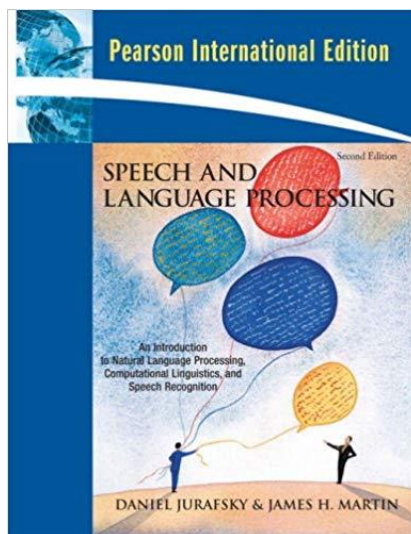# NLP & IR



## Chapter 8

## Sequence Labeling for
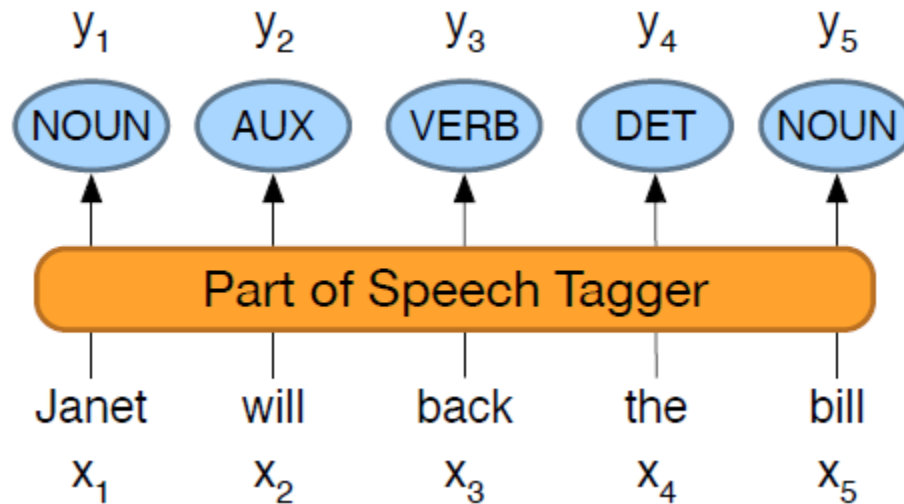## Parts of Speech and Named Entities

Dell Zhang

Birkbeck, University of London

# Sequence Labelling

- Part-of-Speech (POS) Tagging
  - *noun, verb, pronoun, preposition, adverb, conjunction, participle, and article, …*
- Named Entity recognition (NER)
  - *person, location, or organization, …*

# POS Tagging



**Figure 8.3** The task of part-of-speech tagging: mapping from input words $x_1, x_2, ..., x_n$ to output POS tags $y_1, y_2, ..., y_n$.

(8.1) There/PRO/EX are/VERB/VBP 70/NUM/CD children/NOUN/NNS there/ADV/RB ./PUNC/.

(8.2) Preliminary/ADJ/JJ findings/NOUN/NNS were/AUX/VBD reported/VERB/VBN in/ADP/IN today/NOUN/NN 's/PART/POS New/PROPN/NNP England/PROPN/NNP Journal/PROPN/NNP of/ADP/IN Medicine/PROPN/NNP

| | 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 |

**Figure 8.1** The 17 parts of speech in the Universal Dependencies tagset (Nivre et al., 2016a). Features can be added to make finer-grained distinctions (with properties like number, case, definiteness, and so on).

| 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* |

**Figure 8.2**    Penn Treebank part-of-speech tags.

# POS Tagging

- **Closed Class vs Open Class**

    - nouns, verbs, adjectives, adverbs, and interjections

    - function words (like *of*, *it*, *and*, or *you*)

- **Common Nouns vs Proper Nouns**

    - concrete terms (like *cat* and *mango*)
      abstract terms (like *algorithm* and *beauty*)
      verb-like terms (like *pacing*)

    - **named entities** (like *Regina*, *Colorado*, and *IBM*)

# POS Tagging

- Disambiguation
  - For example
    - *book* can be
      a verb (*book that flight*) or
      a noun (*hand me that book*).
    - *that* can be
      a determiner (*Does that flight serve dinner*) or
      a complementizer (*I thought that your flight was earlier*).

earnings growth took a **back/JJ** seat
a small building in the **back/NN**
a clear majority of senators **back/VBP** the bill
Dave began to **back/VB** toward the door
enable the country to buy **back/RP** debt
I was twenty-one **back/RB** then

# POS Tagging

- **Evaluation:** *accuracy*
  - STOA ≈ Human Performance: 97%
  - Most-Frequent-Class Baseline: 92%
    - Assigning each token to the class it occurred in most often in the training set.

| Types: | | WSJ | Brown |
|---|---|---|---|
| Unambiguous | (1 tag) | 44,432 (**86%**) | 45,799 (**85%**) |
| Ambiguous | (2+ tags) | 7,025 (**14%**) | 8,050 (**15%**) |
| **Tokens:** | | | |
| Unambiguous | (1 tag) | 577,421 (**45%**) | 384,349 (**33%**) |
| Ambiguous | (2+ tags) | 711,780 (**55%**) | 786,646 (**67%**) |

**Figure 8.4** Tag ambiguity in the Brown and WSJ corpora (Treebank-3 45-tag tagset).

# NER

- The task of named entity recognition (NER) is to
  - find *spans* of text that constitute proper names and
  - tag the *type* of the entity.

Citing high fuel prices, [ORG United Airlines] said [TIME Friday] it has increased fares by [MONEY $6] per round trip on flights to some cities also served by lower-cost carriers. [ORG American Airlines], a unit of [ORG AMR Corp.], immediately matched the move, spokesman [PER Tim Wagner] said. [ORG United], a unit of [ORG UAL Corp.], said the increase took effect [TIME Thursday] and applies to most routes where it competes against discount carriers, such as [LOC Chicago] to [LOC Dallas] and [LOC Denver] to [LOC San Francisco].

# NER

- Four entity tags are most common

| Type | Tag | Sample Categories | Example sentences |
|---|---|---|---|
| People | PER | people, characters | **Turing** is a giant of computer science. |
| Organization | ORG | companies, sports teams | The **IPCC** warned about the cyclone. |
| Location | LOC | regions, mountains, seas | **Mt. Sanitas** is in **Sunshine Canyon**. |
| Geo-Political Entity | GPE | countries, states | **Palo Alto** is raising the fees for parking. |

**Figure 8.5**   A list of generic named entity types with the kinds of entities they refer to.

# NER

- **Disambiguation**
  - For example,

[PER Washington] was born into slavery on the farm of James Burroughs.
[ORG Washington] went up 2 games to 1 in the four-game series.
Blair arrived in [LOC Washington] for what may well be his last state visit.
In June, [GPE Washington] passed a primary seatbelt law.

**Figure 8.6** Examples of type ambiguities in the use of the name *Washington*.

# NER

- ## BIO Tagging
  - ### **B**egin, **I**nside, and **O**utside

[PER **Jane Villanueva** ] of [ORG **United**] , a unit of [ORG **United Airlines Holding**] , said the fare applies to the [LOC **Chicago** ] route.

| Words | IO Label | BIO Label | BIOES Label |
|---|---|---|---|
| Jane | I-PER | B-PER | B-PER |
| Villanueva | I-PER | I-PER | E-PER |
| of | O | O | O |
| United | I-ORG | B-ORG | B-ORG |
| Airlines | I-ORG | I-ORG | I-ORG |
| Holding | I-ORG | I-ORG | E-ORG |
| discussed | O | O | O |
| the | O | O | O |
| Chicago | I-LOC | B-LOC | S-LOC |
| route | O | O | O |
| . | O | O | O |

**Figure 8.7** NER as a sequence model, showing IO, BIO, and BIOES taggings.

# NER

- Evaluation: $F_1$ measure
  - STOA: 93.39%; Human: 96.95%, 97.60% (MUC-7 in 1998)

- A useful first step in many NLP tasks:
  - Sentiment Analysis
  - Question Answering
  - …