

# STOPWORDS :-

Stopwords are common words in a language that usually do not add important meaning to the text, so they are often removed during the NLP processing.

Example of Stopwords :-

is, am, are, the, a, an, in, on

Example :-

Sentence :- This is a machine learning project

After Removing Stopwords :-

machine learning project

Syntax :-

```
from nltk.corpus import stopwords  
import nltk  
nltk.download('stopwords')
```

]- downloading Stopwords from NLTK.

# we use stopwords.words ('english')

to get stopwords for english language

Main Code :-

paragraph = "God of war --- --- ---  
----- September 2024."

from nltk.tokenize import sent\_tokenize

first we convert paragraph into  
Sentences using sent\_tokenize

Sentences = sent\_tokenize(paragraph)

Storing all sentences in object → 'sentences'

from nltk.tokenize import word\_tokenize

will use word\_tokenize to convert  
all sentences into separate words

from nltk.stem import PorterStemmer

Stemmer-object = PorterStemmer()

So here we are using Porter Stemmer  
also, to reduce words.

from nltk.corpus import stop\_words

import nltk

nltk.download('stopwords')

→ importing  
Stopwords

```
for i in range (len(sentences))  
    words = word_tokenize (sentences[i])  
    # sentences[i] → take one sentence  
    # word_tokenize → split sentence in words  
    words = [stemmer_object.stem(word) for word  
             in words if word not in  
             set(stopwords.words('english'))]
```

# using list comprehension here

for word in words → check each word  
if word not in stopwords → skip common words  
stemmer\_object.stem(word) → shorten the word  
[---] → making a new list

Sentences[i] = " ".join(words)

Joining all words together in one sentence

print (sentences)

Note :- here we are using porter stemmer,  
we can also use Regexp stemme and snowball  
stemmer but to get the best out of all  
we should use lemmatizer.

Main Code Using Lemmatization :-

paragraph = { God of War -- 2024, }

# converting into sentence

```
from nltk.tokenize import sent_tokenize  
sentences = sent_tokenize(paragraph)
```

# import word\_tokenize and lemmatizer

```
from nltk.tokenize import word_tokenize  
from nltk.stem import WordNetLemmatizer  
import  
lemmatizer = WordNetLemmatizer()
```

# import stopwords

```
from nltk.corpus import stopwords  
import nltk  
nltk.download('english_stopwords')
```

for i in range(len(sentences)):

words = word\_tokenize(sentences[i])

words = [lemmatizer.lemmatize(word, pos) for word in words  
if word not in set(stopwords.words('english'))]

Sentences[i] = " ".join(words)

# in this we are using  
lemmatizer. lemmatize(word, pos='v')

# POS-TAG

POS (Part of Speech) tagging means assigning grammatical labels to words. like

- Noun
- Verb
- adjective
- adverb, etc

Pos tagging in NLTK :-

~~#~~ `nltk.pos_tag()`

Example :-

```
import nltk  
from nltk.tokenize import word_tokenize
```

Sentence = "Taj Mahal is a beautiful Monument"  
words = word\_tokenize(sentence)

pos-tag = nltk.pos\_tag(words)

print (pos-tag)

~~#~~ [ ('Taj', 'NNP'), ('Mahal', 'NNP'), ('is', 'VBZ'), ('a', 'DT') ... ]

## Common Pos tags

Tag	Meanings
NN	Noun
NNP	Proper Noun
VB	Verb
VBZ	Verb   3rd person singular)
DT	Determiner

## NAMED-ENTITY RECOGNITION

Named-Entity Recognition is an NLP technique used to identify and classify real world entities in text.

What Named-Entity Recognition detects -

PERSON → name of people

ORG → organisation

GPE / LOCATION → country, city, place

DATE / TIME

MONEY / PERCENT

Example :-

```
import nltk
```

```
from nltk.tokenize import word_tokenize
```

```
from nltk import pos_tag, ne_chunk
```

# ne\_chunk is used for Named Entity Recognition.

Sentence = "Ritik Sharma works at Google in India"

words = word\_tokenize(sentence)

postag = pos\_tag(words)

ner = ne\_chunk(postag)

print(ner)

Output :-

(S

(PERSON RITIK / NNP)

(PERSON Sharma / NNP)

works / VBZ

at / IN

(ORGANIZATION Google / NNP)

in / IN

(CPE India / NNP))

But, we can also draw this using

(.draw) or (ne\_chunk(postag).draw)

continue from there

→ ner = ne-chunk (pos\_tags).draw()  
print (ner)

output :-

