

Text Preparation (NLP)

(1)

Lecture 13

Introduction



1) Lowercasing

→ Whenever you get text, in first step you bring it to lower case e.g.

"Condition of the car is very good, provided that the condition of such old cars".

→ Python language is case-sensitive
df['...'].lower()

2) Remove Unnecessary (e.g. HTML tags, punctuation).

→ By using regex (import re)
→ Online tools available

3) Remove URLs

→ if not needed.
https: & http, www:

4) Remove Punctuations

" ! \$ % & ' () * + , - . : ;

→ We remove these for tokenization.

Buddy! What's up? ⇒ Buddy? What's up

→ In python, import string or for big data use translate.

Chat Word Treatment

→ New trend, where words are being shortened.

→ Asep

→ Lol

→ Gma

→ FYI

Spelling Correction

→ By using different libraries such as textblob.

Removing Stop Words

→ From sentiment analysis & document classification perspective. However you don't do it in POS tagging.

is, are, the, a, my, you, which
These sentence formation words are removed.

→ These are libraries, one such is ~~NLTK~~ NLTK.

Handling Emoji's

→ Either Remove

→ or replace (😊 == happy)

→ Each emoji has unique identifier (U+000024E2)

→ In python import emoji.

Tokenization

→ is breaching of your text data, into smaller parts.

This is Sparta. You fall here

→ Word level tokenization

[This, is, is, Sparta, ...]

→ Sentence level tokenization

[This is Sparta, You fall here]

→ I am new in new delhi
1 2 3 4 3 5
but correct is

→ I am new in new delhi
1 2 3 4 5

→ So one should take care of,
Prefix, suffix, Infix (co-operation), Exceptions (let's u.s.)

\$20 30m

→ In Python

- 1) Split() → on word level
split() → on sentence level

Are you going to peshawar? I am not.
→ on full stop whole two sentences

- 2) Regular Expressions (import re)
- 3) Libraries
↳ NLTK
- 4) The best is Spacy.

Stemming

→ In grammar, inflection is word modification to express different grammatical categories such as tense case, aspect, voice, person numbers, gender & mood.
do → undoable

Eg. Walk Walking
Walked
Walks

→ Stemming is process to reduce inflection in words to their roots. Important in IR (Information Retrieval) systems.

eg. In google search (fish, fishy, fishnet) will show images related to fish.

- NLTK is porter stemmer.
 - ↳ (linguistics make it).
 - Lemmatization
 - ↳ if you want faster stemming.
 - ↳ if you want to show to your user.
- In lemmatization, root will be english word.
story → stor;