

# POS Tagging

(materials mostly from Jurafsky & Martin(2009))

# Parts of Speech

- Perhaps starting with Aristotle in the West (384–322 BCE), there was the idea of having parts of speech
  - a.k.a lexical categories, word classes, “tags”, POS
- It comes from Dionysius Thrax of Alexandria (c. 100 BCE) the idea that is still with us that there are 8 parts of speech

---

# Parts of Speech

- But actually his 8 aren't exactly the ones we are taught today
  - Thrax: noun, verb, article, adverb, preposition, conjunction, participle, pronoun
  - School grammar: noun, verb, adjective, adverb, preposition, conjunction, pronoun, interjection

## Open class (lexical) words

### Nouns

#### Proper

*IBM*  
*Italy*

#### Common

*cat / cats*  
*snow*

### Verbs

#### Main

*see*  
*registered*

Adjectives *old older oldest*

Adverbs *slowly*

#### Numbers

*122,312*  
*one*

*... more*

## Closed class (functional)

Determiners *the some*

Conjunctions *and or*

Pronouns *he its*

#### Modals

*can*  
*had*

Prepositions *to with*

Particles *off up*

*... more*

Interjections *Ow Eh*

# Open vs. Closed classes

- Open vs. Closed classes
  - Closed:
    - determiners: *a, an, the*
    - pronouns: *she, he, I*
    - prepositions: *on, under, over, near, by, ...*
    - Why “closed”?
  - Open:
    - Nouns, Verbs, Adjectives, Adverbs.

# Penn Tagsets

TAG	DESCRIPTION	EXAMPLE
CC	conjunction, coordinating	<i>and, or, but</i>
CD	cardinal number	<i>five, three, 13%</i>
DT	determiner	<i>the, a, these</i>
EX	existential there	<i><u>there</u> were six boys</i>
FW	foreign word	<i>mais</i>
IN	conjunction, subordinating or preposition	<i>of, on, before, unless</i>
JJ	adjective	<i>nice, easy</i>
JJR	adjective, comparative	<i>nicer, easier</i>
JJS	adjective, superlative	<i>nicest, easiest</i>

LS	list item marker	
MD	verb, modal auxillary	<i>may, should</i>
NN	noun, singular or mass	<i>tiger, chair, laughter</i>
NNS	noun, plural	<i>tigers, chairs, insects</i>
NNP	noun, proper singular	<i>Germany, God, Alice</i>
NNPS	noun, proper plural	<i>we met two <u>Christmases</u> ago</i>
PDT	predeterminer	<i><u>both</u> his children</i>
POS	possessive ending	<i>'s</i>
PRP	pronoun, personal	<i>me, you, it</i>
PRP\$	pronoun, possessive	<i>my, your, our</i>

<b>RB</b>	<b>adverb</b>	<i><b>extremely, loudly, hard</b></i>
<b>RBR</b>	<b>adverb, comparative</b>	<i><b>better</b></i>
<b>RBS</b>	<b>adverb, superlative</b>	<i><b>best</b></i>
<b>RP</b>	<b>adverb, particle</b>	<i><b>about, off, up</b></i>
<b>SYM</b>	<b>symbol</b>	<i><b>%</b></i>
<b>TO</b>	<b>infinitival to</b>	<i><b>what <u>to</u> do?</b></i>
<b>UH</b>	<b>interjection</b>	<i><b>oh, oops, gosh</b></i>



<b>VB</b>	verb, base form	<i>think</i>
<b>VBZ</b>	verb, 3rd person singular present	<i>she <u>thinks</u></i>
<b>VBP</b>	verb, non-3rd person singular present	<i>I <u>think</u></i>
<b>VBD</b>	verb, past tense	<i>they <u>thought</u></i>
<b>VCN</b>	verb, past participle	<i>a <u>sunken</u> ship</i>
<b>VBG</b>	verb, gerund or present participle	<i><u>thinking</u> is fun</i>
<b>WDT</b>	<i>wh</i> -determiner	<i>which, whatever, whichever</i>
<b>WP</b>	<i>wh</i> -pronoun, personal	<i>what, who, whom</i>
<b>WP\$</b>	<i>wh</i> -pronoun, possessive	<i>whose, whosever</i>
<b>WRB</b>	<i>wh</i> -adverb	<i>where, when</i>

.	punctuation mark, sentence closer	.;?*
,	punctuation mark, comma	,
:	punctuation mark, colon	:
(	contextual separator, left paren	(
)	contextual separator, right paren	)

# POS Tagging

- 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.

# 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
- Uses:
  - Text-to-speech (how do we pronounce “lead”?)
  - Can write regexps like (Det) Adj\* N+ over the output for phrases, etc.
  - As input to or to speed up a full parser
  - If you know the tag, you can back off to it in other tasks

# 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

# Sources of information

- What are the main sources of information for POS tagging?
  - Knowledge of neighboring words
    - Bill    saw    that   man yesterday
    - NNP NN        DT    NN    NN
    - VB    VB(D) IN     VB    NN
  - Knowledge of word probabilities
    - *man* is rarely used as a verb....
- The latter proves the most useful, but the former also helps



# 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
- Then build a maxent (or whatever) model to predict tag
  - Maxent  $P(t|w)$ :      93.7% overall / 82.6% unknown

# Overview: POS Tagging Accuracies

- Rough accuracies:
  - Most freq tag: ~90% / ~50%
  - Trigram HMM: ~95% / ~55%
  - Maxent  $P(t|w)$ : 93.7% / 82.6%
  - TnT (HMM++): 96.2% / 86.0%
  - MEMM tagger: 96.9% / 86.9%
  - Bidirectional dependencies: 97.2% / 90.0%
  - Upper bound: ~98% (human agreement)

# How to improve supervised results?

- Build better features!

RB

PRP VBD IN RB IN PRP VBD .

They left as soon as he arrived .

- We could fix this with a feature that looked at the next word

JJ

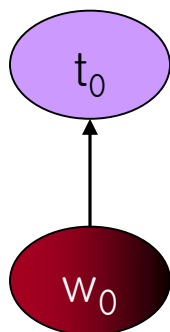
NNP NNS VBD VBN .

Intrinsic flaws remained undetected .

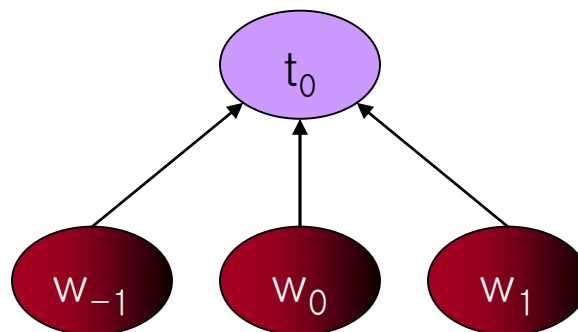
- We could fix this by linking capitalized words to their lowercase versions

# Tagging Without Sequence Information

Baseline



Three Words



Model	Features	Token	Unknown	Sentence
Baseline	56,805	<b>93.69%</b>	82.61%	26.74%
3Words	239,767	<b>96.57%</b>	86.78%	48.27%

Using words only in a straight classifier works as well as a basic (HMM or discriminative) sequence model!!

# POS Tagging with NLTK

## 1. Tokenization

- From `nltk.tokenize` import `word_tokenize`
- `text="Intelligence obtained in recent weeks found that an al Qaeda affiliate was perfecting techniques for hiding explosives in batteries and battery compartments of electronic devices, according to a US official"`
- `tokens=nltk.word_tokenize(text)`

## 2. POS Tagging

- `tagged=nltk.pos_tag(tokens)`

---

# Exercise

- From the example sentence 'text', find all the nouns.
  - Use for-loop
  - Use POS\_tagger