# Natural Launguage Procerring (NLP)

import nltk

nltk.doconload ( )

#note :- Mutiline paragraph can be ured by parting b/w `""""` --- `""""`

* To doconload the paragraph into sentences into we have to ure the below function

    nltk. sent_tokenize (Paragraph)

    Returns a list

* To downlaod the sentences into words

    nltk. wod_tokenize (Sentences)

    Returns a list.

* Instead of tokenizing a paragraph bared on 'space', we tokenize it bared on "." and ", ". Therefore, we get all the different sentences Conristing the paragraph

    Str = " I LOVE DIVYA "

    words = str. split(" ")

    Print (words)

    out: ["I", "LOVE", "NLP"]

For sentence tokenization fust replace str. split(" ") with str. split (".")

## Stemming

→ when Extracting sentences & corpus of sentences

    " John does his Woulc intelligently"
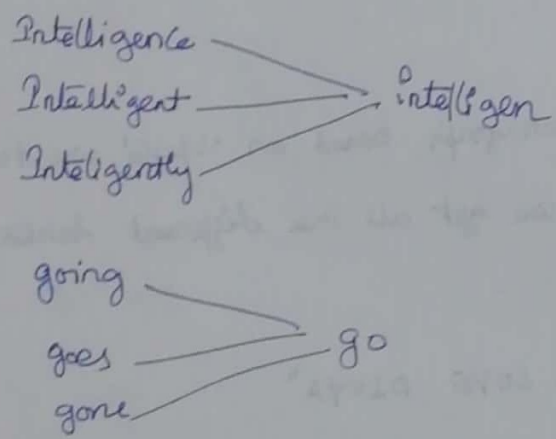
    " John is an intelligent man"

    " John is always woulcing"

tokenize sentence

| Sent 1 | setence 2 | Senten 3 |
|--------|-----------|----------|
| John | John | John |
| does | is | is |
| his | an | always |
| work | intelligent | working |
| Intelligently | man | |

→ We have Consider work and working as same

## Stemming

"Stemming is a process of reducing injected & derived words to their word stem, bare & root form"

Intelligence ⟶ intelligen
Intelligent ⟶ intelligen
Intelligently ⟶ intelligen

going ⟶ go
goes ⟶ go
gone ⟶ go

→ not duplicating the words which gives same meaning.
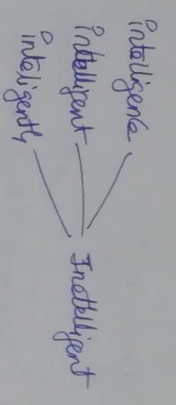
## Problem

Produced intermediate representation of the word may not have any meaning

Ex:- intellijen, fina etc
               ↑
             final

## Lemmatization

Same as stemming but intermediate representation/root form has a meaning.

Intelligence
Intelligent ⟶ Intelligent
intelligently

## Lemmatization | Stemming

* word representations have meaning | ✷ word representation may not have any meaning

* Take more time than stemming | ✷ takes less time

✷ we when meaning of words is important for analysis | ✷ we when meaning of words is not important for analysis

Ex: Question answering application (chatbot) | Ex: spam detection, Text classification, Sentiment analysis.

### for more ( Additional read)
https://nlp.stanford.edu/IR-book/html/htmledition/stemming-and-lemmatization-1.h

### how to implement stemming

```
import nltk
Paragraph = """ --------- """
```

from nltk.stem import PorterStemmer

sentences = nltk.sent_tokenize (Paragraph)
stemmer = PorterStemmer( )

for i in range (len (sentences))

    words = nltk.word_tokenize (sentences [i]);
    words = [stemmer.stem(word) for word in words]
    sentences [i] = ' ' .join (words)

Very ⟶ Veri
academy ⟶ academi
The ⟶ The

# Lemmatization Implementation

Import nltk

from nltk.stem Import WordNetLemmatizer

⋮

(Same code as before)

Except

words = [lemmatizer.lemmatize(word) for word in words]