

# CS371N Lecture 14

## Sequence Modeling, POS Tagging

### Announcements

- A2 solutions posted
- A3 due today

### Recap Language modeling

subword tok

↓  
pre-train an LM on the web

↓  
fine-tune

↓  
classify

↓  
generate

↓  
greedy

↓  
beam

→ nucleus  
sampling

Today Structured prediction:  
sequence modeling  
POS tagging

Part-of-speech tagging

Input sentence:  $x_1, \dots, x_n$

Output: POS tags  $y_1, \dots, y_n$  for  
each word

Different than sentiment class.  
One output per word

What are POS tags (and why?)



## Open-class

(N) Nouns → Proper (Google)  
→ common (cat)  
plural vs. singular

(V) verbs tense, "person"  
(inflection)

(J) adjectives idle, yellow

(RB) adverbs swiftly

## Closed-class

(DT) Determiners: articles (the, a)  
Some, many

(CD) Cardinal: numbers

Prepositions: up, on, in, ...

Partic les: made up

Auxiliaries: had

Modals: could/would/should

① What tags are possible for each word?

Fed raises interest rates 0.5 percent

Fed VBD I fed the cat  
NNP proper noun  
VBN "fed up"

raises NNS plural  
VBZ (3rd person sg present)

interest NN  
VBP present "I interest you"  
VB infinitive "I want NLP  
to interest me"

rates NNS  
VBZ

Interps Standard

0.5 CD

weird

weirder

percent NN

## Two ways to tag

- Using classifiers?
- Using Hidden Markov Models

## Classifiers POS tags $y$

Multiclass LR:

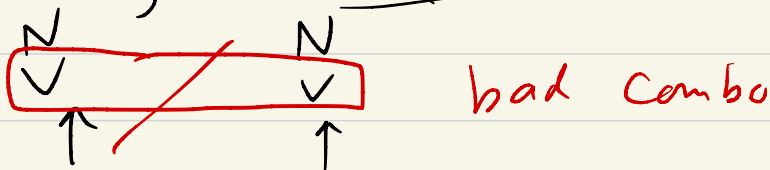
$$P(y_i = t | \bar{x}) \quad \text{iterate over the sequence}$$

letter counting in A3

$$P(\bar{y} | \bar{x}) = \prod_{i=1}^n P(y_i | \bar{x})$$

independent classification of each

Why modeling the sequence?



Fed raises interest rates

$P(y_2 = V | \bar{x})$  is high  $\Rightarrow V$

$P(y_3 = V | \bar{x})$  is high  $\Rightarrow V$

$P(y_2 = V, y_3 = V | \bar{x})$  is low

↓ we don't model this

## Hidden Markov Models

⊕ model the sequence  $(y_3 | y_2)$

⊖ simple generative model

(CRFs: discriminative HMMs,  
Transformer & HMM)

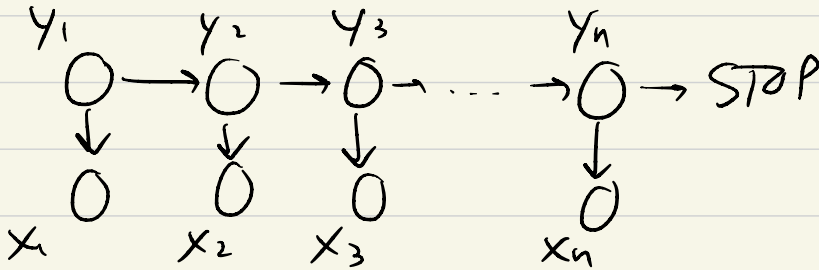
$$y_1 \dots y_n \quad x_1 \dots x_n$$

Generative model  $P(\bar{y} | \bar{x})$

Tags  $y_i \in \mathcal{T}$  words  $x_i \in \mathcal{V}$

HMM:

$$P(\bar{y}, \bar{x}) = P(y_1) P(x_1 | y_1) P(y_2 | y_1) P(x_2 | y_2) P(y_3 | y_2) P(x_3 | y_3) \dots P(\text{STOP} | y_n)$$



Assumptions:

- ① Each  $y$  depends only on the previous  $y$  (Markov)
- ② Each  $x_i$  is indep. of all else given  $y_i$



joint  
↓

Goal: We model  $P(\bar{y}, \bar{x})$  ← conditional  
 but we care about  $P(\bar{y} | \bar{x})$   
 (you give me a sentence  $\bar{x}$  and  
 I give you tags)

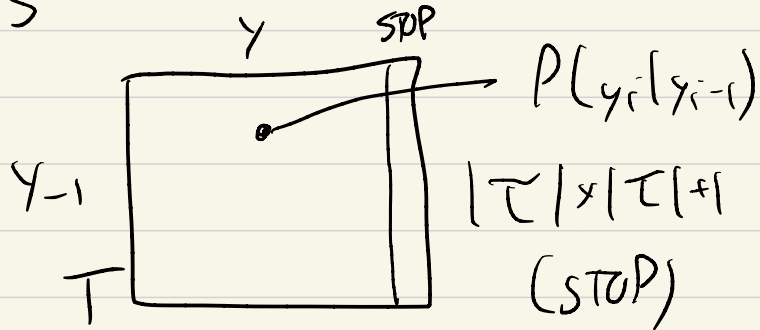
Parameters Three types

$P(y_i)$   
initial dist



$|\tau|$ -len vector  
sums to 1

$P(y_i | y_{i-1})$   
transition probs

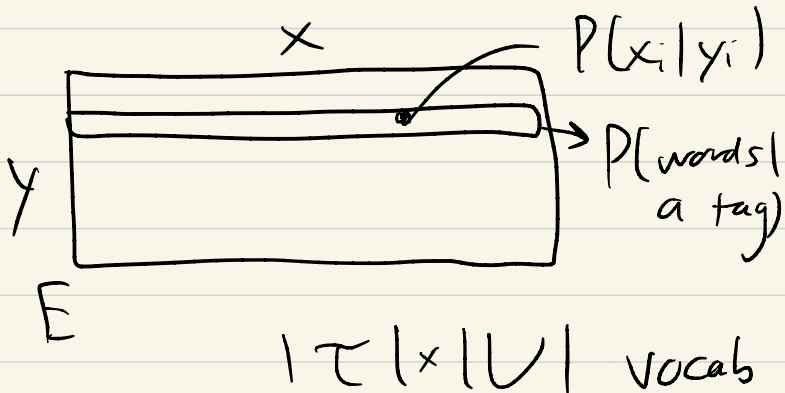


each row is a dist.

$P(\text{rates} | N)$

$P(x_i | y_i)$

emissions



store distributions explicitly