

### Tagging

Steven Bird   Ewan Klein   Edward Loper

University of Melbourne, AUSTRALIA

University of Edinburgh, UK

University of Pennsylvania, USA

May 16, 2007

- How can we predict the behaviour of a previously unseen word?
- 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

#### 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

<b>prepositions</b>	on, under, over, to, with, by
<b>determiners</b>	the, a, an, some
<b>pronouns</b>	she, you, I, who
<b>conjunctions</b>	and, but, or, as, when, if
<b>auxiliary verbs</b>	can, may, are
<b>particles</b>	up, down, at, by
<b>numerals</b>	one, two, first, second

## Open classes

<b>nouns</b>	Proper nouns ( <b>Scotland</b> , <b>BBC</b> ), common nouns: <ul style="list-style-type: none"> <li>count nouns (<b>goat</b>, <b>glass</b>)</li> <li>mass nouns (<b>snow</b>, <b>pacifism</b>)</li> </ul>
<b>verbs</b>	actions and processes ( <b>run</b> , <b>hope</b> ), also auxiliary verbs
<b>adjectives</b>	properties and qualities (age, colour, value)
<b>adverbs</b>	modify verbs, or verb phrases, or other adverbs: <b>Unfortunately John walked home extremely slowly yesterday</b>

## 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>			

## Tagging

- Definition: POS 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 ./.*
- 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)

## Tagging in NLTK

The simplest possible tagger tags everything as a noun:

```
from nltk_lite import tokenize
text = 'There are 11 players in a football team'
text_tokens = list(tokenize.whitespace(text))
# ['There', 'are', '11', 'players', 'in', 'a', 'football', '

from nltk_lite import tag
mytagger = tag.Default('nn')
for t in mytagger.tag(text_tokens):
    print t
# ('There', 'NN')
# ('are', 'NN')
# ...
```

## Information sources for tagging

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

## A regular expression tagger

We can use regular expressions to tag tokens based on regularities in the text, eg numerals:

```
default_pattern = (r'.*', 'NN')
cd_pattern = (r'^[0-9]+(\.[0-9]+)?$', 'CD')
patterns = [cd_pattern, default_pattern]
NN_CD_tagger = tag.Regexp(patterns)
re_tagged = list(NN_CD_tagger.tag(text_tokens))
# [('There', 'NN'), ('are', 'NN'), ('11', 'NN'), ('players', 'NN'), ('in', 'NN'), ('a', 'NN'), ('football', 'NN'), ('team', 'NN')]
```

## A unigram tagger

The NLTK UnigramTagger class implements a tagging algorithm based on a table of unigram probabilities:

$$\text{tag}(w) = \arg \max_{t_i} P(t_i|w)$$

Training a UnigramTagger on the Penn Treebank:

```
from nltk_lite.corpora import treebank
from itertools import islice

# sentences 0-2999
train_sents = list(islice(treebank.tagged(), 3000))
# from sentence 3000 to the end
test_sents = list(islice(treebank.tagged(), 3000, None))

unigram_tagger = tag.Unigram()
unigram_tagger.train(train_sents)
```

## 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%!

## Unigram tagging

```
>>> list(unigram_tagger.tag(tokenize.whitespace("Mr. Jones
the book on the shelf")))
[('Mr.', 'NNP'), ('Jones', 'NNP'), ('saw', 'VBD'), ('the',
('book', 'NN'), ('on', 'IN'), ('the', 'DT'), ('shelf', Non
```

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

We can combine taggers to ensure every word is tagged:

```
>>> unigram_tagger = tag.Unigram(cutoff=0, backoff=NN_CD_tagger)
>>> unigram_tagger.train(train_sents)
>>> list(unigram_tagger.tag(tokenize.whitespace("Mr. Jones
the book on the shelf")))
[('Mr.', 'NNP'), ('Jones', 'NNP'), ('saw', 'VBD'), ('the',
('book', 'VB'), ('on', 'IN'), ('the', 'DT'), ('shelf', 'NN
```

## Evaluating taggers in NLTK

NLTK provides a function `tag.accuracy` to automate evaluation. It needs to be provided with a tagger, together with some text to be tagged and the gold standard tags.

We can make print more prettily:

```
def print_accuracy(tagger, data):
    print '%3.1f%%' % (100 * tag.accuracy(tagger, data))

>>> print_accuracy(NN_CD_tagger, test_sents)
15.0%
>>> print_accuracy(unigram_tagger, train_sents)
93.8%
>>> print_accuracy(unigram_tagger, test_sents)
82.8%
```

## Error analysis

- The % correct score doesn't tell you everything — it is useful to know what is misclassified as what
- *Confusion matrix*: A matrix (ntags x 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

## N-gram taggers

- Basic idea: Choose the tag that maximises:

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

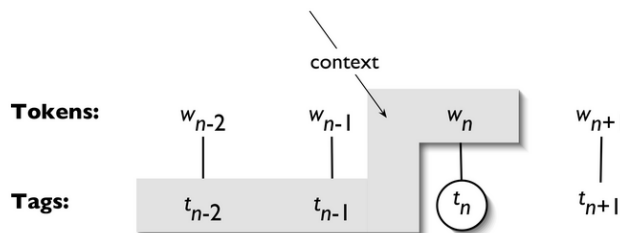
- For a bigram model the best tag at position  $i$  is:

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

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  $t_{i-1}$

## N-gram taggers



## 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:

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

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

## Example (continued)

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

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

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

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

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

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

## Limitation of NLTK n-gram taggers

- Does not find the most likely sequence of tags, simply works left to right always assigning the most probable single tag (given the previous tag assignments)
- Does not cope with zero probability problem well (no smoothing or discounting)
- see module `nltk_lite.tag.hmm`

## Simple bigram tagging in NLTK

```
>>> default_pattern = (r'.*', 'NN')
>>> cd_pattern = (r'^[0-9]+(.[0-9]+)?$', 'CD')
>>> patterns = [cd_pattern, default_pattern]
>>> NN_CD_tagger = tag.Regexp(patterns)
>>> unigram_tagger = tag.Unigram(cutoff=0, backoff=NN_CD_t
>>> unigram_tagger.train(train_sents)
>>> bigram_tagger = tag.Bigram(backoff=unigram_tagger)
>>> bigram_tagger.train(train_sents)

>>> print_accuracy(bigram_tagger, train_sents)
95.6%
>>> print_accuracy(bigram_tagger, test_sents)
84.2%
```

## Brill Tagger

- Problem with n-gram taggers: size
- A rule-based system...
- ...but the rules are learned from a corpus
- Basic approach: start by applying general rules, then successively refine with additional rules that correct the mistakes (painting analogy)
- Learn the rules from a corpus, using a set of rule templates, eg:  
Change tag **a** to **b** when the following word is tagged **z**
- Choose the best rule each iteration

Sentence	Gold	Unigram	Replace nn with vb when the previous word is to	Replace to with in when the next tag is nns
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	vb	
grants	nns	nns		
to	in	to	to	in
states	nns	nns		
for	in	in		
vocational	jj	jj		
rehabilitation	nn	nn		

- **Reading:** Jurafsky and Martin, chapter 8 (esp. sec 8.5); Manning and Schütze, chapter 10;
- Rule-based and statistical tagging
- HMMs and n-grams for statistical tagging
- Operation of a simple bigram tagger
- TnT — an accurate trigram-based tagger