

TDCEN IZAEON

Tokenization  $\rightarrow$  part of text normalization

1. Segmenting/Identifying words
2. Normalizing words format
3. Segmenting sentences

Words

- \* Type : an element of the vocabulary
- \* Token : an instance of that type

! do wh main + mainly } how many  
business tech processing } words?

filled pause

fragments

Titali's cast in the heart  
is different from other casts

} Cast  
vs  
Casts

- \* lemma: same stem, part of speech, rough word sense
  - ↳ cat & cats is same lemma
- \* word form: full inflected surface form
  - ↳ cat & cats is different word forms



## How many words?

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.

tolken  $\rightarrow$  (5) tolken Cor (4) "San", "Francisco"  
er

Type  $\rightarrow$  13 types (or 12) (or 11) "The" "They" "He/She" "San Francisco"

In general,

$N$  = number of tokens

V = Vocabulary

= set of types,

with  $|V|$  is the size of

There are some issues in tokenization:

- use of "what're", "finland's", "isn't"
- use of - "Hawlett-Packard",

the vocabulary

"state-of-the-art"

- "San Francisco", 1 or 2 tokens?
- m.p.h., Ph.D.
- language issues  
(German, French, Japanese, Chinese, etc.)

## NORMALIZATION

"Normalize"

- \* IR → indexed text & query terms must have same form
- \* asymmetric expansion
- \* reduce all letters to lowercase
  - ↳ sometimes "Case" is helpful: sentiment analysis, MT, IE
  - (US and us is important)

Lemmatize → reduce inflection/  
variant forms  
to base form

- \* am, are, is ⇒ be
  - \* car, cars, car's, cars' ⇒ car
- ↳ finding correct dictionary headword form

Looking at parts of words → morphology

↓  
morphemes

↳ small meaningful units that  
make up words

↳ 2 kinds of morphemes:

- stems: core meaning-bearing units
- affixes: pieces that adhere to stems  
↳ often with grammatical function

# SENTENCE SEGMENTATION

usually ends with

“!” “?” “.”

relatively  
unambiguous

quite  
ambiguous

Inc. or Dr.

Numbers like .02 % or 4.3

is it the end of a sentence?

How to detect that? → build a classifier that  
decides whether it's end of sentence /  
not end of sentence.

hand-written /  
regex / ML