

Part-of-Speech Tagging

Chapter 8
(8.1-8.4.6)

Outline

- Parts of speech (POS)
- Tagsets
- POS Tagging
 - Rule-based tagging
 - Probabilistic (HMM) tagging

Garden Path Sentences

- The old dog the footsteps of the young

Parts of Speech

- Traditional parts of speech
 - Noun, verb, adjective, preposition, adverb, article, interjection, pronoun, conjunction, etc
 - Called: parts-of-speech, lexical categories, word classes, morphological classes, lexical tags...
 - Lots of debate within linguistics about the number, nature, and universality of these
 - We'll completely ignore this debate.

Parts of Speech

- Traditional parts of speech
 - ~ 8 of them

Parts of Speech

<p>NOUN</p> <p>Name of a person, place, thing or idea.</p> <p>Examples: Daniel, London, table, hope</p> <p>- Mary uses a blue pen for her notes.</p>	<p>PRONOUN</p> <p>A pronoun is used in place of a noun or noun phrase to avoid repetition.</p> <p>Examples: I, you, it, we, us, them, those</p> <p>- I want her to dance with me.</p>
<p>ADJECTIVE</p> <p>Describes, modifies or gives more information about a noun or pronoun.</p> <p>Examples: cold, happy, young, two, fun</p> <p>- The little girl has a pink hat.</p>	<p>VERB</p> <p>Shows an action or a state of being.</p> <p>Examples: go, speak, eat, live, are, is</p> <p>- I listen to the word and then repeat it.</p>
<p>ADVERB</p> <p>Modifies a verb, an adjective or another adverb. It tells how (often), where, when.</p> <p>Examples: slowly, very, always, well, too</p> <p>- Yesterday, I ate my lunch quickly.</p>	<p>PREPOSITION</p> <p>Shows the relationship of a noun or pronoun to another word.</p> <p>Examples: at, on, in, from, with, about</p> <p>- I left my keys on the table for you.</p>
<p>CONJUNCTION</p> <p>Joins two words, ideas, phrases together and shows how they are connected.</p> <p>Examples: and, or, but, because, yet, so</p> <p>- I was hot and tired but still finished it.</p>	<p>INTERJECTION</p> <p>A word or phrase that expresses a strong emotion. It is a short exclamation.</p> <p>Examples: Ouch! Hey! Oh! Watch out!</p> <p>- Wow! I passed my English exam.</p>

www.grammar4d.com www.woodwardenglish.com www.vocabulary4d.com

5

POS examples

- N noun *chair, bandwidth, pacing*
- V verb *study, debate, munch*
- ADJ adjective *purple, tall, ridiculous*
- ADV adverb *unfortunately, slowly*
- P preposition *of, by, to*
- PRO pronoun *I, me, mine*
- DET determiner *the, a, that, those*

1/21/2020

Speech and Language Processing - Jurafsky and Martin

6

POS Tagging

- The process of assigning a part-of-speech or lexical class marker to each word in a collection.

WORD tag

the
koala
put
the
keys
on
the
table

1/21/2020

Speech and Language Processing - Jurafsky and Martin

7

POS Tagging

- The process of assigning a part-of-speech or lexical class marker to each word in a collection.

WORD tag

the **DET**
koala
put
the
keys
on
the
table

1/21/2020

Speech and Language Processing - Jurafsky and Martin

8

POS Tagging

- The process of assigning a part-of-speech or lexical class marker to each word in a collection.

WORD	tag
the	DET
koala	N
put	
the	
keys	
on	
the	
table	

1/21/2020

Speech and Language Processing - Jurafsky and Martin

9

POS Tagging

- The process of assigning a part-of-speech or lexical class marker to each word in a collection.

WORD	tag
the	DET
koala	N
put	V
the	
keys	
on	
the	
table	

1/21/2020

Speech and Language Processing - Jurafsky and Martin

10

POS Tagging

- The process of assigning a part-of-speech or lexical class marker to each word in a collection.

WORD	tag
the	DET
koala	N
put	V
the	DET
keys	
on	
the	
table	

1/21/2020

Speech and Language Processing - Jurafsky and Martin

11

POS Tagging

- The process of assigning a part-of-speech or lexical class marker to each word in a collection.

WORD	tag
the	DET
koala	N
put	V
the	DET
keys	N
on	
the	
table	

1/21/2020

Speech and Language Processing - Jurafsky and Martin

12

POS Tagging

- The process of assigning a part-of-speech or lexical class marker to each word in a collection.

WORD	tag
the	DET
koala	N
put	V
the	DET
keys	N
on	P
the	
table	

1/21/2020

Speech and Language Processing - Jurafsky and Martin

13

POS Tagging

- The process of assigning a part-of-speech or lexical class marker to each word in a collection.

WORD	tag
the	DET
koala	N
put	V
the	DET
keys	N
on	P
the	DET
table	

1/21/2020

Speech and Language Processing - Jurafsky and Martin

14

POS Tagging

- The process of assigning a part-of-speech or lexical class marker to each word in a collection.

WORD	tag
the	DET
koala	N
put	V
the	DET
keys	N
on	P
the	DET
table	N

1/21/2020

Speech and Language Processing - Jurafsky and Martin

15

Why is POS Tagging Useful?

- First step of many practical tasks, e.g.
- Speech synthesis (aka text to speech)
 - How to pronounce "lead"?
 - OBject obJECT
 - CONtent conTENT
- Parsing
 - Need to know if a word is an N or V before you can parse
- Information extraction
 - Finding names, relations, etc.
- Language modeling
 - Backoff

1/21/2020

Speech and Language Processing - Jurafsky and Martin

16

Why is POS Tagging Difficult?

- Words often have more than one POS:
back
 - The back door = adjective
 - On my back =
 - Win the voters back =
 - Promised to back the bill =

Why is POS Tagging Difficult?

- Words often have more than one POS:
back
 - The back door = adjective
 - On my back = noun
 - Win the voters back =
 - Promised to back the bill =

Why is POS Tagging Difficult?

- Words often have more than one POS:
back
 - The *back* door = adjective
 - On my *back* = noun
 - Win the voters *back* = adverb
 - Promised to *back* the bill =

Why is POS Tagging Difficult?

- Words often have more than one POS:
back
 - The *back* door = adjective
 - On my *back* = noun
 - Win the voters *back* = adverb
 - Promised to *back* the bill = verb
- The POS tagging problem is to determine the POS tag for a particular instance of a word.

POS Tagging

- Input: Plays well with others
- Ambiguity: NNS/VBZ UH/JJ/NN/RB IN NNS
- Output: Plays/VBZ well/RB with/IN others, NNS

Penn
Treebank
POS tags

POS tagging performance

- How many tags are correct? (Tag accuracy)
 - About 97% currently
 - But baseline is already 90%
 - Baseline is performance of stupidest possible method
 - Tag every word with its most frequent tag
 - Tag unknown words as nouns
 - Partly easy because
 - Many words are unambiguous
 - You get points for them (*the*, *a*, etc.) and for punctuation marks!

Deciding on the correct part of speech can be difficult even for people

- Mrs/NNP Shaefer/NNP never/RB got/VBD
around/RP to/TO joining/VBG
- All/DT we/PRP gotta/VBN do/VB is/VBZ
go/VB around/IN the/DT corner/NN
- Chateau/NNP Petrus/NNP costs/VBZ
around/RB 250/CD

How difficult is POS tagging?

- About 11% of the word types in the Brown corpus are ambiguous with regard to part of speech
- But they tend to be very common words.
E.g., *that*
 - I know *that* he is honest = IN
 - Yes, *that* play was nice = DT
 - You can't go *that* far = RB
- 40% of the word tokens are ambiguous

Open vs. Closed Classes

- **Closed class:** *why?*
 - Determiners: a, an, the
 - Prepositions: of, in, by, ...
 - Auxiliaries: may, can, will had, been, ...
 - Pronouns: I, you, she, mine, his, them, ...
 - Usually **function words** (short common words which play a role in grammar)
- **Open class:** *why?*
 - English has 4: Nouns, Verbs, Adjectives, Adverbs
 - Many languages have these 4, but not all!

1/21/2020

Speech and Language Processing - Jurafsky and Martin

25

Open vs. Closed Classes

- **Closed class:** a small fixed membership
 - Determiners: a, an, the
 - Prepositions: of, in, by, ...
 - Auxiliaries: may, can, will had, been, ...
 - Pronouns: I, you, she, mine, his, them, ...
 - Usually **function words** (short common words which play a role in grammar)
- **Open class:** new ones can be created all the time
 - English has 4: Nouns, Verbs, Adjectives, Adverbs
 - Many languages have these 4, but not all!

1/21/2020

Speech and Language Processing - Jurafsky and Martin

26

Open Class Words

▪ Nouns

- Proper nouns (Pittsburgh, Pat Gallagher)
 - English capitalizes these.
- Common nouns (the rest).
- Count nouns and mass nouns
 - Count: have plurals, get counted: goat/goats, one goat, two goats
 - Mass: don't get counted (snow, salt, communism) (*two snows)

▪ Adverbs: tend to modify things

- *Unfortunately*, John walked home *extremely slowly yesterday*
- Directional/locative adverbs (here, home, downhill)
- Degree adverbs (extremely, very, somewhat)
- Manner adverbs (slowly, slinkily, delicately)

▪ Verbs

- In English, have morphological affixes (eat/eats/eaten)

1/21/2020

Speech and Language Processing - Jurafsky and Martin

27

Closed Class Words

Examples:

- prepositions: *on, under, over, ...*
- particles: *up, down, on, off, ...*
- determiners: *a, an, the, ...*
- pronouns: *she, who, I, ..*
- conjunctions: *and, but, or, ...*
- auxiliary verbs: *can, may should, ...*
- numerals: *one, two, three, third, ...*

1/21/2020

Speech and Language Processing - Jurafsky and Martin

28

Prepositions from CELEX

of	540,085	through	14,964	worth	1,563	pace	12
in	331,235	after	13,670	toward	1,390	nigh	9
for	142,421	between	13,275	plus	750	re	4
to	125,691	under	9,525	till	686	mid	3
with	124,965	per	6,515	amongst	525	o'er	2
on	109,129	among	5,090	via	351	but	0
at	100,169	within	5,030	amid	222	ere	0
by	77,794	towards	4,700	underneath	164	less	0
from	74,843	above	3,056	versus	113	midst	0
about	38,428	near	2,026	amidst	67	o'	0
than	20,210	off	1,695	sans	20	thru	0
over	18,071	past	1,575	circa	14	vice	0

1/21/2020

Speech and Language Processing - Jurafsky and Martin

29

POS Tagging Choosing a Tagset

- There are so many parts of speech, potential distinctions we can draw
- To do POS tagging, we need to choose a standard set of tags to work with
- Could pick very coarse tagsets
 - N, V, Adj, Adv.
- More commonly used set is finer grained, the "Penn TreeBank tagset", 45 tags
- Even more fine-grained tagsets exist

1/21/2020

Speech and Language Processing - Jurafsky and Martin

30

Penn TreeBank POS 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>
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>			

1/21/2020

Speech and Language Processing - Jurafsky and Martin

31

Using the Penn Tagset

- The/? grand/? jury/? commmented/? on/? a/? number/? of/? other/? topics/? ./?

1/21/2020

Speech and Language Processing - Jurafsky and Martin

32

Using the Penn Tagset

- The/DT grand/JJ jury/NN commented/VBD on/IN a/DT number/NN of/IN other/JJ topics/NNS ./.

Recall POS Tagging Difficulty

- Words often have more than one POS:
back
 - The *back* door = JJ
 - On my *back* = NN
 - Win the voters *back* = RB
 - Promised to *back* the bill = VB
- The POS tagging problem is to determine the POS tag for a particular instance of a word.

These examples from Dekang Lin

How Hard is POS Tagging? Measuring Ambiguity

	87-tag Original Brown	45-tag Treebank Brown
Unambiguous (1 tag)	44,019	38,857
Ambiguous (2–7 tags)	5,490	8844
Details:		
2 tags	4,967	6,731
3 tags	411	1621
4 tags	91	357
5 tags	17	90
6 tags	2 (<i>well, beat</i>)	32
7 tags	2 (<i>still, down</i>)	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>)

1/21/2020

Speech and Language Processing - Jurafsky and Martin

35

Tagging Whole Sentences with POS is Hard too

- Ambiguous POS contexts
 - E.g., *Time flies like an arrow.*
- Possible POS assignments
 - *Time/[V,N] flies/[V,N] like/[V,Prep] an/Det arrow/N*
 - *Time/N flies/V like/Prep an/Det arrow/N*
 - *Time/V flies/N like/Prep an/Det arrow/N*
 - *Time/N flies/N like/V an/Det arrow/N*
 -

36

How Do We Disambiguate POS?

- Many words have only one POS tag (e.g. **is, Mary, smallest**)
- Others have a single **most likely** tag (e.g. **Dog** is less used as a V)
- Tags also tend to *co-occur* regularly with other tags (e.g. Det, N)
- In addition to conditional probabilities of words $P(w_1|w_{n-1})$, we can look at POS likelihoods $P(t_1|t_{n-1})$ to disambiguate sentences and to assess sentence likelihoods

37

More and Better Features → Feature-based tagger

- Can do surprisingly well just looking at a word by itself:
 - Word the: the → DT
 - Lowercased word Importantly: importantly → RB
 - Prefixes unfathomable: un- → JJ
 - Suffixes Importantly: -ly → RB
 - Capitalization Meridian: CAP → NNP
 - Word shapes 35-year: d-x → JJ

Overview: POS Tagging Accuracies

- Rough accuracies:

- Most freq tag:

~90% / ~50%

Most errors
on unknown
words

- Trigram HMM:

~95% / ~55%

- Maxent P(t|w):

93.7% / 82.6%

- Upper bound:

~98% (human)