

# Text processing

Words and Corpus, text tokenization, stemming, lemmatizing, PoS and NE

Information Retrieval

# IR and NLP

- Text / natural language is everywhere
  - News, emails, clinical reports, finance reports, ...
- Extracting information from documents is challenging
- Understanding user information needs
- Computing an answer for the user information need

# Natural language parsing

1. Word tokenization and sentence delimitation

2. Part of speech tagging

3. Word sense disambiguation

Terms weighting, Vector space models, language models, word embeddings

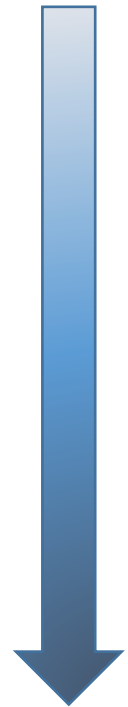
4. Named entities

Linking and relation

5. Subjective attributes

Sentiment, emotion, sarcasm, politeness, ...

**Semantic  
Abstraction**



# Basic Text Processing

## Words and Corpora

# How many words in a sentence?

- "I do uh main- mainly business data processing"
  - Fragments, filled pauses
- "Seuss's **cat** in the hat is different from other **cats**!"
  - **Lemma**: same stem, part of speech, rough word sense
    - **cat** and **cats** = same lemma
  - **Wordform**: the full inflected surface form
    - **cat** and **cats** = different wordforms

# How many words in a sentence?

they lay back on the San Francisco grass and looked at the stars and their

- **Type**: an element of the vocabulary.
- **Token**: an instance of that type in running text.
- How many?
  - 15 tokens (or 14)
  - 13 types (or 12) (or 11?)

# How many words in a corpus?

**$N$**  = number of tokens

**$V$**  = vocabulary = set of types,  **$|V|$**  is size of vocabulary

Heaps Law = Herdan's Law =  $|V| = kN^\beta$  where often  $.67 < \beta < .75$

i.e., vocabulary size grows with  $>$  square root of the number of word tokens

|                                 | Tokens = $N$ | Types = $ V $ |
|---------------------------------|--------------|---------------|
| Switchboard phone conversations | 2.4 million  | 20 thousand   |
| Shakespeare                     | 884,000      | 31 thousand   |
| COCA                            | 440 million  | 2 million     |
| Google N-grams                  | 1 trillion   | 13+ million   |

# Corpora

Words don't appear out of nowhere!

A text is produced by

- a specific writer(s),
- at a specific time,
- in a specific variety,
- of a specific language,
- for a specific function.



# Corpora variations

- **Language:** 7097 languages in the world
- **Variety**, like African American Language varieties.
  - AAE Twitter posts might include forms like "*iont*" (*I don't*)
- **Code switching**, e.g., Spanish/English, Hindi/English:

S/E: Por primera vez veo a @username actually being hateful! It was beautiful:)

*[For the first time I get to see @username actually being hateful! it was beautiful:]*

H/E: dost tha or ra- hega ... dont worry ... but dherya rakhe

*["he was and will remain a friend ... don't worry ... but have faith"]*

- **Genre:** newswire, fiction, scientific articles, Wikipedia
- **Author Demographics:** writer's age, gender, ethnicity, SES

Basic Text Processing

Word tokenization

# Text Normalization

- Every NLP task requires text normalization:
  1. Tokenizing (segmenting) words
  2. Normalizing word formats
  3. Segmenting sentences

# Space-based tokenization

- A very simple way to tokenize
  - For languages that use space characters between words
    - Arabic, Cyrillic, Greek, Latin, etc., based writing systems
  - Segment off a token between instances of spaces
- Unix tools for space-based tokenization
  - The "tr" command
  - Inspired by Ken Church's UNIX for Poets
  - Given a text file, output the word tokens and their frequencies

# Simple Tokenization in UNIX

- (Inspired by Ken Church's UNIX for Poets.)
- Given a text file, output the word tokens and their frequencies

```
tr -sc 'A-Za-z' '\n' < shakes.txt  
    | sort  
    | uniq -c
```

Change all non-alpha to newlines

Sort in alphabetical order

Merge and count each type

```
1945 A  
72 AARON  
19 ABBESS  
5 ABBOT  
... ..
```

```
25 Aaron  
6 Abate  
1 Abates  
5 Abbess  
6 Abbey  
3 Abbot  
.... ..
```

# Scikit-learn Tokenization

```
from sklearn.feature_extraction.text import CountVectorizer

corpus = [
    'This is the first document.',
    'This document is the second document.',
    'And this is the third one.',
    'Is this the first document?',
]

my_stop_words = {'is', 'the'}
# UNIGRAMS
vectorizer = CountVectorizer(ngram_range=(1,1), analyzer='word', stop_words = None)
#vectorizer = CountVectorizer(ngram_range=(1,1), analyzer='word', stop_words = 'english')
#vectorizer = CountVectorizer(ngram_range=(1,1), analyzer='word', stop_words = my_stop_words)

# UNIGRAMS and BIGRAMS
#vectorizer = CountVectorizer(ngram_range=(1,2), analyzer='word')

# Character GRAMS
#vectorizer = CountVectorizer(ngram_range=(3,4), analyzer='char')
X = vectorizer.fit_transform(corpus)

print(vectorizer.get_feature_names())

print(X.todense())
```

# Issues in Tokenization

- Can't just blindly remove punctuation:
  - m.p.h., Ph.D., AT&T, cap'n
  - prices (\$45.55)
  - dates (01/02/06)
  - URLs (<http://www.stanford.edu>)
  - hashtags ([#nlproc](#))
  - email addresses ([someone@cs.colorado.edu](mailto:someone@cs.colorado.edu))
- Clitic: a word that doesn't stand on its own
  - "are" in [we're](#), French "je" in [j'ai](#), "le" in [l'honneur](#)
- When should multiword expressions (MWE) be words?
  - [New York](#), [rock 'n' roll](#)

Basic Text Processing

Data-driven Word Tokenization (BPE)



# Another option for text tokenization

Instead of

- white-space segmentation
- single-character segmentation

**Use the data** to tell us how to tokenize.

**Subword tokenization** (because tokens can be parts of words as well as whole words)

# Subword tokenization

- Three common algorithms:
  - **Byte-Pair Encoding (BPE)** (Sennrich et al., 2016)
  - **Unigram language** modeling tokenization (Kudo, 2018)
  - **WordPiece** (Schuster and Nakajima, 2012)
- All have 2 parts:
  - A token **learner** that takes a raw training corpus and induces a vocabulary (a set of tokens).
  - A token **segmenter** that takes a raw test sentence and tokenizes it according to that vocabulary

# Byte Pair Encoding (BPE) token learner

Let vocabulary be the set of all individual characters

= {A, B, C, D,..., a, b, c, d....}

- Repeat:
  - Choose the two symbols that are most frequently adjacent in the training corpus (say 'A', 'B')
  - Add a new merged symbol 'AB' to the vocabulary
  - Replace every adjacent 'A' 'B' in the corpus with 'AB'.
- Until  $k$  merges have been done.

# BPE token learner algorithm

```
function BYTE-PAIR ENCODING(strings  $C$ , number of merges  $k$ ) returns vocab  $V$ 

 $V \leftarrow$  all unique characters in  $C$            # initial set of tokens is characters
for  $i = 1$  to  $k$  do                           # merge tokens til  $k$  times
     $t_L, t_R \leftarrow$  Most frequent pair of adjacent tokens in  $C$ 
     $t_{NEW} \leftarrow t_L + t_R$                  # make new token by concatenating
     $V \leftarrow V + t_{NEW}$                        # update the vocabulary
    Replace each occurrence of  $t_L, t_R$  in  $C$  with  $t_{NEW}$     # and update the corpus
return  $V$ 
```

# Byte Pair Encoding (BPE) Addendum

Most subword algorithms are run inside space-separated tokens.

So we commonly first add a special end-of-word symbol '\_\_\_\_' before space in training corpus

Next, separate into letters.

# BPE token learner

Original (very fascinating 😊) corpus:

low low low low low lowest lowest newer newer newer newer newer newer wider  
wider wider new new

Add end-of-word tokens, resulting in this vocabulary:

**vocabulary**

—, d, e, i, l, n, o, r, s, t, w

# BPE token learner

## corpus

5 l o w \_  
2 l o w e s t \_  
6 n e w e r \_  
3 w i d e r \_  
2 n e w \_

## vocabulary

\_, d, e, i, l, n, o, r, s, t, w

Merge **e r** to **er**

## corpus

5 l o w \_  
2 l o w e s t \_  
6 n e w er \_  
3 w i d er \_  
2 n e w \_

## vocabulary

\_, d, e, i, l, n, o, r, s, t, w, er

# BPE

## corpus

5 l o w \_  
2 l o w e s t \_  
6 n e w e r \_  
3 w i d e r \_  
2 n e w \_

## vocabulary

\_, d, e, i, l, n, o, r, s, t, w, e r

Merge **er \_** to **er\_**

## corpus

5 l o w \_  
2 l o w e s t \_  
6 n e w e r\_  
3 w i d e r\_  
2 n e w \_

## vocabulary

\_, d, e, i, l, n, o, r, s, t, w, e r, e r\_



# BPE

## corpus

5 l o w \_  
2 l o w e s t \_  
6 n e w er\_  
3 w i d er\_  
2 n e w \_

Merge **n e** to **ne**

## corpus

5 l o w \_  
2 l o w e s t \_  
6 ne w er\_  
3 w i d er\_  
2 ne w \_

## vocabulary

\_, d, e, i, l, n, o, r, s, t, w, er, er\_

## vocabulary

\_, d, e, i, l, n, o, r, s, t, w, er, er\_, ne

# BPE

The next merges are:

| Merge      | Current Vocabulary   |
|------------|--|
| (ne, w)    | —, d, e, i, l, n, o, r, s, t, w, er, er—, ne, new                        |
| (l, o)     | —, d, e, i, l, n, o, r, s, t, w, er, er—, ne, new, lo                    |
| (lo, w)    | —, d, e, i, l, n, o, r, s, t, w, er, er—, ne, new, lo, low               |
| (new, er—) | —, d, e, i, l, n, o, r, s, t, w, er, er—, ne, new, lo, low, newer—       |
| (low, —)   | —, d, e, i, l, n, o, r, s, t, w, er, er—, ne, new, lo, low, newer—, low— |

# BPE token segmenter algorithm

On the test data, run each merge learned from the training data:

- Greedily
- In the order we learned them
- (test frequencies don't play a role)

So: merge every **e r** to **er**, then merge **er \_** to **er\_**, etc.

- Result:
  - Test set "n e w e r \_" would be tokenized as a full word
  - Test set "l o w e r \_" would be two tokens: "low er\_"

# Properties of BPE tokens

Usually include frequent words

And frequent subwords

- Which are often morphemes like *-est* or *-er*

A **morpheme** is the smallest meaning-bearing unit of a language

- *unlikeliest* has 3 morphemes *un-*, *likely*, and *-est*

Basic Text Processing

Word Normalization and other issues

# Character processing and stop-words

- Numbers/dates
- Acronyms
- Multi-language documents
- Stop-words: remove words that are present in all documents
  - *a, and, are, as, at, be, but, by, for, if, in, into, is, it, no, not, of, on, or, such, that, the, their, then, there, these, they, this, to, was, will...*

# Word Normalization

- Putting words/tokens in a standard format
  - U.S.A. or USA
  - uhhuh or uh-huh
  - Fed or fed
  - am, is, be, are

# Case folding

- Applications like IR: reduce all letters to lower case
  - Since users tend to use lower case
  - Possible exception: upper case in mid-sentence?
    - e.g., *General Motors*
    - *Fed* vs. *fed*
    - *SAIL* vs. *sail*
- For sentiment analysis, MT, Information extraction
  - Case is helpful (*US* versus *us* is important)



# Lemmatization

- Represent all words as their lemma, their shared root
  - = dictionary headword form:
    - am, are, is → be
    - car, cars, car's, cars' → car
    - Spanish quiero ('I want'), quieres ('you want')
      - → querer 'want'
- He is reading detective stories
  - → He be read detective story

# Lemmatization is done by Morphological Parsing

- Morphemes:

- The small meaningful units that make up words
- **Stems**: The core meaning-bearing units
- **Affixes**: Parts that adhere to stems, often with grammatical functions

- Morphological Parsers:

- Parse *cats* into two morphemes *cat* and *s*
- Parse Spanish *amaren* ('if in the future they would love') into morpheme *amar* 'to love', and the morphological features *3PL* and *future subjunctive*.

# Dealing with complex morphology is necessary for many languages

- e.g., the Turkish word:
- **Uygarlastiramadiklarimizdanmissinizcasina**
- `(behaving) as if you are among those whom we could not civilize’
- **Uygar** `civilized’ + **las** `become’
  - + **tir** `cause’ + **ama** `not able’
  - + **dik** `past’ + **lar** `plural’
  - + **imiz** `p1pl’ + **dan** `abl’
  - + **mis** `past’ + **siniz** `2pl’ + **casina** `as if’

# Stemming

- Reduce terms to stems, chopping off affixes crudely

This was not the map we found in Billy Bones's chest, but an accurate copy, complete in all things-names and heights and soundings-with the single exception of the red crosses and the written notes.



Thi wa not the map we found in Billi Bone s chest but an accur copi complet in all thing name and height and sound with the singl except of the red cross and the written note

.

# Porter Stemmer

- Based on a series of rewrite rules run in series
  - A cascade, in which output of each pass fed to next pass
- Some sample rules:

ATIONAL → ATE (e.g., relational → relate)

ING → ε if stem contains vowel (e.g., motoring → motor)

SSES → SS (e.g., grasses → grass)

# Sentence Segmentation

!, ? mostly unambiguous but **period** “.” is very ambiguous

- Sentence boundary
- Abbreviations like Inc. or Dr.
- Numbers like .02% or 4.3

Common algorithm: Tokenize first: use rules or ML to classify a period as either (a) part of the word or (b) a sentence-boundary.

- An abbreviation dictionary can help

Sentence segmentation can then often be done by rules based on this tokenization.

# Text processing

## Part of Speech tagging

[https://pt.wikipedia.org/wiki/Análise\\_morfológica](https://pt.wikipedia.org/wiki/Análise_morfológica)

[https://en.wikipedia.org/wiki/Part-of-speech\\_tagging](https://en.wikipedia.org/wiki/Part-of-speech_tagging)

# 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*



## Open class ("content") words

### Nouns

#### Proper

*Janet*  
*Italy*

#### Common

*cat, cats*  
*mango*

### Verbs

#### Main

*eat*  
*went*

### 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*

### Auxiliary

*can*  
*had*

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***).

# "Universal Dependencies" Tagset

Nivre et al. 2016

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

# Sample "Tagged" English sentences

- There/**PRO** were/**VERB** 70/**NUM** children/**NOUN** there/**ADV** ./**PUNC**
- 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 difference

# Text Processing

## Named entities

[https://pt.wikipedia.org/wiki/Análise\\_morfológica](https://pt.wikipedia.org/wiki/Análise_morfológica)

[https://en.wikipedia.org/wiki/Part-of-speech\\_tagging](https://en.wikipedia.org/wiki/Part-of-speech_tagging)

# Named Entities

- **Named entity**, in its core usage, means anything that can be referred to with a proper name. Most common 4 tags:
  - **PER** (Person): “Marie Curie”
  - **LOC** (Location): “New York City”
  - **ORG** (Organization): “Stanford University”
  - **GPE** (Geo-Political Entity): “Boulder, Colorado”
- Often multi-word phrases
- But the term is also extended to things that aren't entities:
  - dates, times, prices

# Named Entity tagging

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



# NER output

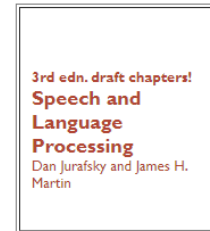
Citing high fuel prices, [ORG United Airlines] said [TIME Friday] it has increased fares by [MONEY \$6] per round trip on flights to some cities also served by lower-cost carriers. [ORG American Airlines], a unit of [ORG AMR Corp.], immediately matched the move, spokesman [PER Tim Wagner] said. [ORG 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].

# Why NER?

- **Sentiment analysis:** consumer's sentiment toward a particular company or person?
- **Question Answering:** answer questions about an entity?
- **Information Extraction:** Extracting facts about entities from text.

# Summary of basic techniques

- Words and Corpus 2.2 + 2.3
- Text normalization 2.4
  - Tokenization 2.4.2 + 2.4.3
  - Stemming, Lemmatization 2.4.4
  - Sentence segmentation 2.4.5
- PoS and NER 8.1, 8.2, 8.4



Chapter 2, 8