**CORK INSTITUTE OF TECHNOLOGY**

INSTITIÚID TEICNEOLAÍOCHTA CHORCAÍ
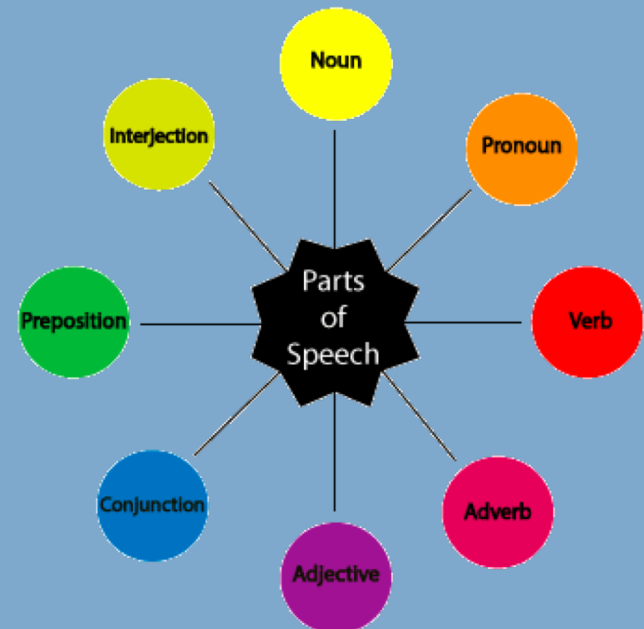
**Computer Science**

**http://www.cit.ie**

# Natural Language Processing
## Week4: Part-of-Speech Tagging

Dr. Haithem Afli

Haithem. afli@cit.ie

@AfliHaithem

2020/2021

# Part-of-Speech Tagging

- *What's the plural of "Part-of-Speech"?*

        *→ Parts-of-Speech*

        *not Part-of-Speeches ☺*


- *Abbreviation: POS*

# Parts of Speech

- 8 (ish) traditional parts of speech

  - Noun, verb, adjective, preposition, adverb, article, interjection, pronoun, conjunction, etc

  - Called: parts-of-speech, lexical categories, word classes, morphological classes, lexical tags...

  - Lots of debate within linguistics about the number, nature, and universality of these

    - We'll completely ignore this debate.

# POS examples

- N     noun
- V     verb
- ADJ   adjective
- ADV   adverb
- P     preposition
- PRO   pronoun
- DET   determiner

# POS examples

- N          noun            *chair, bandwidth, pacing*
- V          verb            *study, debate, munch*
- ADJ      adjective      *purple, tall, ridiculous*
- ADV      adverb          *unfortunately, slowly*
- P          preposition  *of, by, to*
- PRO      pronoun        *I, me, mine*
- DET      determiner    *the, a, that, those*

# POS Tagging

- The process of assigning a part-of-speech or lexical class marker to each word in a collection.

| WORD | tag |
|------|-----|
| the | DET |
| koala | N |
| put | V |
| the | DET |
| keys | N |
| on | P |
| the | DET |
| table | N |

# Why is POS Tagging Useful?

- First step of a vast number of practical tasks

- Speech synthesis

  - How to pronounce "lead"? How about "read"?

  - INsult              inSULT

  - OBject              obJECT

  - CONtent             conTENT

- Parsing

  - Need to know if a word is an N or V before you can parse

- Information extraction

  - Finding names, relations, etc.

- Machine Translation

# Open and Closed Classes

- Closed class: a small fixed membership
  - Prepositions: of, in, by, …
  - Auxiliaries: may, can, will had, been, …
  - Pronouns: I, you, she, mine, his, them, …
  - Usually function words (short common words which play a role in grammar)
- Open class: new ones can be created all the time
  - English has four: Nouns, Verbs, Adjectives, Adverbs
  - Many languages have these four too

# Open Class Words

CIT

- ## Nouns
  - Proper nouns (Boulder, Granby, Eli Manning)
    - English capitalizes these.
  - Common nouns (the rest).
  - Count nouns and mass nouns
    - Count: have plurals, get counted: goat/goats, one goat, two goats
    - Mass: don't get counted (snow, salt, communism) (*two snows)

- ## Adverbs: tend to modify things
  - Unfortunately, John walked home extremely slowly yesterday
  - Directional/locative adverbs (here,home, downhill)
  - Degree adverbs (extremely, very, somewhat)
  - Manner adverbs (slowly, slinkily, delicately)

- ## Verbs
  - In English, have morphological affixes (eat/eats/eaten)

# Closed Class Words

Examples:

- prepositions: *on, under, over, ...*

- particles: *up, down, on, off, ...*

- determiners: *a, an, the, ...*

- pronouns: *she, who, I, ..*

- conjunctions: *and, but, or, ...*

- auxiliary verbs: *can, may should, ...*

- numerals: *one, two, three, third, ...*

# Prepositions from CELEX

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| of | 540,085 | through | 14,964 | worth | 1,563 | pace | 12 |
| in | 331,235 | after | 13,670 | toward | 1,390 | nigh | 9 |
| for | 142,421 | between | 13,275 | plus | 750 | re | 4 |
| to | 125,691 | under | 9,525 | till | 686 | mid | 3 |
| with | 124,965 | per | 6,515 | amongst | 525 | o'er | 2 |
| on | 109,129 | among | 5,090 | via | 351 | but | 0 |
| at | 100,169 | within | 5,030 | amid | 222 | ere | 0 |
| by | 77,794 | towards | 4,700 | underneath | 164 | less | 0 |
| from | 74,843 | above | 3,056 | versus | 113 | midst | 0 |
| about | 38,428 | near | 2,026 | amidst | 67 | o' | 0 |
| than | 20,210 | off | 1,695 | sans | 20 | thru | 0 |
| over | 18,071 | past | 1,575 | circa | 14 | vice | 0 |

# POS Tagging
# Choosing a Tagset

- There are so many parts of speech, potential distinctions we can draw

- To do POS tagging, we need to choose a standard set of tags to work with

- Could pick very coarse tagsets

  - N, V, Adj, Adv.

- More commonly used set is finer grained, the "Penn TreeBank tagset", 45 tags

  - PRP$, WRB, WP$, VBG

- Even more fine-grained tagsets exist

# Penn TreeBank POS Tagset

| Tag | Description | Example | Tag | Description | Example |
|---|---|---|---|---|---|
| CC | coordin. conjunction | *and, but, or* | SYM | symbol | *+,%, &* |
| CD | cardinal number | *one, two, three* | TO | "to" | *to* |
| DT | determiner | *a, the* | UH | interjection | *ah, oops* |
| EX | existential 'there' | *there* | VB | verb, base form | *eat* |
| FW | foreign word | *mea culpa* | VBD | verb, past tense | *ate* |
| IN | preposition/sub-conj | *of, in, by* | VBG | verb, gerund | *eating* |
| JJ | adjective | *yellow* | VBN | verb, past participle | *eaten* |
| JJR | adj., comparative | *bigger* | VBP | verb, non-3sg pres | *eat* |
| JJS | adj., superlative | *wildest* | VBZ | verb, 3sg pres | *eats* |
| LS | list item marker | *1, 2, One* | WDT | wh-determiner | *which, that* |
| MD | modal | *can, should* | WP | wh-pronoun | *what, who* |
| NN | noun, sing. or mass | *llama* | WP$ | possessive wh- | *whose* |
| NNS | noun, plural | *llamas* | WRB | wh-adverb | *how, where* |
| NNP | proper noun, singular | *IBM* | $ | dollar sign | *$* |
| NNPS | proper noun, plural | *Carolinas* | # | pound sign | *#* |
| PDT | predeterminer | *all, both* | " | left quote | *' or "* |
| POS | possessive ending | *'s* | " | right quote | *' or "* |
| PRP | personal pronoun | *I, you, he* | ( | left parenthesis | *[, (, {, <* |
| PRP$ | possessive pronoun | *your, one's* | ) | right parenthesis | *], ), }, >* |
| RB | adverb | *quickly, never* | , | comma | *,* |
| RBR | adverb, comparative | *faster* | . | sentence-final punc | *. ! ?* |
| RBS | adverb, superlative | *fastest* | : | mid-sentence punc | *: ; ... – -* |
| RP | particle | *up, off* | | | |

# 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.

These examples from Dekang Lin

# Two Methods for POS Tagging

1. Rule-based tagging
   - (ENGTWOL)

2. Stochastic
   1. Probabilistic sequence models
      - HMM (Hidden Markov Model) tagging
      - MEMMs (Maximum Entropy Markov Models)

# Discussion

# Thank you

Haithem. afli@cit.ie

@AfliHaithem