

# NLP

DELTA / Page No.

Date

## # Regular Expression

A RegEx, or regular Expression, is a sequence of characters that forms a search pattern.

RegEx can be used to check if a string contains the specified search pattern.

→ Python has a build-in package called re, which can be used to work with Regular Expressions.

Demo ↴

Import re

txt = "The rain in Spain"

x = re.search("The.\*Spain\$", txt)

There is many more function in "re".

- 1) findall — (Return a list containing all matches)
- 2) search — Return a Match object if there is a match anywhere in the string.
- 3) split — Returns a list where the string is at each match
- 4) sub — Replaces one or many matches a string

## Word to vector

DELTA / Page No.

Date

processing natural language text and extract useful information from the given word, a sentence using machine learning and deep learning techniques requires the string/text needs to be converted into a set of real numbers called vector.

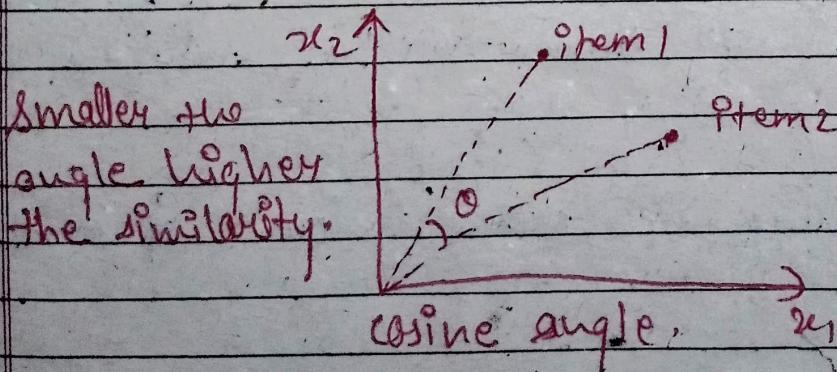
- Now word embedding or word vectorization is a methodology in NLP to map words or phrases from vocabulary to a corresponding vector of real numbers which used to find word predictions, word similarities, semantics.

Cosine similarity approach is used to find the similarity between the documents.

### → Cosine Similarity:

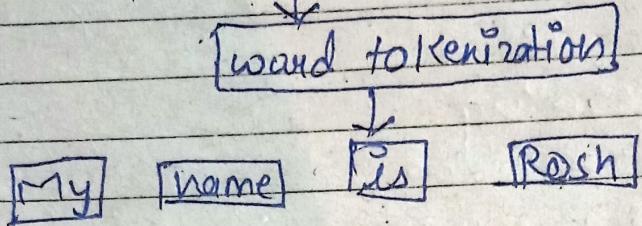
It measures the cosine of the angle b/w two vectors ( $\text{item}_1, \text{item}_2$ ) projected in an N-dimensional vector space.

It predicts the document similarity.



1) **Tokenization:** segmenting the text into a list of tokens.

Ex:- My name is Rosh



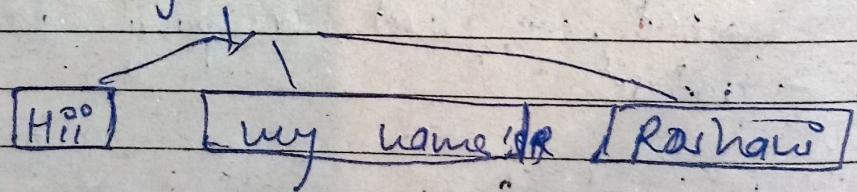
2) **Lowercase!** convert all text into lower case

3) **Stop word removals** it is used to apart the sentences to extract the meaningful explanation.

→ sentence stoppers like " , or . (full stop)"

it is kind off problem in it. ~~for example~~

Hii, my name's Roshni.



for that we use:-

4) **Stemming or Lemmatization:-**

Stemming used to stripping the suffixes from words to get their stem,

→ I [love eating] pizza

→ I. love eat pizza.

→ eating and eat is similar. but not eat and ate.  
lemmatization involve reducing words to their base form based on their part of speech.

example :- eat and ate (both are kinda similar)

### 5) Removing digit (punctuation)-

Removing digits and punctuation from the text.

### 4) feature Engineering:-

feature engineering represent the text in the numeric vector in such a way that the ML algorithm can understand the text attribute.

some of the methods are.

1) TF-IDF vectorizer

2) one hot encoding

3) word Embedding.

### 5) Model Building.

build the model for identify or predict the result.

### 6) Evaluate

Test the model, it's corrections, faults.

If it is able to correct. Then, we must move to the text processing step.

## # Spacy & NLTK

Spacy is object oriented. we need to create an object of that string then work on it.

It is user friendly. It is a kind of in-built feature which provides best algorithm and gives best result.

It provides most efficient NLP algorithm for a given task. Hence if you care about the end result, go with spacy.

→ Code :-

→ import spacy

nlp = spacy.load("en\_core\_web\_sm")

doc = nlp("I am Rosh. This is my code.")

for sentence in doc.sents:  
    print(sentence)

Output → I am Rosh.

This is my code.

for sentence in doc.sents:

    for word in sentence:  
        print(word)

Output → I              This

am

in

Rosh

my

code

NLTK is a string processing library. we directly pass the string and perform text preprocessing.

NLTK is also user friendly but probably less user friendly compared to spaCy.

provides access to many algorithms. if you care about specific algo and customizations go with NLTK

~~NLTK~~ NLTK use - by mostly researchers.

Code :-

```
import nltk  
nltk.download("punkt")  
from nltk.tokenize import sent_tokenize  
sent_tokenize("I am Rosh. This is my python")
```

Output :- I am Rosh.

This is my python.

## # NLP Language processing pipeline in spaCy:

- 1) NLP pipeline is a series of steps that transform a raw text data into a desired output like :- label, summary, or a response

2) This pipeline performs a specific function, such as tokenization, lemmatization, sentiment analysis, named entity recognition.

→ Design a Pipeline.

To design a pipeline, define goal (output)  
your data  
your tasks  
your model.

We can either import blank pipeline (with no functions) or we can import pre-trained pipeline (with all functions).

→ blank pipeline

nlp = spacy.blank("En")

nlp.pipe\_names

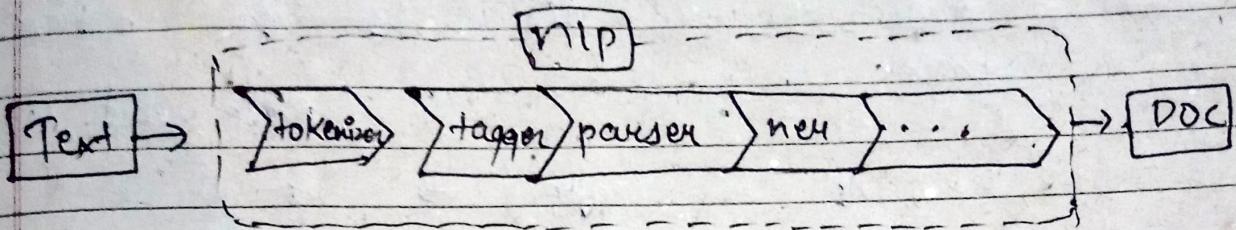
output ⇒ []

→ pre-trained pipeline.

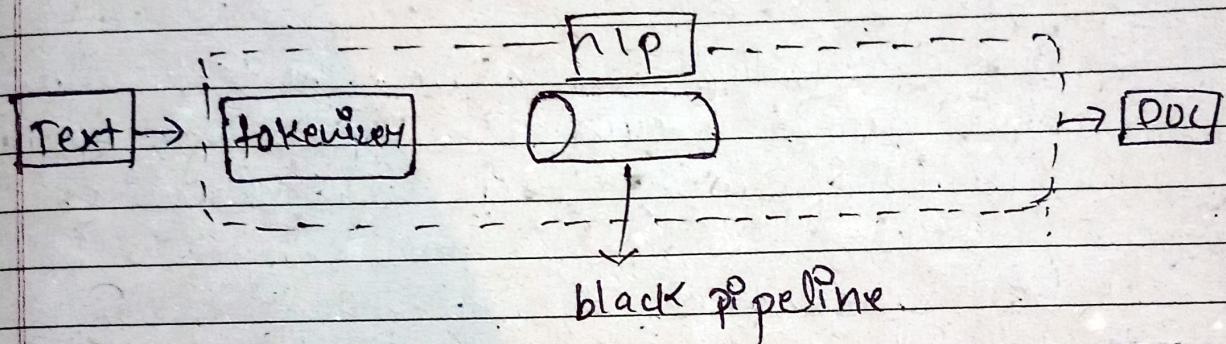
nlp = spacy.load("En\\_core\\_web\\_sm")

nlp.pipe\_names

output ⇒ ['tok2vec', 'tagger', 'parser', 'attribute\_ruler', 'lemmatizer', 'ner']



dig.: nlp pipeline with functions.



→ add component in black pipeline.

nlp = spacy.blank ("en")

nlp.add\_pipe("ner"); ~~source = source-nlp~~  
nlp.pipe\_names.

output → ['ner']

we can add functions individually in pipeline.

# stemming and lemmatization function in pipeline

stemming and lemmatization in nlp is play a crucial role.

→ It reduce words to their base or root forms. (Base word).

example → ate → eat

talking → talk

adjustable → adjust

→ This technique normalize the text, allowing far more accurate analysis, information retrieval, and lang. understanding.

let's understand it deep:-

### ⇒ stemming (stem)

- 1) It involve reducing words to their base or root form, known as stem.
- 2) It remove prefixes, suffixes, and inflections from words.
- 3) help in reducing vocabulary size, normalizing text, and improving text analysis task.

### ⇒ lemmatization (lemma)

- 1) It is a technique that involves reducing words to their base or dictionary form, known as lemma.
- 2) Lemmatization takes into account the word's part of speech (POS) and context to deter-

wined its canonical form.

ex:-

[lemma. of [ORG]]

tesla | ORG | like - company, organization,  
industrial group.

⇒ stemming and lemmatizing in nlp

import nlp

from nlp import PorterStemmer  
stemmer = PorterStemmer()

→ stemmer.stem("helping")

output ⇒ 'help'

⇒ In spacy

import spacy

nlp = spacy.load("en-core-web-sm")

→ doc = nlp("helping helps helped")

for token in doc:

print(token, "||", token.lemma\_)

Output ⇒ helping || help  
helps || help  
helped || helped.

# modification in lemmatization dictionary :-

→ attribute\_ruler :-

It is used to add an addition word as according to your preference.

→ code :-

```
att = nlp.get_pipe("attribute_ruler")
```

```
att.add([{"text": "helping"}, {"text": "helped"}],
```

```
{ "LEMMA": "help" } ]
```

```
doc = nlp("helping helps helped")
```

```
for token in doc:
```

```
    print(token, "||", token.lemma_)
```

Output:-

helping || help

helps || help

helped || help

## # Part of Speech (POS).

It is a process of converting a sentence to forms - list of words, list of tuples (where each tuple is having a form (word, tag))

gt signifies whether the word is a noun, adjective, verb, and so on.

POS	Tag
Noun	n
Verb	v
Adjective	a
adverb	t

⇒ Some of the function in POS

1) token-pos -

provide the type of token it is (noun, verb, pronoun).

2) spacy-explain (explained-junction-provided).

explain the function.

3) token-tag -

it provide the tag given to token (VBP, NN, CC)

4) doc-count-by (spacy. attrs. POS)

count the no. of different token in sentence.

## # NER (Named Entity Recognition)

It is one of the most popular data processing.

It identify the key information in the text and classified it in a predefined categories.

Some of the categories are:-

- 1) Person
- 2) Organization
- 3) place/ location
- 4) date / time
- 5) expression
- 6) money, percent, weight, etc)
- 7) E-mail address;

Code :-

```
nlp = spacy.load("en_core_web_sm")
```

```
doc = nlp("Twitter Inc. is worth of $48 million")
```

for ent in doc.ents:

```
print(ent.text, "ent.label_ ", spacy.explain(ent.label_))
```

→ output

Twitter Inc. Arg.

\$48 million Money.

## ⇒ feature Engineering section

### # Text representation:-

Text representation is a crucial aspect of NLP that involves converting raw text data into machine-readable form using one-hot encoding, Bag of words, TF-IDF, word embedding.

→ Limitation of one-hot-encoding / Label coding.

- 1) No fixed length representation.
- 2) out of vocabulary, if new word is added
- 3) consume too much memory & compute resources
- 4) similar words do not have similar representation.

⇒ due to this limitation we use Bag of words, TF-IDF, word embedding techniques.

### 1) Bag of words (BOW)

Each word in a text considered as a feature, and the number of times that word appears in text represent the importance of the word.

Disregarding grammar and word order but

Keeping track of the frequency of each word.

ex:- the cat in the hat

the dog in the house

dog	cat	bird	in	house	hat	the
The cat in the hat	0	1	0	1	0	1
the dog in the house	1	0	0	1	1	0

Raw text



Number vector



Machine learning

1) CountVectorizer(), (Count the no. of words appearing)

2) Pipeline

→ Code

From sklearn.pipeline import Pipeline

pipe = Pipeline([

① [ ] ('vectorizer', CountVectorizer()),  
 ('nb:', MultinomialNB())  
 7)

② - pipe model .fit(x\_train, y\_train)

1) ~~train~~ train model using pipeline for text training. (NLP)

2) Train dataset using pipe-model function by method "fit".

## # Stop words.

It is a set of commonly used words in a language.

Ex:- "a", "the", "is", "are", etc.

Commonly used in text mining and NLP to eliminate words that are so widely used that they carry very little useful information.

Wikipedia is great example for this, it provides keyword related to person or organization using this feature.

→ Limitation :-

chatbot, AI can't use this feature cuz  
it remove the meaning of the sentence.

ex:- I am using laptop.

after using feature.

output → Laptop { This is useless and meaning less }

We can't use stop word feature in chatbot;  
① in system, language translation etc

⇒ Code

## ⇒ N-GRAMS in BOW

It is a set of words, so that word vectorization has more meaningful.

Useful in many text analytics applications where sequence of words are relevant, such as in sentiment analysis, text classification, and text generation.

ex:- Hello, this is my python code.

1-gram. 'Hello', 'this', 'is', 'my', 'python', 'code'

2-gram 'Hello, this', 'this is', 'is my', 'my python', 'python code'.

## ⇒ TF-IDF (Term frequency - Inverse Document freq)

- determines the importance of word in a document
- identify key terms and concepts in text.
- Improve text classification and clustering.

### 1) term frequency (TF):

how frequently a term appears within a document.

$$TF(t, d) = \frac{\text{Number of times term 't' appears in doc}}{\text{Total terms in document 'd'}}$$

### 2) Inverse Document Frequency (IDF):

Measure how rare a term is across the entire collection of document

$$\text{IDF}(t) = \log \left( \frac{\text{Total number of doc}}{\text{No. of doc. with term}} \right)$$

### 3) TF-IDF

→ combine TF and IDF to assign a weight to each term in a document.

⇒

$$\rightarrow \text{TF-IDF}(t, d) = \text{TF}(t, d) \times \text{IDF}(t)$$

→ Higher TF-IDF scores indicate more important terms.

## # FastText Embedding :-

FastText is a word embedding that provides embedding to the character n-grams. It is the extension of the word2vec model.

word  $\Rightarrow$  cat

1-gram = "C", "a", "t"

2-gram = "Ca", "at"

3-gram = "cat"

→ It's advantages:-

- i) Better handling of rare and OOV words :  
It can generate vector representations for

words it has never seen before by simply summing the vectors (n-gram).

- OOV → out of vocabulary.
- Fasttext is often a first choice when you want to train custom word embeddings for your domain.
- It is a technique as well as a library.
- fasttext is pre-trained model on wikipedia.