

①

## Natural Language processing (NLP)

↳ To make machine learn  
natural language.

→ Real World Applications of NLP

- Contextual Advertisement.
- Email Clients - spam filtering,  
Smart reply.
- Social Media - removing adult  
content, opinion mining.
- Search Engine
- Chatbots.

→ Common NLP Tasks

- Text or Document Classification.
- Sentiment Analysis.
- Information Retrieval.
- Parts of speech Tagging.
- Language Detection and  
machine translation.
- Conversational Agents.
- Knowledge graph and QA  
Systems.
- Text Summarization.
- Topic Modelling.
- Text Generation.
- Spell checking and grammar  
Correction.
- Text Parsing and Speech to text

→ Approaches to NLP:

↳ Heuristic Methods.

→ ML based Models or Methods.

→ DL based Methods.

→ Heuristic Approach

↳ Examples:

→ Regular Expression

Can help remove html tags from scrapped text.

or it can also help to search some specific words.

→ Wordnet

↳ Stores the relationship of word with other words.

↳ i.e.: Run ↔ jogg

↑  
Shoe ⇒ (you need shoe while running).

→ Open mind Common Sense

↳ you can store Common Sense facts from all over the world here.

→ Advantages:

↳ You get quick response

→ Less error

→ NLP (Pipeline)

↳ Set of steps followed to build an end to end NLP software.

NLP software consists of following steps:

↳ Data Acquisition.

↳ Text preparation.

↳ Text cleanup

→ Basic Preprocessing

→ Advance Preprocessing

→ Feature Engineering

↳ Converting text into numbers

→ Modelling

↳ Applying actual Algorithm

↳ Model Building.

→ Evaluation.

→ Deployment

↳ Deploy

↳ Monitoring

↳ Model Update

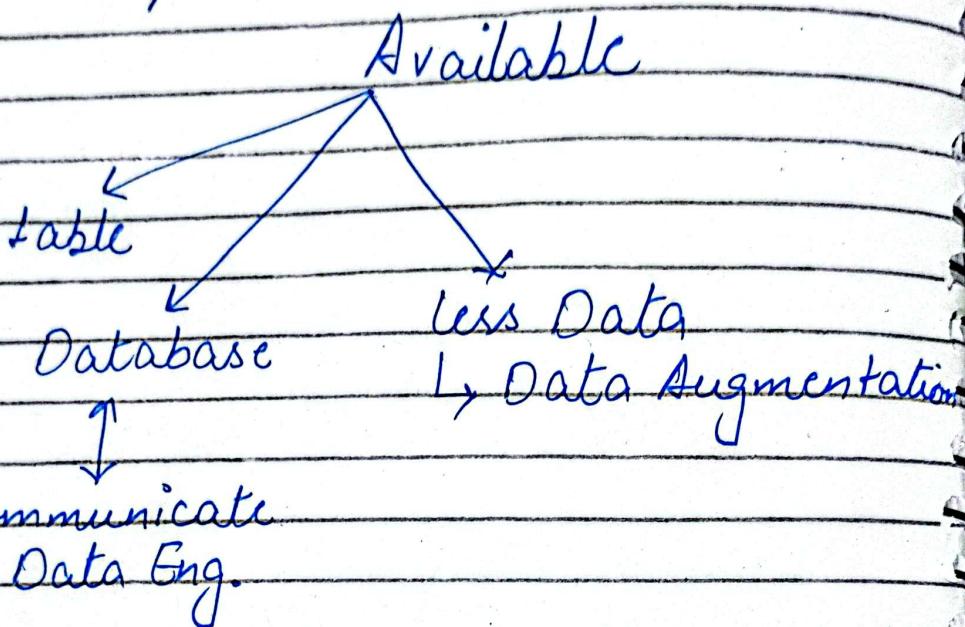
## Data Acquisition:

→ Available:

↳ on the table

→ Others (External Data).

→ Nobody.



### Others

→ Public Dataset (Kaggle)

→ Web Scrapping (Beautiful Soup).

→ API → Rapid API

↓ request  
+

Data in JSON

→ Pdf → use Python libraries.

→ image → use Python libraries

→ Audio → audio to text using Python libraries.

## Text Preparation:

- ↳ Cleaning
- ↳ Basic preprocessing
- ↳ Advance preprocessing

### → Cleaning text:

- ↳ Cleaning Emojis is called unicodenormalization
- ↳ Converting emojis into machine language

### → Checking spelling:

incorrect = " - - - "

```
from textblob import TextBlob  
txt_blo = TextBlob("incorrect")  
txt_blo.correct()
```

— .

### → Basic preprocessing

↳ Tokenization ↳ Basic

↳ Sentence

→ Words

→ OPTIONAL

↳ Stopword removal

→ Stemming

→ Removing punctuations, digits

→ Lowercasing or uppercasing

→ Language detection.

→ Tokenization:

↳ Converting text into tokens  
it could be a sentence or word.

```
from nltk.tokenize import  
    sent_tokenize,  
    word_tokenize
```

→ Optional:

↳ Stop word removal:

We use some words in English whose task is to form sentence but they don't have any contribution in meaning i.e.: a, an, for, the etc.

So if we are making NLP application from where we want to derive a context or meaning so there these stop words will not help and we remove them.

→ Stemming

If we have different words for the same meaning i.e.: dance, dancing and danced, so we can bring these words in root form.

↓  
their

Diff b/w Stemming and Lemmatization?

→ Lowercasing or uppercasing the txt:  
↳ So that our text remains consistent.

→ Advance Preprocessing  
↳ POS tagging (Parts of speech)  
→ Parsing  
→ Co-reference resolution

⇒ POS tagging can only be done if we have not used "word removal"

Note: There is a technique in Deep learning that is mostly used nowadays called Transfer learning.

↳ you take advance DL model that is trained on very Big data. i.e.: BERT Transformer.

↓  
it is trained with 40GB of data already.

(3)

## → Text Pre-processing (with Code implementation)

### ① Lowercasing:

Suppose if we have data then  
↳ data.iloc[3, 0].lower()

### ② Remove HTML Tags:

↳ Do it by using regex

```
def remove_html(text):  
    Pattern = re.compile('<.*?>')  
    return Pattern.sub(r'', text)
```

```
clean_text = lower_case.apply(remove_html)
```

```
Print(clean_text).
```

### ③ Remove URL's:

```
def remove_url(text):  
    Pattern = re.compile(r'^(https?:\/\/)?(www\.)?([A-z0-9]+(\.[A-z0-9]+)*(\.[A-z]{2,})?)$')  
    return Pattern.sub(r'', text)
```

## ④ Chat word treatment:

chat words = {

'ASAP': 'As Soon as Possible',  
'FYI': 'For Your Information',

}

def chat conver(text):

new\_text = []

for w in text.split():

if w.upper() in chat words:

new\_text.append(chatwords[w.upper()])

else:

new\_text.append(w)

return " ".join(new\_text).

## ⑤ Removing Stop words

↳ Do it if you are doing

↳ Document

→ Use NLTK library

Analysis

→ Sentiment  
Analysis

a the of are

But you don't remove them if you want to do POS (Parts of speech tagging).

```
def remove_stop(text):
```

```
    new_text = []
```

```
    for w in text.split():
```

```
        if w in stopwords.words('english'):
```

```
            new_text.append('')
```

```
        else:
```

```
            new_text.append(w)
```

```
x = new_text[:]
```

```
new_text.clear()
```

```
return " ".join(x).
```

## ⑥ Handling Emojis

→ 2 options

↳ remove emoji

↳ replace emoji with its meaning

## ⑦ Tokenization:

We can do it by using:

→ Split function

→ regular expression

But they might have some problems.

→ Best way is to use libraries  
use NLTK:

```
from nltk.tokenize import  
    word_tokenize,  
    sent_tokenize
```

These 2 functions have internal algorithms that are ready to handle situations like:

- Prefix
- Suffix
- Infix
- Exception

→ Another Best tokenization technique is given by Spacy library it is better than nLTK.

```
import spacy  
nlp = spacy.load('en_core_web_sm')
```

Sent = "We're here to help! mail us at nks@gmail.com".

```
d = nlp(Sent)
```

```
for i in d:  
    print(i)
```

## ⑧ Stemming:

{ walk → stem  
walking  
walks  
walked

Majorly used when we are making I R (information retrieval) system.

→ Use nLTK library

i.e: Google → we'll use porter stemmer in this library. ↓

Snowball Stemmer

Algorithm

→ for other language

## ⑨ Lemmatization

It searches (slow)

↳ Same as Stemming But it will take only specific language word as root word.

Lemmatization is slower than stemming as it tries to find the actual word of a language. (Use it if you want to show the output to user.)

We use Wordnet lemmatizer

↳ It has stored the English language words that are related to each other.

(4)

## Feature Extraction (In NLP)

↳ we extract features from text after converting text into numbers.

→ Core Idea? (Goal)

↳ whenever you are converting from text to numbers then this conversion or set of numbers should tell the meaning (Semantic) of something.

↳ hidden meaning should be convey through numbers

→ Techniques (for converting txt to num)  
(with semantic meaning)

→ One hot Encoding

→ Bag of words

→ n grams

→ Tf Idf

→ Custom features

→ Word 2 Vec {Embedding}

→ Common Terms

↳ Corpus

→ Vocabulary

→ Document

→ Word

Corpus  $\Rightarrow$  Combination every word in whole dataset. (C)

Vocabulary  $\Rightarrow$  Combination of every unique word in corpus. (V)

Document  $\Rightarrow$  Individual block in a feature or column. (D)

Word  $\Rightarrow$  Individual word in a document.

$\rightarrow$  One-hot Encoding

D<sub>1</sub> People watch movie

D<sub>2</sub> rafay watch movie

D<sub>3</sub> People write letter

D<sub>4</sub> rafay write letter

Corpus  $\Rightarrow$  People watch movie rafay  
watch movie people write  
letter rafay write letter

Vocabulary  $\Rightarrow$  Corpus (List of words in) (~ L)

People watch movie rafay letter  
 $V = 5$

Now OHE will make 5 features of every unique word.

write

People watch movie rafay letter

D <sub>1</sub>	1	1	0	1	0	0
D <sub>2</sub>	0	1	0	1	1	0
D <sub>3</sub>	1	0	1	0	0	1
D <sub>4</sub>	0	0	1	0	1	1

$$D_1 = [[1 \ 0 \ 0 \ 0 \ 0 \ 0], [0 \ 1 \ 0 \ 0 \ 0 \ 0], [0 \ 0 \ 0 \ 1 \ 0 \ 0]]$$

$$D_2 = [[0 \ 0 \ 0 \ 0 \ 1 \ 0], [0 \ 1 \ 0 \ 0 \ 0 \ 0], [0 \ 0 \ 0 \ 1 \ 0 \ 0]]$$

$$D_3 = [[1 \ 0 \ 0 \ 0 \ 0 \ 0], [0 \ 0 \ 1 \ 0 \ 0 \ 0], [0 \ 0 \ 0 \ 0 \ 0 \ 1]]$$

$$D_4 = [[0 \ 0 \ 0 \ 0 \ 1 \ 0], [0 \ 0 \ 1 \ 0 \ 0 \ 0], [0 \ 0 \ 0 \ 0 \ 0 \ 1]]$$

One-hot encoding

Pros  $\Rightarrow$  Intuitive

$\rightarrow$  Easy to implement  
~~(Efficient)~~

Flaws  $\Rightarrow$  Sparsity  $\Rightarrow$  Sparse array  
causes overfitting.

- No fixed size.
- OOV (out of vocabulary)
- No capturing of semantic meaning.

### → Bag of Words

↳ Specially used when we are working on text classification problem. i.e., News Document

↳ Dividing it in some categories like: sport news, Political news etc.

- D<sub>1</sub> People watch rafay
- D<sub>2</sub> rafay watch rafay
- D<sub>3</sub> People write comment
- D<sub>4</sub> rafay write comment

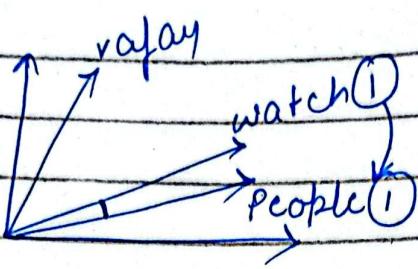
BOW:

People watch rafay write comment

D <sub>1</sub>	1	1	1	0	0
D <sub>2</sub>	0	1	2	0	0
D <sub>3</sub>	1	0	0	1	1
D <sub>4</sub>	0	0	1	1	1

Suppose if we have only 3 features then:

People watch rafay



Giving Sentence a vector form  
then checking the angle with the  
closest one.

→ Pros

→ Simple and intuitive

→ Disadvantages

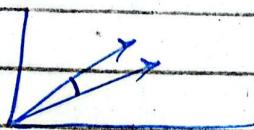
→ Sparsity

→ OOV

→ Un Order

→ This is a good movie

This is [not] a good movie



→ Bag of n-grams

Pros

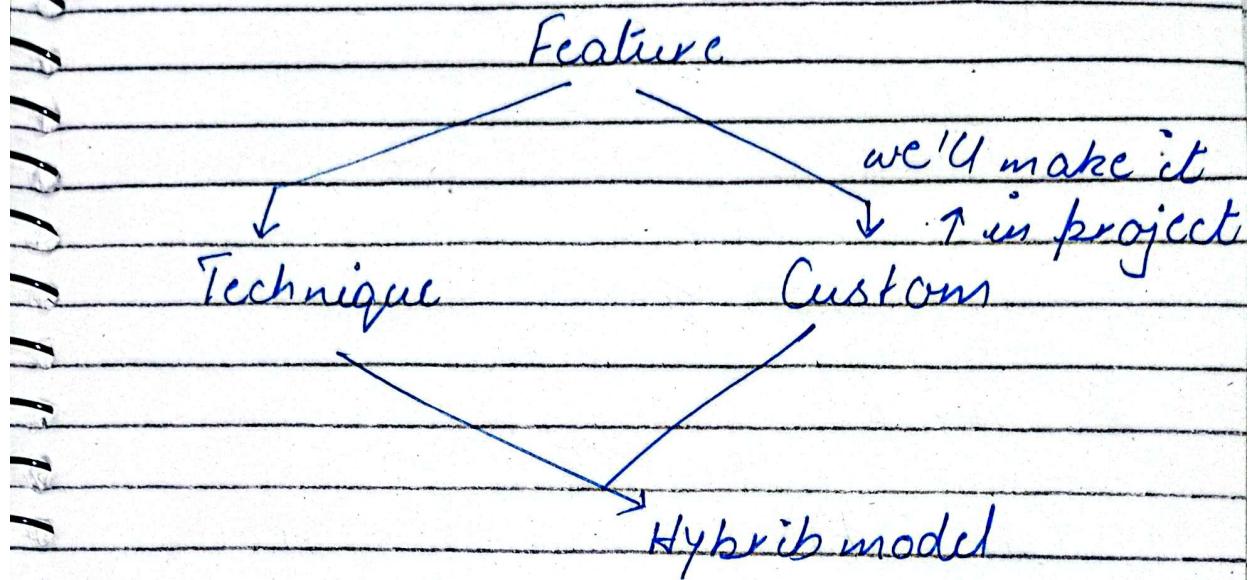
→ Able to capture Semantic meaning.

→ Easy implementation → OOV

Cons

→ As we increase n-grams the dimension of vector also increases due to which machine takes time to learn.

→ Custom Feature:



→ Word Embeddings

↳ BOW

↳ grams  
→ Tf Tdf

Frequency

Base

In NLP, word embedding is a term used for the representation of words for text analysis, typically in the form of a real valued vector that encodes the meaning of words that are closer in vector space are expected to be similar in meaning.

↳ In Short ⇒ Converting words into vector

→ 2 types ⇒ Frequency based  
Prediction based

✓ 5

Word2Vec (uses DL)

↳ Prediction base word embedding

## Similarity

Can see (difference) between words that have same meaning

i.e.: happy  
joy

## Advantages

- Can Capture Semantic meaning.
  - Improve sparsity (Compute fast).
  - Dense vector. (most words are non zero).  $\Sigma \otimes \Sigma \otimes \Sigma$

	King	Queen	Man	Woman	Monkey
King	1	0	1	1	1
Queen	0	1	0.3	0.3	0
Man	1	0.4	0.2	0.2	0
Woman	0	0.4	0.6	0.5	0.3
Monkey	1	1	1	1	1
<hr/>					
King =	1	1	1	0.8	1
King - man + Queen	-1	+0	-0.3	+0.3	=0
Man	1	-0.3	+0.3	-1	0
Queen	0	+0.2	-0.2	=1	0
Monkey	1	-0.1	+0.5	=0.7	1
<hr/>					
close to					
Queen					

→ Types of Word2Vec

CBOW

Skip gram

→ CBOW (Continuous Bag of words)

fake problem.

↳ Solve

↳ Vector

↳ by product

Watch ML for Data Science

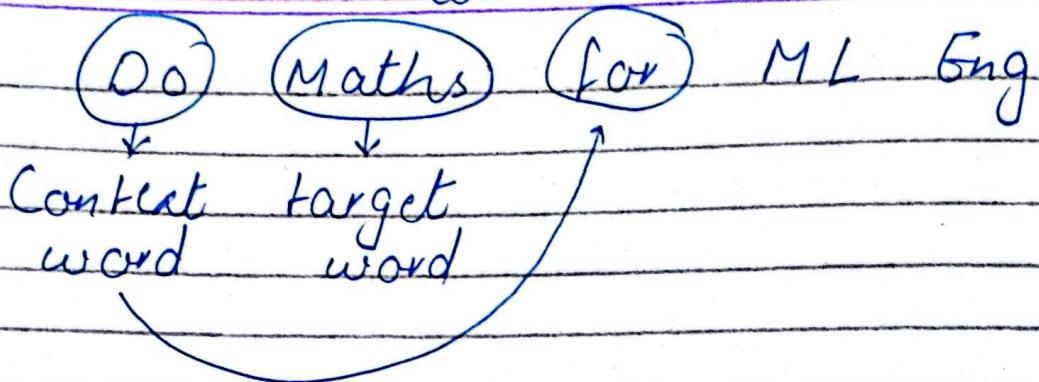
Context target 3 words  
word word

X for Y  
watch, (ML) ML  
ML, Data for  
for, Science Data

→ Skip gram (Opposite of CBOW)

↳ given target words you have  
to predict the context word (Dummy  
Problem).

window = 3



X	Y
Maths	Do, for
for	Maths, ML
ML	for, Eng

use NN:

Maths I/P

0		
1		
0		
0		
0		

5x3

3x5

3x5

0	1
0.2	0
0	0
1	0
0	0

for O/P

0	0
0.3	0
1	1
1	0
0	0

loss

→ When to use CBOW / Skip gram

Small data

↳ CBOW

Large data

↳ Skip gram

→ How to improve Word 2 Vec

↳ Increase training data.

→ increase dimension of hL vector

→ increase window size for words.

→ Training our own model

(Using Word 2 Vec)

↳ by gensim lib

## → Text Classification:

### ↳ Types

↳ Binary Classification

→ multi class classification.

→ multi label classification

↳ For multiple input

you'll have multiple output.

↳ example → one particular

news can come under multiple  
news

## → Applications

↳ Gmail spam classification.

→ Customer Support.

→ Sentiment Analysis.

→ Language Detection.

→ Fake news Detection.

## → Pipeline

→ Data Acquisition

→ Text Preprocessing

→ Feature Extracting

→ Modelling

→ Evaluation

→ Different approaches to do text classification

↳ Heuristic approach

→ APIs

→ ML

→ DL

→ APIs:

↳ Amazon and Google have their own cloud platform that have built solutions for most frequently occurring problem.