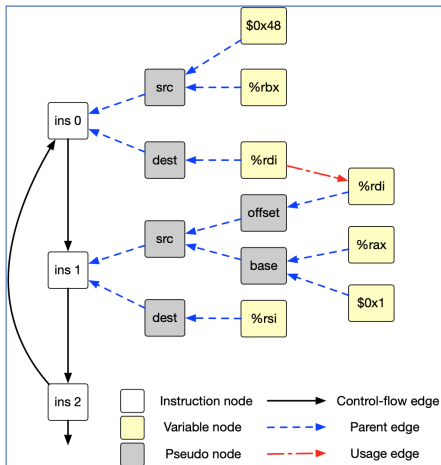
# Interpretability and Adversarial Learning

Projects from last year:

- ▶ Incorporating Textual Data from Reviews for Explainable Recommendation
- ▶ Label Flipping Attacks on Sentiment Analysis Systems
- ▶ Topic-based Interpretability Improvement on Citation Recommendation System

# Sequential Modeling for Computer Architecture

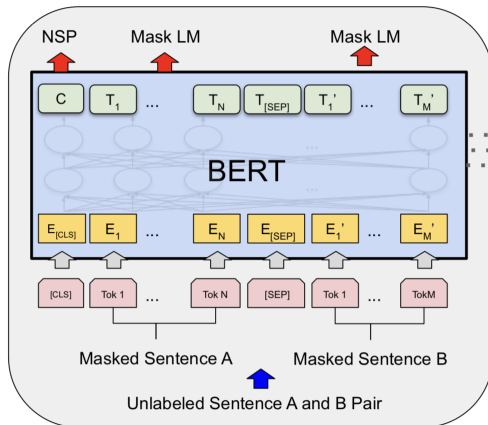
- Using neural network and NLP models to analyze machine code



Ashish Venkat  
(venkat@virginia.edu)



# Running BERT on Edge Devices



Marzieh Lenjani  
(marzieh.lenjani@gmail.com)

# Problem Definition

---

# Part of Speech

- ▶ A way to categorize words with similar *grammatical* properties
- ▶ Common English POS tags
  - ▶ NOUN: used to name persons, things, animals, places etc.  
e.g., Tom Hanks, yesterday, Grounds
  - ▶ VERB: show an action or state  
e.g., fight, was
  - ▶ PRONOUN: replacement of nouns  
e.g., she, his, it, theirs
  - ▶ ADJECTIVE: used to describe a noun or a pronoun  
e.g., large, beautiful

# Part of Speech (II)

- ▶ Common English POS tags (cont.)
  - ▶ ADVERB: used to describe adjectives, verbs, or another adverb  
e.g., gracefully, yesterday, very
  - ▶ PREPOSITION: specify location or a location in time  
e.g., above, near, since
  - ▶ CONJUNCTION: join words, phrases, or clauses together  
e.g., and, for
  - ▶ INTERJECTION: convey strong emotions  
e.g., Ouch, Hey

## Example

Teacher Strikes Idle Children

- ▶ Teacher<sub>Noun</sub> Strikes<sub>Noun</sub> Idle<sub>Verb</sub> Children<sub>Noun</sub>
- ▶ Teacher<sub>Noun</sub> Strikes<sub>Verb</sub> Idle<sub>Adj</sub> Children<sub>Noun</sub>

[Eisenstein, 2018, Chap 8]

# Sequence Labeling

From a training set, to learn a mapping  $p(\mathbf{y} \mid \mathbf{x})$ ,

$$\hat{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmax}} p(\mathbf{y} \mid \mathbf{x}) \quad (1)$$

where

- ▶  $\mathbf{x}$ : a sentence (a sequence of words)
- ▶  $\mathbf{y}$ : the POS tag sequence of  $\mathbf{x}$  (a **sequence** of POS tags)

## Example

$\mathbf{x}$	$x_1$	$x_2$	$x_3$	$x_4$
	Teacher	Strikes	Idle	Children
$\mathbf{y}$	$y_1$	$y_2$	$y_3$	$y_4$
	NOUN	NOUN	VERB	NOUN

# Label Classification

$$\hat{y} = \underset{y}{\operatorname{argmax}} p(y \mid x) \quad (2)$$

- ▶  $x$ : entire sentence only one token
- ▶  $y$ : entire sequence only the corresponding tag
- ▶  $P(y|x)P(y_i|x_i)$

## Example

$x$	$x_1$	$x_2$	$x_3$	$x_4$
	Teacher	Strikes	Idle	Children
$y$	$y_1$	$y_2$	$y_3$	$y_4$
	NOUN	NOUN	VERB	NOUN

# Label Classification

## Example

$x$	$x_1$	$x_2$	$x_3$	$x_4$
	Teacher	Strikes	Idle	Children
$y$	$y_1$	$y_2$	$y_3$	$y_4$
	NOUN	NOUN	VERB	NOUN
	NOUN	VERB	ADJ	NOUN

Limitations [TODO: revise the following description]

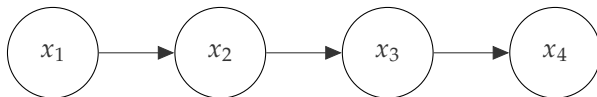
- ▶ No constraint from the previous POS tag
  - ▶ Solution: sequence labeling (e.g., hidden Markov models, conditional random fields)
- ▶ No information from the surrounding words
  - ▶ Solution: conditional random fields



# Markov Chain

Modeling the dependency between  $\{y_i\}$  as

$$p(\mathbf{y}) = p(y_1)p(y_2 | y_1)p(y_3 | y_2)p(y_4 | y_3) \quad (3)$$



## Question

How to merge the conditional dependence from  $p(\mathbf{y})$  into  $p(\mathbf{y} | \mathbf{x})$ ?

# Hidden Markov Models

---

# Generative Models

- ▶ Observation  $x$
- ▶ Target variable  $y$

Bayes rule

$$\begin{aligned} p(y|x) &= \frac{p(x|y) \cdot p(y)}{p(x)} \\ &\approx \underbrace{p(y)}_{\text{prior}} \cdot \underbrace{p(x|y)}_{\text{likelihood}} \end{aligned} \tag{4}$$

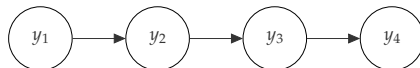
# Prior $p(\mathbf{y})$

$$p(\mathbf{y} \mid \mathbf{x}) \approx p(\mathbf{y}) \cdot p(\mathbf{x} \mid \mathbf{y}) \quad (5)$$

Factorization

$$p(\mathbf{y}) = \prod_{i=1} \underbrace{p(y_i \mid y_{i-1})}_{\text{Transition probability}} \quad (6)$$

Graphical model



$$p(\mathbf{y}) = p(y_1) \cdot p(y_2 \mid y_1) \cdot p(y_3 \mid y_2) \cdot p(y_4 \mid y_3)$$

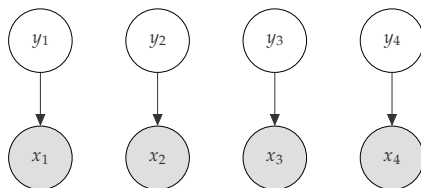
# Likelihood $P(\mathbf{x}|\mathbf{y})$

$$p(\mathbf{y} | \mathbf{x}) = p(\mathbf{y}) \cdot p(\mathbf{x} | \mathbf{y}) \quad (7)$$

Factorization

$$p(\mathbf{x} | \mathbf{y}) = \prod_{i=1} \underbrace{p(x_i | y_i)}_{\text{Emission probability}} \quad (8)$$

Graphical model

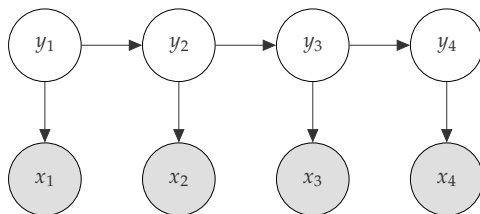


$$p(\mathbf{x} | \mathbf{y}) = p(x_1 | y_1) \cdot p(x_2 | y_2) \cdot p(x_3 | y_3) \cdot p(x_4 | y_4)$$

# Hidden Markov Models

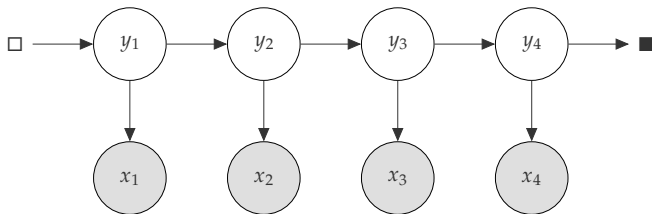
$$P(\mathbf{x}, \mathbf{y}) = \prod_{i=1} \left\{ P(y_i | y_{i-1}) P(x_i | y_i) \right\} \quad (9)$$

Graphical model



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

# Two Special Tokens



$$P(\mathbf{x}, \mathbf{y}) = P(y_1 \mid \square) \prod_{i=1}^4 \left\{ P(y_i \mid y_{i-1}) P(x_i \mid y_i) \right\} p(\blacksquare \mid y_4) \quad (10)$$

# Two Questions

- ▶ Learning: parameter estimation
  - ▶  $p(y_n \mid y_{n-1}) = ?$
  - ▶  $p(x_n \mid y_n) = ?$
- ▶ Prediction: inference/decoding
  - ▶  $\hat{y} = \operatorname{argmax}_y p(y \mid x)$



# Parameter Estimation

---

# Training Corpus

## Training corpus

- ▶ they<sub>PRON</sub> can<sub>VERB</sub> fish<sub>NOUN</sub>
- ▶ teacher<sub>NOUN</sub> strikes<sub>VERB</sub> idle<sub>ADJ</sub> children<sub>NOUN</sub>
- ▶ ...

How to estimate the following probabilities?

$$\begin{aligned}P(x_i|y_i) &=? \\ P(y_i|y_{i-1}) &=?\end{aligned}\tag{11}$$

Transition probability

$$P(y_i|y_{i-1}) = \frac{\#(y_i, y_{i-1})}{\#(y_{i-1})} \quad (12)$$

Emission probability

$$P(x_i|y_i) = \frac{\#(x_i, y_i)}{\#(y_i)} \quad (13)$$

# Viterbi Decoding

---

# Decoding by Brute-force Search

Given a sentence

The dog barks

and the possible POS tags {D,N,V}. A brute-force algorithm will try every possible combination as

D	D	D
D	D	N
D	D	V
D	N	D
D	N	N
D	N	N
⋮	⋮	⋮

There are  $3^3 = 27$  possible sequences in this case.

# Decoding $y_i$

$$\hat{y} = \operatorname{argmax}_y p(y \mid x) \quad (14)$$

$$= \operatorname{argmax}_y p(x, y) \quad (15)$$

$$(16)$$

With conditional dependency

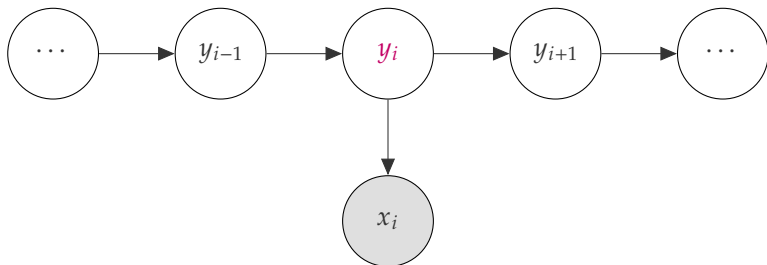
$$p(x, y) = \prod_{i=1} \left\{ p(y_i \mid y_{i-1}) p(x_i \mid y_i) \right\} \quad (17)$$

$$= \cdots \underbrace{p(y_i \mid y_{i-1}) \cdot p(y_{i+1} \mid y_i) \cdot p(x_i \mid y_i)}_{\text{items related to } y_i} \cdots \quad (18)$$

# Decoding $y_i$ (II)

$$p(\mathbf{x}, \mathbf{y}) = \prod_{i=1} \left\{ p(\mathbf{y}_i \mid y_{i-1}) p(x_i \mid \mathbf{y}_i) \right\} \quad (19)$$

$$= \cdots \underbrace{p(y_i \mid y_{i-1}) \cdot p(y_{i+1} \mid y_i) \cdot p(x_i \mid y_i)}_{\text{items related to } y_i} \cdots \quad (20)$$



# Factorization

Factorize  $p(\mathbf{x}, \mathbf{y})$  with respect to  $(x_i, y_i)$

$$\begin{aligned} p(\mathbf{x}, \mathbf{y}) &= p(\mathbf{x}_{\leq i-1}, \mathbf{y}_{\leq i-1}) \cdot p(x_i, y_i | \mathbf{y}_{\leq i-1}) \cdot p(\mathbf{x}_{\geq i+1}, \mathbf{y}_{\geq i+1} | y_i) \\ &= p(\mathbf{x}_{\leq i-1}, \mathbf{y}_{\leq i-1}) \cdot p(x_i | y_i) \cdot p(y_i | y_{i-1}) \end{aligned} \quad (21)$$

$$\cdot p(\mathbf{x}_{\geq i+1}, \mathbf{y}_{\geq i+1} | y_i) \quad (22)$$

Three components

$$\underbrace{p(\mathbf{x}_{\leq i-1}, \mathbf{y}_{\leq i-1})}_{\leq i-1} \cdot \underbrace{p(x_i | y_i) \cdot p(y_i | y_{i-1})}_i \cdot \underbrace{p(\mathbf{x}_{\geq i+1}, \mathbf{y}_{\geq i+1} | y_i)}_{> i} \quad (23)$$



# Principle of Optimality

Assume we have the optimal decoded sequence  $\hat{\mathbf{y}}$ , such that

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} p(\mathbf{x}, \mathbf{y}), \quad (24)$$

then,  $\mathbf{y}_{\geq i}$  is the also the optimal subsequence of the following subproblem

$$\hat{\mathbf{y}}_{\geq i} = \operatorname{argmax}_{\mathbf{y}_{\geq i}} p(\mathbf{x}_{\geq i}, \mathbf{y}_{\geq i} \mid \hat{\mathbf{y}}_i) \quad (25)$$

Justification:

$$\begin{aligned} p(\mathbf{x}, \mathbf{y}) &= p(\mathbf{x}_{\leq i-1}, \hat{\mathbf{y}}_{\leq i-1}) \\ &\quad \cdot p(x_i \mid \mathbf{y}_i) \cdot p(y_i \mid \hat{\mathbf{y}}_{i-1}) \\ &\quad \cdot p(\mathbf{x}_{\geq i+1}, \mathbf{y}_{\geq i+1} \mid \mathbf{y}_i) \end{aligned} \quad (26)$$

# Basic Idea of Decoding

$$\underbrace{p(\mathbf{x}_{\leq i-1}, \mathbf{y}_{\leq i-1})}_{\leq i-1} \cdot \underbrace{p(x_i | y_i) \cdot p(y_i | y_{i-1})}_i \cdot \underbrace{p(\mathbf{x}_{\geq i+1}, \mathbf{y}_{\geq i+1} | y_i)}_{> i} \quad (27)$$

- ▶ Forward computation: start from  $y_1$ , for **every** possible value of  $y_i$ , from the best path from  $y_{i-1}$

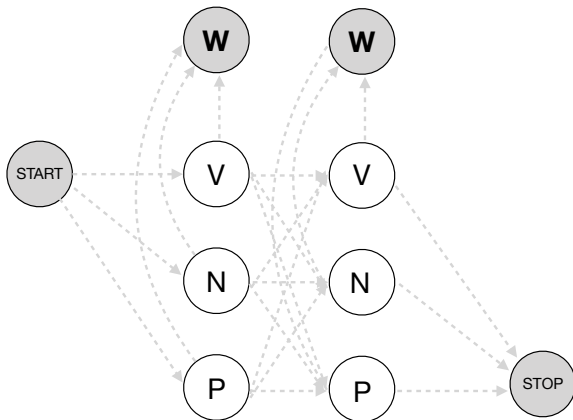
$$\max_{y_{i-1}} p(\mathbf{x}_{\leq i-1}, \mathbf{y}_{\leq i-1}) \cdot p(x_i | y_i) \cdot p(y_i | y_{i-1}),$$

depending on past and present states  $\{\mathbf{y}_{\leq i}\}$

- ▶ Backward tracing: start from  $y_T = \blacksquare$ , for a **given**  $y_{i+1}$  find the best  $y_i$

# Example

$$\max_{y_{i-1}} p(\mathbf{x}_{\leq i-1}, \mathbf{y}_{\leq i-1}) \cdot p(x_i | y_i) \cdot p(y_i | y_{i-1}),$$



# Reference



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