

***CS4442/9542b: Artificial Intelligence II***  
***Prof. Olga Veksler***

***Lecture 10***  
***NLP: Part of Speech Tagging (POS)***

Many slides from: L. Kosseim, M. Hearst, K. McCoy,  
Yair Halevi

1

***Outline***

---

- What is POS and POS tagging
- Why we need POS tagging
- Different Approaches to POS
  1. rule-based tagging
  2. statistical tagging

2

## ***What is a Part of Speech ?***

---

- Words that somehow 'behave' alike:
  - Appear in similar contexts
  - Perform similar functions in sentences
  - Undergo similar transformations
- Terminology
  - POS (part-of-speech tag) are also called grammatical tag, grammatical category, syntactic class, word class

3

## ***Substitution Test***

---

- Two words belong to the same part of speech if replacing one with another does not change the grammaticality of a sentence.
  - The {sad, big, fat, green, ...} dog is barking.

4

## ***How many word classes are there?***

---

- A basic set:
  - N(oun), V(erb), Adj(ective), Adv(erb), Prep(osition), Det(erminer), Aux(iliaries), Part(icle), Conj(unction)
- A simple division: open/content vs. closed/function
  - Open: N, V, Adj, Adv
    - New members are added frequently
  - Closed: Prep, Det, Aux, Part, Conj, Num
    - New members are added rarely
- Many subclasses, e.g.
  - eats/V ⇒ eat/VB, eat/VBP, eats/VBZ, ate/VBD, eaten/VBN, eating/VBG, ...

5

## ***POS tagging***

---

- Goal: assign the right part of speech tag (noun, verb, ...) to words in a text

*"The /AT representative /NN put /VBD chairs /NNS on /IN the /AT table /NN."*
- What set of parts of speech do we use?
  - There are various standard tagsets to choose from; some have a lot more tags than others
  - The choice of tagset is based on the application
  - Accurate tagging can be done with even large tagsets

6

## ***What does Tagging do?***

---

1. Collapses distinctions
  - E.g., all personal pronouns tagged as **PRP**
  - Lexical identity may be completely discarded
2. Introduces distinctions (by reducing ambiguity)
  - E.g., “deal” tagged with **NN** or **VB**

7

## ***Why do POS Tagging?***

---

- **word sense disambiguation (semantics)**
  - Limits the range of meanings, “deal” as noun vs. “deal” as verb
- **speech recognition / synthesis (better accuracy)**
  - how to recognize/pronounce a word
  - *CON*tent/noun VS *con*TENT/adj
- **stemming**
  - which morphological affixes the word can take
  - adverb - *ly* = noun (*friendly* - *ly* = friend)
- **question answering**
  - analyzing a query to understand what type of entity the user is looking for and how it is related to other noun phrases mentioned in the question
- **partial parsing/chunking**
  - to find noun phrases/verb phrases
- **information extraction**
  - tagging and partial parsing help identify useful terms and relationships between them

8

## ***Tag sets***

---

- Different tag sets, depends on the purpose of the application
  - 45 tags in Penn Treebank
  - 62 tags in CLAWS with BNC corpus
  - 79 tags in Church (1991)
  - 87 tags in Brown corpus
  - 147 tags in C7 tagset
  - 258 tags in Tzoukermann and Radev (1995)

9

## ***Penn Treebank***

---

- First syntactically annotated corpus
- 1 million words from Wall Street Journal
- Part of speech tags and syntax trees

## Important Penn Treebank tags

### ■ 45 tags

IN	preposition or subordinating conjunct.
JJ	adjective or numeral, ordinal
JJR	adjective, comparative
NN	noun, common, singular or mass
NNP	noun, proper, singular
NNS	noun, common, plural
TO	"to" as preposition or infinitive marker
VB	verb, base form
VBD	verb, past tense
VBG	verb, present participle or gerund
VBN	verb, past participle
VBP	verb, present tense, not 3rd p. singular
VBZ	verb, present tense, 3rd p. singular
...	

Slide modified from Massimo  
Poesio's

11

## Verb inflection tags

VBP	base present	<i>take</i>
VB	infinitive	<i>take</i>
VBD	past	<i>took</i>
VBG	present participle	<i>taking</i>
VBN	past participle	<i>taken</i>
VBZ	present 3sg	<i>takes</i>
MD	modal	<i>can, would</i>

Slide modified from Massimo  
Poesio's

12

## The entire Penn Treebank tagset

Tag	Description	Example	Tag	Description	Example
CC	Coordin. Conjunction	<i>and, but, or</i>	SYM	Symbol	<i>+, %, &amp;</i>
CD	Cardinal number	<i>one, two, three</i>	TO	“to”	<i>to</i>
DT	Determiner	<i>a, the</i>	UH	Interjection	<i>ah, oops</i>
EX	Existential ‘there’	<i>there</i>	VB	Verb, base form	<i>eat</i>
FW	Foreign word	<i>mea culpa</i>	VBD	Verb, past tense	<i>ate</i>
IN	Preposition/sub-conj	<i>of, in, by</i>	VBG	Verb, gerund	<i>eating</i>
JJ	Adjective	<i>yellow</i>	VCN	Verb, past participle	<i>eaten</i>
JJR	Adj., comparative	<i>bigger</i>	VBP	Verb, non-3sg pres	<i>eat</i>
JJS	Adj., superlative	<i>wildest</i>	VBZ	Verb, 3sg pres	<i>eats</i>
LS	List item marker	<i>1, 2, One</i>	WDT	Wh-determiner	<i>which, that</i>
MD	Modal	<i>can, should</i>	WP	Wh-pronoun	<i>what, who</i>
NN	Noun, sing. or mass	<i>llama</i>	WP\$	Possessive wh-	<i>whose</i>
NNS	Noun, plural	<i>llamas</i>	WRB	Wh-adverb	<i>how, where</i>
NNP	Proper noun, singular	<i>IBM</i>	\$	Dollar sign	<i>\$</i>
NNPS	Proper noun, plural	<i>Carolinas</i>	#	Pound sign	<i>#</i>
PDT	Predeterminer	<i>all, both</i>	“	Left quote	<i>( ‘ or “</i>
POS	Possessive ending	<i>’s</i>	”	Right quote	<i>( ‘ or ”</i>
PP	Personal pronoun	<i>I, you, he</i>	(	Left parenthesis	<i>( [ , { , &lt;</i>
PP\$	Possessive pronoun	<i>your, one’s</i>	)	Right parenthesis	<i>( ], }, , &gt;</i>
RB	Adverb	<i>quickly, never</i>	,	Comma	<i>,</i>
RBR	Adverb, comparative	<i>faster</i>	.	Sentence-final punc	<i>( . ! ?</i>
RBS	Adverb, superlative	<i>fastest</i>	:	Mid-sentence punc	<i>( ; ... - -)</i>
RP	Particle	<i>up, off</i>			

Slide modified from Massimo Poesio's

13

## Terminology

- Suppose we have text “The cat decided to jump on the couch to play with another cat”
- Terminology
  - Word type
    - Distinct words in the text (vocabulary), the text above has 10 word types: “the cat decided to jump on couch play with another”
  - Word token
    - any word occurring in the text
    - The text above has 13 word tokens

14

## ***Distribution of Tags***

- Parts of speech follow the usual frequency-based distributional behavior
  - Most word types have only one part of speech
  - Of the rest, most have two
  - A small number of word types have lots of parts of speech
  - Unfortunately, the word types with lots of parts of speech occur with high frequency (and words that occur most frequently tend to have multiple tags)

15

## ***Most word types are not ambiguous but...***

- but most word types are rare...
- Brown corpus (Francis&Kucera, 1982):
  - 11.5% word types are ambiguous (>1 tag)
  - 40% word tokens are ambiguous (>1 tag)

	num. word types
<b>Unambiguous (1 tag)</b>	<b>35 340</b>
<b>Ambiguous (&gt;1 tag)</b>	<b>4 100</b>
2 tags	3760
3 tags	264
4 tags	61
5 tags	12
6 tags	2
7 tags	1 <i>"still"</i>

16



## ***Why is Tagging Hard?***

---

- Tagging is a type of disambiguation
- Examples:
  1. Book/**VB** that/**DT** flight/**NN**
    - “book” can also be **NN**
    - Can I read a **book** on this flight?
  2. Does/**VBZ** that/**DT** flight/**NN** serve/**VB** dinner/**NN** ?
    - “that” can also be a **complementizer**
    - My travel agent said **that** there would be a meal on this flight.

17

## ***Potential Sources of Disambiguation***

---

1. **Lexical information:**
  - look up all possible POS for a word in a dictionary
  - “table”: {noun, verb} but not a {adj, prep,...}
  - “rose”: {noun, adj, verb} but not {prep, ...}
2. **Syntagmatic information:**
  - some tag sequences are more probable than others:
  - **DET + N** occur frequently but **DET+V** never occurs
  - **ART+ADJ+N** is more probable than **ART+ADJ+VB**
  - We can find the syntagmatic information
    - by talking to the experts
    - or, better, from training corups

## Syntagmatic information from Corpus

- For a is a sequence of tags  $t_1, t_2, \dots, t_k$  compute  $P(t_1, t_2, \dots, t_k)$ , which will tell us how likely this tag sequence is
  - we have done something very similar before, i.e. we computed probability of a sequence of words
  - will make similar approximations as before,

$$P(t_n | t_1, t_2, \dots, t_{n-1}) = P(t_n | t_{n-k} \dots t_{n-1})$$

- in fact, for computational efficiency, the assumption will be

$$P(t_n | t_1, t_2, \dots, t_{n-1}) = P(t_n | t_{n-1})$$

## Techniques to POS tagging

1. rule-based tagging
  - uses hand-written rules
2. statistical tagging
  - uses probabilities computed from training corpus
    - Charniak
    - Markov Model based

## ***Rule-based POS Tagging***

---

- Step 1: Assign each word with all possible tags
  - use dictionary
- Step 2: Use if-then rules to identify the correct tag in context (disambiguation rules)

21

## ***Rule-based POS Tagging: Sample rules***

---

- **N-IP rule:**  
A tag N (noun) cannot be followed by a tag IP (interrogative pronoun)  
*... man who ...*
  - man: {N}
  - who: {RP, IP} --> {RP} relative pronoun
- **ART-V rule:**  
A tag ART (article) cannot be followed by a tag V (verb)  
*...the book...*
  - the: {ART}
  - book: {N, V} --> {N}

22

## ***Rule-based Tagger***

---

- using only syntagmatic patterns
  - Green & Rubin (1971)
  - accuracy of 77%
- In addition, it is very time consuming to come up with the rules and need an expert in English to come up with the rules

23

## ***Statistical POSTagger: Charniak 1993***

---

- Simplest statistical tagger
- Use corpus to calculate most probable tag for each *word*
  - that is the one maximizing
$$\text{count}(\textit{word} \text{ has tag } t) / \text{count}(\textit{word})$$
- Charniak tagging assigns POS tag to each word separately
- Given a *word* to tag,
  1. for each possible tag *t* for this *word*, compute
$$\text{count}(\textit{word} \text{ has tag } t)$$
  2. choose tag *t* that maximizes the above

24

## **Statistical POSTagger: Charniak 1993**

- Accuracy of 90%
  - contrast with 77% accuracy of the rule-based tagger!
  - Another evidence of the power of statistical methods over rule-based methods
  - MUCH better than rule based, but not very good...
    - 1 mistake every 10 words
  - funny fact: every word will have only one POS assigned to it (book will always be assigned the noun tag)
- This tagger is used mostly as baseline for evaluation
- How do we improve it?
  - tag of a word should depend on tags of other words around it, i.e. have to take “context” in the account
  - in other words, some sequence of tags are much more likely than others

25

## **Statistical Tagger: Markov Model Based**

- Suppose we want to tag sentence of words  $w_1 w_2 \dots w_n$
- Let  $t_1 t_2 \dots t_n$  be a possible sequence of tags corresponding to the sentence above
  - That is  $t_i$  is a tag for word  $w_i$
- Let  $w_{1..n}$  be shorthand for sentence  $w_1 w_2 \dots w_n$  and  $t_{1..n}$  be shorthand for its tagging  $t_1 t_2 \dots t_n$
- We want to find the “best” tagging  $t_{1..n}$  out of all possible taggings
- We have 2 sources of information in our corpus:
  1. given that the previous word tag is  $t_i$ , we can find how likely the tag of the next word is  $t_{i+1}$ , namely  $P(t_{i+1}|t_i)$
  2. we can find how likely is each word for each tag, namely  $P(w_i|t_i)$ 
    - tells us how likely part of speech  $t_i$  will “generate” a word  $w_i$
    - For example, if we know that that tag of the word is  $t_i = \text{noun}$ , what is the probability of the word to be “book”
    - $P(\text{book}|\text{verb}) > P(\text{book}|\text{noun})$ , because there are many more nouns than verbs

26

## Statistical Tagger: Markov Model Based

- Suppose we are given tag assignment  $t_{1:n}$  for sentence  $w_{1:n}$

- Using Bayes law:

$$P(t_{1:n} | w_{1:n}) = \frac{P(w_{1:n} | t_{1:n})P(t_{1:n})}{P(w_{1:n})}$$

- This says:

$$P\left( \begin{array}{|c|} \hline \text{word}_1 \\ \hline \text{has tag}_1 \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{word}_2 \\ \hline \text{has tag}_2 \\ \hline \end{array} \quad \dots \quad \begin{array}{|c|} \hline \text{word}_n \\ \hline \text{has tag}_n \\ \hline \end{array} \right) =$$

$$= \frac{P\left( \begin{array}{|c|} \hline \text{tag}_1 \text{ "gives"} \\ \hline \text{word}_1 \\ \hline \end{array} \quad \dots \quad \begin{array}{|c|} \hline \text{tag}_n \text{ "gives"} \\ \hline \text{word}_n \\ \hline \end{array} \right) P\left( \begin{array}{|c|} \hline \text{tag}_1 \\ \hline \end{array} \quad \dots \quad \begin{array}{|c|} \hline \text{tag}_n \\ \hline \end{array} \right)}{P\left( \begin{array}{|c|} \hline \text{word}_1 \\ \hline \end{array} \quad \dots \quad \begin{array}{|c|} \hline \text{word}_n \\ \hline \end{array} \right)}$$

27

## Statistical Tagger: Markov Model Based

$$P(t_{1:n} | w_{1:n}) = \frac{P(w_{1:n} | t_{1:n})P(t_{1:n})}{P(w_{1:n})}$$

- We will make two simplifying assumptions

  - Given a POS tag, probability of a word is independent of the tags of other words in a sentence:

$$P(w_{1:n} | t_{1:n}) = \prod_{i=1}^n P(w_i | t_i)$$

$$P\left( \begin{array}{|c|} \hline \text{tag}_1 \text{ "gives"} \\ \hline \text{word}_1 \\ \hline \end{array} \quad \dots \quad \begin{array}{|c|} \hline \text{tag}_n \text{ "gives"} \\ \hline \text{word}_n \\ \hline \end{array} \right) =$$

$$= P\left( \begin{array}{|c|} \hline \text{tag}_1 \text{ "gives"} \\ \hline \text{word}_1 \\ \hline \end{array} \right) \dots P\left( \begin{array}{|c|} \hline \text{tag}_n \text{ "gives"} \\ \hline \text{word}_n \\ \hline \end{array} \right)$$

- $P(w_i | t_k)$  can be approximated from the tagged corpus as:  $C(w_i \text{ has tag } t_k) / C(t_k)$ , that is how many times word  $w_i$  has tag  $t_k$  divided by how often tag  $t_k$  occurs in the corpus

## Statistical Tagger: Markov Model Based

$$P(t_{1,n} | w_{1,n}) = \frac{P(w_{1,n} | t_{1,n})P(t_{1,n})}{P(w_{1,n})}$$

2. Each tag is only dependent only on one previous tag:

$$P(t_{1,n}) = \prod_{i=1}^n P(t_i | t_{i-1})$$

$$\begin{aligned} P(\text{tag}_1 \dots \text{tag}_n) &= \\ &= P(\text{tag}_1 | \text{tag}_0) \dots P(\text{tag}_n | \text{tag}_{(n-1)}) \end{aligned}$$

- this is **Markov** assumption that we have seen before in language modeling
- Recall that  $P(\text{tag}_1 | \text{tag}_2)$  can be approximated by  $C(\text{tag}_2 \text{tag}_1) / C(\text{tag}_2)$
- Here  $P(\text{tag}_1 | \text{tag}_0)$  stands for  $P(\text{tag}_1)$  and is approximated by  $C(\text{tag}_1) / (\text{number of tokens in the training corpus})$

## Statistical Tagger: Markov Model Based

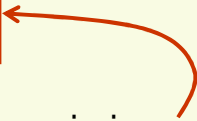
- Using these 2 simplifications, we get

$$P(t_{1,n} | w_{1,n}) = \frac{\prod_{i=1}^n P(w_i | t_i) P(t_i | t_{i-1})}{P(w_{1,n})}$$

- Given a possible sequence of tags  $t_{1,n}$  for sentence  $w_{1,n}$ , we can evaluate how good this tag/sequence assignment is using  $P(t_{1,n} | w_{1,n})$
- Algorithm: go over all possible tag assignments and choose the tag assignment which gives highest  $P(t_{1,n} | w_{1,n})$ 
  - Notice that  $P(w_{1,n})$  **does not** effect maximization so we do not have to compute it

## Statistical Tagger: Markov Model Based

- Algorithm: given sentence  $w_{1,n}$  go over all possible tag assignments  $t_{1,n}$  and compute

$$\prod_{i=1}^n P(w_i | t_i) P(t_i | t_{i-1})$$


- Choose final tagging  $t_{1,n}$  which maximizes
- Efficiency
  - For each word  $w_i$ , try only the tags given by the dictionary (lexical information)
  - Ex: for “fly”, possible tags are NOUN, VERB and also ADJECTIVE (meaning “keen” or “artful”, mainly in England)

31

## Statistical Tagger: Markov Model Based

- Side note: Markov tagger becomes Charniak’s tagger if tags are assumed independent, that is if  $P(t_i | t_{i-1}) = P(t_i)$ , then

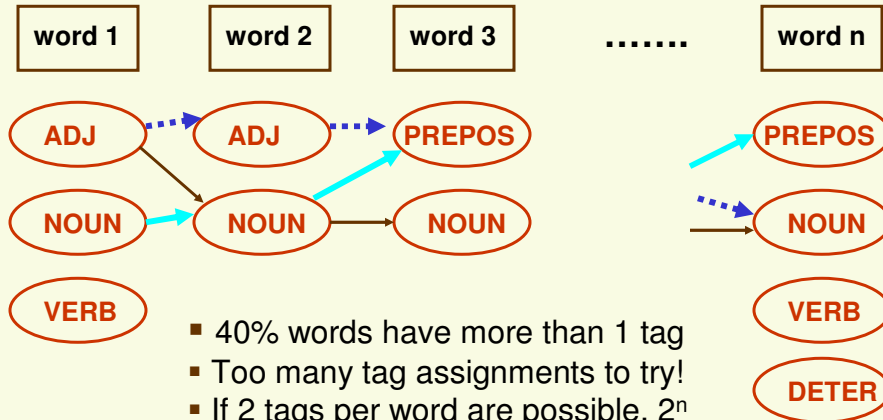
$$\begin{aligned} \prod_{i=1}^n P(w_i | t_i) P(t_i | t_{i-1}) &= \prod_{i=1}^n P(w_i | t_i) P(t_i) = \prod_{i=1}^n \frac{P(w_i, t_i)}{P(t_i)} P(t_i) \\ &= \prod_{i=1}^n P(w_i, t_i) \end{aligned}$$

32



## Statistical Tagger: Markov Model Based

- Algorithm: given sentence  $w_{1,n}$  go over all possible tag assignments  $t_{1,n}$  and compute  $\prod_{i=1}^n P(w_i | t_i) P(t_i | t_{i-1})$



- 40% words have more than 1 tag
- Too many tag assignments to try!
- If 2 tags per word are possible,  $2^n$  tag assignments are possible
  - exhaustive search is exponential

33

## MM based Tagger: Viterbi Algorithm

- Fortunately there is a very useful algorithm (Viterbi)
- If there are  $k$  tags per word and  $n$  words, can find best tagging in time  $k^2n$ 
  - There are  $k^n$  different tag assignments possible
- First, to avoid floating point underflows, take logarithm of

$$\prod_{i=1}^n P(w_i | t_i) P(t_i | t_{i-1})$$

$$\log \left[ \prod_{i=1}^n P(w_i | t_i) P(t_i | t_{i-1}) \right] = \sum_{i=1}^n (\log P(w_i | t_i) + \log P(t_i | t_{i-1}))$$

- Now we want to maximize:

$$\sum_{i=1}^n [\log P(w_i | t_i) + \log P(t_i | t_{i-1})] = \sum_{i=1}^n \underbrace{\log P(w_i | t_i)}_{\text{how likely word } w_i \text{ is for tag } t_i} + \sum_{i=1}^n \underbrace{\log P(t_i | t_{i-1})}_{\text{how likely tag } t_i \text{ to follow tag } t_{i-1}}$$

## ***MM based Tagger: Viterbi Algorithm***

---

- Maximizing:

$$\sum_{i=1}^n \log P(w_i | t_i) + \sum_{i=1}^n \log P(t_i | t_{i-1})$$

- Is equivalent to minimizing

$$-\sum_{i=1}^n \log P(w_i | t_i) - \sum_{i=1}^n \log P(t_i | t_{i-1})$$

- It is more convenient to minimize the expression above

## ***MM based Tagger: Viterbi Algorithm***

---

- So we need to find a sequence of tags  $t_1, t_2, \dots, t_n$  to minimize:

$$-\sum_{i=1}^n \log P(w_i | t_i) - \sum_{i=1}^n \log P(t_i | t_{i-1})$$

- To simplify notation, will write

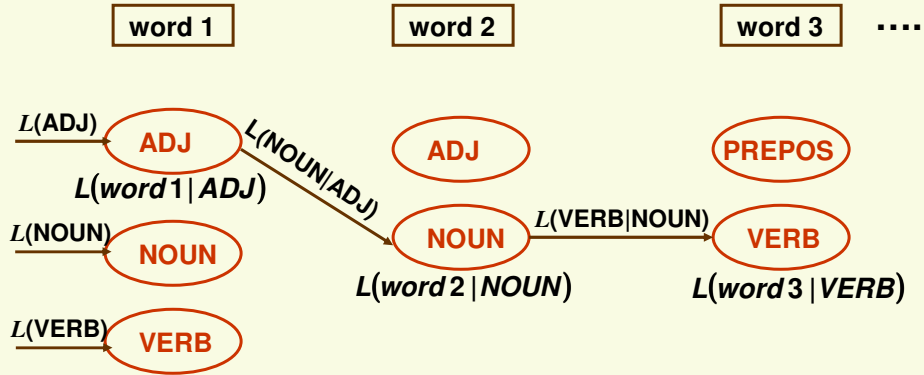
- $L(w|t)$  instead of  $-\log[P(w|t)]$
- $L(t_i|t_{i-1})$  instead of  $-\log[P(t_i|t_{i-1})]$

- In new notation, we need to find a sequence of tags  $t_1, t_2, \dots, t_n$  to minimize:

$$\sum_{i=1}^n [L(w_i | t_i) + L(t_i | t_{i-1})]$$

## MM based Tagger: Viterbi Algorithm

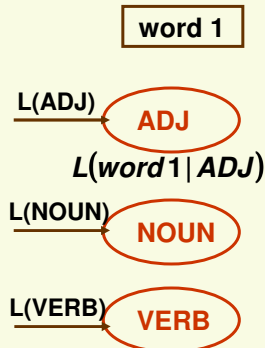
$$\sum_{i=1}^n [L(w_i | t_i) + L(t_i | t_{i-1})]$$



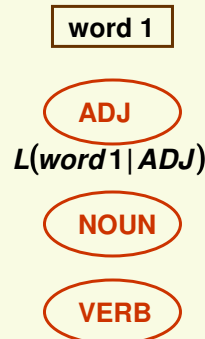
## MM based Tagger: Viterbi Algorithm

- Let's change notation slightly for the first word
  - $L(\text{word1}|\text{tag})$  will stand for  $-\log[P(\text{word1}|\text{tag})] - \log [P(\text{tag})]$
  - This will simplify the pictures:

*instead of this:*



*we will picture this:*

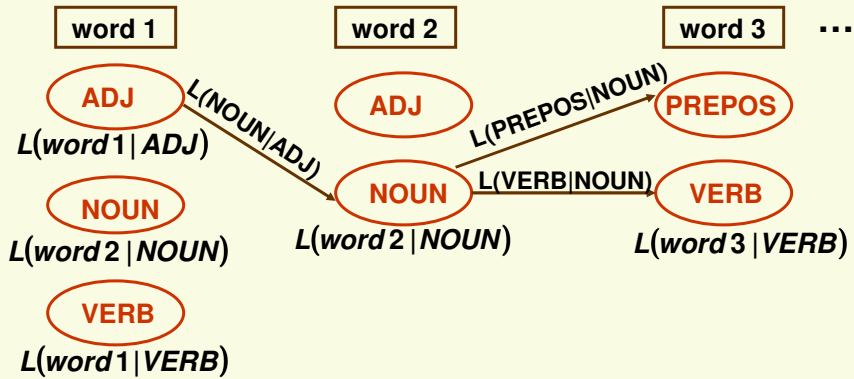


$$\sum_{i=1}^n \log P(w_i | t_i) + \sum_{i=2}^n \log P(t_i | t_{i-1})$$

## MM based Tagger: Viterbi Algorithm

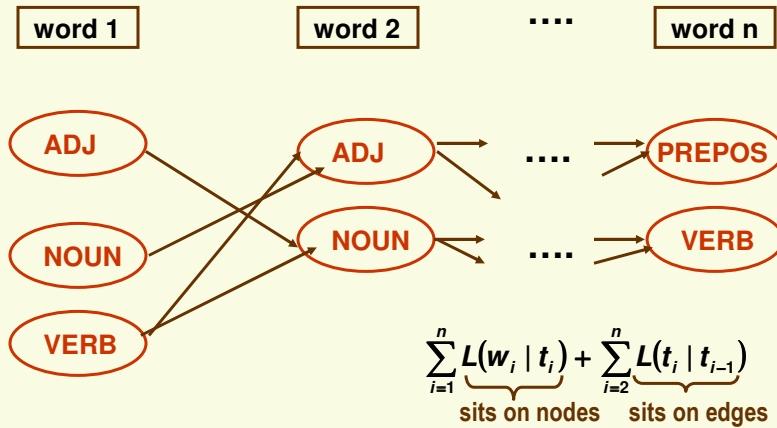
- Each node has cost  $L(\text{word}|\text{tag})$
- Each link between nodes has cost  $L(\text{tag 1} | \text{tag 2})$
- Cost of path, summing up node costs and edge costs is:

$$\sum_{i=1}^n L(w_i | t_i) + \sum_{i=2}^n L(t_i | t_{i-1})$$



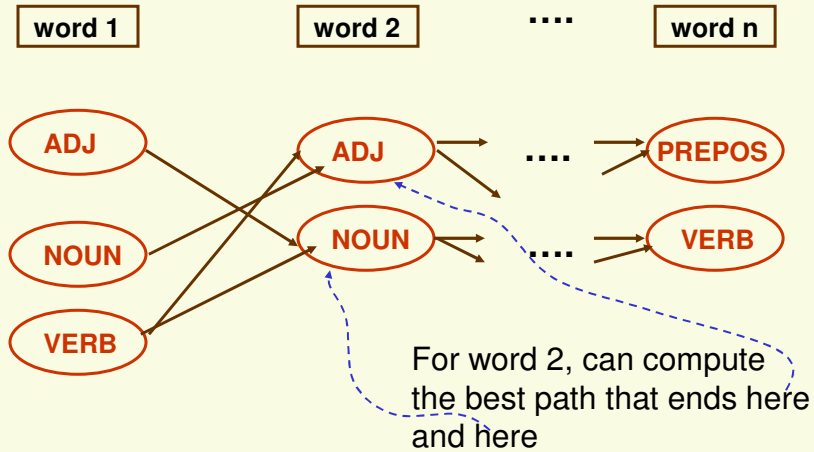
## MM based Tagger: Viterbi Algorithm

- So we need to find the path with smallest cost that starts at some node corresponding to word 1 and ends at some node corresponding to word n



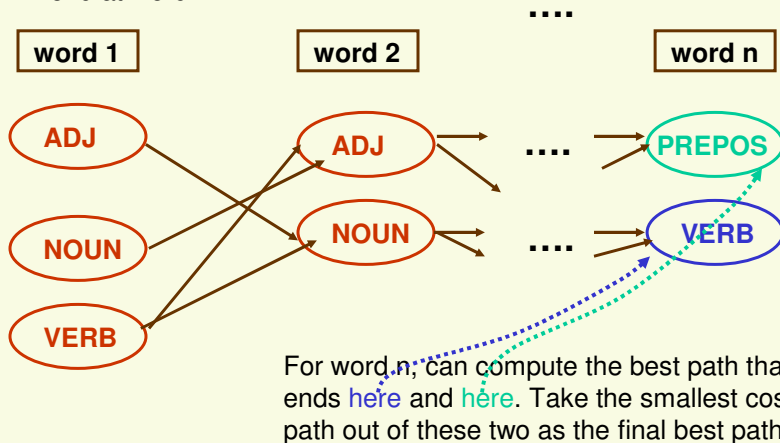
## MM based Tagger: Viterbi Algorithm

- Idea: for every node corresponding to word  $i$ , we can efficiently find the best (smallest cost) path that ends at it (and starts at any node corresponding to word 1)



## MM based Tagger: Viterbi Algorithm

- First compute the best path that ends at any node for word 1
- Then compute the best path that ends at any node for word 2
- ....
- Finally compute the best path that ends at any node for word  $n$ 
  - The best path overall is the smallest cost path out of those paths that end at word  $n$



## MM based Tagger: Viterbi Algorithm

- First compute the best path that ends at any node for word 1
- Trivial, since the path has just 1 node
- $C(w_1, \text{tag}) = L(w_1 | \text{tag})$ , holds the cost of the best path ending at  $(w_1, \text{tag})$
- $P(w_1, \text{tag}) = \text{null}$ , holds the parent node on the best path ending at  $(w_1, \text{tag})$

word 1

ADJ

$L(\text{word1} | \text{ADJ})$

NOUN

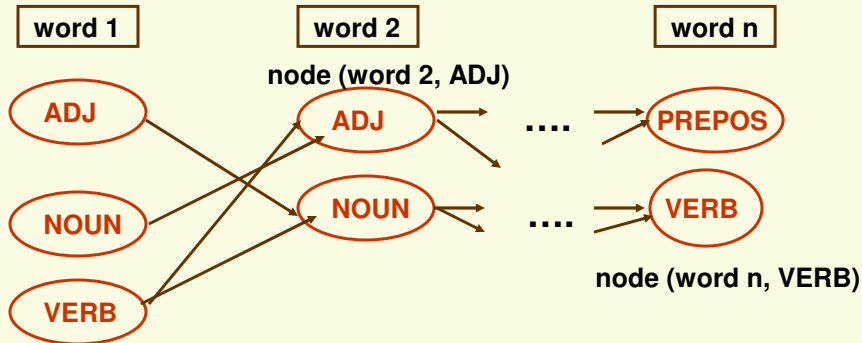
$L(\text{word1} | \text{NOUN})$

VERB

$L(\text{word1} | \text{VERB})$

## MM based Tagger: Viterbi Algorithm

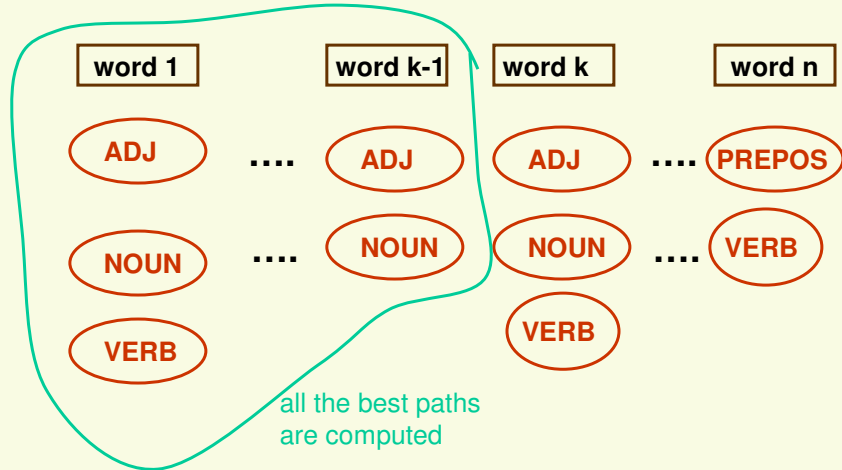
- In general, any node is specified by the word and the tag (word  $i$ , tag)



- Let  $C(\text{word } k, \text{tag})$  stand for the cost of the best path starting at any node for word 1 and ending at node (word  $k$ , tag)
- Let  $P(\text{word } k, \text{tag})$  stand for the parent on the best cost path starting at any node for word 1 and ending at node (word  $k$ , tag). Note that the parent must be the node (word  $k-1$ , tag')
- After all  $C(w, t)$  are computed, the best cost path overall is given by the minimum over all  $t$  of  $C(\text{word } n, t)$

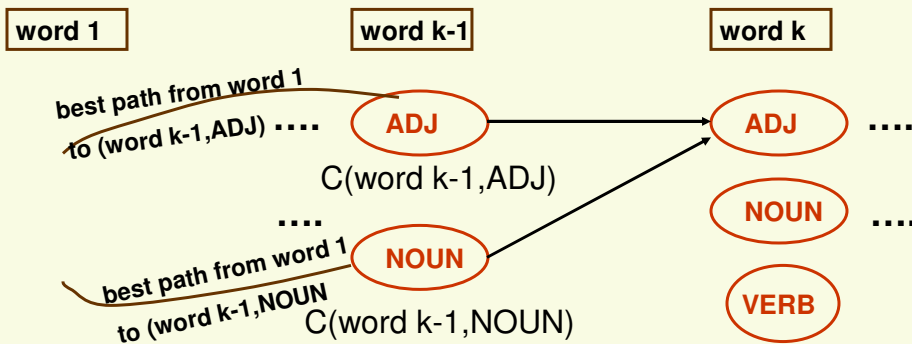
## MM based Tagger: Viterbi Algorithm

- We saw that for all possible values of tag, computing  $C(\text{word } 1, \text{tag})$  and  $P(\text{word } 1, \text{tag})$  is trivial
- Suppose we have computed  $C(\text{word } i, \text{tag})$  and  $P(\text{word } i, \text{tag})$  for all tags and all  $i < k$



## MM based Tagger: Viterbi Algorithm

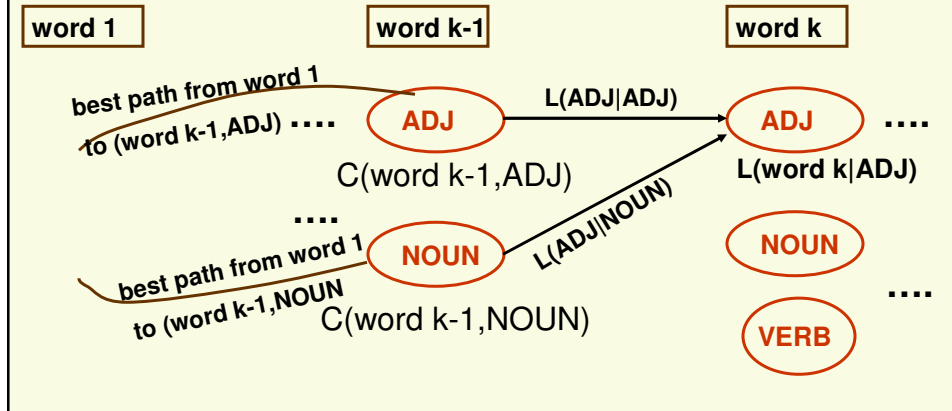
- Suppose we have computed  $C(\text{word } i, \text{tag})$  for all tags and all  $i < k$
- Need to compute  $C(\text{word } k, \text{tag})$  and  $P(\text{word } k, \text{tag})$  for all possible tags of word k



- Consider node (word k, ADJ). Let P be the best path from the word 1 to node (word k, ADJ). Path P will go either through either
  - node (word k-1, ADJ). In this case P follows the best path from word 1 to node (word k-1, ADJ)
  - or through node (word k-1, NOUN). In this case P follows the best path from word 1 to node (word k-1, NOUN)
  - we using property that a subpath of the best path is a best path itself

## MM based Tagger: Viterbi Algorithm

- Therefore  $C(\text{word } k, \text{ADJ})$  is the smaller of 2 quantities
  - $C(\text{word } k-1, \text{ADJ}) + L(\text{ADJ}|\text{ADJ}) + L(\text{word } k-1|\text{ADJ})$ 
    - In this case,  $P(\text{word } k, \text{ADJ}) = (\text{word } k-1, \text{ADJ})$
  - $C(\text{word } k-1, \text{NOUN}) + L(\text{ADJ}|\text{NOUN}) + L(\text{word } k-1|\text{ADJ})$ 
    - In this case,  $P(\text{word } k, \text{ADJ}) = (\text{word } k-1, \text{NOUN})$



## MM based Tagger: Viterbi Algorithm

- In general,  $C(\text{word } k, \text{tag})$  is computed as follows:

$C(\text{word } k, \text{tag}) =$

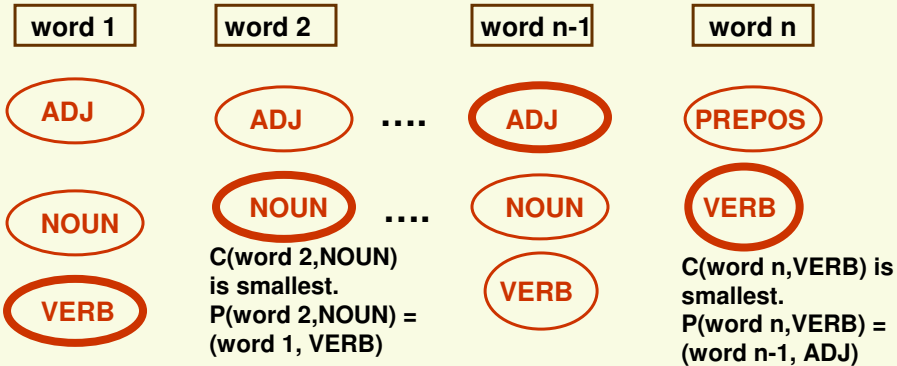
$$= \min_{t \in T(\text{word } k-1)} \{ \underbrace{C(\text{word } k-1, t)}_{\text{cost of best path from first word to node (word } k-1, t)} + \underbrace{L(\text{tag} | t)}_{\text{cost of going through node (word } k, \text{tag)}} \} + \underbrace{L(\text{word } k | \text{tag})}_{\text{cost of going between nodes (word } k-1, t) \text{ and (word } k, \text{tag)}}$$

- $P(\text{word } k, \text{tag}) = (\text{word } k-1, t^*)$  where  $t^*$  is the tag for word k-1 minimizing the expression above



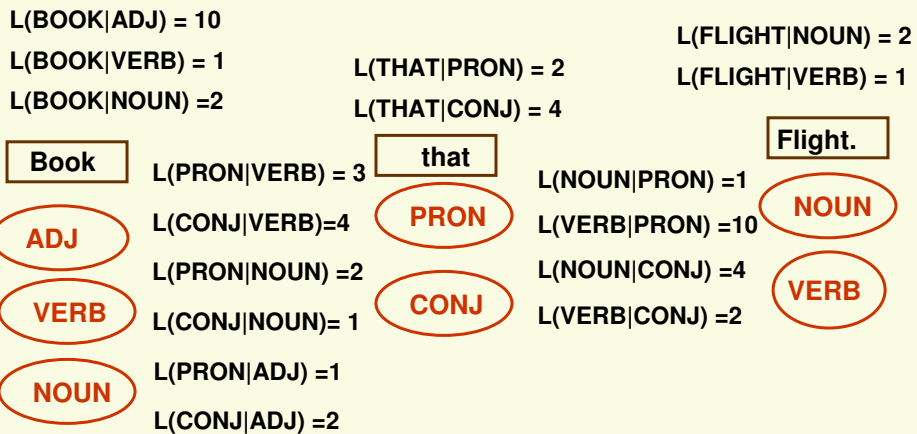
## MM based Tagger: Viterbi Algorithm

- After we computed  $C(\text{word } i, t)$  for all  $i$  and  $t$ , the best cost path is found as the maximum of  $C(\text{word } n, t)$  over all tags  $t$  that word  $n$  can have
- The parents on the path can be traced back using the computed  $P(\text{word } i, t)$



- Final tagging is: VERB NOUN ... ADJ VERB

## MM based Tagger: Small Example



- Notice that I made these log probabilities up

## MM based Tagger: Small Example

$L(\text{BOOK}|\text{ADJ}) = 10$

$L(\text{BOOK}|\text{VERB}) = 1$

$L(\text{BOOK}|\text{NOUN}) = 2$

Book

ADJ

VERB

NOUN

- Iteration 1:
  - $C(\text{book}, \text{adj}) = 10$ ,  $P(\text{book}, \text{adj}) = \text{null}$
  - $C(\text{book}, \text{verb}) = 1$ ,  $P(\text{book}, \text{verb}) = \text{null}$
  - $C(\text{book}, \text{noun}) = 2$ ,  $P(\text{book}, \text{noun}) = \text{null}$

## MM based Tagger: Small Example

$L(\text{PRON}|\text{ADJ}) = 1$

$L(\text{PRON}|\text{VERB}) = 3$

$L(\text{PRON}|\text{NOUN}) = 2$

$L(\text{THAT}|\text{PRON}) = 2$

$L(\text{THAT}|\text{CONJ}) = 4$

Book

that

ADJ

PRON

VERB

CONJ

NOUN

$C(\text{book}, \text{adj}) + L(\text{pron}|\text{adj}) + L(\text{that}|\text{pron}) = 13$

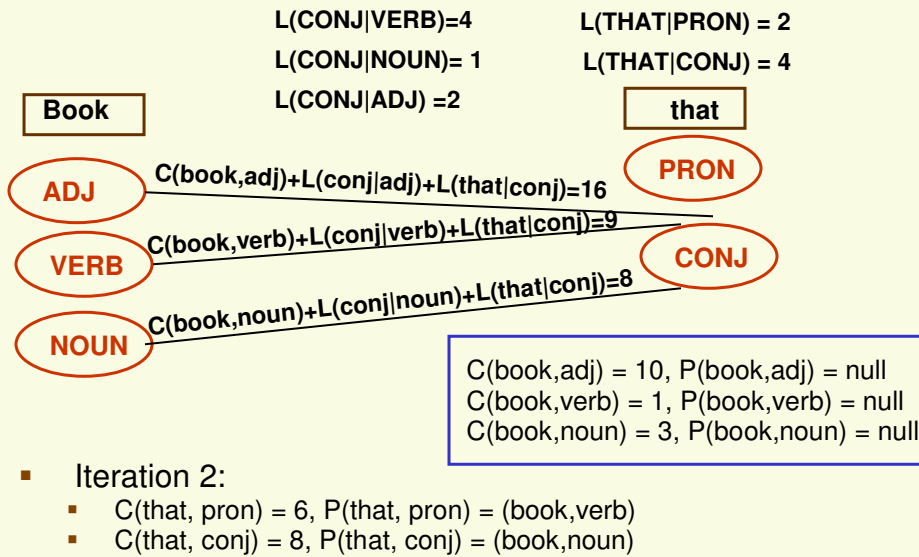
$C(\text{book}, \text{verb}) + L(\text{pron}|\text{verb}) + L(\text{that}|\text{pron}) = 6$

$C(\text{book}, \text{noun}) + L(\text{pron}|\text{noun}) + L(\text{that}|\text{pron}) = 7$

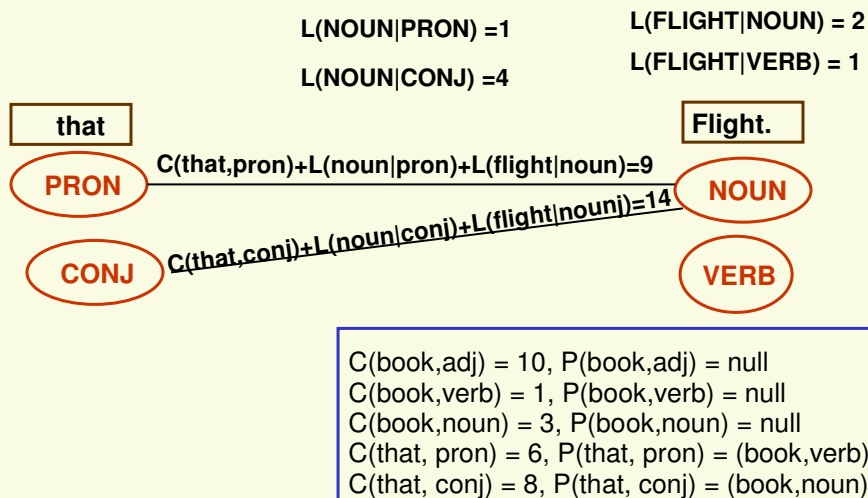
$C(\text{book}, \text{adj}) = 10$ ,  $P(\text{book}, \text{adj}) = \text{null}$   
 $C(\text{book}, \text{verb}) = 1$ ,  $P(\text{book}, \text{verb}) = \text{null}$   
 $C(\text{book}, \text{noun}) = 3$ ,  $P(\text{book}, \text{noun}) = \text{null}$

- Iteration 2:
  - $C(\text{that}, \text{pron}) = 6$ ,  $P(\text{that}, \text{pron}) = (\text{book}, \text{verb})$

## MM based Tagger: Small Example



## MM based Tagger: Small Example



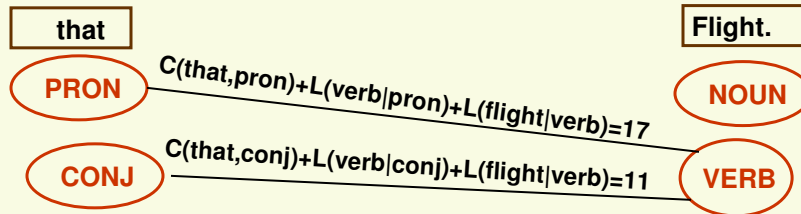
## MM based Tagger: Small Example

$L(\text{VERB}|\text{PRON}) = 10$

$L(\text{FLIGHT}|\text{NOUN}) = 2$

$L(\text{VERB}|\text{CONJ}) = 2$

$L(\text{FLIGHT}|\text{VERB}) = 1$

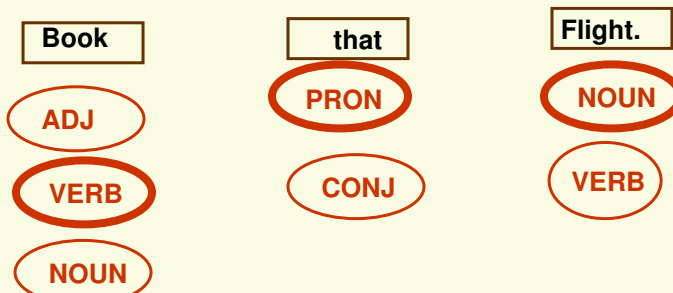


$C(\text{book,adj}) = 10$ ,  $P(\text{book,adj}) = \text{null}$   
 $C(\text{book,verb}) = 1$ ,  $P(\text{book,verb}) = \text{null}$   
 $C(\text{book,noun}) = 3$ ,  $P(\text{book,noun}) = \text{null}$   
 $C(\text{that,pron}) = 6$ ,  $P(\text{that,pron}) = (\text{book,verb})$   
 $C(\text{that,conj}) = 8$ ,  $P(\text{that,conj}) = (\text{book,noun})$

### Iteration 3:

- $C(\text{flight,noun}) = 9$ ,  $P(\text{flight,noun}) = (\text{that,pron})$
- $C(\text{flight,verb}) = 11$ ,  $P(\text{flight,verb}) = (\text{that,conj})$

## MM based Tagger: Small Example



$C(\text{book,adj}) = 10$ ,  $P(\text{book,adj}) = \text{null}$   
 $C(\text{book,verb}) = 1$ ,  $P(\text{book,verb}) = \text{null}$   
 $C(\text{book,noun}) = 3$ ,  $P(\text{book,noun}) = \text{null}$   
 $C(\text{that,pron}) = 6$ ,  $P(\text{that,pron}) = (\text{book,verb})$   
 $C(\text{that,conj}) = 8$ ,  $P(\text{that,conj}) = (\text{book,noun})$   
 **$C(\text{flight,noun}) = 9$ ,  $P(\text{flight,noun}) = (\text{that,pron})$**   
 $C(\text{flight,verb}) = 11$ ,  $P(\text{flight,verb}) = (\text{that,pron})$

**Final Tagging is: Book<verb> that <pron> flight<noun>**

## Viterbi Algorithm

```
for each  $t \in \text{Tags}(\text{word } 1)$  do  
     $C(\text{word } 1, t) = L(\text{word } 1 | t)$ ,  $P(\text{word } 1, t) = \text{null}$   
for  $i \leftarrow 2$  to  $n$  do  
    for each  $t \in \text{Tag}(\text{word } i)$  do  
         $C(\text{word } i, t) = -\infty$   
        for each  $t' \in \text{Tag}(\text{word } i - 1)$  do  
             $\text{nextCost} = C(\text{word } i - 1, t') + L(t|t') + L(\text{word } i|t)$   
            if  $\text{nextCost} < \text{cost}(\text{word } i, t)$  do  
                 $C(\text{word } i, t) = \text{nextCost}$   
                 $P(\text{word } i, t) = t'$ 
```

Note:  $\text{Tags}(\text{word } i)$  is the set of all possible tags for word  $i$

## Unknown Words

- Simplest method: assume an unknown word could belong to any tag; unknown words are assigned the distribution over POS over the whole lexicon
  - $P(\text{verb} | \text{"karumbula"}) = P(\text{noun} | \text{"karumbula"}) = P(\text{adjective} | \text{"karumbula"}) = \dots$  etc
- Some tags are more common than others
  - for example a new word can be most likely a verb, a noun etc. but not a preposition or an article
- Use features of the word (morphological and other cues, for example words ending in *-ed* are likely to be past tense forms or past participles)

## ***Tagging Accuracy***

---

- Ranges from 96%-97%
- Depends on:
  - Amount of training data available
  - The tag set
  - Difference between training corpus and dictionary and the corpus of application
  - Unknown words in the corpus of application
- A change in any of these factors can have a dramatic effect on tagging accuracy – often much more stronger than the choice of tagging method