# Tokenization in NLP

① Topics

    ① Corpus  →  Paragraphs

    ② Documents  →  Sentences

    ③ vocabulary  →  unique words

    ④ words  →  All words present in corpus.

② Tokenization :-

" My name is prathamesh I have intrest in teaching in ML, NLP & DL." I am also a "Youtber".

$$\Downarrow$$

Tokens { Sentences }

$$\downarrow$$

1) my name is prathamesh & I have intrest in teaching ML, NLP & DL.

2) I am also a Youtuber.

$$\Downarrow$$

Tokenisation

③ Another example :- "I like to drink Apple Juice my friends like mango Jucie".

$$\downarrow$$ Tokenization.

① So Corpus :- "I like to drink Apple Juice & my friend like mango Juice."

② Documents :-

   ① I like to drink Apple Juice.

   ② my friend like mango Juice.

③ Total words :- 11 words.

④ Vocabulary :- 9 words unique words.

④ 4) Practicle session NLP

→ Diffrence Between NLTK & SPacy.

| Parameters | SPacy | NLTK |
|---|---|---|
| Use | High Performance & Producho ready application focousing Speed & efficiency. | NLTK is more towords researc & education, offer wider range of algorithum & flextibility. |
| Performance | High | low |
| Approach | Object-orient and have rich information about text | Primary work on strings require manual approch or work. |
| focous | in production high accuracy & Performance | research & experiments algorithms. |

(5) Practice session

→ Tokenization.

10. What we learnt ?

Text preprocessing

①

↑

| DATASET |
|---|

②

↑

| Text prepres |
|---|
| Preprocessin |

→ lower case word

→ ~~Fster~~ lower case word

→ Tokenization

→ regular expression

③

↑

| Text |
|---|
| Preprocessing -2 |

→ Stemming

→ lemmatization

→ Stopwords

⑤

↑

| ML |
|---|
| Algo |
| train with |

→ Gensim

④

↑

| Text |
|---|
| to |
| vectors |

→ OHE

→ BOW

→ TF-IDF

→ word2vec

→ Avg word2vec