# CS 371 Lecture 14
## Sequence Modeling I: Part of speech
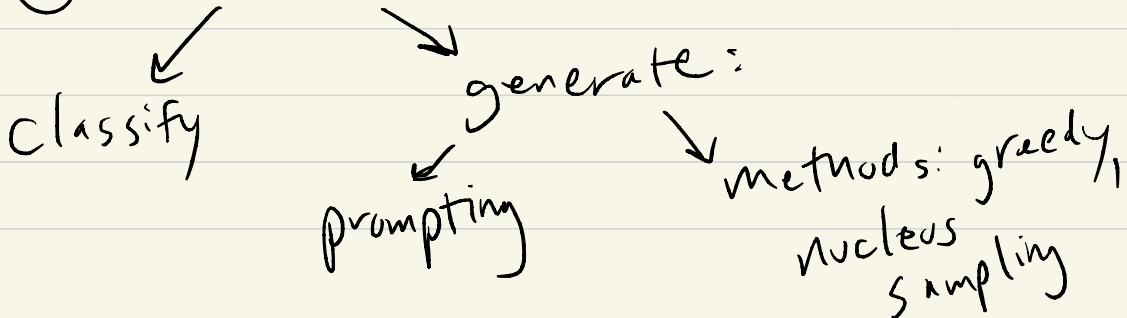
## Recap   Language modeling

① Tokenization (subword) / featurization
② Pre-training phase: skip-gram
   language modeling
③ Fine-tuning phase: train something like
   in A2
   or fine-tune BERT/GPT/...

④ Inference

classify       generate:

prompting       methods: greedy,
                  nucleus
                    sampling

Today   Structured prediction
        — sequence modeling: part of speech
        — syntactic parsing         tagging
        → today: POS and Hidden Markov Models

## Part-of-speech tagging

Input: sentence $X_1 \ldots X_n$

Output: POS tags $y_1 \ldots y_n$ for each word

What are POS tags?

| N | N | V | N |
|---|---|---|---|
| N | V | ADJ | N |
| teacher | strikes | idle | kids |

Predicting POS $\Longleftrightarrow$ interpreting the sentence

Text-to-speech: record   verb or
                              noun

# POS tags

open-class: new words with these tags are always emerging

Closed-class: known set

## Open-Class:

(N) Nouns:
- → Proper (Google)    NNP
- → Common (shoe) ← NN
- → plural vs. singular
  (NNS)

(V) Verbs: features like tense, person (1st or 3rd)

In standard datasets: VBZ

VBD: past tense    "3rd person present singular verb"

(J) Adjectives → yellow, idle

(RB) adverbs → swiftly

## Closed-class

(DT) Determiner: articles (the, a)
     some, many

(CD) cardinals: numbers

(IN) prepositions : up, on, in,

(RP) particles : made _up_

Modals (could/would/should), auxiliary
verbs (had)

① what tags are possible for each word?
② what sentences make sense?

Fed raises interest rates 0.5 percent

Fed {
  NNP    Federal Reserve
  VBD    "They Fed me"
  VBN    "I was fed up"
}

raises {
  NNS    $ raises
  VBZ    "she raises"
}

interest {
  NN
  VBP    "Pyramids interest me"
  VB     "I want NLP to interest me"
}

rates {
  NNS
  VBZ
}

0.5     CD
percent NN          Sentences: standard

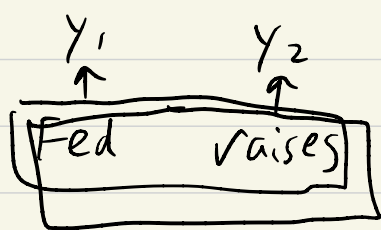                              alt 1        alt 2
                                    alt 3

# Methods for POS tagging

(later) Hidden Markov Models

(now) Classifiers

## Classifier   POS tags $y$

MC class.   $P(y \mid \bar{x}) \in Y$

For seqs :   $P(y_i = t \mid \bar{x}, i)$

$y_1$      $y_2$          Run classifier twice

↑          ↑

| Fed    raises |            BoW ✗     no
                                    position info

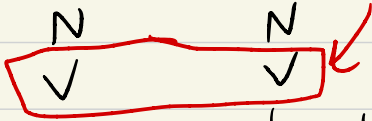Position-sensitive BoW:  "Unigram = Fed &
                                offset = −1 "

          predicting $y_2$  ↙

$$P(\bar{y} \mid \bar{x}) = \prod_{i=1}^{n} P(y_i \mid \bar{x}, i)$$          independent
                                                    classifier

$$N \qquad N$$
$$\vee \qquad \vee$$

Fed   raises   interest   rates

$(y_2, y_3)$   should not be   $(V, V)$
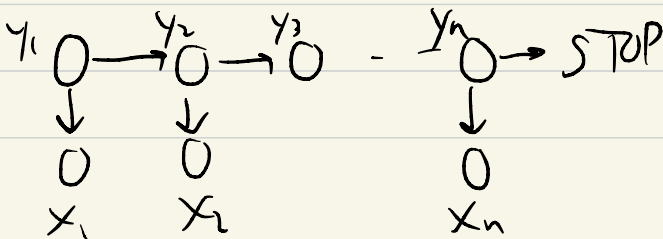
Instead we want a <u>Sequence model</u>
really model $P(\bar{y} | \bar{x})$

↗
the whole sequence

<u>HMMs</u>   models of sequences, can capture
$$P(y_i | y_{i-1}) = \text{transitions}$$

<u>generative</u> models   $P(\bar{x}, \bar{y})$

HMM: $P(\bar{y}, \bar{x}) =$

$P(y_1) P(x_1 | y_1) \; P(y_2 | y_1) \, P(x_2 | y_2) \; P(y_3 | y_2) \cdots$
$$P(STOP | y_n)$$

$y_1$ O $\xrightarrow{y_2}$ O $\xrightarrow{y_3}$ O — $\overset{y_n}{O} \rightarrow$ STOP
↓ ↓ ↓
O O O
$x_1$ $x_2$ $x_n$

# Assumptions

(1) $\bar{y}$s are modeled with a "bigram LM"
(Markov property: $y_i$ is conditionally
independent of $y_1 \dots y_{i-2}, x_1 \dots x_{i-1}$
given $y_{i-1}$)

(2) Each $x_i$ is indep. of everything else
given $y_i$

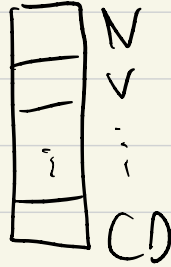Generative story: (1) Pick $y_1$ (2) Pick $x_1 | y_1$
(3) Pick $y_2 | y_1$ (4) Pick $x_2 | y_2$ --

Goal: model $P(\bar{x}, \bar{y})$, but ultimately
we want $P(\bar{y} | \bar{x})$
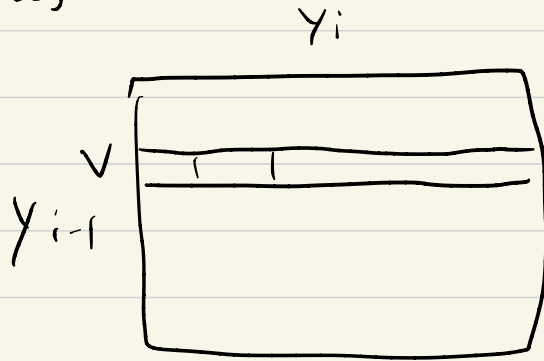
$\mathcal{V}$ vocab, $\tau$ tags

Three types of parameters:

$P(y_1)$ initial distribution

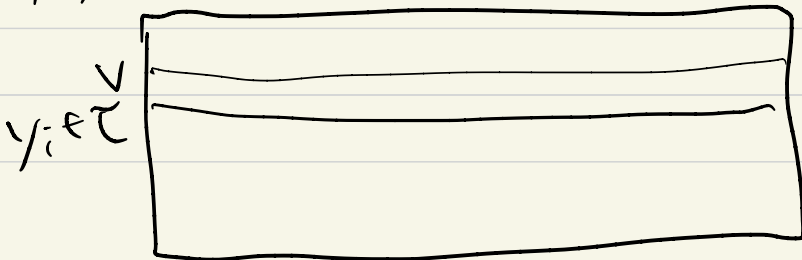$|\tau|$-len vector, adds to 1

Transition probs

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

$y_i$

$y_{i-1}$

$\rightarrow P(y_i | V)$

" 70% N

0% V

: "

$|\tau| \times (|\tau|+1)$ matrix

↖ STOP

Emission probs

$P(x_i | y_i)$

$x_i \in \mathcal{V}$

$y_i \in \tau$

$P(x_i | V)$

" 5% go

5% eat

:

"