

Day 2

① Preprocessing

① Tokenization → Sentence Tokenizer

→ Word Tokenizer

↳ ① Word Tokenizer

↳ ② White space

→ Regex Tokenizer

② Normalization - lower case

③ stopwords Removal - ① language specific
↳ ② Domain specific
stopwords

④ Remove punctuation

⑤ Contraction mapping

I didn't like the movie -ve
I like the movie +ve

didn't = did not, never

don't = do not

⑥ stemming & lemmatization

↓
stem

(root word)

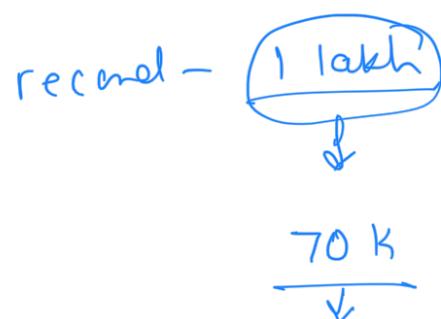
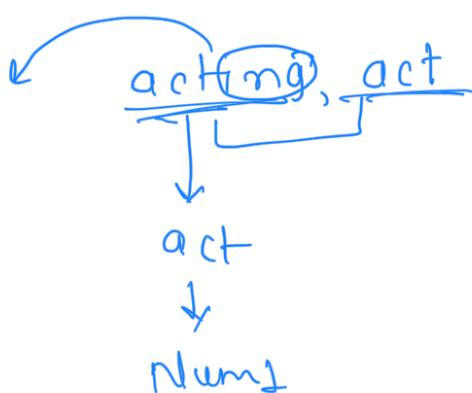
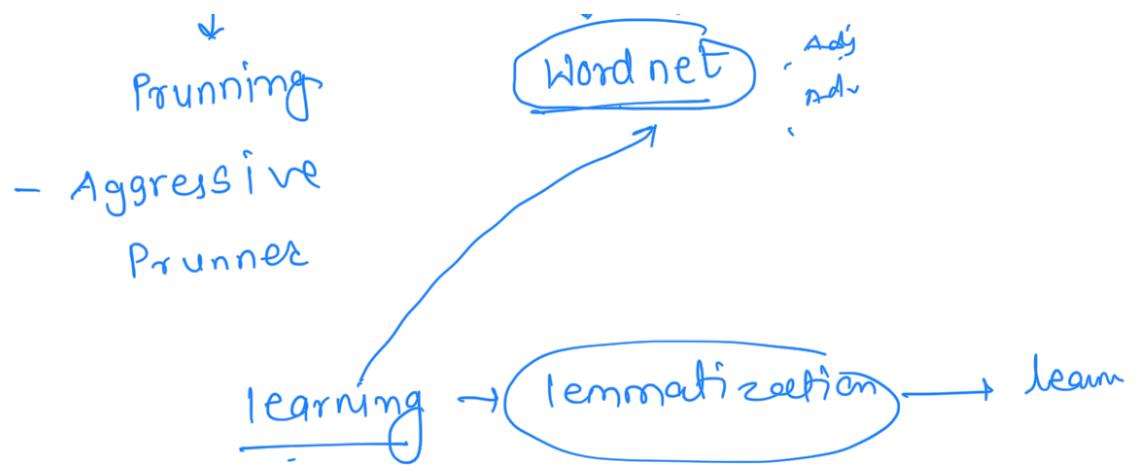
| ...

↓

lemma

(meaningful root word)

↓ .



⑦ Handling accented characters

ä, ö, ÿ

apple → apple

ßat → bat

⑧ cleaning → regex -
- "html, >> ? a" b

C

I like the

product "

* ⑨ Remove digits

* ⑯ Auto correction → *

* Feature Engineering

Text → Numerical → Word embeddings

Word Embedding

① Frequency Based

- 1) One hot encoding
- 2) Count vectorizer
- 3) TF IDF

② Prediction Based

- 1) Word2Vec
- 2) Doc2Vec

① Count vectorizer

Sent ① We are learning nlp

Sent ② NLP is a part of data science

Sent ③ machine learning and deep learning
are imp in data science

	learning	nlp,	data	science	part	machine	deep	imp	ora
①	1	1	0	0	0	0	0	0	0
②	0	1	1	1	1	0	0	0	0
③	2	0	1	1	1	0	1	1	1
test	0	0	1	1	1	0	0	0	2

test ×

test = data science is interesting field

Drawbacks

① Curse Of Dimensionality

② It will not give meaning of word

* Creating vocabulary

1 lakh → 10 k

max-df = 0.95, 0.92

most frequent

100 doc

95 >

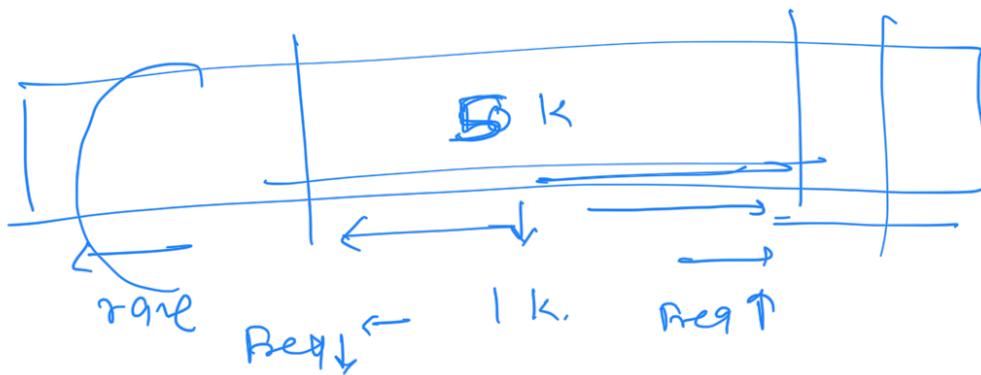
X

$1k\downarrow$

$$\min-df = \underline{1}, \underline{5}$$

100 doc \rightarrow 25 \times rare words

$$= 1k$$



OOV = Out of vocabulary

Count vectorizer \rightarrow train \rightarrow columns \rightarrow $\frac{100}{col}$

test \rightarrow C.V. \rightarrow OOV
3 \rightarrow + 3

100 words \rightarrow 100 column + 1 OOV

TF-IDF

TF = Term frequency

= $\frac{\text{Count of word in a doc } d}{\text{Total number of words in } d}$

IDF = Inverse Documents Frequency

$DF = \frac{\text{number of doc. containing term } t}{\text{total no. of doc}}$

$$IDF = \log \left(\frac{1}{DF} \right)$$

Sent ① = "NLP part data science
NLP huge domain"

$$TF(\text{nlp}) = \frac{2}{7}$$

Sent ① = "learning nlp"

Sent ② = "NLP part data science"

Sent ③ = "machine learning
part data science" & nlp
deep learning

	learning	nlp	part	data	science	machine	deep
①	$\frac{1}{2} \times \log(3)$	$\frac{1}{2} \times \log(3)$	0	0	0	0	0

	$U \cup \bar{v}$	$U - \bar{v}$					
②	0	$\frac{1}{4} \log\left(\frac{3}{2}\right)$	$\frac{1}{4} \log\left(\frac{3}{2}\right)$	$\frac{1}{4} \log\left(\frac{3}{2}\right)$	$\frac{1}{4} \log\left(\frac{3}{2}\right)$	0	0
③	$\frac{2}{7} \log\left(\frac{3}{2}\right)$	0	$\frac{1}{7} \log\left(\frac{3}{2}\right)$	$\frac{1}{7} \log\left(\frac{3}{2}\right)$	$\frac{1}{7} \log\left(\frac{3}{2}\right)$	$\frac{1}{7} \log\left(\frac{3}{2}\right)$	$\frac{1}{7} \log(3)$

Freq ↑ → weightage ↓

Drawbacks

① Curse Of dimensionality

100 → movie
 → not

②