## Hidden Markov Models
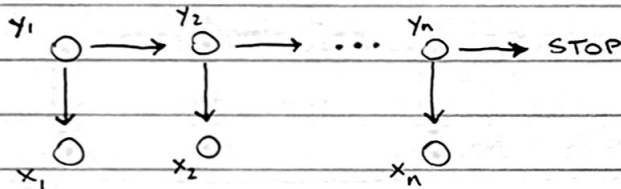
- Generative Sequence Model

Tags $y_i \in T$ (tags)     Words $x_i \in V$ (vocabulary)

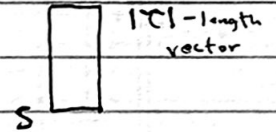$$P(\bar{y}, \bar{x}) = P(y_1) \, P(x_1 | y_1) \, P(y_2 | y_1) \, P(x_2 | y_2) \, \cdots \, P(stop | y_n)$$

$y_1 \;\bigcirc \longrightarrow\; y_2 \;\bigcirc \longrightarrow\; \cdots \; y_n \;\bigcirc \longrightarrow\; STOP$

$\bigcirc$ $\qquad$ $\bigcirc$ $\qquad$ $\bigcirc$

$x_1$ $\qquad\qquad$ $x_2$ $\qquad\qquad$ $x_n$
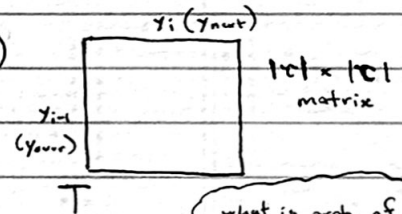
* $y$'s form a <u>markov process</u>:
$y_i$ is conditionally independent
of $y_1 \cdots y_{i-2}$ given $y_{i-1}$.

### Parameters

$P(y_1)$
initial distribution

$\;\;\;\;$ $|T|$-length vector

S

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

$y_i \,(y_{next})$

$y_{i-1}$ $(y_{curr})$

T

$|T| \times |T|$ matrix
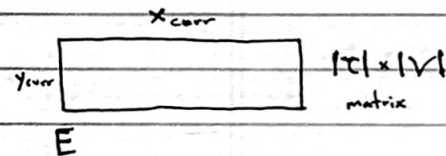
what is prob. of $y_{next}$ given $y_{curr}$?

- Two steps to use model
  1. Parameter Estimation
  2. Inference

$P(x_i | y_i)$
Emissions

$x_{curr}$

$y_{curr}$

E

$|T| \times |V|$ matrix

Given a tag (NN), what is prob. of seeing each word in vocab? (fixed vocab)

## HMM's : Parameter Estimation

- Labeled data : $\left(\bar{x}^{(i)}, \bar{y}^{(i)}\right)_{i=1}^{D}$

  unlike logistic regression
  which is $P(\gamma | x)$

- Maximize $\longrightarrow \sum_{i} \log P(\bar{y}^{(i)}, \bar{x}^{(i)})$  generative (joint) likelihood

$$= \sum_{i} \log P(y_1^{(i)}) + \sum_{i} \sum_{j} \log P(x_j^{(i)} | y_j^{(i)}) + \sum_{i} \sum_{j} \log P(y_j^{(i)} | y_{j-1}^{(i)})$$

log prob. $y_1$ for
each sequence in
training data

$i$ = sum over training data
$j$ = sum over sentence index
* Accumulate prob. of seeing
that $x$ given that $y$

Accumulate probability of
transition to $y_j$ given
$y_{j-1}$

- MLE with frequency counts

  Biased coin w/ probability $P$ of heads
  - observe:  HHHT
  - what is maximum likelihood probability $P$ for this coin? ( value of $P$ that maximizes data likelihood )
  - $3/4$ at first glance $\searrow$

  $\text{argmax}(3 \log P + \log(1-P)) = 3/4$ ✓ correct

  * HMM parameter estimation doesn't involve gradient descent.
    You can estimate parameters by counting and normalizing

Ex.

$$\Upsilon = \{N, V, STOP\} \qquad V = \{They, can, fish\}$$

Data:

| | N | V | STOP |
|---|---|---|---|
| | They | can | |

| | N | V | STOP |
|---|---|---|---|
| | They | fish | |

$S = $

| | |
|---|---|
| N | 2 |
| V | 0 |

$\downarrow$ normalize

| | |
|---|---|
| N | 1 |
| V | 0 |

$T = $

| | N | V | STOP |
|---|---|---|---|
| N | 0 | 2 | 0 |
| V | 0 | 0 | 2 |

$\downarrow$ normalize

| | N | V | STOP |
|---|---|---|---|
| N | 0 | 1 | 0 |
| V | 0 | 0 | 1 |

$E = $

| | They | can | fish |
|---|---|---|---|
| N | 2 | 0 | 0 |
| V | 0 | 1 | 1 |

$\longrightarrow$

| | T | c | F |
|---|---|---|---|
| N | 1 | 0 | 0 |
| V | 0 | 1/2 | 1/2 |

$P(\text{fish} | V)$

- **Smoothing**

    Add counts (fake data) to avoid 0's

$$T \longrightarrow \begin{array}{|c|c|c|} \hline 1 & 3 & 1 \\ \hline 1 & 1 & 3 \\ \hline \end{array} \xrightarrow{\text{normal.}} \begin{array}{ccc} 1/5 & 3/5 & 1/5 \\ 1/5 & 1/5 & 3/5 \end{array}$$

Ex)  $\underset{\text{They}}{N} \quad \underset{\text{can}}{V} \quad \underset{\text{fish}}{V} \Longrightarrow P(\bar{y}, \bar{x}) = \begin{array}{cccc} y_1 & y_2|y_1 & y_3|y_2 & s|y_3 \\ x_1|y_1 & x_2|y_2 & x_3|y_3 \end{array}$

$$\downarrow$$

$$\begin{array}{cccc} 1 & 3/5 & 1/5 & 3/5 \\ 1 & 1/2 & 1/2 \end{array}$$

$\boxed{\text{HMM's : Viterbi Algorithm}}$

- HMM's : model of $P(\bar{y}, \bar{x}) = P(y_1) P(x_1 | y_1) P(y_2 | y_1) \ldots$
- Inference : $\underset{\bar{y}}{\text{argmax }} P(\bar{y} | \bar{x}) \longrightarrow$ Given a sentence, what is most likely POS tag sequence that could've produced that sentence

$$\underset{\bar{y}}{\text{argmax }} P(\bar{y}|\bar{x}) = \underset{\bar{y}}{\text{argmax }} \frac{P(\bar{y}, \bar{x})}{P(\bar{x})} = \underset{\bar{y}}{\text{argmax }} \log P(\bar{y}, \bar{x})$$

$$\uparrow$$
$$\text{constant w.r.t. } \bar{y}$$

$$= \underset{\tilde{y}_1 \ldots \tilde{y}_n}{\text{argmax }} \log P(\tilde{y}_1) + \log P(x_1 | \tilde{y}_1) + \log(\tilde{y}_2 | \tilde{y}_1) + \ldots$$

- <u>Viterbi Dynamic Program</u>

Define $v_i(\tilde{y}) = n \times |\mathcal{T}|$  $\quad n = $ sentence length  $\quad |\mathcal{T}| = $ number of tags

score of best path ending in $\tilde{y}$ at time $i$

Base : $v_1(\tilde{y}) = \log P(x_1 | \tilde{y}) + \log P(\tilde{y})$

Recurrence : $v_i(\tilde{y}) = \log P(x_i | \tilde{y}) + \underset{\tilde{y}_{prev}}{\max} \log P(\tilde{y} | \tilde{y}_{prev}) + v_{i-1}(\tilde{y}_{prev})$

Viterbi    for $i = 1 \ldots n$ :
　　　　　for $\tilde{y}$ in $|\mathcal{T}|$ :
　　　　　　　compute $v_i(\tilde{y})$
　　　Compute $v_{n+1}(\text{STOP})$, this $= \underset{\bar{y}}{\max} \log P(\bar{x}, \bar{y})$
　　　Track "backpointers"

$\boxed{\text{Ex.}}$   $S = \begin{array}{c|c} N & -1 \\ \hline V & -1 \end{array}$

$T = \begin{array}{c|c|c|c} & N & V & STOP \\ \hline N & -2 & -1 & -1 \\ \hline V & -1 & -1 & -2 \end{array}$  (log probs)

$E = \begin{array}{c|c|c|c} & They & can & fish \\ \hline N & -1 & -3 & -1 \\ \hline V & -3 & -1 & -1 \end{array}$

$v_i(\tilde{y})$

| | they | can | can | fish | STOP |
|---|---|---|---|---|---|
| N | $-2$ | $-2-2-3 : \boxed{-7}$ | | | |
| | | $-4-1-3 : -8$ | | | |
| V | $-4$ | $-2-1-1 : \boxed{-4}$ | | | |

$-2 - 2 - 3 = -7$　　　　　　　　　　　　　　　　STOP
prev  tr  em