

Sequence Tagging

- **Today**
 - Part-of-speech tagging
 - Introduction

Part of speech tagging

“There are 10 parts of speech, and they are all troublesome.”

-Mark Twain

- POS tags are also known as word classes, morphological classes, or lexical tags.
- Typically much larger than Twain's 10:
 - Penn Treebank: 45
 - Brown corpus: 87
 - C7 tagset: 146

Part of speech tagging

- **Assign the correct part of speech (word class) to each word/token in a document**

“The/DT planet/NN Jupiter/NNP and/CC its/PPS moons/NNS are/VBP in/IN effect/NN a/DT mini-solar/JJ system/NN ,/, and/CC Jupiter/NNP itself/PRP is/VBZ often/RB called/VBN a/DT star/NN that/IN never/RB caught/VBN fire/NN ./.”

- **Needed as an initial processing step for a number of language technology applications**
 - Answer extraction in Question Answering systems
 - Base step in identifying syntactic phrases for IR systems
 - Critical for word-sense disambiguation
 - Information extraction
 - ...

Why is p-o-s tagging hard?

- **Ambiguity**
 - He will **race**/VB the car.
 - When will the **race**/NOUN end?
 - The boat **floated**/ VBD.
 - The boat **floated**/ VBD down Fall Creek.
 - The boat **floated**/**VBN**down Fall Creek sank.
- **Average of ~2 parts of speech for each word**
- **The number of tags used by different systems varies a lot. Some systems use < 20 tags, while others use > 400.**

Hard for Humans

- **particle vs. preposition**
 - He talked *over* the deal.
 - He talked *over* the telephone.
- **past tense vs. past participle**
 - The horse *walked* past the barn.
 - The horse *walked* past the barn fell.
- **noun vs. adjective?**
 - The *executive* decision.
- **noun vs. present participle**
 - *Fishing* can be fun.

To obtain gold standards for evaluation, annotators rely on a set of tagging guidelines.

From Ralph Grishman, NYU

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>	WPS	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>			

Let's give it a try...

P-o-s tagging exercise

1. It is a nice night.

It/~~PRP~~ is/VBZ a/DT nice/JJ night/NN ./.

5. . . . I am sitting in Mindy's restaurant putting on the gefillte fish, which is a dish I am very fond of, . . .

. . . I/~~PRP~~ am/VBP sitting/VBG in/IN Mindy/NNP 's/POS

restaurant/NN putting/VBG on/RP the/DT gefillte/NN

fish/NN ,/, which/WDT is/VBZ a/DT dish/NN I/~~PRP~~ am/VBP

very/RB fond/JJ of/RP ,/, . . .

Think buffalo

buffalo buffalo buffalo buffalo buffalo buffalo
buffalo.

Buffalo buffalo Buffalo buffalo buffalo buffalo Buffalo
buffalo.

Buffalo buffalo, Buffalo buffalo buffalo, buffalo Buffalo
buffalo.



Think buffalo

n1. the city of Buffalo, NY

n2. an animal...the American bison

v. to bully, confuse, deceive, or intimidate

Buffaloⁿ¹ buffaloⁿ² Buffaloⁿ¹ buffaloⁿ² buffalo^v buffalo^v Buffaloⁿ¹ buffaloⁿ².

[Those] (Buffalo buffalo) [whom] (Buffalo buffalo) buffalo, buffalo (Buffalo buffalo).

[Those] buffalo(es) from Buffalo [that are intimidated by] buffalo(es) from Buffalo intimidate buffalo(es) from Buffalo.

Bison from Buffalo, New York, who are intimidated by other bison in their community, also happen to intimidate other bison in their community.

THE buffalo FROM Buffalo WHO ARE buffaloes BY buffalo FROM Buffalo, buffalo (verb) OTHER buffalo FROM Buffalo.

Among easiest of NLP problems

- **State-of-the-art methods achieve ~97% accuracy.**
- **Simple heuristics can go a long way.**
 - ~90% accuracy just by choosing the most frequent tag for a word (MLE)
 - To improve reliability: *need to use some of the local context.*
- **But defining the rules for special cases can be time-consuming, difficult, and prone to errors and omissions**

Approaches

1. **rule-based**: involve a large database of hand-written disambiguation rules, e.g. that specify that an ambiguous word is a noun rather than a verb if it follows a determiner.
2. **learning-based**: resolve tagging ambiguities by using a training corpus to compute the probability of a given word having a given tag in a given context.
 - HMM tagger
3. **hybrid ML-/rule-based**: E.g. transformation-based tagger (Brill tagger); learns symbolic rules based on a corpus.
4. **ensemble methods**: combine the results of multiple taggers.

- **Today**

- Part-of-speech tagging

- HMM' s for p-o-s tagging

HMM p-o-s Tagger

Given $W = w_1, \dots, w_n$, find $T = t_1, \dots, t_n$ that maximizes

$$P(t_1, \dots, t_n | w_1, \dots, w_n)$$

Restate using Bayes' rule:

$$(P(t_1, \dots, t_n) * P(w_1, \dots, w_n | t_1, \dots, t_n)) / P(w_1, \dots, w_n)$$

Ignore denominator...

Make independence assumptions...

$P(t_1, \dots, t_n)$: approximate using n-gram model

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

trigram $\prod_{i=1, n} P(t_i | t_{i-2}t_{i-1})$

$P(w_1, \dots, w_n | t_1, \dots, t_n)$: approximate by assuming that a word appears in a category independent of its neighbors

$$\prod_{i=1, n} P(w_i | t_i)$$

Assuming bigram model:

$$P(t_1, \dots, t_n) * P(w_1, \dots, w_n | t_1, \dots, t_n) \approx$$

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

↑
transition
probabilities

↙
lexical generation
probabilities