



Statistical Natural Language Processing

Lecture 8: Part of Speech Tagging

Dr. Momtazi

Amirkabir University of Technology

Outline

2

- 1 Part of Speech Tagging
- 2 Sequential Modeling
- 3 Evaluation

Outline

3

1 Part of Speech Tagging

2 Sequential Modeling

3 Evaluation

Parts Of Speech (POS)

4

- 8 Parts of speech are traditionally used to summarize the linguistic knowledge
 - Noun, Verb, Preposition, Adverb, Article, Interjection, Pronoun, Conjunction
- The modified list is currently used
 - Noun, Verb, Auxiliary, Preposition, Adjective, Adverb, Number, Determiner, Interjection, Pronoun, Conjunction, Particle
- Known as:
 - Parts of speech
 - Lexical categories
 - Word classes
 - Morphological classes
 - Lexical tags

POS Examples

5

Noun	book/books, sugar, Germany, Sony
Verb	eat, wrote
Auxiliary	can, should, have
Adjective	new, newer, newest
Adverb	well, urgently
Numbers	872, two, first
Determiner	the, some
Conjunction	and, or
Pronoun	he, my
Preposition	to, in
Particle	off, up
Interjection	Ow, Eh

Open vs. Closed Classes

6

- Closed (limited number of words, do not grow usually)
 - Determiners: the, some, a, an, ...
 - Pronouns: she, he, I, ...
 - Prepositions: to, in, on, under, over, by, ...
 - Auxiliaries: can, should, have, had, ...
 - Conjunctions: and, or
 - Particles: off, up
 - Interjections: Ow, Eh

- Open (unlimited number of words)
 - Nouns
 - Verbs
 - Adjectives
 - Adverbs

Applications

7

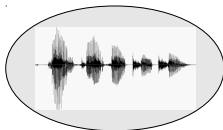
- Speech Synthesis
- Parsing
- Machine Translation
- Information Extraction

Applications

8

■ Speech Synthesis

How to pronounce “*lead*” ?



Applications

9

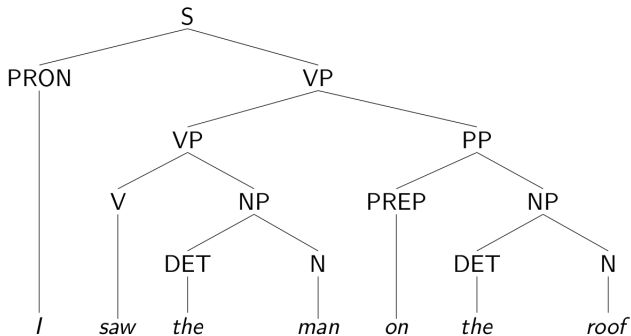
- Machine Translation

"I like ..." ⇒

Applications

10

■ Parsing



POS Tagset

- There are so many parts of speech tagsets we can draw
- Choosing a standard tagset is essential
- Tag types
 - Coarse-grained
 - noun
 - verb
 - adjective
 - ...
 - Fine-grained
 - noun-proper-singular, noun-proper-plural, noun-common-mass, ..
 - verb-past, verb-present-3rd, verb-base, ...
 - adjective-simple, adjective-comparative, ...
 - ...

Penn TreeBank

A large annotated corpus of English
tagset: 45 tags

Penn TreeBank Tagset

12

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>	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>
PRP	personal pronoun	<i>I, you, he</i>	(left parenthesis	<i>[, (, {, <</i>
PRP\$	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>			

Ambiguity

13

- Definition
 - The process of assigning a part of speech to each word in a text
- Challenge
 - Words often have more than one POS

On my back

The back door

Pay the money back

Promised to back the bill

Ambiguity

14

- Definition
 - The process of assigning a part of speech to each word in a text
- Challenge
 - Words often have more than one POS

On my back_[NN]

The back_[JJ] door

Pay the money back_[RB]

Promised to back_[VB] the bill

Distribution of Ambiguities

15

45-tag Treebank Brown		
Unambiguous (1 tag)		38,857
Ambiguous (2–7 tags)		8844
Details:	2 tags	6,731
	3 tags	1621
	4 tags	357
	5 tags	90
	6 tags	32
	7 tags	6 (<i>well, set, round, open, fit, down</i>)
	8 tags	4 (<i>'s, half, back, a</i>)
	9 tags	3 (<i>that, more, in</i>)

Distribution of Ambiguities

- The frequency of ambiguous words are relatively high
 - 11.5% of word types
 - 40% of word tokens

Goal

- Using a set of labeled data to train a model
- Using the trained model to predict the POS tag of the unseen words

POS Tagging

18

Plays well with others

Plays	NNS/VBZ
well	UH/JJ/NN/RB
with	IN
others	NNS

Plays_[VBZ] well_[RB] with_[IN] others_[NNS]

Performance

■ Baseline model

- Tagging unambiguous words with the correct label
- Tagging ambiguous words with their most frequent label
- Tagging unknown words as a noun

Already performs around 90%

Outline

20

① Part of Speech Tagging

② Sequential Modeling

③ Evaluation

Task

21

- Similar to a normal classification task
 - Feature Selection
 - Algorithm

POS Tagging

22

■ Features

Word

the: the → DT

Prefixes

unbelievable: un- → JJ

Suffixes

slowly: -ly → RB

Lowercased word

Importantly: importantly → RB

Capitalization

Stefan: [CAP] → NNP

Word shapes

35-year: d-x → JJ

■ Model

- Maximum Entropy
 $P(t|w)$

Data	Performance
Overall	93.7
Unknown	82.6

POS Tagging

23

- More Features?

*They*_[PRP] *left*_[VBD] *as*_[IN] *soon*_[RB] *as*_[IN] *he*_[PRP] *arrived*_[VBD]

- Better Algorithm

- Using Sequence Modeling

Sequence Modeling

24

- Many of the NLP techniques should deal with data represented as sequence of items
 - Characters, Words, Phrases, Lines, ...

警察枪杀了那个逃

B I B I B B B B I

*I*_[PRP] *saw*_[VBP] *the*_[DT] *man*_[NN] *on*_[IN] *the*_[DT] *roof*_[NN].

Sequence Modeling

25

- Two types of information
 - Local
 - Contextual

Sequence Modeling

26

■ Making a decision based on the

□ Current Observation

- Word (W_0)
- Prefix
- Suffix
- Lowercased word
- Capitalization
- Word shape

□ Surrounding observations

- W_{+1}
- W_{-1}

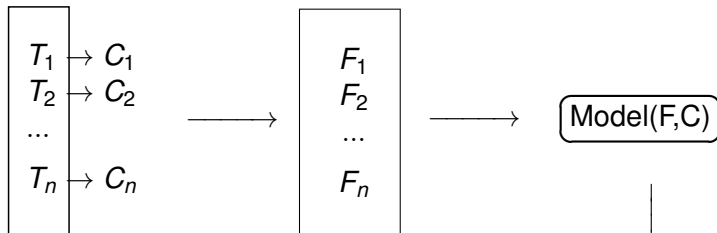
□ Previous decisions

- T_{-1}
- T_{-2}

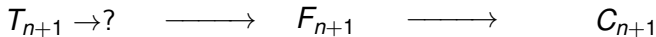
Learning Model

27

Training



Testing



Sequence Modeling

28

■ Greedy inference

- Starting from the beginning of the sequence
- Assigning a label to each item using the classifier in that position
- Using previous decisions as well as the observed data

■ Beam inference

- Keeping the top k labels in each position
- Extending each sequence in each local way
- Finding the best k labels for the next position

Hidden Markov Model (HMM)

29

- Finding the best sequence of tags ($t_1 \dots t_n$) that corresponds to the sequence of observations ($w_1 \dots w_n$)
- Probabilistic View
 - Considering all possible sequences of tags
 - Choosing the tag sequence from this universe of sequences, which is most probable given the observation sequence

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(t_1^n | w_1^n)$$

Using Bayes Rule

30

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(t_1^n | w_1^n)$$

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

$$P(t_1^n | w_1^n) = \frac{P(w_1^n | t_1^n) \cdot P(t_1^n)}{P(w_1^n)}$$



$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(w_1^n | t_1^n) \cdot P(t_1^n)$$

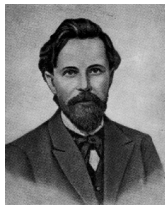
Using Markov Assumption

31

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(w_1^n | t_1^n) \cdot P(t_1^n)$$

$$P(w_1^n | t_1^n) \approx \prod_{i=1}^n P(w_i | t_i)$$

$$P(t_1^n) \approx \prod_{i=1}^n P(t_i | t_{i-1})$$



$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} \prod_{i=1}^n P(w_i | t_i) \cdot P(t_i | t_{i-1})$$

Two Probabilities

32

- The tag transition probabilities: $P(t_i|t_{i-1})$
 - Finding the likelihood of a tag to proceed by another tag
 - Similar to the normal bigram model

$$P(t_i|t_{i-1}) = \frac{C(t_{i-1}, t_i)}{C(t_{i-1})}$$

Two Probabilities

33

- The word likelihood probabilities: $P(w_i|t_i)$
 - Finding the likelihood of a word to appear given a tag

$$P(w_i|t_i) = \frac{C(t_i, w_i)}{C(t_i)}$$

Two Probabilities

34

- Zero probability problem
 - Solution: similar to language modelling, use the smoothing method for both probabilities

Two Probabilities

35

I_[PRP] saw_[VBP] the_[DT] man_[NN?] on the roof.

$$P([NN] | [DT]) = \frac{C([DT], [NN])}{C([DT])}$$

$$P(man | [NN]) = \frac{C([NN], man)}{C([NN])}$$

Ambiguity

36

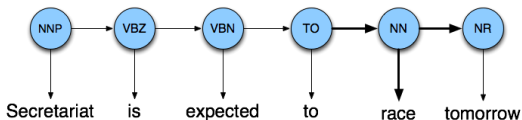
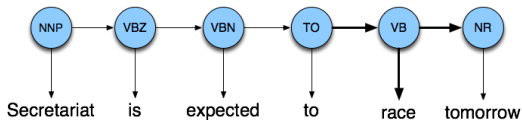
Secretariat_[NNP] is_[VBZ] expected_[VBN] to_[TO] **race**_[VB] tomorrow_[NR].

People_[NNS] inquire_[VB] the_[DT] reason_[NN] for_[IN] the_[DT] **race**_[NN].

Ambiguity

37

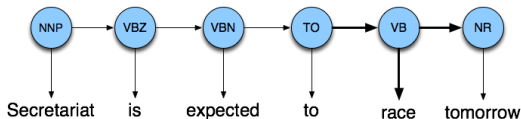
*Secretariat*_[NNP] *is*_[VBZ] *expected*_[VBN] *to*_[TO] ***race***_[VB] *tomorrow*_[NR].



Ambiguity

38

*Secretariat*_[NNP] *is*_[VBZ] *expected*_[VBN] *to*_[TO] ***race***_[VB] *tomorrow*_[NR].



$$P(VB|TO) = 0.83$$

$$P(\text{race}|VB) = 0.00012$$

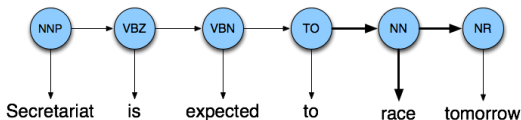
$$P(NR|VB) = 0.0027$$

$$P(VB|TO)P(NR|VB)P(\text{race}|VB) = 0.00000027$$

Ambiguity

39

*Secretariat*_[NNP] *is*_[VBZ] *expected*_[VBN] *to*_[TO] ***race***_[VB] *tomorrow*_[NR].



$$P(NN|TO) = 0.00047$$

$$P(race|NN) = 0.00057$$

$$P(NR|NN) = 0.0012$$

$$P(NN|TO)P(NR|NN)P(race|NN) = 0.00000000032$$

Performance

40

- Model

- Maximum Entropy
 $P(t|w)$

Data	Performance
Overall	93.7
Unknown	82.6

- HMM

Data	Performance
Overall	96.2
Unknown	86.0

- Upper bound (human agreement): $\sim 98\%$

Hidden Markov Model (HMM)

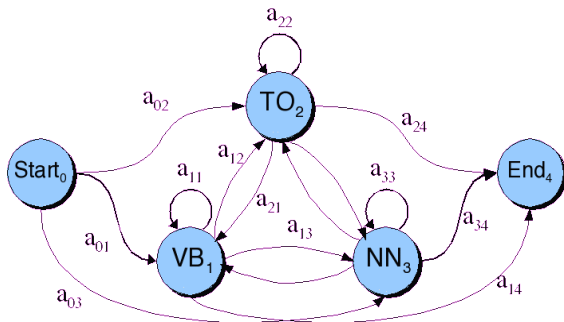
41

- A weighted finite-state automaton adds probabilities to the arcs
 - The probabilities leaving any arc must sum to one
- An HMM is an extension of a Markov chain in which the input symbols are not the same as the states
- We do not know which state we are in
 - The output symbols are words
 - The hidden states are POS tags

Hidden Markov Model (HMM)

42

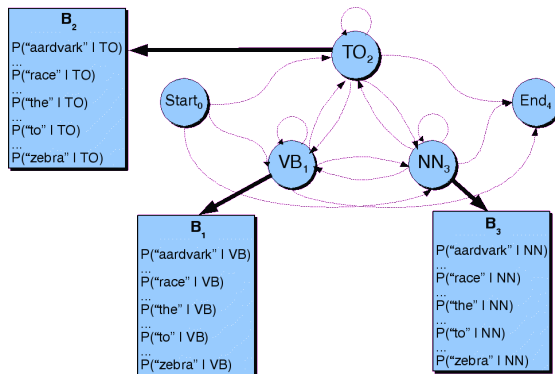
- Transition probabilities



Hidden Markov Model (HMM)

43

■ Word likelihood probabilities



The Viterbi Algorithm

44

- Viterbi inference
 - Memorizing the model using dynamic programming
 - Considering the small window of previous decisions

The Viterbi Algorithm

45

- Creating an array
 - Columns corresponding to inputs
 - Rows corresponding to possible states
- Sweeping through the array in one pass filling the columns left to right using the transition probabilities and observation probabilities
- Storing the max probability path to each cell (not all paths) using dynamic programming

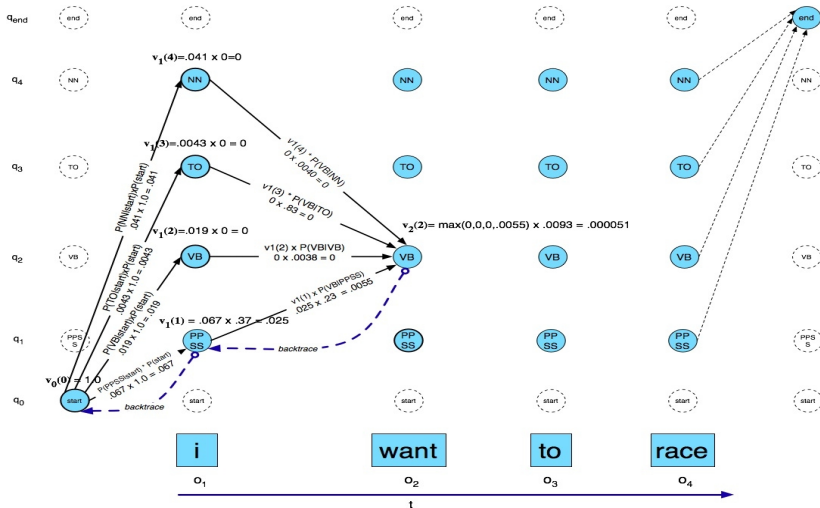
The Viterbi Algorithm

46

- Basic idea behind the algorithm:
the recursive definition for finding the maximum probability

The Viterbi Algorithm

47



Hidden Markov Model (HMM)

48

■ Proc and Cons

- HMM taggers are very simple to train
 - Just need to compile counts from the training corpus
- Perform relatively well
- Main difficulty is modelling $P(\text{word}|\text{tag})$ specially for complex words

Outline

49

① Part of Speech Tagging

② Sequential Modeling

③ Evaluation

Evaluating POS Taggers

50

- Comparing the output of a tagger with a human-labelled gold standard
- Accuracy:

$$Accuracy = \frac{\text{\#correctly tagged words}}{\text{\#total word token}}$$

Evaluating POS Taggers

51

- Accuracy:

$$Accuracy = \frac{tp}{N}$$

$$Accuracy = \frac{\sum_c^C tp_c}{N}$$

Evaluating POS Taggers

52

- The accuracy score doesn't show everything
- It is useful to know what is misclassified as what
- Solution: providing a confusion matrix
 - A matrix ($\# \text{ tags} \times \# \text{ tags}$): the rows correspond to the correct tags and the columns correspond to the tagger output
 - $Cell(i, j)$ gives the count of the number of times tag i was classified as tag j
 - The leading diagonal elements correspond to correct classifications
 - Off diagonal elements correspond to misclassifications
- A good approach for error analysis

Further Reading

53

- Speech and Language Processing
 - Chapter 5: POS Tagging
 - Chapter 6: MaxEnt & HMM