

Day 2, 10 Dec 2023

## Natural Language Processing (NLP)

- by mentor Dr. Ayan Debnath, IIT Delhi + Harvard Alumni

- Recap
- ① What is NLP
  - ② NLP → application
  - ③ Roadmap
  - ④ Level-1

## Tokenization

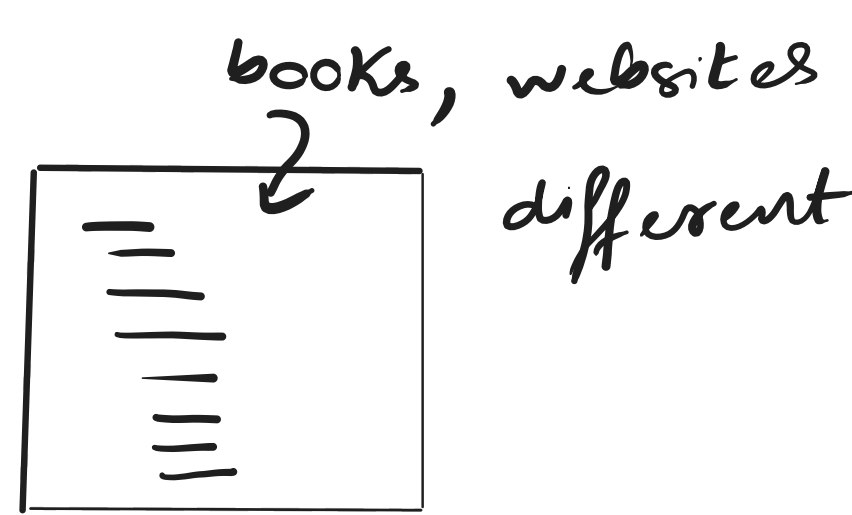
terminology

Corpus - paragraph

Document - sentences

Vocabulary - unique words

words -



tokenization → splitting text into smaller chunks/parts

paragraph → " sent1. sent2. .... "

eg " I Love NLP. It is very interesting topic. .... "

### Sentence Tokenization

paragraph → sentences

sentence → I Love NLP

word Tokenization

sentence to words

Why Tokenization is needed?

easy computation

easy processing

## Stopwords

I am a boy and I love to eat chicken → boy love eat chicken

Q) Why we remove stopwords?

easy <sup>①</sup> computation, faster <sup>②</sup> processing, Redundant <sup>③</sup>

text → no.

10 ----- ⇒ ----- ⇒ Computation resource ↑  
3 ---- ⇒ ---- ⇒ " " ↓

Stemming → Process of converting inflected word to stem/base word

- eg ① history → histori  
historical → histori stem word / base word
- ② loving → love  
love → love  
loveable → love
- ③ final → fina  
finalization → fina  
finally → fina

Lemmatization → Convert the word to base word with some meaning to it.

eg history → history  
historical → history

final → final  
finalization → final  
finally → final

Stemming	vs	Lemmatization
① Convert to base word		① ... but with meaningful word
② Faster process		② Slower process wot stemming
③ Application		③ Application
Sentiment analysis		chatbot / Q&A
Restaurant review		
→ ~+ve/-ve		

$w_1, w_2, w_3$  Sent1  
2  
3

if  $w_1$  not in set(stopwords('english'))  
for  $w_1$  in word  
rootword = stemmer.stem( $w_1$ )

Tokenization, Stopwords, Stemming, Lemmatization

Text processing: I Love NLP  
① ↓  
love NLP  
↓  
numerical <sup>②</sup> 12 256  
representation

Bows, TFIDF

Bag of words (Bow)

Sent1: He is intelligent boy  
Sent2: she is intelligent girl  
Sent3: Both of them are intelligent boy and girl

	intelligent	boy	girl
Sent1	✓	✓	
" 2	✓		✓
" 3	✓	✓	✓

test ✓ ✓ → model