

# NLP Applications

## Lecture 5: Tagging with NLTK

Claire Gardent

CNRS/LORIA  
Campus Scientifique,  
BP 239,  
F-54 506 Vandœuvre-lès-Nancy, France

2007/2008

### Parts of speech

- ▶ Words can be divided into classes that behave similarly.
- ▶ Traditionally eight parts of speech: noun, verb, pronoun, preposition, adverb, conjunction, adjective and article.
- ▶ More recently larger sets have been used: eg Penn Treebank (45 tags), Susanne (353 tags).

### Parts of Speech

Constructing simple taggers in NLTK

Evaluating taggers

Other approaches to POS Tagging

The Brill Tagger

Summary

### Parts of Speech

What use are parts of speech?

They tell us a lot about a word (and the words near it).

- ▶ Tell us what words are likely to occur in the neighbourhood (eg adjectives often followed by nouns, personal pronouns often followed by verbs, possessive pronouns by nouns)
- ▶ Pronunciations can be dependent on part of speech, eg object, content, discount (useful for speech synthesis and speech recognition)
- ▶ Can help information retrieval and extraction (stemming, partial parsing)
- ▶ Useful component in many NLP systems

## Closed and open classes

- ▶ Parts of speech may be categorised as *open* or *closed* classes
- ▶ Closed classes have a fixed membership of words (more or less), eg determiners, pronouns, prepositions
- ▶ Closed class words are usually *function words* — frequently occurring, grammatically important, often short (eg of,it,the,in)
- ▶ The major open classes are *nouns*, *verbs*, *adjectives* and *adverbs*

## Closed classes in English

**prepositions** on, under, over, to, with, by

**determiners** the, a, an, some

**pronouns** she, you, I, who

**conjunctions** and, but, or, as, when, if

**auxiliary verbs** can, may, are

**particles** up, down, at, by

**numerals** one, two, first, second

## Open classes

**nouns** Proper nouns (Scotland, BBC),  
common nouns:

- ▶ count nouns (goat, glass)
- ▶ mass nouns (snow, pacifism)

**verbs** actions and processes (run, hope), also auxiliary verbs

**adjectives** properties and qualities (age, colour, value)

**adverbs** modify verbs, or verb phrases, or other adverbs:  
*Unfortunately John walked home extremely slowly yesterday*

## The Penn Treebank tagset (1)

CC	Coord Conjunction	<i>and, but, or</i>	NN	Noun, sing. or mass	<i>dog</i>
CD	Cardinal number	<i>one, two</i>	NNS	Noun, plural	<i>dogs</i>
DT	Determiner	<i>the, some</i>	NNP	Proper noun, sing.	<i>Edinburgh</i>
EX	Existential there	<i>there</i>	NNPS	Proper noun, plural	<i>Orkneys</i>
FW	Foreign Word	<i>mon dieu</i>	PDT	Predeterminer	<i>all, both</i>
IN	Preposition	<i>of, in, by</i>	POS	Possessive ending	<i>'s</i>
JJ	Adjective	<i>big</i>	PP	Personal pronoun	<i>I, you, she</i>
JJR	Adj., comparative	<i>bigger</i>	PP\$	Possessive pronoun	<i>my, one's</i>
JJS	Adj., superlative	<i>biggest</i>	RB	Adverb	<i>quickly</i>
LS	List item marker	<i>1, One</i>	RBR	Adverb, comparative	<i>faster</i>
MD	Modal	<i>can, should</i>	RBS	Adverb, superlative	<i>fastest</i>

## The Penn Treebank tagset (2)

RP	Particle	<i>up, off</i>	WP\$	Possessive-Wh	<i>whose</i>
SYM	Symbol	<i>+, %, &amp;</i>	WRB	Wh-adverb	<i>how, where</i>
TO	"to"	<i>to</i>	\$	Dollar sign	<i>\$</i>
UH	Interjection	<i>oh, oops</i>	#	Pound sign	<i>#</i>
VB	verb, base form	<i>eat</i>	"	Left quote	<i>' , "</i>
VBD	verb, past tense	<i>ate</i>	"	Right quote	<i> ; "</i>
VBG	verb, gerund	<i>eating</i>	(	Left paren	<i>(</i>
VBN	verb, past part	<i>eaten</i>	)	Right paren	<i>)</i>
VBP	Verb, non-3sg, pres	<i>eat</i>	,	Comma	<i>,</i>
VBZ	Verb, 3sg, pres	<i>eats</i>	.	Sent-final punct	<i>. ! ?</i>
WDT	Wh-determiner	<i>which, that</i>	:	Mid-sent punct.	<i>: ; — ...</i>
WP	Wh-pronoun	<i>what, who</i>			

## The Brown tagset

ap	determiner-pronoun, post-determiner	many other next
at	article	the an no a every
cc	conjunction, coordinating	and or but plus
cs	conjunction, subordinating	that as after whether
in	preposition	in for by
md	modal auxiliary	may might will
pn	pronoun, nominal	something, none
ppl	pronoun, singular, reflexive	himself myself
pp\$	determiner, possessive	its his their my
pp\$\$	pronoun, possessive	mine his hers theirs yours
pps	pronoun, personal, nom, 3rd pers sng	he she
ppss	pronoun, personal, nom, not 3rd pers sng	we I you
wdt	WH-determiner	what whatever
wps	WH-pronoun, nominative	who whoever

## Tagging

- ▶ Definition: Tagging is the assignment of a single part-of-speech tag to each word (and punctuation marker) in a corpus.
- ▶ For example:  
 “/“ The/DT guys/NNS that/WDT make/VBP traditional/JJ hardware/NN are/VBP really/RB being/VBG obsoleted/VBN by/IN microprocessor-based/JJ machines/NNS ,/, ”/” said/VBD Mr./NNP Benton/NNP ./.

## Tagging

- ▶ Non-trivial: POS tagging must resolve ambiguities since the same word can have different tags in different contexts
- ▶ In the Brown corpus 11.5% of word types and 40% of word tokens are ambiguous
- ▶ In many cases one tag is much more likely for a given word than any other
- ▶ Limited scope: only supplying a tag for each word, no larger structures created (eg prepositional phrase attachment)

What information can help decide the correct PoS tag for a word?

**Other PoS tags** Even though the PoS tags of other words may be uncertain too, we can use information that some tag sequences are more likely than others (eg *the/AT red/JJ drink/NN* vs *the/AT red/JJ drink/VBP*). Using *only* information about the most likely PoS tag sequence does not result in an accurate tagger (about 77% correct)

**The word identity** Many words can have multiple possible tags, but some are more likely than others (eg *fall/VBP* vs *fall/NN*)  
Tagging each word with its most common tag results in a tagger with about 90% accuracy

NLTK provides several means of developing a tagger:

- ▶ Assigning a default tag
- ▶ Using regular expressions on strings (the shape of the string helps guessing the associated PoS tag)
- ▶ Unigram tagging: assigning the most probable tag
- ▶ Bigram tagging: assigning the most probable tag given a left-adjacent PoS
- ▶ Brill tagging: transformation based learning of tags

The simplest possible tagger tags everything as a noun:

```
>>> tokens = 'John saw 3 polar bears .'.split()
>>> default_tagger = nltk.DefaultTagger('NN')
>>> default_tagger.tag(tokens)
[('John', 'NN'), ('saw', 'NN'), ('3', 'NN'), ('polar', 'NN'), ('bears', 'NN'), ('.', 'NN')]
```

The regular expression tagger assigns tags to tokens on the basis of matching patterns in the token's text. For instance, the following tagger assigns "cd" to cardinal numbers, and "nn" to everything else:

```
>>> default_pattern = (r'.*', 'NN')
>>> cd_pattern = (r'^[0-9]+(.[0-9]+)?$', 'CD')
>>> patterns = [cd_pattern, default_pattern]
>>> tokens = 'There are 11 players in a football team'.split()
['There', 'are', '11', 'players', 'in', 'a', 'football', 'team']
>>> NN_CD_tagger = nltk.RegexpTagger(patterns)
>>> NN_CD_tagger.tag(tokens)
[('There', 'NN'), ('are', 'NN'), ('11', 'CD'), ('players', 'NN'), ('in', 'NN'), ('a', 'NN'), ('football', 'NN'), ('team', 'NN')]
```

## A more sophisticated regular expression tagger

Guess the correct tag for words based on the presence of certain prefix or suffix strings.

```
>>> patterns = [
...     (r'.*ing$', 'VBG'),           # gerunds
...     (r'.*ed$', 'VBD'),           # simple past
...     (r'.*es$', 'VBZ'),           # 3rd singular present
...     (r'.*ould$', 'MD'),          # modals
...     (r'.*\ 's$', 'NN$'),         # possessive nouns
...     (r'.*s$', 'NNS'),            # plural nouns
...     (r'^-?[0-9]+(\.[0-9]+)?$', 'CD'), # cardinal numbers
...     (r'.*', 'NN')                # nouns (default)
... ]
>>> regexp_tagger = nltk.RegexpTagger(patterns)
>>> regexp_tagger.tag(nltk.corpus.brown.sents('a')[3])
[('‘‘', 'NN'), ('Only', 'NN'), ('a', 'NN'), ('relative', 'NN'),
 ('handful', 'NN'), ('of', 'NN'), ('such', 'NN'), ('reports', 'NN'),
 ('was', 'NNS'), ('received', 'VBD'), ('"', 'NN'), (',', 'NN'),
 ('the', 'NN'), ('jury', 'NN'), ('said', 'NN'), (',', 'NN'), ('‘‘',
 ('considering', 'VBG'), ('the', 'NN'), ('widespread', 'NN'),
```

17 / 38

## Unigram taggers

- ▶ A UnigramTagger assigns to each token its most likely tag
- ▶ A UnigramTagger is based on a table of unigram probabilities
- ▶ This table is derived from an annotated corpus
- ▶ Before a “UnigramTagger“ can be used to tag data, it must be trained on a \*training corpus\*. This corpus is used to determine which tags are most common for each word.
- ▶ Annotated data can also be used to evaluate the tagger

18 / 38

## A unigram tagger

The NLTK UnigramTagger class implements a unigram tagging algorithm

```
>>> brown_a = nltk.corpus.brown.tagged_sents('a')
>>> unigram_tagger = nltk.UnigramTagger(brown_a)
>>> sent = nltk.corpus.brown.sents('a')[2007]
>>> unigram_tagger.tag(sent)
[('Various', None), ('of', 'IN'), ('the', 'AT'), ('apartments',
 ('of', 'IN'), ('the', 'AT'), ('terrace', 'NN'), ('type', 'NN'),
 ('being', 'BEG'), ('on', 'IN'), ('the', 'AT'), ('ground', 'NN'),
 ('so', 'QL'), ('that', 'CS'), ('entrance', 'NN'), ('is', 'BEZ'),
>>> nltk.tag.accuracy(unigram_tagger, brown_a)
0.8550331165343994
```

The UnigramTagger assigns the default tag None to words that are not in the training data (eg *Various*)

19 / 38

## Using a backoff

We can combine taggers to ensure every word is tagged:

```
>>> t0 = nltk.DefaultTagger('NN')
>>> t1 = nltk.UnigramTagger(brown_a, backoff=t0)
```

20 / 38

## Evaluating taggers

Sentence	Gold Standard	Unigram Tagger
The	at	at
President	nn-tl	nn-tl
said	vbd	vbd
he	pps	pps
will	md	md
ask	vb	vb
Congress	np	np
to	to	to
increase	vb	*nn*
grants	nns	nns
to	in	*to*
states	nns	nns
.	.	.

The tagger correctly tagged 10 out of 12 words.

Accuracy = 10/12 or 83%

## Evaluating taggers

- ▶ Basic idea: compare the output of a tagger with a human-labelled *gold standard*
- ▶ Need to compare how well an automatic method does with the agreement between people
- ▶ The best automatic methods have an accuracy of about 96-97% when using the (small) Penn treebank tagset (but this is still an average of one error every couple of sentences...)
- ▶ Inter-annotator agreement is also only about 97%
- ▶ A good unigram baseline (with smoothing) can obtain 90-91%!

## Evaluating taggers in NLTK

NLTK provides a function `nltk.tag.accuracy` to automate evaluation.

```
>>> acc = nltk.tag.accuracy(unigram_tagger, train_sents)
>>> print 'Accuracy = %4.1f%%' % (100 * acc)
Accuracy = 81.8%
```

## Evaluating taggers in NLTK

Training and testing corpus should be distinct:

```
>>> nltk.tag.accuracy(unigram_tagger, train_sents)
93.8%
>>> nltk.tag.accuracy(unigram_tagger, test_sents)
82.8%
```

- ▶ The % correct score doesn't tell you everything — it is useful know what is misclassified as what
- ▶ *Confusion matrix*: A matrix (ntags × ntags) where 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
- ▶ Thus a confusion matrix gives information on the major problems facing a tagger (eg NNP vs. NN vs. JJ)
- ▶ See section 3 of the NLTK tutorial on Tagging

- ▶ Rule-based (briefly discussed in Jurafsky and Martin,, sec 8.4)
- ▶ N-gram / HMM based tagging (eg TnT)
- ▶ Transformation-based learning (the Brill tagger)

## n-gram based tagging

- ▶ Also referred to as hidden Markov model (HMM) tagging
- ▶ Basic idea: Choose the tag that maximises:

$$P(\text{word}|\text{tag}) \cdot P(\text{tag}|\text{previous } n \text{ tags})$$

- ▶ For a bigram model the best tag at position  $i$  is:

$$t_i = \arg \max_{t_j} P(w_i|t_j)P(t_j|t_{i-1})$$

Assuming that you know the previous tag,  $t_{i-1}$ .

- ▶ Interpretation: choose the tag  $t_i$  that is most likely to generate word  $w_i$  given that the previous tag was  $t_{i-1}$

## Example (J+M, p304)

Secretariat/NNP is/VBZ expected/VBZ to/TO race/VB tomorrow/NN  
 People/NNS continue/VBP to/TO inquire/VB the/DT reason/NN  
 for/IN the/DT race/NN for/IN outer/JJ space/NN

- ▶ “race” is a verb in the first, a noun in the second.
- ▶ Assume that race is the only untagged word, so we can assume the tags of the others.
- ▶ Probabilities of “race” being a verb, or race being a noun in the first example:

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

$$P(\text{race is } NN) = P(NN|TO)P(\text{race}|NN)$$

$$P(NN|TO) = 0.021$$

$$P(VB|TO) = 0.34$$

$$P(\text{race}|NN) = 0.00041$$

$$P(\text{race}|VB) = 0.00003$$

$$\begin{aligned} P(\text{race is } VB) &= P(VB|TO)P(\text{race}|VB) \\ &= 0.34 \times 0.00003 = 0.00001 \end{aligned}$$

$$\begin{aligned} P(\text{race is } NN) &= P(NN|TO)P(\text{race}|NN) \\ &= 0.021 \times 0.00041 = 0.000007 \end{aligned}$$

## Brill tagging

- ▶ A rule-based system...
- ▶ ...but the rules are learned from a corpus
- ▶ Transformation-based tagging (TBL)

## Limitation of n-gram taggers

- ▶ Size of n-gram table (*Language model*): hundred of millions of entries
- ▶ Context limited to tags; Information about preceding or neighbouring word is often useful

## Transformation-based tagging

Basic idea:

1. Label every word with its most likely tag (Unigram tagging)
2. Generate finite set of transformation rules (*In context C, replace T1 with T2*)
3. Select best rule i.e., rule that most improve tagging
4. Re-tag the data using this rule

Repeat 1-3 until stopping criteria is met (e.g., insufficient improvement over previous pass).

Output: ordered list of transformations that can be applied to a new corpus; tagging procedure.



- ▶ Parts of speech and tagsets
- ▶ Tagging
- ▶ Constructing simple taggers in NLTK
- ▶ Evaluating taggers
- ▶ Other types of taggers (n-gram, TBL)

- ▶ Jurafsky and Martin, chapter 8 (esp. sec 8.5);
- ▶ Manning and Schütze, chapter 10  
Chris Manning and Hinrich Schütze, Foundations of Statistical Natural Language Processing, MIT Press. Cambridge, MA: May 1999.
- ▶ T. Brants (2000). "TnT – a statistical part-of-speech tagger". In Proceedings of the 6th Applied NLP Conference, ANLP-2000.  
<http://uk.arxiv.org/abs/cs.CL/0003055>
- ▶ Schuetze, H. Distributional Part-of-Speech Tagging. Proceedings of the seventh conference on European chapter of the Association for Computational Linguistics, 1995.  
<ftp://csli.stanford.edu/pub/prosit/DisPosTag.ps>