# Natural Language Processing

## Part of Speech Tagging

Edwin Puertas, Ph.D(c).
epuerta@utb.edu.co

# Parts of Speech

- From the earliest linguistic traditions (Yaska and Panini 5ᵗʰ C. BCE, Aristotle 4ᵗʰ C. BCE), the idea that words can be classified into grammatical categories.

- part of speech, word classes, POS, POS tags

- 8 parts of speech attributed to Dionysius Thrax of Alexandria (c. 1ˢᵗ C. BCE):

- noun, verb, pronoun, preposition, adverb, conjunction, participle, article

- These categories are relevant for NLP today.

# Two classes of words: Open vs. Closed

- Closed class words
  - Relatively fixed membership
  - Usually **function** words: short, frequent words with grammatical function
    - determiners: ***a, an, the***
    - pronouns: ***she, he, I***
    - prepositions: ***on, under, over, near, by, …***
- Open class words
  - Usually **content** words: Nouns, Verbs, Adjectives, Adverbs
    - Plus interjections: **oh, ouch, uh-huh, yes, hello**
  - New nouns and verbs like *iPhone* or *to fax*

**Open class** ("content") words

Nouns
- Proper

  *Janet*

  *Italy*
- Common

  *cat, cats*

  *mango*

Verbs
- Main

  *eat*

  *went*
- Auxiliary

  *can*

  *had*

Adjectives  *old  green  tasty*

Adverbs  *slowly yesterday*

Numbers

  *122,312*

  *one*

Interjections  *Ow  hello*

*… more*

**Closed class** ("function")

Determiners *the some*

Conjunctions  *and or*

Pronouns  *they its*

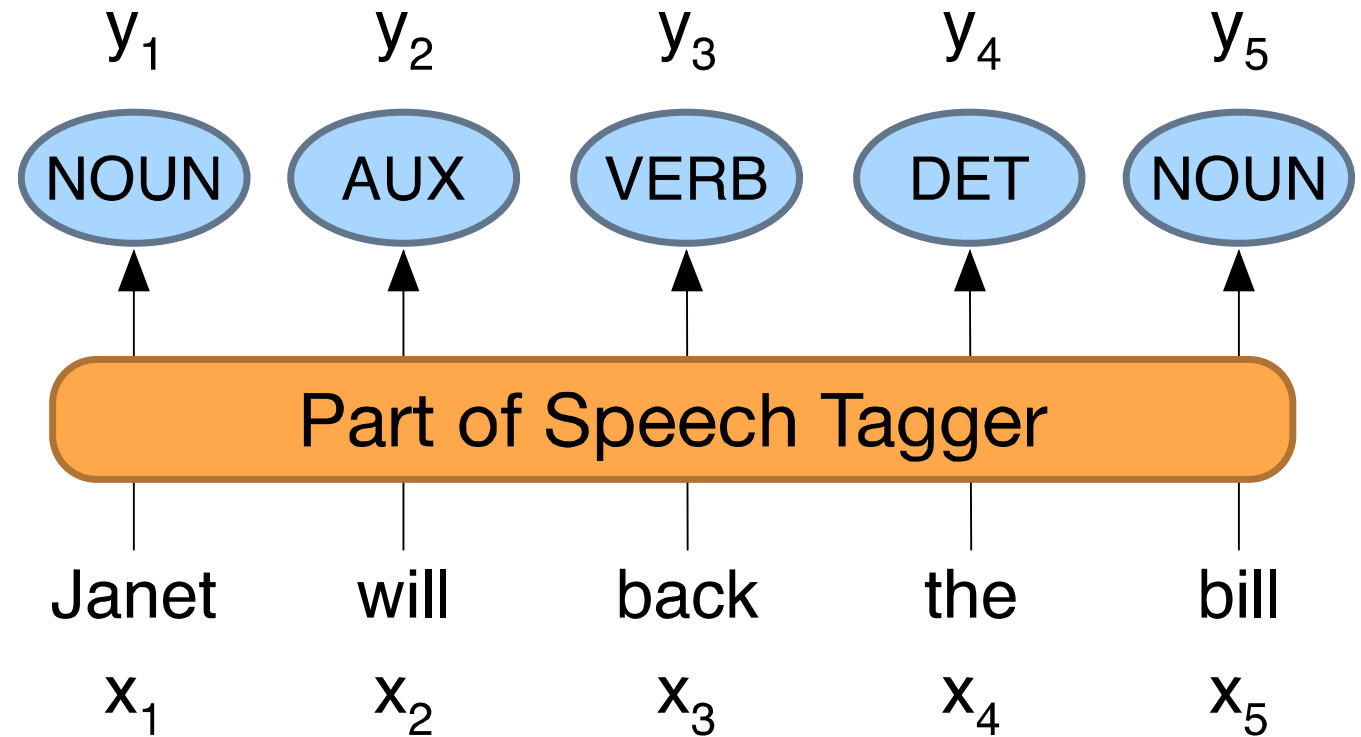Prepositions  *to with*

Particles  *off  up*

*… more*

# Part-of-Speech Tagging

- Assigning a part-of-speech to each word in a text.
- Words often have more than one POS.
- **book**:
    - VERB: (***Book*** *that flight*)
    - NOUN: (*Hand me that **book***).

# Part-of-Speech Tagging

Map from sequence $x_1,...,x_n$ of words to $y_1,...,y_n$ of POS tags

# "Universal Dependencies" Tagset

| | 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 |

- Nivre et al. 2016

# Sample "Tagged" English sentences

There/PRO were/VERB 70/NUM children/NOUN there/ADV ./PUNC  [ENG]

Había / AUX 70 / NUM niños / NOUN allí / ADV. / PUNC   [SPA]


Preliminary/ADJ findings/NOUN were/AUX reported/VERB in/ADP today/NOUN 's/PART New/PROPN England/PROPN Journal/PROPN of/ADP Medicine/PROPN

# Why Part of Speech Tagging?

- Can be useful for other NLP tasks
  - Parsing: POS tagging can improve syntactic parsing
  - MT: reordering of adjectives and nouns (say from Spanish to English)
  - Sentiment or affective tasks: may want to distinguish adjectives or other POS
  - Text-to-speech (how do we pronounce "lead" or "object"?)
- Or linguistic or language-analytic computational tasks
  - Need to control for POS when studying linguistic change like creation of new words, or meaning shift
  - Or control for POS in measuring meaning similarity or differenc

# How difficult is POS tagging in English?

- Roughly 15% of word types are ambiguous
- Hence 85% of word types are unambiguous
- *Janet* is always PROPN, *hesitantly* is always ADV
- But those 15% tend to be very common.
- So ~60% of word tokens are ambiguous
- E.g., *back*
    - earnings growth took a back/ADJ seat
    - a small building in the back/NOUN
    - a clear majority of senators back/VERB the bill
    - enable the country to buy back/PART debt
    - I was twenty-one back/ADV then

# POS tagging performance in English

- How many tags are correct?  (Tag accuracy)
  - About 97%
    - Hasn't changed in the last 10+ years
    - HMMs, CRFs, BERT perform similarly .
    - Human accuracy about the same
- But baseline is 92%!
  - Baseline is performance of stupidest possible method
    - "Most frequent class baseline" is an important baseline for many tasks
      - Tag every word with its most frequent tag
      - (and tag unknown words as nouns)
  - Partly easy because
    - Many words are unambiguous

# Sources of information for POS tagging

- `Janet will back the bill`
    AUX/NOUN/VERB?    NOUN/VERB?

- Prior probabilities of word/tag
    - "will" is usually an AUX
- Identity of neighboring words
    - "the" means the next word is probably not a verb
- Morphology and wordshape:
    - Prefixes            unable:        un- → ADJ
    - Suffixes            importantly:   -ly → ADJ
    - Capitalization      Janet:         CAP → PROPN

# Standard algorithms for POS tagging

- Supervised Machine Learning Algorithms:
- Hidden Markov Models
- Conditional Random Fields (CRF)/ Maximum Entropy Markov Models (MEMM)
- Neural sequence models (RNNs or Transformers)
- Large Language Models (like BERT), finetuned
- All required a hand-labeled training set, all about equal performance (97% on English)
- All make use of information sources we discussed
- Via human created features: HMMs and CRFs
- Via representation learning:  Neural LMs