

CS4442/9542b
Artificial Intelligence II
prof. Olga Veksler

Lecture 16

Natural Language Processing

Part of Speech Tagging

Many slides from: Joshua Goodman, L. Kosseim, D. Klein, D. Jurafsky, M. Hearst, K. McCoy, Y. Halevi, C. Manning, M. Poesio

Outline

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

What is a Part of Speech ?

- Words that behave alike
 - appear in similar contexts
 - perform similar functions in sentences
 - undergo similar transformations
- Terminology
 - **POS** (part-of-speech tag)
 - also called
 - grammatical tag
 - grammatical category
 - syntactic word class

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, green, ...} dog is barking.

Origin

- Perhaps started with Aristotle (384–322 BCE)
- From Dionysius Thrax of Alexandria (c. 100 BCE) the idea that is still with us
 - 8 main parts of speech
- Those 8 are not exactly the ones taught today
 - **Thrax:** noun, verb, article, adverb, preposition, conjunction, participle, pronoun
 - **School grammar:** noun, verb, adjective, adverb, preposition, conjunction, pronoun, interjection

How Many POS are there?

- A basic set:
 - N(oun), V(erb), Adj(ective), Adv(erb), Prep(osition), Det(erminer), Aux(ilaries), 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 \Rightarrow eat/VB, eat/VBP, eats/VBZ, ate/VBD, eaten/VBN, eating/VBG, ...

POS tagging

- Goal: assign POS tag (noun, verb, ...) to text

The/AT girl/NN put/VBD chairs/NNS on/IN the/AT table/NN.

- What set of parts of speech do we use?
 - various standard tagsets to choose from, some have a lot more tags than others
 - choice of tagset is based on application
 - accurate tagging possible with even large tagsets

Why do POS Tagging?

- Word sense disambiguation (semantics)
 - limits the range of meanings: *deal* as noun vs. *deal* as verb
- Speech recognition and synthesis
 - how to recognize/pronounce a word:
 - *content/noun* vs. *content/adj*
- Stemming: which morphological affixes word can take
 - adverb - *ly* = noun: *friendly* - *ly* = friend
 - cannot apply to adjectives, example: *sly*
- Partial parsing/chunking
 - to find noun phrases/verb phrases
- Information extraction
 - helps identify useful terms and relationships between them

Common Tagged Datasets

- 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)

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 total

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
VCN	verb, past participle
VBP	verb, present tense, not 3rd p. singular
VBZ	verb, present tense, 3rd p. singular

...

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>

The entire Penn Treebank tagset

Tag	Description	Example	Tag	Description	Example
CC	Coordin. Conjunction	<i>and, but, or</i>	SYM	Symbol	<i>+, %, &</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>	VBN	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>([, ({ , <)</i>
PP\$	Possessive pronoun	<i>your, one's</i>)	Right parenthesis	<i>(],), }, >)</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>			

Terminology

- Given text

The cat decided to jump on the couch to play with another cat

- Terminology

- **Word type**

- Distinct words in the text (vocabulary)

- text above has 10 word types

- *the, cat, decided, to, jump, on, couch, play, with, another*

- **Word token**

- any word occurring in the text

- text above has 13 word tokens

Distribution of Tags

- POS follow typical frequency-based behavior
 - most word types have only one part of speech
 - of the rest, most have two
 - only a small number of word types have lots of parts of speech
 - but these occur with high frequency

Most Word Types not Ambiguous but

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

- 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)

Tagging is a Type of Disambiguation

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.

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**
- 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 $\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_k$ compute

$$\mathbf{P}(\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_k)$$

- tells us how likely this tag sequence is
- similar to computing probability of a sequence of words $\mathbf{P}(\mathbf{w})$
- make the same approximation as before

$$\mathbf{P}(\mathbf{t}_n | \mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_{n-1}) = \mathbf{P}(\mathbf{t}_n | \mathbf{t}_{n-k} \dots \mathbf{t}_{n-1})$$

- for computational efficiency, our assumption is

$$\mathbf{P}(\mathbf{t}_n | \mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_{n-1}) = \mathbf{P}(\mathbf{t}_n | \mathbf{t}_{n-1})$$

POS Tagging Techniques

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)

Rule-based POS Tagging: Sample rules

- **ART-V rule:**

tag ART (article) cannot be followed by a tag V (verb)

...the book...

- the: {ART}
- book: {N, V} --> {N}

- **N-IP rule:**

tag N (noun) cannot be followed by tag IP (interrogative pronoun)

... man who ...

- man: {N}
- who: {RP, IP} --> {RP} relative pronoun

Rule-based Tagger

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

Statistical POS Tagger: Charniak 1993

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

Statistical POS Tagger: Charniak 1993

- Accuracy of 90%
 - contrast with 77% accuracy of the rule-based tagger!
 - evidence of power of statistical 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?
 - take the context of the surrounding words into account
 - some sequence of tags are much more likely than others

Statistical Tagger: Markov Model Based

- Tag sentence of words, t_i is a tag for word w_i

$$\mathbf{w}_{1,n} = w_1 w_2 \dots w_n$$

$$\mathbf{t}_{1,n} = t_1 t_2 \dots t_n$$

- Find best tagging $\mathbf{t}_{1,n}$
- Define best tagging?
- Use statistical principle, maximize

$$P(\mathbf{t}_{1,n} \mid \mathbf{w}_{1,n})$$

Markov Model Tagger

- The best tagging maximizes

$$P(\mathbf{t}_{1,n} \mid \mathbf{w}_{1,n})$$

- Hard to estimate directly
- Bayes law

$$P(\mathbf{t}_{1,n} \mid \mathbf{w}_{1,n}) = \frac{P(\mathbf{w}_{1,n} \mid \mathbf{t}_{1,n})P(\mathbf{t}_{1,n})}{P(\mathbf{w}_{1,n})}$$

- Bottom does not effect maximization,
 - constant over all possible taggings $\mathbf{t}_{1,n}$
- Find tagging that maximizes

$$P(\mathbf{w}_{1,n} \mid \mathbf{t}_{1,n})P(\mathbf{t}_{1,n})$$

Markov Model Tagger: First Assumption

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

- $P(w_1 = \text{get}, w_2 = \text{the}, w_3 = \text{pen} | t_1 = \text{VERB}, t_2 = \text{DET}, t_3 = \text{NOUN}) =$
 $P(\text{get}, \text{the}, \text{pen} | \text{VERB}, \text{DET}, \text{NOUN})$

- Knowing its own tag is **VERB** is helpful

- $P(\text{get} | \text{VERB}), P(\text{ask} | \text{VERB}), P(\text{place} | \text{VERB})$ should be high
- $P(\text{an} | \text{VERB}), P(\text{that} | \text{VERB}), P(\text{one} | \text{VERB})$, should be very low

- Knowing other tags are not very helpful

- $P(\text{get} | \text{VERB}, \text{DET}, \text{NOUN}), P(\text{ask} | \text{VERB}, \text{DET}, \text{NOUN})$ should still be high
- $P(\text{an} | \text{VERB}, \text{DET}, \text{NOUN}), P(\text{that} | \text{VERB}, \text{DET}, \text{NOUN})$ should still be low

- **Given its own tag, probability of word is independent of other tags**

$$P(\text{get}, \text{the}, \text{pen} | \text{VERB}, \text{DET}, \text{NOUN}) = P(\text{get} | \text{VERB}) P(\text{the} | \text{DET}) P(\text{pen} | \text{NOUN})$$

Markov Model Tagger: First Assumption

$$P(\mathbf{w}_{1,n} | \mathbf{t}_{1,n}) = \prod_{i=1}^n P(\mathbf{w}_i | \mathbf{t}_i) = P(\mathbf{w}_1 | \mathbf{t}_1) P(\mathbf{w}_2 | \mathbf{t}_2) \dots P(\mathbf{w}_n | \mathbf{t}_n)$$

- $P(\mathbf{w}_i | \mathbf{t}_i)$ estimated from tagged corpus:

$$\frac{C(\mathbf{w}_i \text{ has tag } \mathbf{t}_i)}{C(\mathbf{t}_i)}$$

- example: $P(\text{book} | \text{verb})$ is count of how many times **book** has tag **verb** divided by how many times tag **verb** occurs in corpus
- $P(\text{book} | \text{verb}) > P(\text{book} | \text{noun})$
 - there are many more nouns than verbs
 - say 1,000 verbs and 10,000 nouns

Markov Model Tagger: Second Assumption

$$P(\mathbf{w}_{1,n} | \mathbf{t}_{1,n}) P(\mathbf{t}_{1,n})$$

2. Each tag depends only on one previous tag

$$P(\mathbf{t}_{1,n}) = \prod_{i=1}^n P(\mathbf{t}_i | \mathbf{t}_{i-1}) = P(\mathbf{t}_1 | \mathbf{t}_0) P(\mathbf{t}_2 | \mathbf{t}_1) \dots P(\mathbf{t}_n | \mathbf{t}_{n-1})$$

- this is **Markov** assumption we saw in language modeling
- estimate as in language modeling

$$P(\mathbf{t}_i | \mathbf{t}_{i-1}) = \frac{C(\mathbf{t}_{i-1} \mathbf{t}_i)}{C(\mathbf{t}_{i-1})}$$

- $P(\mathbf{t}_1 | \mathbf{t}_0)$ stands for $P(\mathbf{t}_1)$, estimated by $P(\mathbf{t}_1) = \frac{C(\mathbf{t}_1)}{N}$

Markov Model Tagger

- Using these 2 assumptions, find tagging that maximizes

$$\prod_{i=1}^n P(\mathbf{w}_i | \mathbf{t}_i) P(\mathbf{t}_i | \mathbf{t}_{i-1}) \quad (1)$$

- Naïve algorithm: given sentence $\mathbf{w}_{1,n}$ go over all possible tag assignments $\mathbf{t}_{1,n}$ and compute (1)
- Choose final tagging $\mathbf{t}_{1,n}$ which maximizes (1)
 - efficiency: for each word try only tags given by the dictionary
 - example: for **fly**, possible tags are **noun**, **verb** and also **adjective** (meaning keen or artful, mainly in England)

Markov Model Tagger

- Markov tagger becomes Charniak's tagger if assume independent tags

$$P(\mathbf{t}_i | \mathbf{t}_{i-1}) = P(\mathbf{t}_i)$$

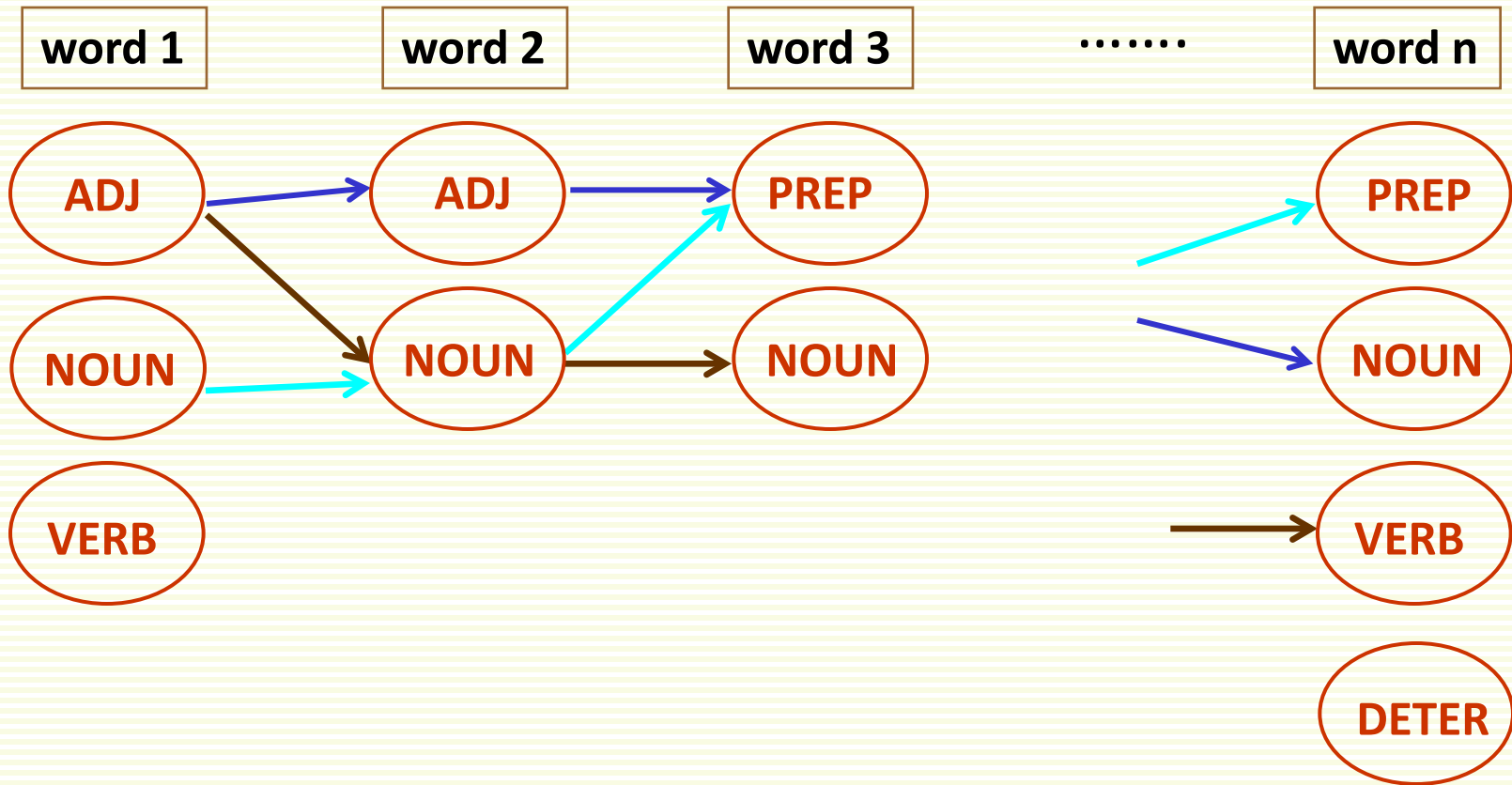
$$\begin{aligned} \prod_{i=1}^n P(\mathbf{w}_i | \mathbf{t}_i) P(\mathbf{t}_i | \mathbf{t}_{i-1}) &= \prod_{i=1}^n P(\mathbf{w}_i | \mathbf{t}_i) P(\mathbf{t}_i) \\ &= \prod_{i=1}^n \frac{P(\mathbf{w}_i, \mathbf{t}_i)}{P(\mathbf{t}_i)} P(\mathbf{t}_i) \\ &= \prod_{i=1}^n P(\mathbf{w}_i, \mathbf{t}_i) \end{aligned}$$

Markov Model Tagger

- Naïve algorithm: given sentence $\mathbf{w}_{1,n}$ go over all possible tag assignments $\mathbf{t}_{1,n}$
- 40 % words have more than 1 tag
- too many tag assignments to try
- if 2 tags per word, then 2^n possible assignments
- exhaustive search is exponential

Markov Model Tagger

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



Markov Model Tagger: DP

- DP (dynamic programming) to find optimal tagging efficiently
 - if k tags per word and n words, can find best tagging in $O(k^2n)$
 - also called Viterbi algorithm
- To avoid floating point underflows, take logarithms

$$\log \left[\prod_{i=1}^n \mathbf{P}(\mathbf{w}_i | \mathbf{t}_i) \mathbf{P}(\mathbf{t}_i | \mathbf{t}_{i-1}) \right] = \sum_{i=1}^n \left(\underbrace{\log \mathbf{P}(\mathbf{w}_i | \mathbf{t}_i)}_{\text{how likely word } \mathbf{w}_i \text{ is for tag } \mathbf{t}_i} + \underbrace{\log \mathbf{P}(\mathbf{t}_i | \mathbf{t}_{i-1})}_{\text{how likely tag } \mathbf{t}_i \text{ to follow tag } \mathbf{t}_{i-1}} \right)$$

- Break into two sums

$$\sum_{i=1}^n \log \mathbf{P}(\mathbf{w}_i | \mathbf{t}_i) + \sum_{i=1}^n \log \mathbf{P}(\mathbf{t}_i | \mathbf{t}_{i-1})$$

- Turn maximization into equivalent minimization

$$- \sum_{i=1}^n \log \mathbf{P}(\mathbf{w}_i | \mathbf{t}_i) - \sum_{i=1}^n \log \mathbf{P}(\mathbf{t}_i | \mathbf{t}_{i-1})$$

Markov Model Tagger: DP

- Find a sequence of tags $\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_n$ to minimize

$$\sum_{i=1}^n -\log P(\mathbf{w}_i | \mathbf{t}_i) + \sum_{i=1}^n -\log P(\mathbf{t}_i | \mathbf{t}_{i-1})$$

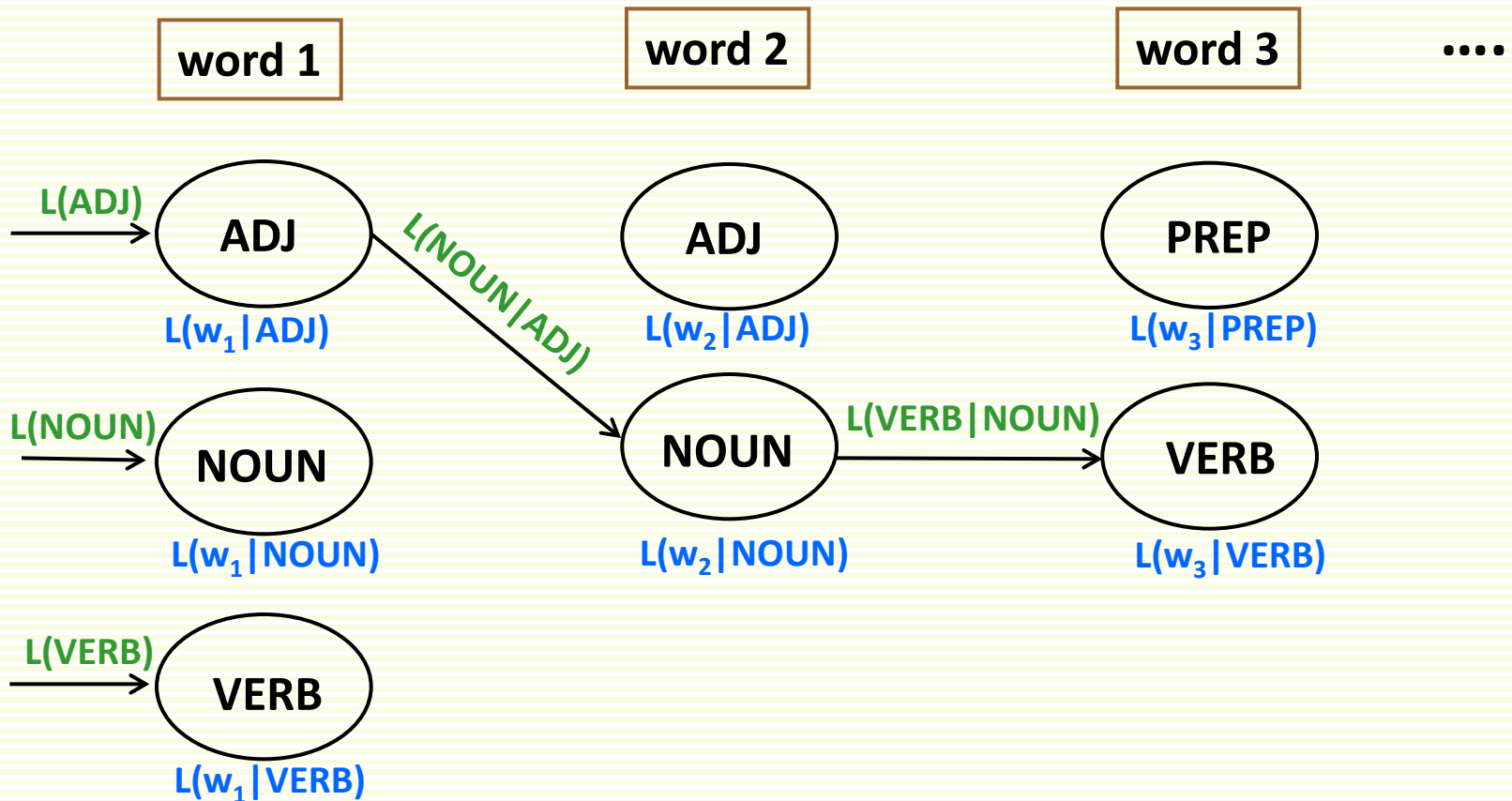
$L(\mathbf{w}_i | \mathbf{t}_i)$ $L(\mathbf{t}_i | \mathbf{t}_{i-1})$

- Using new notation, find tags $\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_n$ to minimize

$$\sum_{i=1}^n L(\mathbf{w}_i | \mathbf{t}_i) + \sum_{i=1}^n L(\mathbf{t}_i | \mathbf{t}_{i-1})$$

Markov Model Tagger: DP

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

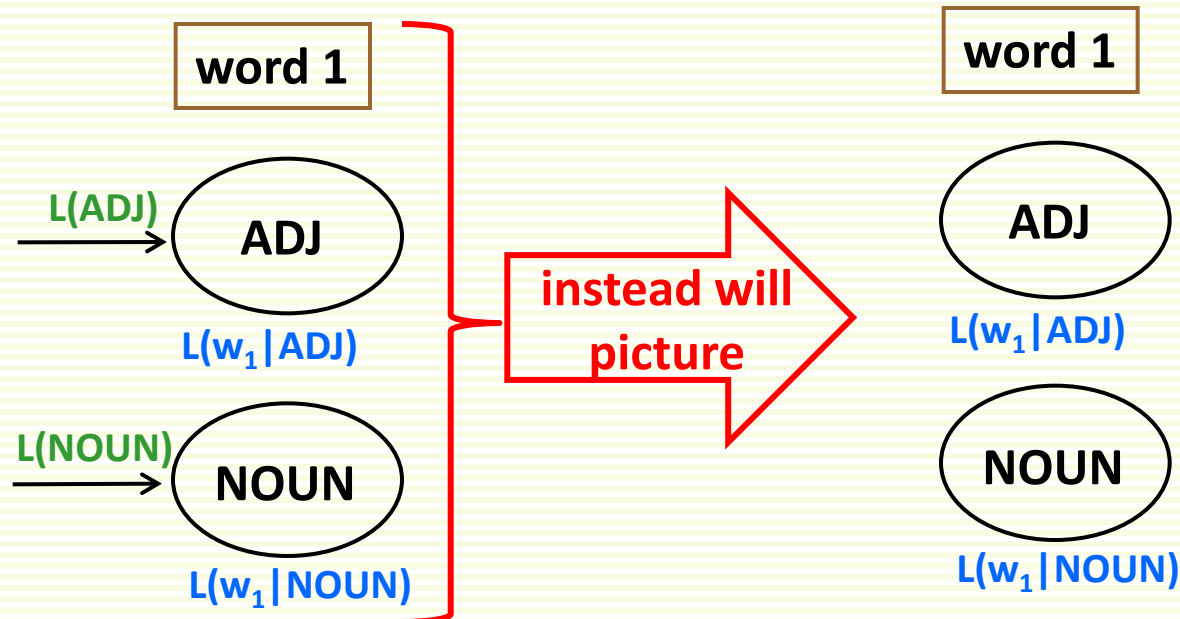


Markov Model Tagger: DP

- Change notation for the first word

$$L(\mathbf{w}_1 | \mathbf{t}_1) = -\log[P(\mathbf{w}_1 | \mathbf{t}_1)] - \log [P(\mathbf{t}_1 | \mathbf{t}_0)]$$

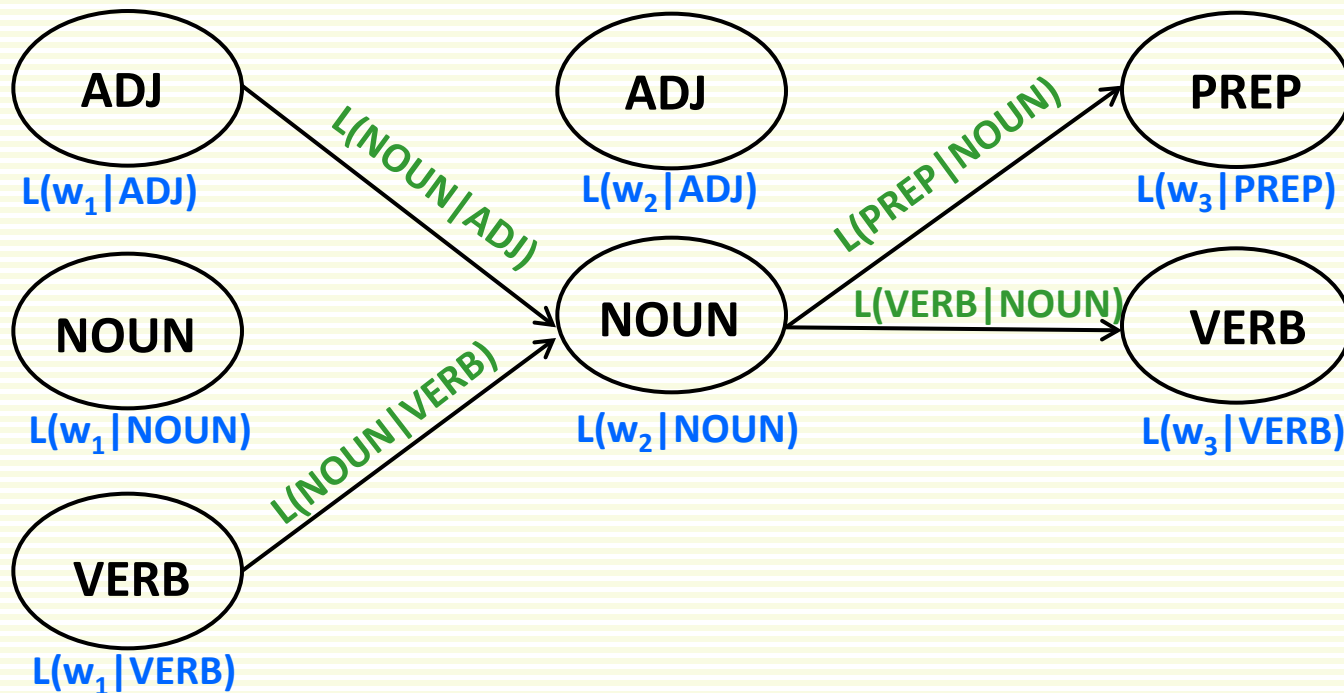
$$\sum_{i=1}^n L(\mathbf{w}_i | \mathbf{t}_i) + \sum_{i=1}^n L(\mathbf{t}_i | \mathbf{t}_{i-1}) \Rightarrow \sum_{i=1}^n L(\mathbf{w}_i | \mathbf{t}_i) + \sum_{i=2}^n L(\mathbf{t}_i | \mathbf{t}_{i-1})$$



Markov Model Tagger: DP

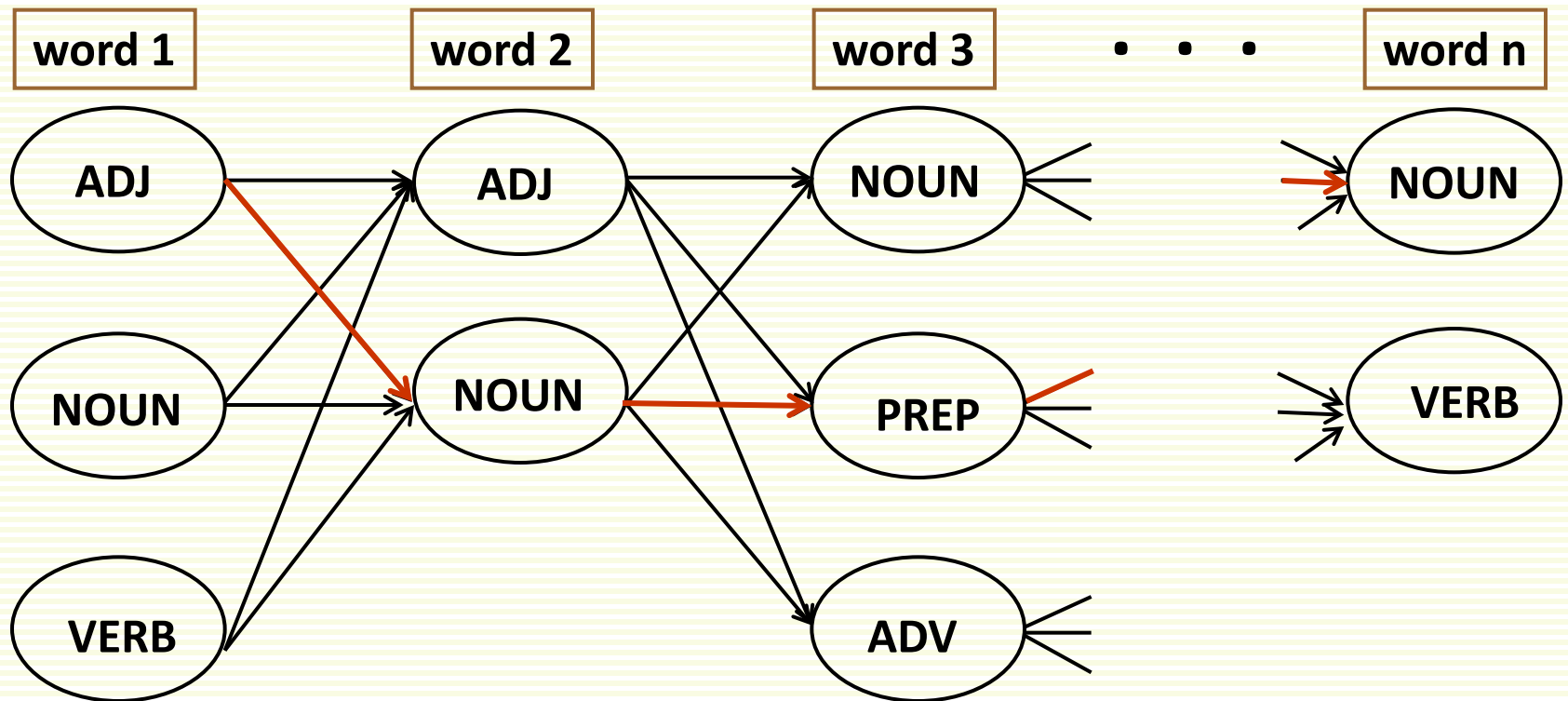
- Each node has cost $L(w_i | t_i)$
- Each edge has cost $L(t_i | t_{i-1})$

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



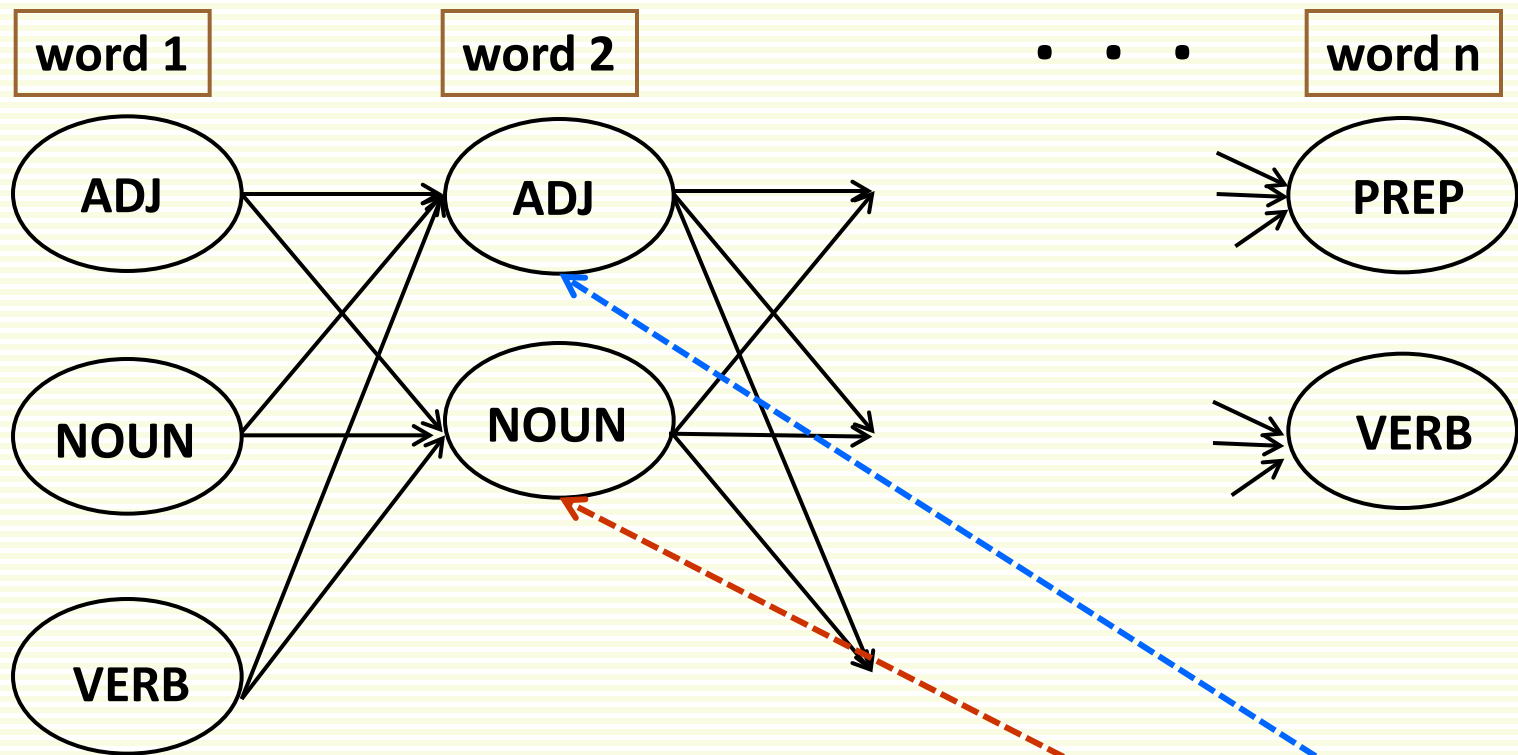
Markov Model Tagger: DP Graph

- Find minimum cost path that starts at **word 1** node and ends at **word n** node



Markov Model Tagger: DP Iteration Step

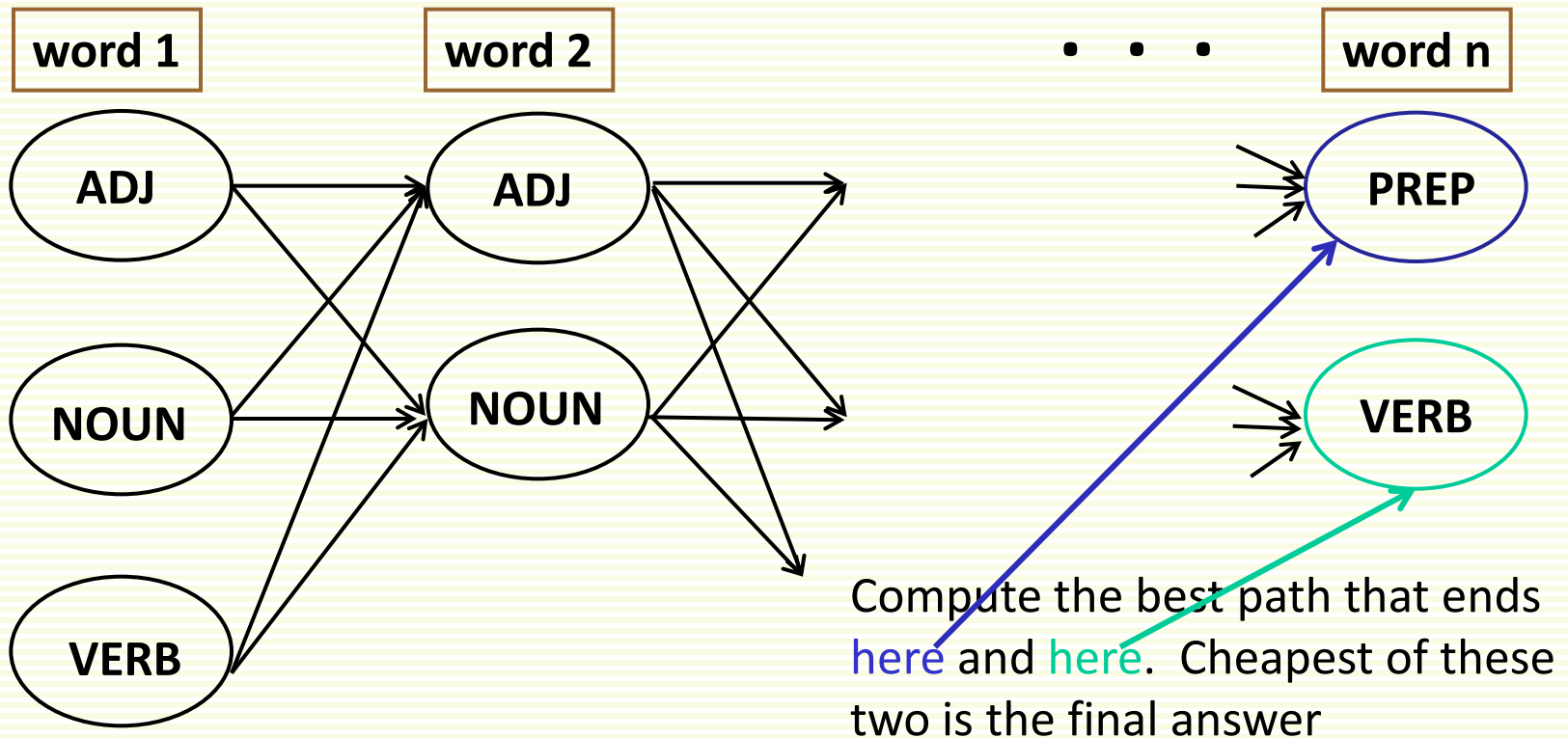
- Iteration i : for each word w_i node, find smallest cost path that ends there. This path starts at any word w_1 node.



iteration 2: compute best path that ends **here** and **here**

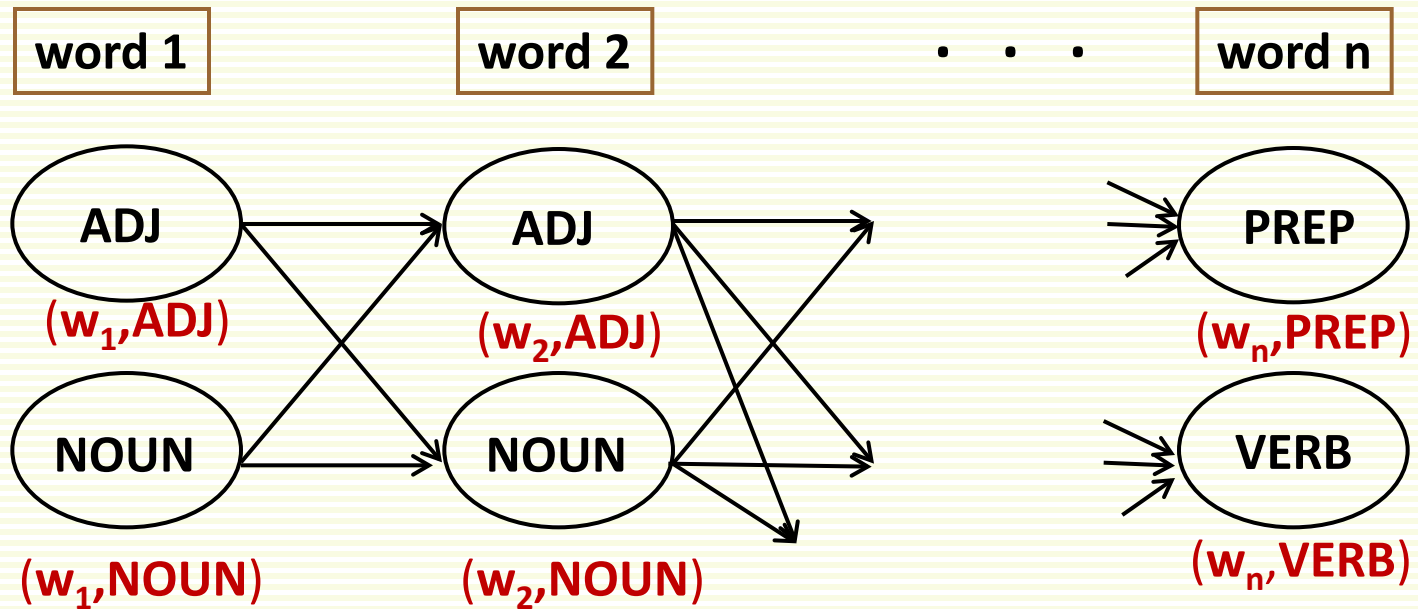
Markov Model Tagger: DP Overview

- All paths start at any w_1 node
- 1. Find smallest cost path that ends at w_1 node for each w_1 node
- 2. Find smallest cost path that ends at w_2 node for each w_2 node
-
- n. Find smallest cost path that ends at w_n node for each w_n node
- The overall best path is smallest cost path that ends at w_n node



Markov Model Tagger: DP Variables

- For word w_i tag t node is (w_i, t)



- $C(w_i, t)$ cost of best path that starts at any (w_1, t) and ends at (w_i, t)
- $P(w_i, t)$ is parent of node (w_i, t) on this path
- After all $C(w_i, t)$ computed, min of $C(w_n, t)$ over all t gives best path

Markov Model Tagger: DP Initialization

- First compute the best path that ends at any node for w_1
 - trivial, since the path has just one node
- For all tags of the first word t :

$$C(w_1, t) = L(w_1 | t)$$

$$P(w_1, t) = \text{null}$$

word 1

ADJ

$L(w_1 | \text{ADJ})$

NOUN

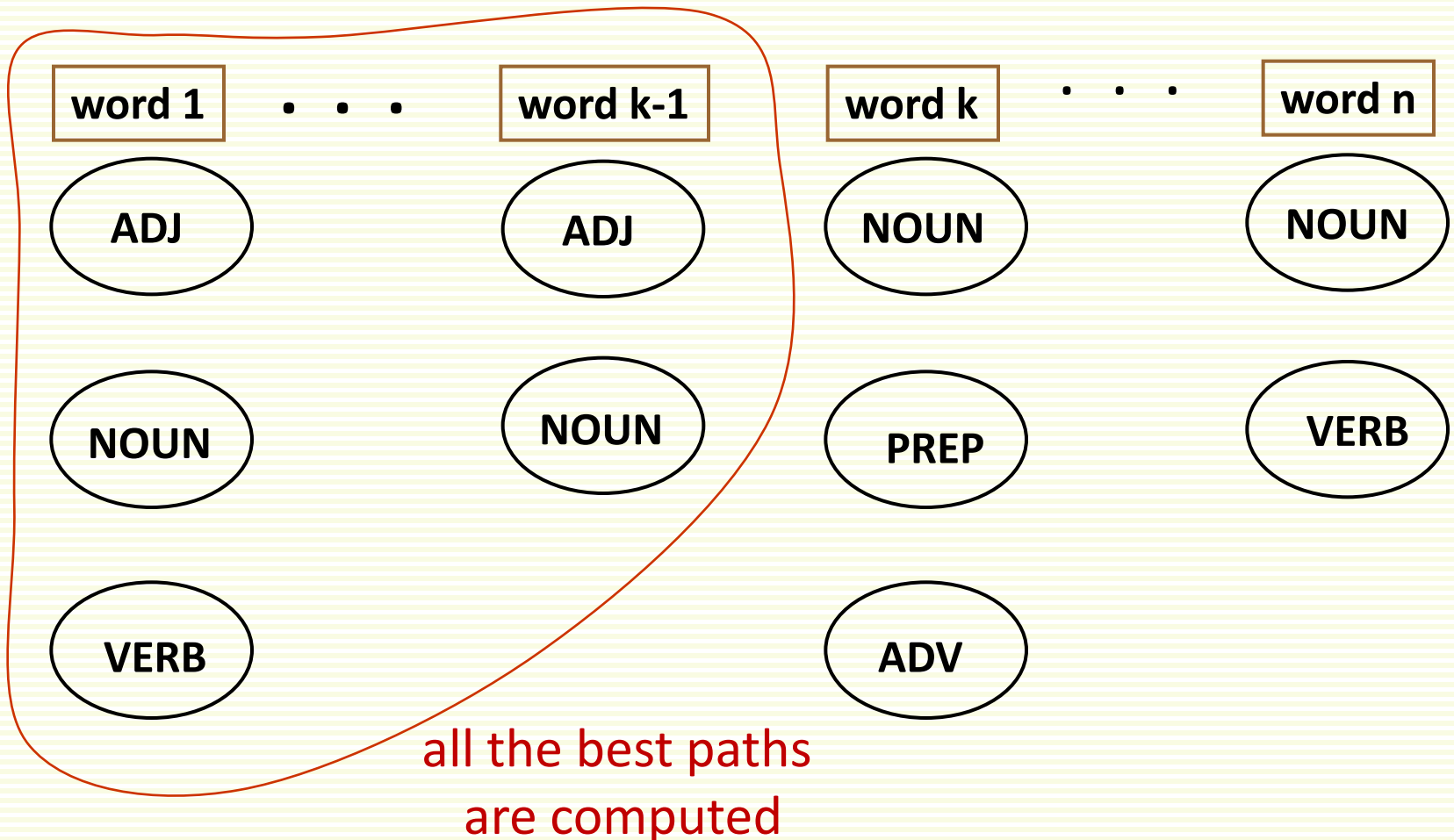
$L(w_1 | \text{NOUN})$

VERB

$L(w_1 | \text{VERB})$

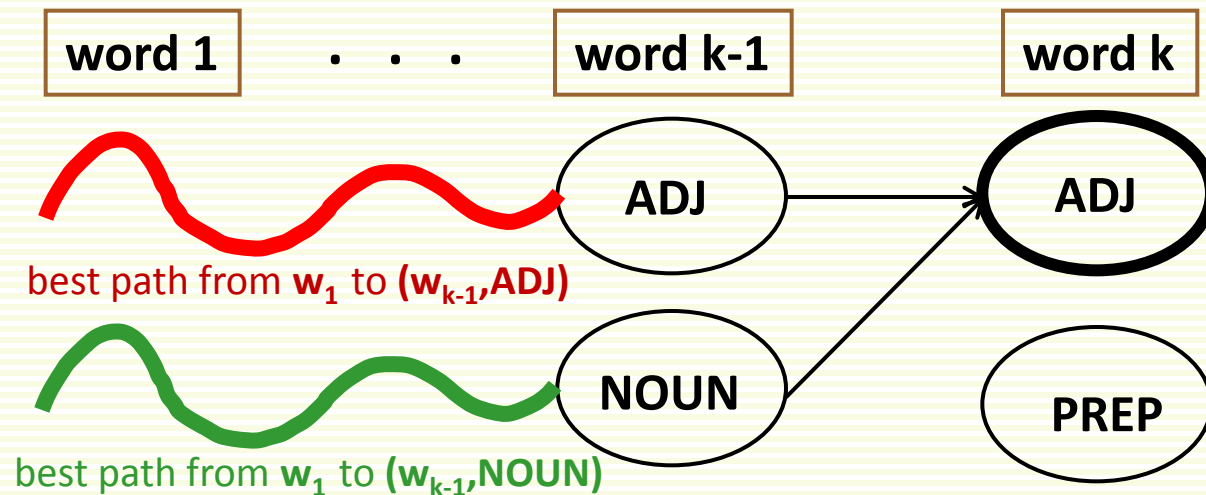
Markov Model Tagger: DP Iteration

- Computed $C(w_i, t)$ and $P(w_i, t)$ for all tags t and $i < k$



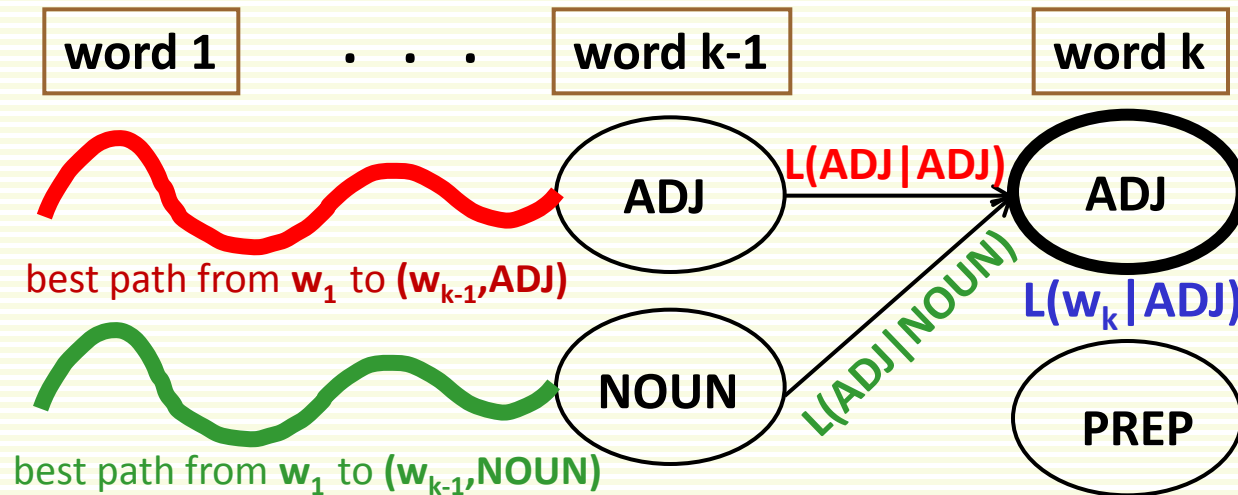
Markov Model Tagger: DP Iteration

- Now compute $C(\mathbf{w}_k, \mathbf{t})$ and $P(\mathbf{w}_k, \mathbf{t})$ for k
- Consider node $(\mathbf{w}_k, \text{ADJ})$



- The best path from \mathbf{w}_1 to $(\mathbf{w}_k, \text{ADJ})$ goes through either
 1. $(\mathbf{w}_{k-1}, \text{ADJ})$: then it follows **best path from \mathbf{w}_1 to $(\mathbf{w}_{k-1}, \text{ADJ})$**
 2. $(\mathbf{w}_{k-1}, \text{NOUN})$: then it follows **best path from \mathbf{w}_1 to $(\mathbf{w}_{k-1}, \text{NOUN})$**
 - because a sub-path of the best path is a best path itself

Markov Model Tagger: DP Iteration



- $C(w_k, \text{ADJ})$ is the smaller of two quantities
 1. $C(w_{k-1}, \text{ADJ}) + L(\text{ADJ}|\text{ADJ}) + L(w_k|\text{ADJ})$
 - then $P(w_k, \text{ADJ}) = (w_{k-1}, \text{ADJ})$
 2. $C(w_{k-1}, \text{NOUN}) + L(\text{ADJ}|\text{NOUN}) + L(w_k|\text{ADJ})$
 - then $P(w_k, \text{ADJ}) = (w_{k-1}, \text{NOUN})$

Markov Model Tagger: DP Iteration

- In general, $C(\mathbf{w}_k, \mathbf{t})$ is computed as follows:

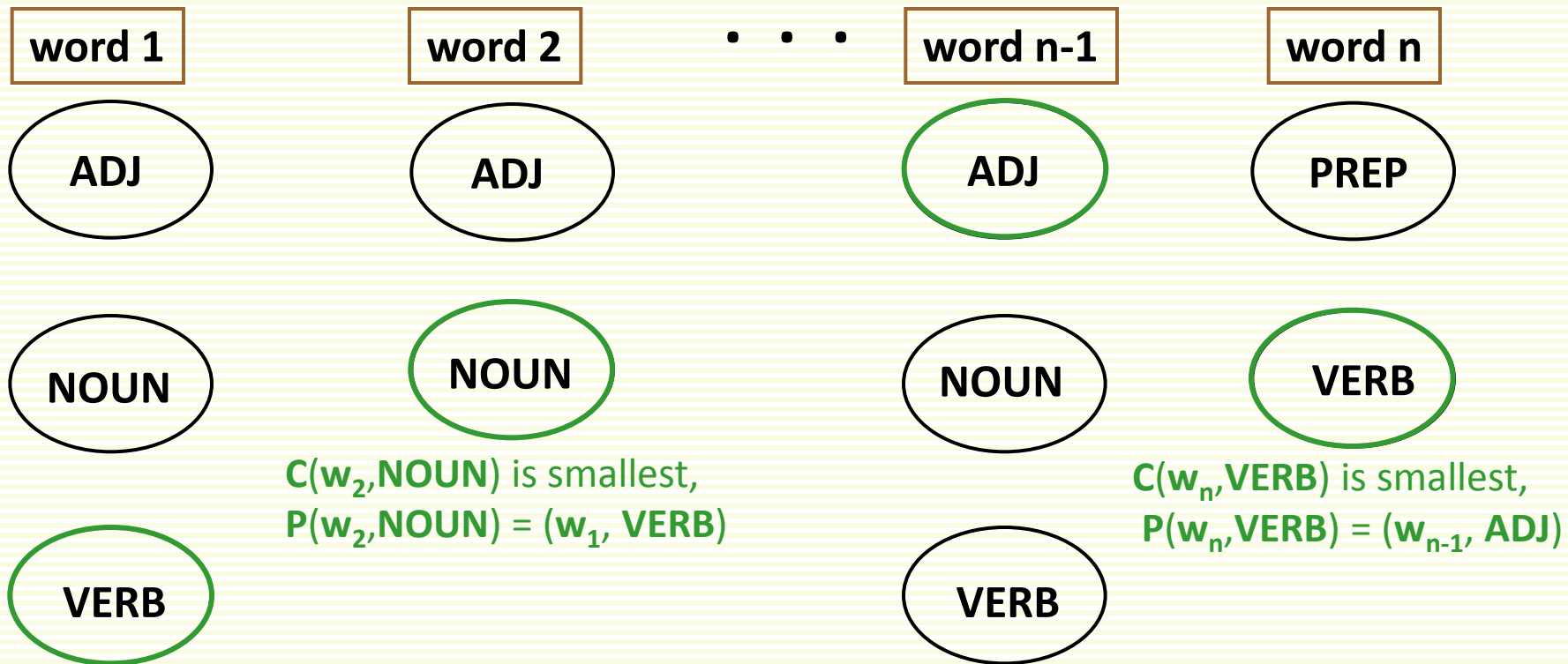
$$C(\mathbf{w}_k, \mathbf{t}) = \min_{\mathbf{t}' \in T(\mathbf{w}_{k-1})} \left\{ \underbrace{C(\mathbf{w}_{k-1}, \mathbf{t}')}_{\text{cost of best path from first word to node (word k-1, t')}} + \underbrace{L(\mathbf{t} | \mathbf{t}')}_{\text{cost of going between nodes (w}_{k-1}, \mathbf{t}') \text{ and (w}_k, \mathbf{t})} \right\} + \underbrace{L(\mathbf{w}_k | \mathbf{t})}_{\text{cost of going through node (w}_k, \mathbf{t})}$$

search over all tags \mathbf{t}' for word $k-1$

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

Markov Model Tagger: DP Termination

- After computed all $C(w_i, t)$ best cost path is found as the minimum of $C(w_n, t)$ over all tags t
- Parents on the path traced back using $P(w_i, t)$



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

MMT Example

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

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

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

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

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

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

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

book

ADJ

VERB

NOUN

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

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

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

$L(\text{CONJ}|\text{NOUN}) = 1$

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

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

that

PRON

CONJ

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

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

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

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

flight

NOUN

VERB

MMT 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:
 - $C(\text{book},\text{ADJ}) = 10$, $P(\text{book},\text{ADJ}) = \text{null}$
 - $C(\text{book},\text{VERB}) = 1$, $P(\text{book},\text{VERB}) = \text{null}$
 - $C(\text{book},\text{NOUN}) = 2$, $P(\text{book},\text{NOUN}) = \text{null}$

MMT Example

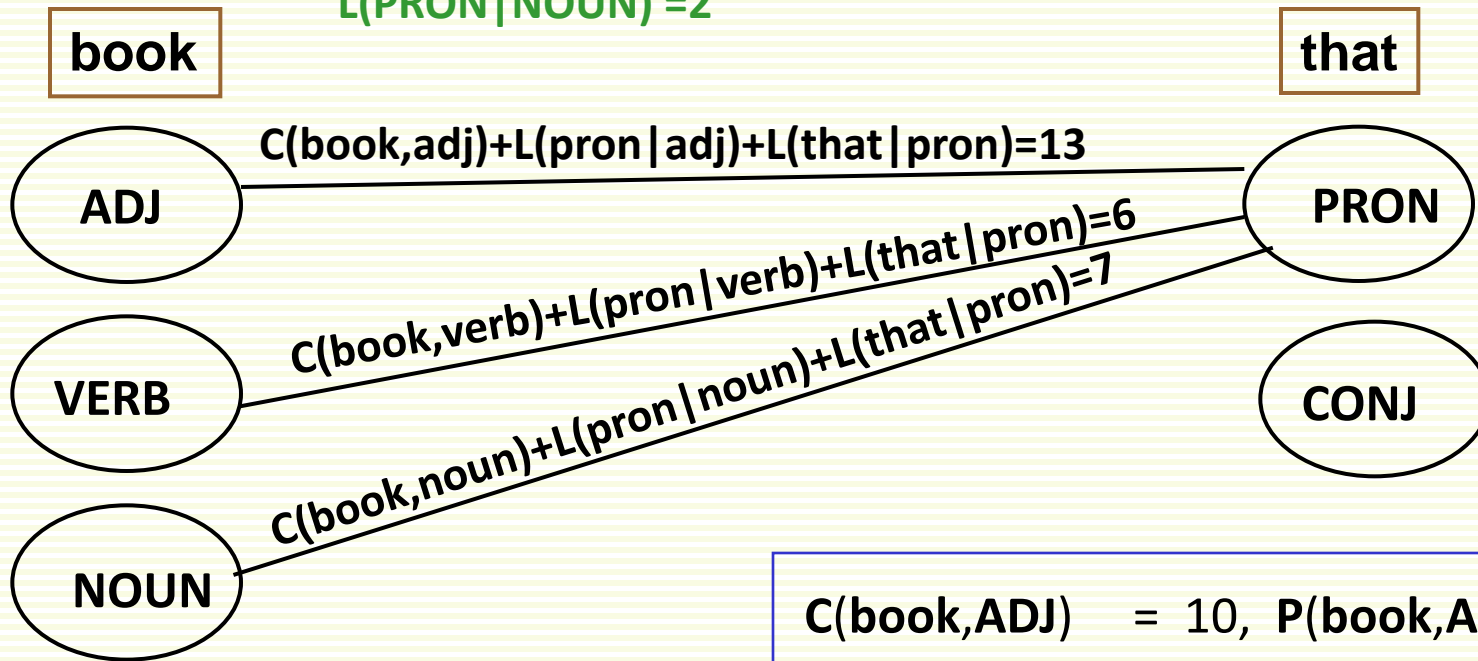
$$L(\text{PRON} | \text{ADJ}) = 1$$

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

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

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

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



$C(\text{book}, \text{ADJ}) = 10$	$P(\text{book}, \text{ADJ}) = \text{null}$
$C(\text{book}, \text{VERB}) = 1$	$P(\text{book}, \text{VERB}) = \text{null}$
$C(\text{book}, \text{NOUN}) = 2$	$P(\text{book}, \text{NOUN}) = \text{null}$

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

MMT Example

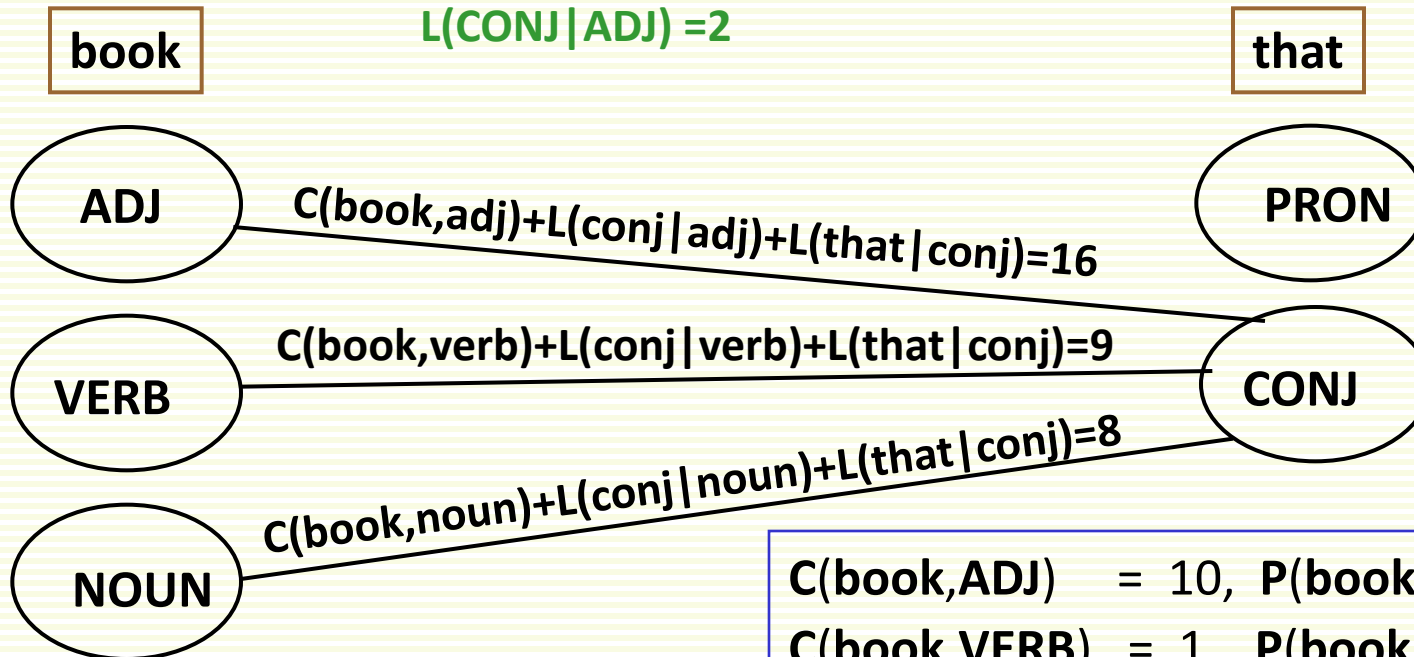
$$L(\text{CONJ}|\text{VERB})=4$$

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

$$L(\text{CONJ}|\text{NOUN})= 1$$

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

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



$$\begin{aligned} C(\text{book,ADJ}) &= 10, & P(\text{book,ADJ}) &= \text{null} \\ C(\text{book,VERB}) &= 1, & P(\text{book,VERB}) &= \text{null} \\ C(\text{book,NOUN}) &= 2, & P(\text{book,NOUN}) &= \text{null} \end{aligned}$$

- Iteration 2:

- $C(\text{that, CONJ}) = 8, P(\text{that, CONJ}) = (\text{book,NOUN})$

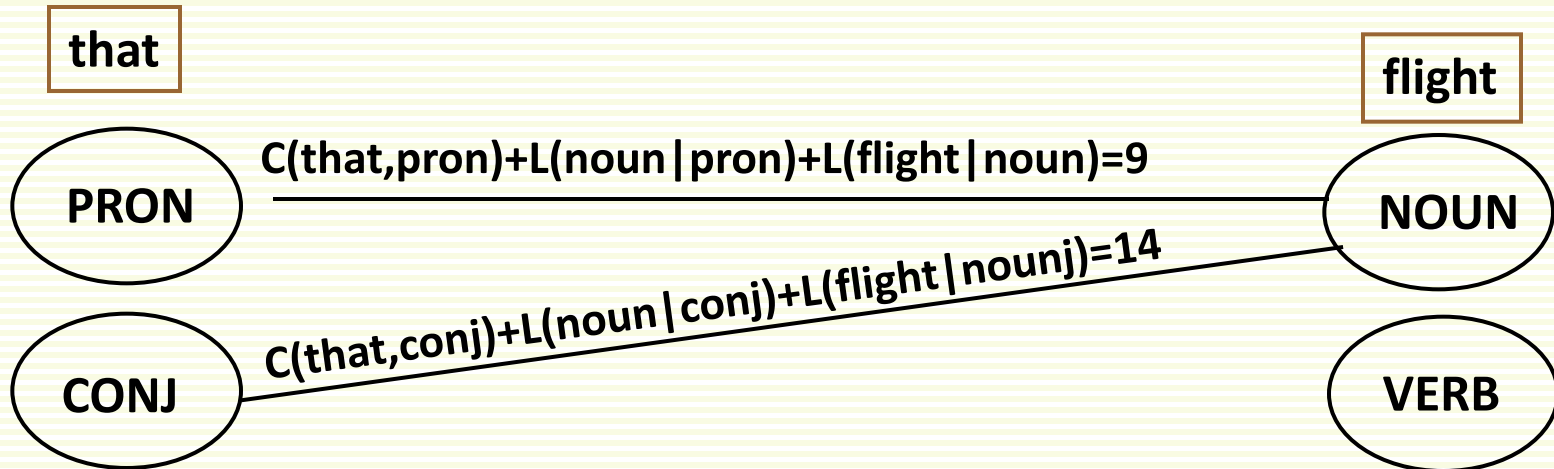
MMT Example

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

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

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

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



$C(\text{book,ADJ}) = 10, P(\text{book,ADJ}) = \text{null}$

$C(\text{book,VERB}) = 1, P(\text{book,VERB}) = \text{null}$

$C(\text{book,NOUN}) = 2, P(\text{book,NOUN}) = \text{null}$

$C(\text{that, PRON}) = 6, P(\text{that, PRON}) = (\text{book,VERB})$

$C(\text{that, CONJ}) = 8, P(\text{that, CONJ}) = (\text{book,NOUN})$

- Iteration 3:

- $C(\text{flight, NOUN}) = 9, P(\text{flight, NOUN}) = (\text{that,PRON})$

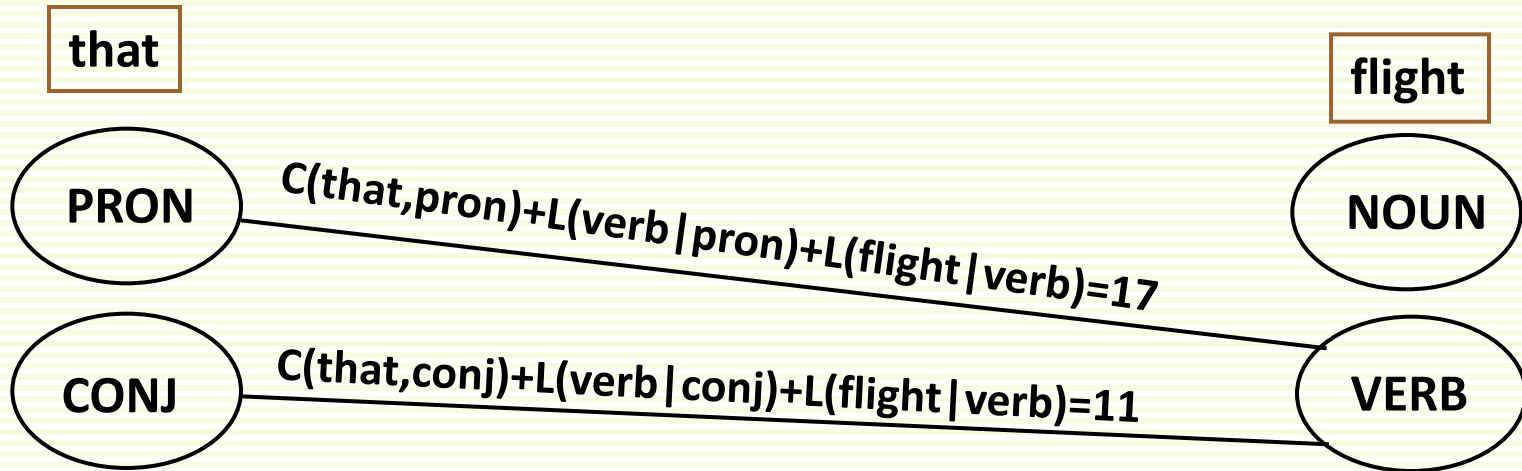
MMT Example

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

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

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

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



$C(\text{book,ADJ}) = 10, P(\text{book,ADJ}) = \text{null}$

$C(\text{book,VERB}) = 1, P(\text{book,VERB}) = \text{null}$

$C(\text{book,NOUN}) = 2, P(\text{book,NOUN}) = \text{null}$

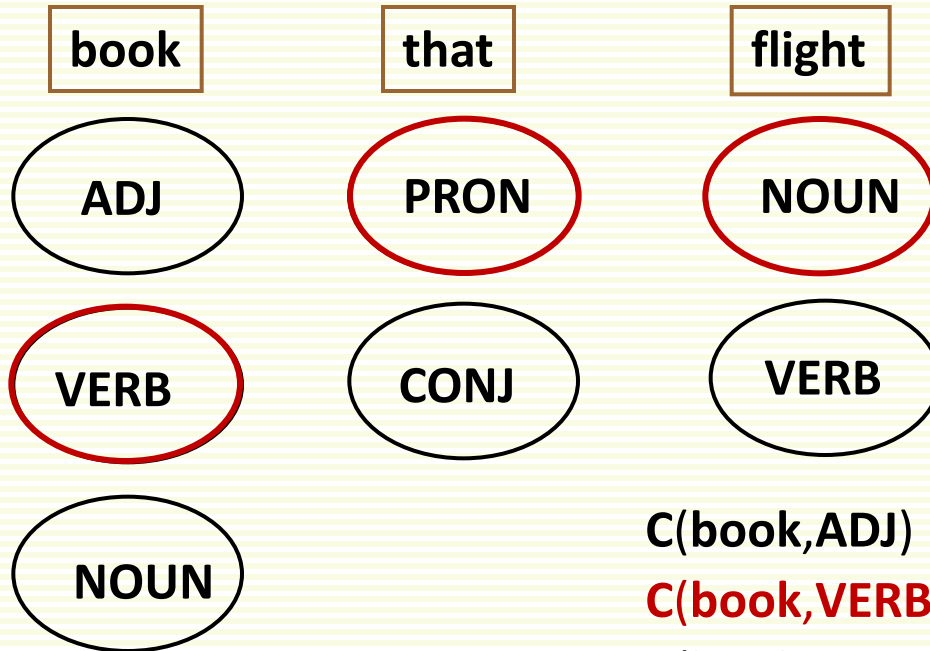
$C(\text{that, PRON}) = 6, P(\text{that, PRON}) = (\text{book,VERB})$

$C(\text{that, CONJ}) = 8, P(\text{that, CONJ}) = (\text{book,NOUN})$

- Iteration 3:

- $C(\text{flight, VERB}) = 11, P(\text{flight, VERB}) = (\text{that,CONJ})$

MMT Example



$C(\text{book}, \text{ADJ}) = 10$, $P(\text{book}, \text{ADJ}) = \text{null}$

$C(\text{book}, \text{VERB}) = 1$, $P(\text{book}, \text{VERB}) = \text{null}$

$C(\text{book}, \text{NOUN}) = 2$, $P(\text{book}, \text{NOUN}) = \text{null}$

$C(\text{that}, \text{PRON}) = 6$, $P(\text{that}, \text{PRON}) = (\text{book}, \text{VERB})$

$C(\text{that}, \text{CONJ}) = 8$, $P(\text{that}, \text{CONJ}) = (\text{book}, \text{NOUN})$

$C(\text{flight}, \text{NOUN}) = 9$, $P(\text{flight}, \text{NOUN}) = (\text{that}, \text{PRON})$

$C(\text{flight}, \text{VERB}) = 11$, $P(\text{flight}, \text{VERB}) = (\text{that}, \text{PRON})$

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

MMT: Pseudo Code for DP

- $\text{Tags}(\mathbf{w}_i)$ is the set of all possible tags for \mathbf{w}_i

for each $t \in \text{Tags}(\mathbf{w}_1)$ do

$\mathbf{C}(\mathbf{w}_1, t) = \mathbf{L}(\mathbf{w}_1 | t), \mathbf{P}(\mathbf{w}_1, t) = \text{null}$

for $i \leftarrow 2$ to n do

for each $t \in \text{Tag}(\mathbf{w}_i)$ do

$\mathbf{C}(\mathbf{w}_i, t) = -\infty$

for each $t' \in \text{Tag}(\mathbf{w}_{i-1})$ do

$\text{nextCost} = \mathbf{C}(\mathbf{w}_{i-1}, t') + \mathbf{L}(t | t') + \mathbf{L}(\mathbf{w}_i | t)$

if $\text{nextCost} < \text{cost}(\mathbf{w}_i, t)$ do

$\mathbf{C}(\mathbf{w}_i, t) = \text{nextCost}$

$\mathbf{P}(\mathbf{w}_i, t) = t'$

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
 - $P(\text{"karumbula"} | \text{verb}) = P(\text{"karumbula"} | \text{noun}) = P(\text{"karumbula"} | \text{adjective}) = \dots \text{ etc}$
- 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 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