

## Natural Language Processing Tokenization

↳ Importing NLTK library

↳ `import nltk` `nltk.download()` `nltk.download('punkt')`

↳ Paragraph 1 = "" I have three vision for India. In 3000 years of our history, people from all over the world have come invaded us, Captured our lands, Enslaved our minds. """

↳ `sentences = nltk.sent_tokenize(paragraph)`

↳ `Sentences = nltk.sent_tokenize(paragraph)`  $\rightarrow$  Sentence in the form of List

↳ Words = `nltk.word_tokenize(paragraph)`  $\rightarrow$  word

↳ Counts how many words

↳ What sentence wrote at appropriate place in it.

↳ Tokenization: Breaks into individual words.

- Taking the paragraph or corpus and then will be converting those paragraph into sentences.

## Stemming And Lemmatization

Stemming  $\rightarrow$  Lemmatization

history  $\rightarrow$  histon  
historical  $\rightarrow$  historical  
finally  $\rightarrow$  finally  
final  $\rightarrow$  final  
finalized  $\rightarrow$  finalized

Processes of Reducing Inflected words to their

Word STEM

history  $\rightarrow$  history  
historical  $\rightarrow$  historical

final  $\rightarrow$  final  
finalized  $\rightarrow$  finalized

going  $\rightarrow$  go  
goes  $\rightarrow$  goes

Tokenerization = Is the process of breaking text into pieces called tokens.

Date: \_\_\_\_\_  
Page: \_\_\_\_\_

etiology: Hemorrhage → primary generalized hemarthrosis

→ Use Car like Sentiment Analysis, Restaurant Review  
Speech classifier. Some of the word representation will  
not have the meaning. It does not take time to  
find out the word.

a) Lemmatization

## a) Linnearizations

↳ Basically it creates more "meaning forward"

~~We need~~ that can help us understand

~~an~~ writing, history → history

historical and no blots)

→ It takes time for process

→ use code :- chartbot ; Q & A application.

\* Stemming (Stop words), other : changed

Stopwords { we apply stopwords to remove these kind of words like that don't ~~has~~ plays much meaningful  $\Rightarrow$  is, of, this, I'll, ... etc (this kind of words are called articles and conjunction) helps to remove these type of words import nltk

→ from nikkō stem import Porter stemmers

→ from Nitro-Cropus impact afterwards

stomach (is there any?) present in middle part

metabolism / protein synthesis → protein synthesis

→ Paragraph 2 (cont'd. from p. I) --- " " "

> Sentences = nltk.sent\_tokenize (processing)

Der mit  $\rightarrow$  stammmer = PorterStammer()

~~positive with~~ ~~# for p, step by step~~ purpose is to  
~~p~~

## #1 Stemmetry

~~for i in range(len(sentences)):~~

`words = nltk.word_tokenize(sentences[i])`

words = [stemmer, stem (word)] for word in words

if word not in set(stopwords.words('english'))]

Sentences [i] = 1. join (words)  can be  
any language

Set helps us to take only the unique stopwords in the English language.

## Problem

→ Produces an intermediate representation of the work

~~may not have any meaning.~~

Sg :- Intelligent  $\Rightarrow$  Intelligent

final → final, etc

## \* Lemmatization

$\delta = \text{loop } (ii) \text{ from}$

$\delta = \text{path}(\text{iii})$

5 Learning objectives

?import nltk

```
from nltk.stem import WordNetLemmatizer
```

from nltk.corpus import stopwords

Then i good news start of print factory

3. Sentences = ntkc sent\_tokenize (paragraph)

Lemmatizer = WordNetLemmatizer() this is responsible

> for  $\frac{w}{n}$  in range ( $\log(\text{sentence})$ ):

Word = nicht-kontrahierbare Sätze zu Konzern (sentences [i])

`word = Lemmatizer.lemmatize(word)` for word in

- words if word not in set (stop words). word[eyn]

Sentences [i] = 1. join (words)

## \* Bag of Words

↳ ((Example 2) no2) approach in ai nlp

((i) brow, nose, mouth, eye, etc) are not words

show Sentence 1 → He is a good boy = sentence (good) (boy)

((i) brow, nose, mouth, eye, etc) are not words

Sentence 2 → She is a good girl = sentence (good) (girl)

Sentence 3 → Boy & girl = sentence (boy) (girl)

Creating a histogram  
basically what is the frequency, what is the total word count

	good	boy	girl	op
freq	3	2	2	0
vector	→	→	→	→
Bow	sent 1	sent 2	sent 3	independent
	1	1	1	0

Count (i) good = 3      \* for word

(ii) boy = 2

(iii) girl = 2

After training

→ terms like boy, girl, mouth, nose, etc

• Binary Vector → form: expression of term

Disadvantages - My values - are - either "1" or "0", but one

important thing to note over here is that

the good value = 1, the boy value = 1, both

are having equal representation. The semantics

of this particular words are almost same,

we are not able to determine which word

(i) is much more important. like brow

→ In Sentiment Analysis the major thing

is how to decide which word is more important

is that we ~~should~~ know that which word is more important like in this particular case Good should be more important than Intelligent.

e.g. - He is Intelligent.

Note:- To avoid this disadvantage of Bow we use TF-IDF

- Bow can be used in Sentiment Analysis
- But, If you have a huge Dataset Don't go with Bow, use word2vec to solve.

## \* Bag of Words

() document representation

(disjoint) bag of words

[] = bag of words

Count: (estimator) no. appearance of i-th word

(B) matrix of the form  $[S \cdot A \cdot D^T]$

() word-document matrix

() document-term matrix

() matrix of BoW, not (bow) matrix

() document-term matrix

## Bag of words ("Code")

import nltk  
from sklearn.feature\_extraction.text import CountVectorizer

Document matrix all together

Paragraph h = " " " - - - - - " "

# cleaning the texts

import re import regular expression for lowercasing

from nltk.corpus import stopwords

porter stemmer import PorterStemmer

wordnet lemmatizer import WordNetLemmatizer

Ps = PorterStemmer() # for stemming

Wordnet = WordNetLemmatizer()

Sentences = nltk.sent\_tokenize(paragraph)

Corpus = [ ]

for not

for i in range(len(sentences)):

review = re.sub('[^a-zA-Z]', ' ', sentences[i])

review = review.lower()

removing all the

commas, punctuation

space, question

marks  
not required

review = review.split() # list of words

review = [Ps.stem(word) for word in review]

if word not in set(stopwords.words('english'))]

corpus.append(review) # lemmatize word

Lemmatizer

# Creating the Bag of words model

from sklearn.feature\_extraction.text import CountVectorizer

cv = CountVectorizer(max\_features=1500)

X = cv.fit\_transform(corpus).toarray()  
creating a matrix

## \* TF-IDF

Sent 1 → He is good boy

after applying  
Stemming,

Sent 1 → good boy

Sent 2 → She is good girl

Lemmatization

Sent 2 → good girl

Stop words: P.

Sent 3 → Boy & Girl are

good

Extract S + P + V + Adj good

(Sent 1) → he good boy

she is good girl

(Sent 2) → she good girl

(Sent 3) → boy girl good

Sent 3 → boy girl good

Term Frequency (TF) = No. of repetition of words

→ word freq word in sentence

→  $\frac{\text{No. of words}}{\text{No. of words in sentence}}$

Inverse Document frequency (IDF)

$IDF = \log \left( \frac{\text{No. of sentences}}{\text{No. of sentences containing word}} \right)$

TF-IDF =  $TF * IDF$

Finally TF-IDF =  $TF * IDF$

TF-IDF =  $TF * IDF$

TF-IDF =  $TF * IDF$

words frequency

TF-IDF \*

good 3

bad boy 2 - good boy is often → freq.

girl 2 →

bad girl 2 → good girl is not → freq.

TF-IDF

vectors

word	TF	IDF	TF-IDF
good	1/2	1/2	1/3
boy	1/2	0	1/3
girl	0	1/2	1/3

good boy girl → 8 freq.

Set 1	good	boy	girl
sent 1	1/2 * log(3/3)	1/2 * log(3/2)	0 * log(3/2)
Set 2	1/2 * log(3/3)	0 * log(3/2)	1/2 * log(3/2)
(e.g. sent 3) 0.5 * 1/3 * log(3/2)	1/2 * log(3/2)	1/2 * log(3/2)	

Note: In Bow we have values in 1 and 0

In TF-IDF after computing feature (independent role)  
we will not get values like 1 and 0,  
we will get values in Decimal like 1.32, 4.7

\* TF-IDF ("code")

~~→ import ntk~~

(operator (cogn) confident file) and  
with a

? paragraph = " " - - - - - - - - -

> cleaning the text

import re

from ntk.corpus import stopwords

from ntk. stem -perfer imperfer Porter Stemmer

from nltk.stem import WordNetLemmatizer

~~> PS = PorterStemmer() # Stemming~~

> wordnet = WordNetLemmatizer() # Lemmatization

> Sentence = nltk.sent\_tokenize (paragraph)

```
> corpus = [ ]
```

? for i in range(len(questions)):

```
review = re.sub('^\[^a-zA-Z\]', '', sentences[i])
```

review = review + lower

`service = service.split`

~~stemming~~ review = review · split  
→ review = [ps. stem (word) for word in review]

if not word in set (stopwords.words('cng'))]

→ interview = [wordnet, lemmatize (word) for word in

answers if next word in set (stopwords.words('english'))

review = 1 ' . join (review)

~~the~~ Corpus opened (review)

# creating the Bag of words model  
 from sklearn.feature\_extraction.text import TfidfVectorizer

cv = TfidfVectorizer()

x = cv.fit\_transform(corpus).toarray()

print(x)

print(x[0])

print(x[1])

print(x[2])

print(x[3])

print(x[4])

print(x[5])

print(x[6])

print(x[7])

print(x[8])

print(x[9])

print(x[10])

print(x[11])

print(x[12])

print(x[13])

print(x[14])

print(x[15])

print(x[16])

print(x[17])

print(x[18])

print(x[19])

print(x[20])

print(x[21])

print(x[22])

print(x[23])

print(x[24])

## N-Gram

→ n-Gram is a way to help machines understand a word in the context to get a better understanding of words. N-gram is a neighboring sequence of n-items from a given sample of text. N-item means we can have two items, three items and so on. So, it is a contiguous sequence of some items. It helps to predicting the next words in a sentence. Things will be characters, words, sentences.

When n is 2 then we can call it as bigrams.  
When n is 3 then we can call it as trigrams.  
Based on sentence we can change the value of "n".

## Unigram Model

$$\text{Uni} = 1$$

Unigram = 1 word in a sentence

Ex:- "he tried to find a position as a teacher but had no luck at first"

Unigram:

he

tried

to

find

a

position

Unigram model  $\Rightarrow$

$P(w_t)$

## Bigram Model

$M_i = 2$

Bigram = 2 consecutive words in a sentence

Ex :- "he tried to find a position as a teacher  
but had no luck at first"

Bigram :

he tried

tried to

to find

find a

a position

⋮  
⋮

Bigram model :-  $P(w_t | w_{t-1})$

## Trigram Model

$M_i = 3$

Trigram = 3 consecutive words in a sentence

Ex :- "he tried to find a position as a teacher  
but had no luck at first"

Trigram:

he tried to

tried to find

to find a

find a position

a position as

⋮

⋮

⋮ Trigram model =  $P(w_t | w_{t-1}, w_{t-2})$

## \* Word2Vec

- Bag of words and TF-IDF (Problems)
  - Both BoW and TF-IDF approach semantic information is not stored. TF-IDF gives importance to uncommon words.
  - There is definitely chance of overfitting.

~~Learn - Does BoW & TF-IDF work well?~~

Solution → Word2Vec

- In this specific model, each word is basically represented as a vector of 32 or more dimension instead of a single number.
- Here the semantic information and relations between different words is also preserved.

### • Visual Representation - Word2Vec

(2D) (x,y) Since women are basically inter-related

There is semantic information present between them.

Man (3,6)

women (3,2,6,2)

Dimension 2  
(King, woman)

(King + woman) = Queen

(roughly) = King

Word (King + woman) = Queen

Play may not be selected for this kind of word (like, man, woman), so, there is no difference.

Dimension 1

King - Man + woman = Queen

- steps to create word2vec
  - Tokenization of the sentences
  - Create histograms if too many words
  - Take most frequent words
  - Create a matrix with all the unique words. It also represent the occurrence relation between the words.
- Gensim library → Word2Vec model.

```

> import nltk
> from nltk.corpus import stopwords
> from gensim.models import Word2Vec
> from nltk.corpus import stopwords
> import re

```

? Paragraph = "Mr. Shrek - considerate - kind  
 (P)(O&)

### > # Preprocessing the data

To remove space & special character. Some words

text = re.sub(r'\[\d-\d\]\*\]', '', Paragraph)

text = re.sub(r'\st+', '', text)

text = text.lower()

text = re.sub(r'\d', ' ', text)

text = re.sub(r'\s+', ' ', text)

{ remove extra spaces }

Output = remove extra and

## > # Preparing the dataset

sentences = nltk.sent\_tokenize(text)

sentences = [nltk.word\_tokenize(sentence) for sentence

in sentences] # converting all the sentences into words  
for removing all the stopwords

> for i in range(len(sentences)):

    sentences[i] = [word for word in sentences[i] if word  
        not in stopwords.words('eng')])]

## > # Training the Word2Vec model

,model = Word2Vec(sentences, min\_count=1)

if the word is present less than 1, i am going to  
just skip that particular word

? words = model.wv.vocab

for each one every word we will get vector, dimensions, sparse,  
it is going to represent itself as vectors

## > # Finding word Vectors

? vector = model.wv["word"]

## # Most ~~similar~~ similar words

? similar = model.wv.most\_similar("word")

From "freedom"