

## Natural Language Processing (NLP)

input	output
I like to watch movies	+ve → 1
I lost my watch	-ve → 0

→ Machine can't understand the data, till now we deal

with the models which are having numerical no text &

no images.

→ So, we need to convert this need to numerical data.

then we can predict the output (0 or 1) → (Binary classification).

→ Converting Text to numerical algorithms are

1. BOW → Bag of words.
2. TF-IDF
3. Word2Vectorizer
4. BERT

Spacy & NLTK are the libraries for Text data  
SKlearn library for ML Models.

My dataset is having Reviews & Sentiments

Reviews (Input)	Sentiments (Output)
- - - - -	+ve (1)
- - - - -	-ve (0)
- - - - -	+ve (1)

Document      Document      Document

→ Each Row is a "document"

→ Collection of Documents is also known as Corpus.

1st step: Collect Data:-

Ex:-

Input
It was the best of times it
It was the worst of times
It was the age of wisdom
It was the age of foolishness.

} each row is a document.

1

2nd step:- Learn Vocabulary

→ First we need to understand the Vocabulary

↳ It's nothing but learning the ~~Vocab~~ unique words in the Corpus.

it	It	was	the	best	of	Times	worst	age	Wisdom
----	----	-----	-----	------	----	-------	-------	-----	--------

→ Vocabulary

→ Now update the Input data with "Count" with the Vocabulary.

Lowcase  
Uppercase  
(different)

bit	IT	Was	the	best	of	times	Worst	age	Wisdom	Foolishness.
2	0	1	1	1	1	1	0	0	0	0
1	1	1	1	0	1	1	1	0	0	0
0	1	1	1	0	1	0	0	1	1	0
0	1	1	1	0	1	0	0	1	0	1

00  
Count  
ent  
Vector

↳ Now, Here we are showing the Count • "Count" of the Document Term Matrix

each word.

↳ And we converted it into Numerical Form.

↳ Each row is called Document Vector

(Converted Document to Document Vector)  
 ↓ text                          ↓ Numerical

↳ Collection of Document Vector is called Document Term Matrix.

→ In the above example, we have 1st two columns i.e IT, if

↳ Both are same, Unnecessarily increased the Dimensionality (more no. of Columns)  
 (IT, IT)

↳ We can convert it into a Standard case then  
 (Upper (or) Lower)

then the no. of columns are reduced.

## Stop Words:

→ If we look into the Dataset (each document)

↳ we have unnecessary words which are used for sentence formation. These words are called stop words.

Formation. These words are called stop words.



This words won't make any impact on the the Dataset (~~dataset~~, Document Term Matrix)

↳ best, times, worst, age, freedom, wisdom, foolishness

all these ↓  
↓ are not stopwords.

This words will define my input

If, was, on, of, at, ...

↳

all these are stopwords

↳ It won't have any impact.

→ We can Remove these words. otherwise if we keep this

Stop words then My dimensionality will increase.

→ And, I

Stemming:-

↳ Converting the Word into root word.

→ 1. I am the happiest person right now

2. I am happy person right now

happy & happiest  
↓      ↓  
            Happ (trimmed)

Both are of same words.

↳ If we trim the words then it will be

good

↳ In a single word (happ)

So, my dimensionality will reduce.

Lemmization & Stemming → both are same

↳ But the words are converted into meaningful form in Lemmitization

↳ In above example Lemmitization → is → Happ.

Remove special characters.

↳ sometimes in the Corpus we have ; @, #, ... ; ? , !

↳ all these are special words, it won't make any impact to the Corpus.

↳ These special characters are for our purpose like

tagging some particular location, we use @<sup>hyderabad</sup>, ~~ben~~

↳ we remove @ & mentioned hyderabad is also same.

→ So, we can remove special characters.

Step-1 ↳ 1) Remove special characters

→ ↴

2) Convert to lower case

3) Remove the stop words.

4) Stemming / Lemmatization

Preprocessing steps

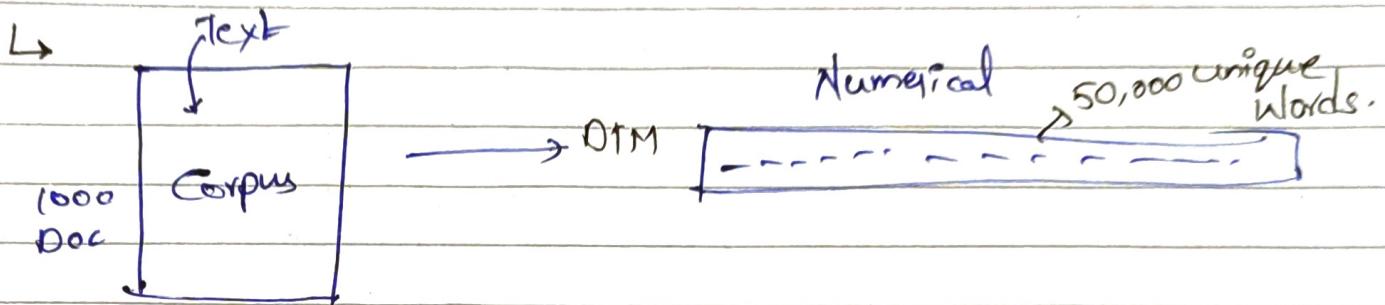
for Text preprocessing

(Cleaning the Data)

Step-2: Learn the Vocabulary

Step-3: → Convert Corpus to Document Term Matrix

↳ Suppose if my Corpus is having 1000 Doc. and it has 50,000 unique words.



↳ 50,000 columns is Very Very Very high & we can't run ML Model on this

↳ our accuracy will gonna reduce a lot.

↳ That's why we need to do Text preprocessing.

(Note:- But the issue is this BowL is not preserving the order)

Ex:- 1) I do have a house?  
2) Do I have a house?

I	do	have	a	house
1	1	1	1	1
1	1	1	1	1



Here both the statements are different, but my Document Vector is same.

→ So, we need preserve the sequential order

→ In this <sup>case</sup> we uses 1-gram (1word)