

Stop words are function words.

words w/ little lexical meaning
BUT
are imp't. to sentence structure

* ~ funcⁿ word = content word.

Type vs token

unique word

all words

Lexical richness = avg. count of how many times each unique word appears in text.

$\frac{1}{LR}$ = type token ratio

→ measure of vocabulary variation

TTR = $\frac{\text{type count}}{\text{token count}}$.

Tokenization → "process — segmenting — string(char) → words"

→ starting + ending pt?

Binary classifier — simple decision tree.

→ Also handle: a) handling hyphens

b) named entities

c) word segmentation (greedy algo) — match words w/ lexicon.

But why normalization for tokenizⁿ?

Equivalence classes

Case folding

all in lower case

Normalization strategies:

A. Lemmatization

- variant → base form

- dictionary head word req.

} use morphology.

B. Stemming

- rule based — not dict. dependent

- remove suffixes (crudely)

- language dependent.

} set of 'if-then-else'

- doesn't always produce words.

Algorithmic Stemmers

↓ Algorithmic Stemmers

Eg:

1/p
Hoped

Ship 'ed'
hop

Repair
hope (add e if short)
hop (delete one if doubled)

- + fast + lean
- + more efficient to not use dict.
- + ignore irregular forms rather than complicating algo.

What if we've spelling errors?

Minimum edit distance - Levenshtein. Nuh dry.

→ Valid changes?

- insert char
- delete char
- change char.
- transpose.

Usual normalization flow:

1. Convert to lowercase
2. Remove numb. + whitespace (regex)
3. Tokenization
4. Remove stopwords.
5. Stemming / lemmatization
6. POS tagging
7. Chunking / shallow parsing
8. Named entity recog.
9. Coreference resolution.