

## Regular Expressions

- Formally, a regular expression is an algebraic notation of characterizing a set of strings. particularly used for searching in texts, when we have a pattern to search for and a corpus of texts to search through
- Before we take processing words, we need to decide what counts as a word??
- Corpus ((pl: corpora)): a computer readable collection of text or speech. it is the ~~dataset~~ we use for NLP
- Token: Each individual occurrence of a word in a text, including punctuation is counted ((depending on context))
- Type: unique words in corpus ((vocabulary size))  
capitalization and punctuation may or may not count as separate tokens depending on NLP tasks.
- Disfluency in speech:
  - Disfluency: unplanned pause or corrections in speech.
    - broken-off words: main-
    - Fillers: uh, um

## ● Lemma vs word forms:

- lemma: the base form of a word, representing all its inflected word forms.
- word form: a fully inflected or derived form  
eg: "cat" is lemma, "cats" is a word form.

## ● Vocabulary size:

Type ( $|V|$ ): size of corpus

Tokens ( $N$ ): total running words

## ● Zipf's Law / Heaps Law

- observation: as a corpus gets larger, new word types keep appearing, but at a decreasing rate.

Formula:  $|V| = KN^\beta$        $|V| = \text{types}, N = \text{tokens}$   
 $K, \beta = \text{constants}$

$0 < \beta < 1$  usually 0.65-0.75

## ● Corpora and variations in language:

① context of words

② language varieties

- american mainstream english
- afro-american english

③ code switching

- using more than one language in a stream.  
like hindi words in english

④ genre variation

- different genres, scientific, poetic, etc.



- tokenization: splitting text into words or meaningful units.
- word normalization: standardizing word forms.
- sentence segmentation: splitting text into sentences.

### word tokenization:

#### - punctuation separation:

"Hello, world!" → ["Hello", ",", "world"]

#### - clitic expression:

"we're" → ["we", "are"]

#### - special cases:

"Ph.D." → ["Ph.D"]

#### - numbers:

"\$45.55" → ["\$45.55"]

### Subword Tokenization (BPE)



## ■ Lemma and Senses:

■ Lemma; (citation form): the base dictionary form of a word

— sing is lemma for sing, sang, sung

■ Word form: any specific inflected form of a lemma

— sung, carpets, duermes

■ Sense: a distinct meaning of a lemma

— mouse = rodent, mouse = computer device

■ Lemmas can be polysemous (multiple senses), this will create interpretation problems.

■ Word senses disambiguation (wsd):

Task of deciding which sense applies in context

## ■ Synonymy:

Two senses are synonyms if they can substitute for each other in any sentence without changing its truth conditions.

No two words are truly identical

→ Differences in linguistic form correspond to some to some difference in meaning.



## Word Similarity

- not the same as synonymy:

-- cat and dog are not synonyms but are similar

## Word Relatedness:

- words can be related without being similar. these words often have a semantic relation with each other.

## Semantic Frames and Roles:

- semantic frame: a set of words representing roles in an event.

eg: commercial transaction

verbs: buy, sell, pay  
noun: buyer  
roles: buyers, sellers, money

## Connotation:

- emotional part of the meaning, positive or negative.

- words with similar meaning, can differ in connotation

eg: copy      copy seems positive/neutral  
fake      seems negative.