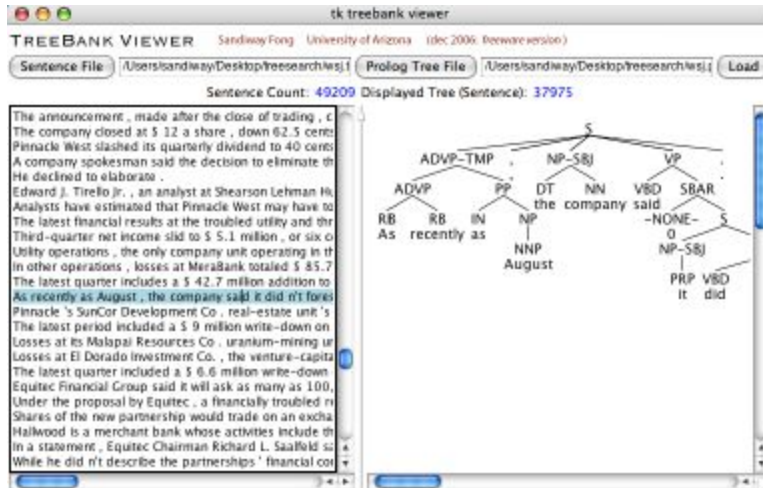# Part-of-speech tagging



A simple but useful form of linguistic analysis

## Open class (lexical) words

### Nouns

**Proper**

*IBM*
*Italy*

**Common**

*cat / cats*
*snow*

### Verbs

**Main**

*see*
*registered*

**Modals**

*can*
*had*

### Adjectives   *old   older   oldest*

### Adverbs   *slowly*

### Numbers

*122,312*
*one*

*… more*

## Closed class (functional)

**Determiners** *the some*

**Conjunctions**   *and or*

**Pronouns**   *he its*

**Prepositions**   *to with*

**Particles**   *off   up*

**Interjections**   *Ow   Eh*

*… more*

# 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.

# 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.

| | Tag | Description | Example |
|---|---|---|---|
| **Open Class** | **ADJ** | Adjective: noun modifiers describing properties | *red, young, awesome* |
| | **ADV** | Adverb: verb modifiers of time, place, manner | *very, slowly, home, yesterday* |
| | **NOUN** | words for persons, places, things, etc. | *algorithm, cat, mango, beauty* |
| | **VERB** | words for actions and processes | *draw, provide, go* |
| | **PROPN** | Proper noun: name of a person, organization, place, etc.. | *Regina, IBM, Colorado* |
| | **INTJ** | Interjection: exclamation, greeting, yes/no response, etc. | *oh, um, yes, hello* |
| **Closed Class Words** | **ADP** | Adposition (Preposition/Postposition): marks a noun's spacial, temporal, or other relation | *in, on, by, under* |
| | **AUX** | Auxiliary: helping verb marking tense, aspect, mood, etc., | *can, may, should, are* |
| | **CCONJ** | Coordinating Conjunction: joins two phrases/clauses | *and, or, but* |
| | **DET** | Determiner: marks noun phrase properties | *a, an, the, this* |
| | **NUM** | Numeral | *one, two, first, second* |
| | **PART** | Particle: a preposition-like form used together with a verb | *up, down, on, off, in, out, at, by* |
| | **PRON** | Pronoun: a shorthand for referring to an entity or event | *she, who, I, others* |
| | **SCONJ** | Subordinating Conjunction: joins a main clause with a subordinate clause such as a sentential complement | *that, which* |
| **Other** | **PUNCT** | Punctuation | *, , ()* |
| | **SYM** | Symbols like $ or emoji | *$, %* |
| | **X** | Other | asdf, qwfg |

5

| Tag | Description | Example | Tag | Description | Example | Tag | Description | Example |
|-----|-------------|---------|-----|-------------|---------|-----|-------------|---------|
| CC | coord. conj. | *and, but, or* | NNP | proper noun, sing. | *IBM* | TO | "to" | *to* |
| CD | cardinal number | *one, two* | NNPS | proper noun, plu. | *Carolinas* | UH | interjection | *ah, oops* |
| DT | determiner | *a, the* | NNS | noun, plural | *llamas* | VB | verb base | *eat* |
| EX | existential 'there' | *there* | PDT | predeterminer | *all, both* | VBD | verb past tense | *ate* |
| FW | foreign word | *mea culpa* | POS | possessive ending | *'s* | VBG | verb gerund | *eating* |
| IN | preposition/ subordin-conj | *of, in, by* | PRP | personal pronoun | *I, you, he* | VBN | verb past participle | *eaten* |
| JJ | adjective | *yellow* | PRP$ | possess. pronoun | *your, one's* | VBP | verb non-3sg-pr | *eat* |
| JJR | comparative adj | *bigger* | RB | adverb | *quickly* | VBZ | verb 3sg pres | *eats* |
| JJS | superlative adj | *wildest* | RBR | comparative adv | *faster* | WDT | wh-determ. | *which, that* |
| LS | list item marker | *1, 2, One* | RBS | superlatv. adv | *fastest* | WP | wh-pronoun | *what, who* |
| MD | modal | *can, should* | RP | particle | *up, off* | WP$ | wh-possess. | *whose* |
| NN | sing or mass noun | *llama* | SYM | symbol | *+,%, &* | WRB | wh-adverb | *how, where* |

# 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

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

# 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%

Most errors on unknown words