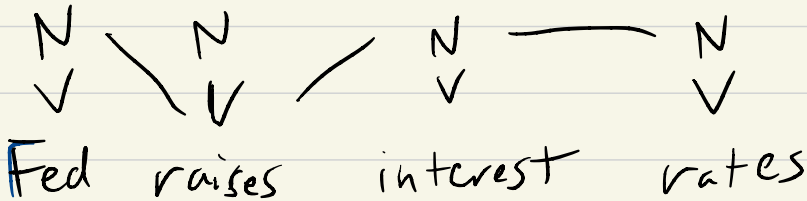# CS371N Lecture 15
# HMMs, Viterbi

<u>Announcements</u>

- A4 due in a week
- Midterm due next Thurs
- OPTIONAL: Independent FP proposals due
  after midterm

<u>Recap</u> POS tagging

N   N   N —— N
V   V   V   V
Fed raises interest rates
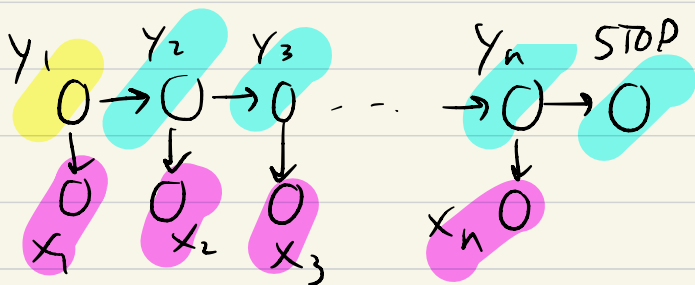
There are constraints on what makes

a well-formed tag sequence (e.g., rare
                              to have V-V)

(discrim.: $P(\bar{y}|\bar{x})$)

# HMMs Generative model of sequences

$$P(\bar{y}, \bar{x}) = P(y_1) P(x_1|y_1) P(y_2|y_1) P(x_2|y_2) \cdots$$



Parameters:

Initial $P(y_1)$ vector

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

Emissions $P(x_i|y_i)$ matrix

$P(go|V) = 0.2$

$P(is|V) = 0.2$

$P(eat|V) = 0.1$

Goal: compute $P(\bar{y}|\bar{x})$ given a sequence $x$

$\underline{Ex}$  $T = \{N, V, STOP\}$
$V = \{they, can, fish\}$

Initial $P(y) = \begin{cases} 1.0 & N \\ 0 & V \\ 0 & STOP \end{cases}$

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

|   | N | V | STOP |
|---|---|---|------|
| N | 1/5 | 3/5 | 1/5 |
| V | 1/5 | 1/5 | 3/5 |

Emissions $P(x_i | y_i) =$

|   | they | can | fish |
|---|------|-----|------|
| N | 1 | 0 | 0 |
| V | 0 | 1/2 | 1/2 |

① Compute the probability of

$\begin{pmatrix} N & V & V & STOP \\ they & can & fish & \end{pmatrix}$  $P(STOP|V)$

$P(y_1 = N)$  .  $P(y_2 = V | y_1 = N) \cdot P(V|V)$

$P(x_1 = they | y_1 = N)$  $P(x_2 = can | y_2 = V)$  $P(fish|V)$

$$\frac{1.0 \cdot 3/5 \cdot 1/5 \cdot 3/5}{1.0 \quad 1/2 \quad 1/2} = \frac{9}{500}$$

② Is there a higher-scoring tag sequence for "they can fish"?

## Goal of HMMs

HMMs model $P(\bar{y}, \bar{x})$

They are not good generative models of
  text

What we use them for is $P(\bar{y} \mid \bar{x})$

$$P(\bar{y} \mid \bar{x}) = \frac{P(y \mid \bar{x}) \cdot P(\bar{x})}{P(\bar{x})} = \frac{P(\bar{y}, \bar{x})}{P(\bar{x})}$$

$$\sum_{Y} P(\bar{y}, \bar{x}) \qquad P(\bar{y} \mid \bar{x}) \propto P(\bar{y}, \bar{x})$$

proportional
  to

What this means:

$$\underset{\tilde{y}}{argmax} \; P(\tilde{y} \mid \bar{x}) = \underset{\tilde{y}}{argmax} \; P(\hat{\tilde{y}}, \bar{x})$$

$$= \underset{\tilde{y}}{argmax} \; \log P(\tilde{y}, \bar{x})$$

$$\tilde{y} = \tilde{y}_1, \tilde{y}_2, \tilde{y}_3 \cdots \tilde{y}_n$$

$$= \underset{\tilde{y}_1, \tilde{y}_2 \cdots \tilde{y}_n}{argmax} \; \log P(\tilde{y}_1) + \log P(x_1 \mid \tilde{y}_1)$$
$$+ \log P(\tilde{y}_2 \mid \tilde{y}_1) +$$
$$\log P(x_2 \mid \tilde{y}_2) + \cdots$$

# Viterbi Algorithm
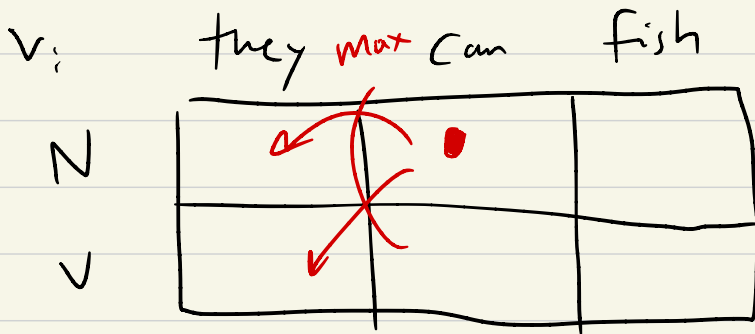
Define $V_i(\tilde{y}_i)$ as the chart

$i$ is index from 1 to $n$

$\tilde{y}_i$ is a tag in $T$

$n \times |T|$ matrix

$V_i(\tilde{y}_i) = \log$ prob of the best sequence
of tags ending in $\tilde{y}$ at index $i$

$V_i$    they max can    fish



$V_2(N)$

max over

$\tilde{y}_i \in \{N, V\}$

compute $V_i \longrightarrow$

N

$\vee$

## Initial

$$V_1(\tilde{y}_1) = \log \underbrace{P(x_1|\tilde{y}_1)}_{\text{emission}} + \log \underbrace{P(\tilde{y}_1)}_{\text{initial}}$$

## Recurrent Compute $V_i$ using $V_{i-1}$

$$V_i(\tilde{y}_i) = \log \underbrace{P(x_i|\tilde{y}_i)}_{\text{emission}}$$

$$+ \max_{\tilde{y}_{prev}} \left[ \log P(\tilde{y}_i|\tilde{y}_{prev}) + V_{i-1}(\tilde{y}_{prev}) \right]$$

$$V_2(\tilde{y}_2) = \log P(x_2|\tilde{y}_2)$$

$$+ \max_{\tilde{y}_1} \left[ \log P(\tilde{y}_2|\tilde{y}_1) + V_1(\tilde{y}_1) \right]$$

## End  $V_n(\tilde{y}_n) = $ recurrent formula +

$$\log P(STOP|\tilde{y}_n)$$

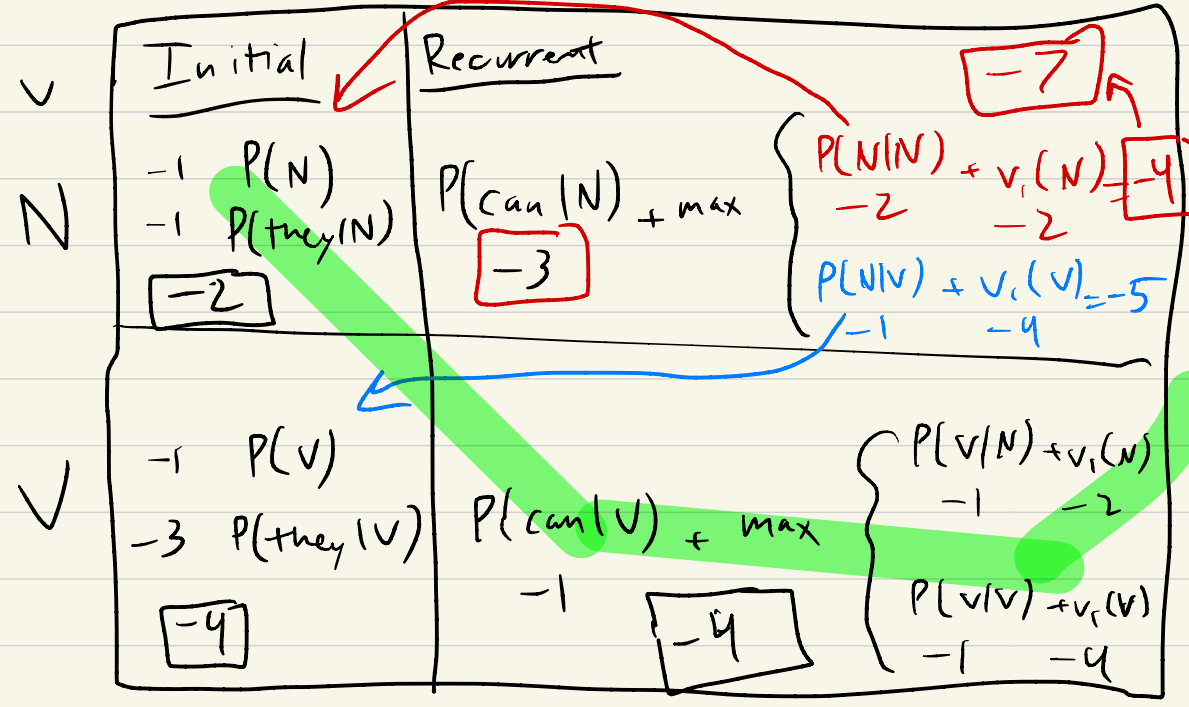Given $V$ chart, extract best sequence
with backpointers

$$\log P(y_1) = \begin{cases} N & -1 \\ V & -1 \end{cases} \leftarrow \text{log probs}$$

$$\log P(y_i | y_{i-1}) = \begin{array}{c|ccc} & N & V & STOP \\ N & -2 & -1 & -1 \\ V & -1 & -1 & -2 \end{array}$$

$$\log P(x_i | y_i) = \begin{array}{c|ccccc} & they & fish & can & dog & a \\ N & -1 & -1 & -3 & & \\ V & -3 & -1 & -1 & & \end{array}$$

$\bar{x} =$ they can fish



|  | 1 Initial | 2 Recurrent |  |
|---|---|---|---|
| N | $-1$ $P(N)$ $-1$ $P(they|N)$ $\boxed{-2}$ | $P(can|N) + max$ $\boxed{-3}$ | $\begin{cases} P(N|N) + v_1(N) = \boxed{-4} \\ -2 \quad\quad -2 \\ P(N|V) + v_1(V) = -5 \\ -1 \quad\quad -4 \end{cases}$ $\boxed{-7}$ |
| V | $-1$ $P(V)$ $-3$ $P(they|V)$ $\boxed{-4}$ | $P(can|V) + max$ $-1$ $\boxed{-4}$ | $\begin{cases} P(V|N) + v_1(N) \\ -1 \quad\quad -2 \\ P(V|V) + v_1(V) \\ -1 \quad\quad -4 \end{cases}$ |

3 fish                    ↙ STOP

$$N \begin{vmatrix} & \\ & \end{vmatrix}$$

-6 $\Rightarrow$ -7 w/ STOP

$Y_{prev} = V$

─────────────────

-6 $\Rightarrow$ -8 w/ STOP

$Y_{prev} = V$

Sequence: N V N    $\log P = -7$

Markov property allowed us to do this efficiently!

# Parameter Estimation

Suppose we have labeled data

$\overline{y}^{(1)}, \overline{x}^{(1)}$     Estimate params by

$\overline{y}^{(2)}, \overline{x}^{(2)}$     counting + normalizing

$\quad \vdots$

$\overline{y}^{(N)}, \overline{x}^{(N)}$

Initial prob $(N) = \dfrac{\text{number of } \overline{y}^{(i)} \text{ w/ } Y_1 = N}{\text{total num exs.}}$

Transition prob $(N \to V) = \dfrac{\text{number of times we saw } N \to V}{\text{number of times we saw } N}$

These maximize $\displaystyle\sum_{i=1}^{N} \log P(\overline{y}^{(i)}, \overline{x}^{(i)})$

<u>Data</u>: English Penn Treebank

44 tags

Assign each word its most frequent
tag = 90% acc.

Trigram HMM tagger: 95%

BERT: 97.5%