

# Parts of Speech Taggings (POS)

Ex. POS tag: NN (Noun, singular), NNP (Noun, plural), VB (Verb), ... . In this case: State is a tag

## I. Data Source

• Training data: WSJ\_02\_21.pos; hmm-vocab.txt

word	Tag
In	IN
an	DT
Oct.	NNP

word	Tag
-- unk --	
-- unk-adj --	
Bello	
Belo	

Note: in WSJ  
 ≈ 87% token is unambiguous  
 (one token has one tag)  
 ≈ 23% token is ambiguous  
 (one token has more than one tag)  
 word index of this word in a sorted list.  
 Change to Vocab = {  
 -- unk -- : 9,  
 -- unk-adj -- : 10,  
 'Bello' : 11, ... }  
 'Belo': 12, ... }

## II. POS tagging

### Training

- Transition count: Number of times (#) tag  $t_{i-1}$  happened next to tag  $t_i$ . Ex. DT  $\rightarrow$  NNP

- Emission count: # a word  $W_i$  given tag  $t_i$ . Ex. How many words given tag NNP

- Tag count: # each tag appears

def create\_dics (training\_corpus, vocab):

transition\_count = defaultdict(int) # access a key not exist, return value=0

Ex: transition\_count = { ('--1--', 'IN') : 5050,  
 ('DT', 'NNP') : 9049, ... }

prev\_tag  $\xrightarrow{\text{tag}}$

emission\_count = { ('DT', 'ang') : 722,

('NN', 'decrease') : 7, ... }

start state (a start of a sentence)

tag\_count = { 'IN' : 90000, '--1--' : 39800, ... }

## III. Hidden Markov Models (HMM) for POS

state is a POS which is hidden.

• Markov Model utilizes a transition matrix A

• HMM adds an Emission matrix B which describes the probability of a visible observation when we are in a particular state.

is a word in the corpus curr state

• Transition Probability Matrix A. Ex.

	NN	VB	0
NN	0.5	0.2	0.3
VB	0.4	0.2	0.4
0	0.1	0.1	0.8

$$P(t_i | t_{i-1}) = \frac{C(t_{i-1}, t_i) + \alpha}{\sum_{j=1}^N C(t_{i-1}, t_j) + \alpha N}$$

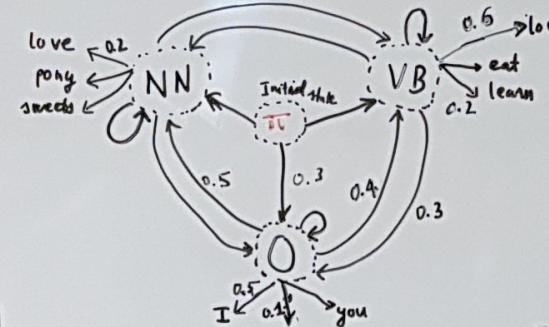
tag (state)  $t_{i-1} \rightarrow t_i$

( $t_{i-1}, t_i$ ): transition count from  $t_{i-1}, t_i$

N: number of tags (states)

$\alpha$  and  $\alpha \times N$  are used for smoothing

$$C(t_{i-1}) = \sum_{j=1}^N C(t_{i-1}, t_j)$$



. Emission Probability Matrix B<sub>N × M</sub>

$$P(W_i | t_i) = \frac{C(t_i, Word_i) + \alpha}{\sum_{j=1}^M C(t_i, word_j) + \alpha \cdot M}$$

C( $t_i, word_i$ ): count number of time of word<sub>i</sub> given by tag  $t_i$

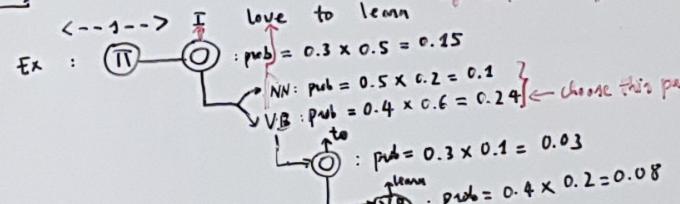
M: total number of words

$\alpha \cdot M$ : smoothing

Note: love can be in states NN or VB

## IV. Viterbi Algo

Ex. I love to learn → Find the sequence of hidden states (tags) that has the highest probability for this sequence.



choose the path:  
 $\langle \dots \rangle \xrightarrow{\text{I}} \langle \text{VB} \rangle \xrightarrow{\text{to}} \langle \text{VB} \rangle \xrightarrow{\text{learn}}$   
 with prob =  $0.15 \times 0.24 \times 0.03 \times 0.08$   
 highest

Algo: Need 2 matrix

	W <sub>1</sub>	W <sub>2</sub>	...	Word/M
t <sub>1</sub>	C <sub>1,1</sub>			
t <sub>2</sub>				
⋮				
t <sub>N</sub>	C <sub>N,1</sub>			

to store intermediate optimal probability

D =  
 NxM  
 to store indices of visited state

	W <sub>1</sub>	W <sub>2</sub>	...	Word/M
t <sub>1</sub>				
t <sub>2</sub>				
⋮				
t <sub>N</sub>				

(1) Initialization step (first word in a sentence)

$$C_{i,1} : \langle \dots \rangle \xrightarrow{\text{I}} \langle t_i \rangle \xrightarrow{\text{word}_1}$$

$$d_{i,1} = 0$$

Joint prob so we make log to if  $A('---1---', t_i) = 0$ :  $C_{i,1} = \text{float}('inf')$

become sum else  $C_{i,1} = \log A('---1---', t_i) + \log B[t_i, \text{first-word-at-corpus}]$

for i in range (1, M): # for second word

for j in range (N): best-d<sub>i,j</sub> = float('inf'); best-t<sub>i,j</sub> = None

for k in range (N):

$$\{ c = C_{k,i-1} + \log A[k,j] + \log B[j, \text{word}_i]$$

if best-d<sub>i,j</sub> < c: best-d<sub>i,j</sub> = c; best-t<sub>i,j</sub> = k

$$C_{j,i} = best-d_{i,j}; d_{j,i} = best-d_{i,j}$$

(2) Forward Step (populate C, D matrix, column by column  $\equiv$  word by word)

(3) Backward Path

	W <sub>1</sub>	W <sub>2</sub>	W <sub>3</sub>	W <sub>4</sub>	W <sub>5</sub>
t <sub>1</sub>	0.35	0.125	0.025	0.0125	0.00625
t <sub>2</sub>	0.1	0.25	0.05	0.02	0.008
t <sub>3</sub>	0.3	0.05	0.25	0.02	0.008
t <sub>4</sub>	0.2	0.1	0.008	0.0025	0.0008

$\langle \dots \rangle \xrightarrow{\text{t}_1} \langle \text{t}_2 \rangle \xrightarrow{\text{t}_3} \langle \text{t}_4 \rangle \xrightarrow{\text{t}_5}$

• Hop 1:  $j = \arg\max(C_{i,M})$

• Hop 2: look for last column at D with  $i=3 \Rightarrow$  it is  $t_2$ .

• Hop 3: at D move it back word.

i index: 1 to N states (tags)

j index: 1 to M words

t<sub>i</sub> states = ['---1---', 'IN', 'DT',

'NN', 'NNP', ... ] having 46

1 tags in this case.

i.e.  $t_1$  to  $t_5$  etc etc

$\pi_{t_1} \rightarrow \pi_{t_2} \rightarrow \pi_{t_3} \rightarrow \pi_{t_4} \rightarrow \pi_{t_5}$

$\pi_{t_1} = \text{argmax}(C_{1,M})$

$\pi_{t_2} = \text{argmax}(C_{2,M})$

$\pi_{t_3} = \text{argmax}(C_{3,M})$

$\pi_{t_4} = \text{argmax}(C_{4,M})$

$\pi_{t_5} = \text{argmax}(C_{5,M})$