

Natural Language Processing

16/11/22

1) Basic NLP

→ i) Text pre-processing

ii) text - vector

→ techniques - BOW

TF, IDF

Python - library - NLTK

↓
Natural Language Tool Kit

Coding

```
import nltk
```

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

```
sent = "biryani is"
```

whenever you are using paragraph & splitting by sent-tokenization always the sentence should have space before it starts or after ending with symbols.

stem → is an attribute

Parts of speech also matters

POS tagging → then using this we can use lemmatization correctly.

Sklearn - Library - scientific
used for all ML algorithms

CountVectorizer - is a class → initialize the class,
↓ same as BOW returns object.

whenever you want to train the data,
use fit → This learning is called fit

fit → fit & transforming
learning from the data
↓
whatever learnt from data
transform it to machine
vectors.

unique words will automatically be saved
in class not to the assigned variable.


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

for a sentence → use word_tokenize
for a paragraph → use sent_tokenize

after ending the sentence with any character like, sentence ending symbols, ex- ., ?, ! the next sentence should start with space.

```
from nltk.corpus import stopwords
then assign it to a variable
```

stop = stopwords.words("english")

⇒ stemming & lemmatization +

```
from nltk.stem import PorterStemmer, LancasterStemmer, SnowballStemmer, WordNetLemmatizer
```

⇒ ps = PorterStemmer()
ps.stem("give the word here")

Same for others as well.

Text pre-processing

```
tf = []
```

```
for sent in data["text"]:
```

```
    l = []
```

```
    ls = sent.lower() — (Ist step) *
```

```
    ls = re.sub("[^a-zA-Z]", "", ls): (IInd)
```

```
    for word in word_tokenize(ls): (IIIrd)
```

```
        if word in stop: (stop words)
```

```
            pass
```

```
        else:
```

```
            stems = PorterStemmer().stem(word)
```

```
            l.append(stems)
```

(stemming)

```
tf.append(" ".join(l))
```



```
import sklearn
```

```
from sklearn sklearn.feature_extraction.text
```

```
import CountVectorizer
```

```
cv = CountVectorizer()
```

```
CountVectorizer(binary=True)
```



This will only 0's & 1's



no word



word yes.

⇒ cv.fit_transform(data["review"])

this will give sparse-vector matrix

which contains more zeros.

17/11/22

to correct spelling.

import TextBlob

!st → pip install textblob

from textblob import TextBlob

pd.DataFrame(news_cv.toarray(), columns=sorted(cv.vocabulary_))

18/11/22

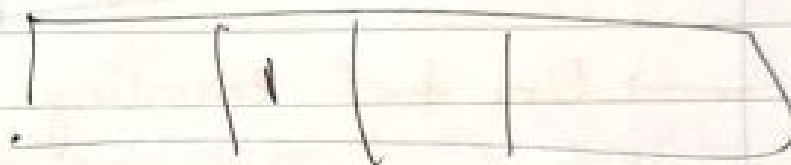
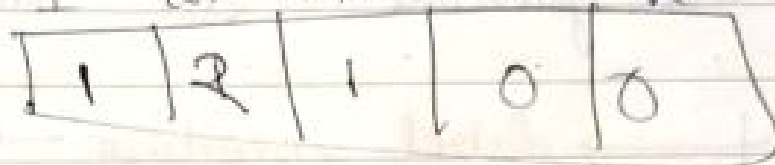
→ when you have reviews in different languages use Translator.

from translate import Translator
translator = Translator(from_lang="hindi", to_lang="en")

Q1) I am happy am

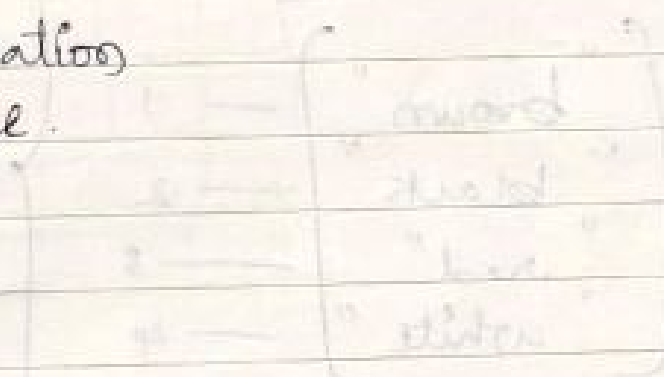
I will be happy

I am not will be



translation = translator.translate("am I")

⇒ output → translation
Bye-Bye



18/11/22

whenever your categorical column is in text format - we have two techniques

- **label encoding** → simply mapping
- **One hot encoding** → we use whenever your column already have an order

ex:- Hair colour

"brown"	— 1
"black"	— 2
"red"	— 3
"white"	— 4

} this is called mapping technique

this is not correct as this column has no order.

this is leading to → machine thinks "white" is more important and follows.
by giving this
to get rid of this we use
One hot encoding

whenever our column is categorical
 & is having no order or nominal use
One hot encoding

How does it work -

- 1) Count
 Unique categories — and the length
 with this we know dimension
- 2) vector of d-dim \rightarrow from length.

3)	"black"	1	0	0	0
	"brown"	0	1	0	0
	"red"	0	0	1	0
	"white"	0	0	0	1

when dimensions \uparrow ses. curse of dimensionality \uparrow ses

Task

Country	ht
India	
Aust	
Afri	

x_i y_i

PS \rightarrow predict ht.

take country and
 do caught of persons
 from that country &
 replace it wherever
 that country is there

ex1-	ht	Country		ht	Country
	10	Ind (10)		10	10
	15	Ameri (16)		16	16
	14	Aust (14)		14	14
	13	UK (13)	\Rightarrow	13	13
	11	Ind (10)		11	10
	16	Ameri (16)		16	16
	9	Ind (10)		9	10

will consider all the ht from country and do average & replace it at that particular column

$$\underline{\text{Ind}} \rightarrow \frac{10 + 11 + 9}{3} = \frac{30}{3} = 10$$

$$\text{Ameri} \rightarrow \frac{15 + 16}{2} = \frac{32}{2} = 16$$

to get rid of curse of dimensionality
do feature engineering based on PS.

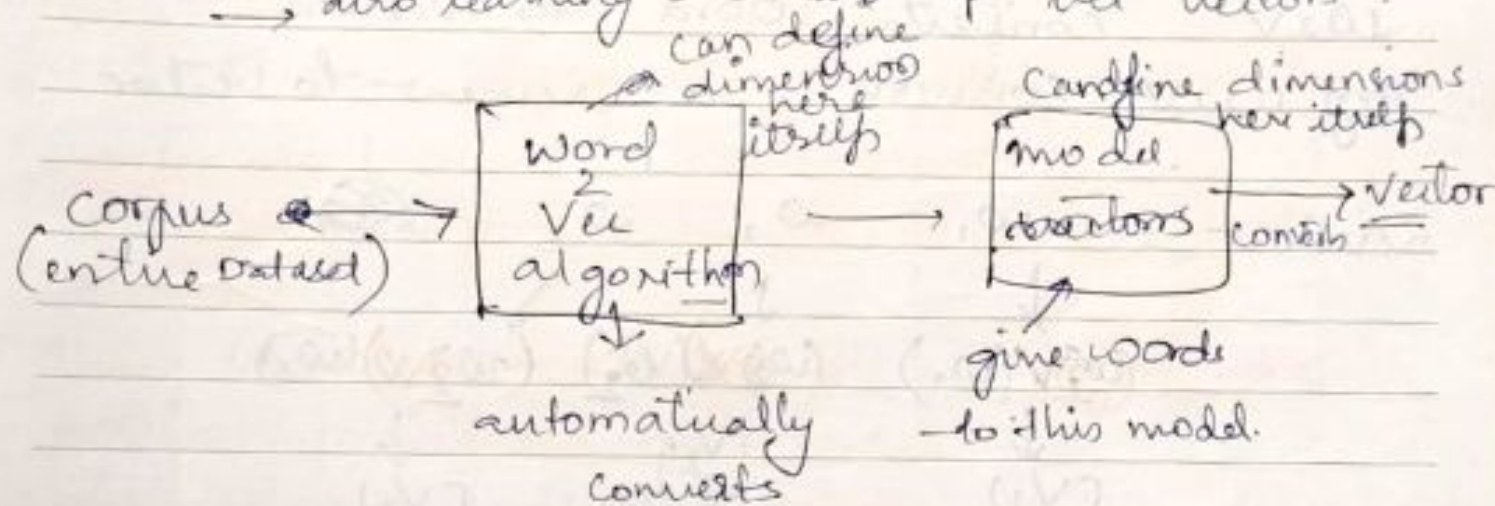
= also having good domain knowledge.

↓
using this we can create new features.

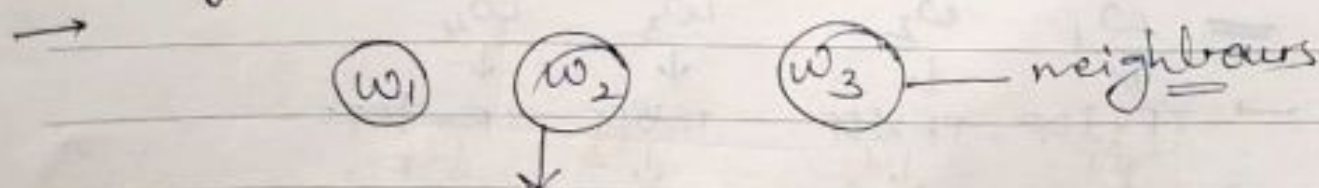
state of art of technique :- Word 2 Vec

Deep learning Technique

- Word 2 Vec → considers semantic meaning
 - It says give me a word I will create a vector.
 - also learning relationship betⁿ vectors.



- It works best for large Dataset
- the vectors max^m dimension would be 500
- this will be in dense matrix, where there are less zeros.
- as dimensionality increase the vector will be strong.



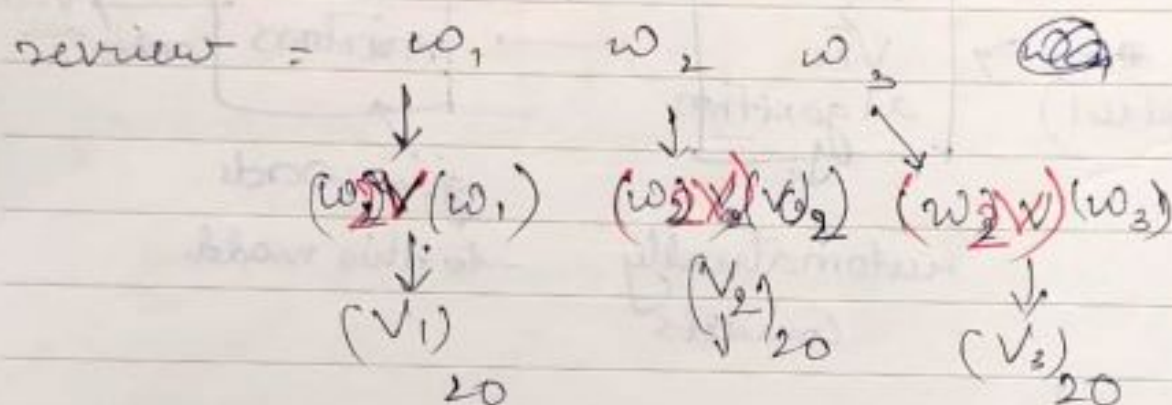
$$\text{If } v_1 \parallel v_2 \parallel v_3 \parallel v_4$$

There is a semantic relation & it is parallel to each other.

Average W2V

W2V converts word to vectors

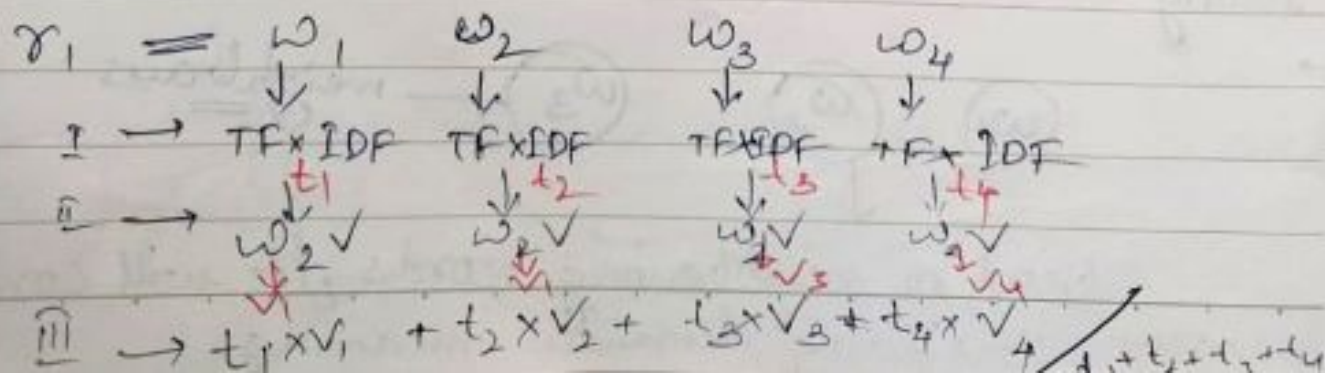
Avg W2V converts a review to vector



then add vectors & take avg.

$$\frac{1}{n} (V_1 + V_2 + V_3 + V_4)$$

TF IDF - W2V



$\left. \begin{array}{l} \text{Avg } w_2 V \\ \text{TFIDF } w_2 V \end{array} \right\} \text{convert a sentence to vector}$

$w_2 V \rightarrow$ Convert word to vector

One hot encoder / Label Encoder

import sklearn

from sklearn.preprocessing import LabelEncoder,

OneHotEncoder

LabelEncoder

→ this class object \rightarrow has to be assigned to a variable.

→ this variable must be fit-transform the particular/mentioned column.

→ this column will be given labels based on alphabetical order.

→ One hot —

array should always be ad

→ here, considering the rank it gives output ad , ad ~~array~~ array object.

based on ranks — 0001 \rightarrow this means the 1000

ad row 1 with number 1 means the object is present in that place.

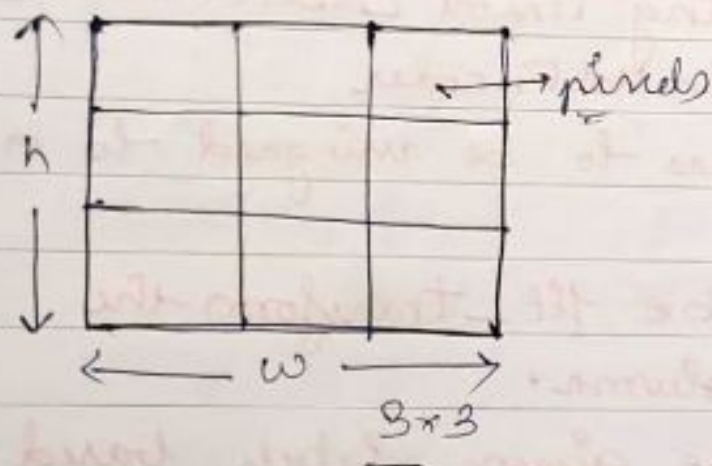
19/11/22

Image:—

2-dimensional representation of a visible spectrum

- representation of visible light.

more the pixels more the clarity of image



if someone say your image has 2 level -
it means Black, white i.e., 0, 1

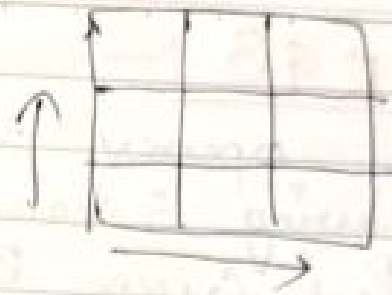
if your image is Gray scale image
there values range from 0 - 255
0 → Black
1 → white

remaining 244 - grey.

256 → in total → 2^8 → 8 bit value.
start with black ends with white.

Color Image :-

we have three values




Red (0-255) Green (0-255) Blue (0-255)

Black & white and Gray scale image \rightarrow stored in machine with (length, ht) 2-dimensional array
 Colour image \rightarrow 3-dimensional array (R, G, B)

21/11/22

Colour space \rightarrow It is a method by using which we can store the image in a machine.
 * We can also define as storing the colour.

1) Black and White \rightarrow If Image stores 0, 1 in each pixel.

0 \rightarrow black, 1 \rightarrow white. $2^1 \rightarrow 1$ bit 
 \rightarrow stored in black & white space.

2) Gray scale image :- (0 \rightarrow 255)

between 0 to 255 it is levels of gray scale

$2^8 \rightarrow 8$ bit (each pixel stored in) stored in Gray scale space

→ 3) RGB —

opencv says that each pixel is having 3 values belongs to Three Colours Red, Green, Blue.

Red → 0 to 255 shades of red.

↓ high intense red.

The image is stored in RGB colour space.

4) HSV — Hue, Saturation, Value.

stands for colour

↓
states
Intensity of
Colour

↓
brightness of colour

opencv → range Hue → 0 → 180 (or 255)

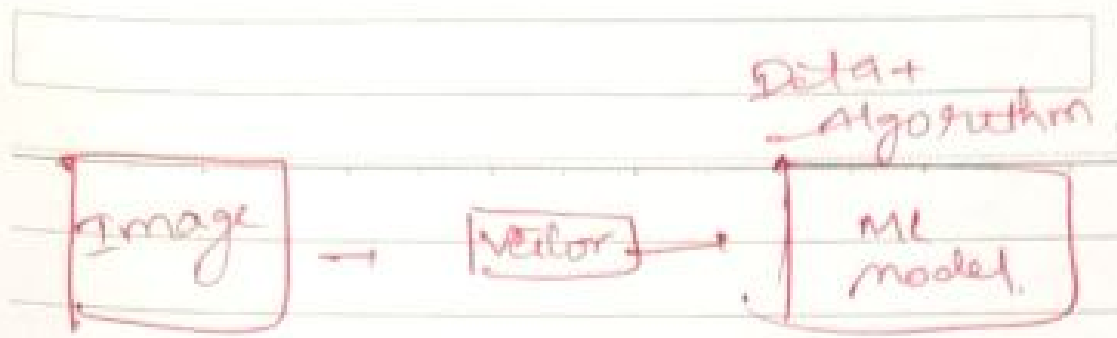
Saturation → 0 - 255

Value → 0 - 255

Colour space used here is HSV.

Open CV → Computer Vision

python library which contains lot of tools & methods we can apply them on images & videos.



How to convert image to vectors +

Copy image path & read

`cv2.imread (path)` → convert image to array.

→ to see image +

`imshow ("any msg." img array)`

this converts into 2-d image.

we can give waitkey → if waitkey = 0
at time it will wait till you close.

`cv2.imshow (`

`cv2.waitKey (0)`

`cv2.destroyAllWindows`

if waitkey (10) → 10 milli secs.

your image automatically gone

1) Bag of words :-

It is a technique by which we can convert text to vectors

- 1) Sentence → group of words/tokens.
- 2) Documents → can be a single sentence, paragraph, files
collection of sentences, paragraphs, files.
- 3) Corpus :- Collection of all the documents

4) Tokens :- breaking a huge thing into sub-groups is called tokens
ex:- Today is ML class

breaking it to sub words.
today
is
ML
class } tokens

→ tokens → can be words, sentence, paragraph etc.

5) tokenization :- process by using which we can get the tokens.

There are two types of tokenization

- 1) sentence tokenization
- 2) word tokenization

→ Vocabulary → set of all the words.

Bag of words (BOW) :-

let us say, we have ³ documents.

d_1 :- $w_1, w_2, w_3, w_4, w_5, w_6$

d_2 :- $w_1, w_2, w_4, w_8, w_5, w_6$

d_3 :- $w_1, w_9, w_3, w_7, w_5, w_6$

① create a dictionary :-

↳ not python dictionary

→ In this we have all the unique words from corpus.

the length of .

② no. of unique words indicates → no. of dimension

ex - 9 dimensions i.e. 9 unique words.
each word is called dimension.

③ $V_1 =$

w_1	w_6	w_3	w_7	w_4	w_8	w_9	w_5	w_2
1	1	1	0	1	0	0	1	1

no sequential order.

→ How many no. of times a word is occurring

$$V_2 =$$

w_1	w_6	w_3	w_7	w_4	w_8	w_9	w_5	w_2
1	1			1			1	1

features = 0

$T_1 \rightarrow V_1 \rightarrow P_1$	words	-	-	-
$T_2 \rightarrow V_2 \rightarrow P_2$	-	-	-	-
$T_3 \rightarrow V_3 \rightarrow P_3$	-	-	-	-

Technique bag of words

ex1- d_1 :- Biryani is very good & less
 d_2 :- " " " bad & cheap
 d_3 " " " good " affordable.

unique words = biryani is very good, bad, less, cheap
 & affordable

$V_1 =$

1	1	1	1	0	1	0	1	0
---	---	---	---	---	---	---	---	---

 $\Rightarrow P_1$
 bir is very good bad less cheap and affor

$V_2 =$

1	1	1	0	1	0	1	1	0
---	---	---	---	---	---	---	---	---

 $\Rightarrow P_2$

$V_3 =$

1	1	1	1	0	0	0	1	1
---	---	---	---	---	---	---	---	---

 $\Rightarrow P_3$

~~$V_1 - V_3$~~

$$P_1 - P_3 = \sqrt{(1-1)^2 + (1-1)^2 + (1-1)^2 + (1-1)^2 + (1-0)^2 + (0-0)^2 + (0-0)^2 + (1-1)^2 + (1-1)^2}$$

$$= \sqrt{1+1} = \sqrt{2}$$

$$V_1 - V_3 = \sqrt{2}$$

$$S(x_1, x_3) > Sim(x_1, x_2)$$

$$\therefore d(V_1 - V_3) < d(V_1 - V_2)$$

$$\underline{V_1 - V_2} :-$$

$$= (1-1)^2 + (1-1)^2 + (1-1)^2 + (1-1)^2 + (1-1)^2 + (1-1)^2 + (1-1)^2 + (0)^2$$

$$\boxed{V_1 - V_2 = \sqrt{4}}$$

$$S(V_1 - V_3) > S(V_1 - V_2) \quad \checkmark \quad \text{Hence proved}$$

Binary bag of words :-

- ① creating a dictionary \rightarrow set of all unique words.
- ② the number is 1 or 0, if the word occurs give 1, if doesn't occur give '0'

Disadvantage :-

bag of words.

- \rightarrow doesn't consider semantic meaning.
- \rightarrow No sequential order.

Sparse \rightarrow used for n-dim

In a vector, if more no. of elements are '0' it is sparse vectors.

as your corpus size \uparrow dimensionality \uparrow

before converting text to vector we do pre-processing on text called text pre-processing.

this is mainly done to check whether the dimensionality can be decreased.

Text pre-processing

- 1) Convert all your text to lower case letters. ^{upper}
- 2) Remove all the stopwords (***) ^{second step}

^{first step}

those words which will not give meaning to a ~~text~~ sentence.

ex-1 - 1) this is python class.
python class.

2) this Biryani is very good
Biryani very good

→ very } remove
→ not } them without changing the meaning

3) Remove all unnecessary symbols :-

using Reg X

4) Stemming & Lemmatization :-

process by using which we will get the root of the word.

this will try to remove suffixes of the word and give the root word.

Stemming
ex - tasty, tastier, tasteful.

- give root word \Rightarrow taste

Stemming \rightarrow give root word may not be actual english word.

Lemmatization \rightarrow give root word which have actual english word.

lemmatization - tasty, tastier, tasteful

\rightarrow give root word \Rightarrow taste

15/11/22

Types of stemming lemmatization algorithms

- 1) Porter stemmer
 - 2) Lancaster stemmer
 - 3) Snowball stemmer
- root word
↓
called as stem

= Porter, Lancaster are english stemmer i.e., work only on english

= Snowball stemmer → for other languages i.e., non-english.

Porter stemmer — (oldest)

has set of five rules. $C[CV]^m V$

$C \Rightarrow$ Consonant

$V =$ vowels.

It will take a word & apply rules.

→ after stemming the last letter it is not checking for the word.

Lancaster stemmer — (new stemmer) Iterative algorithm.

we have 120 rules.

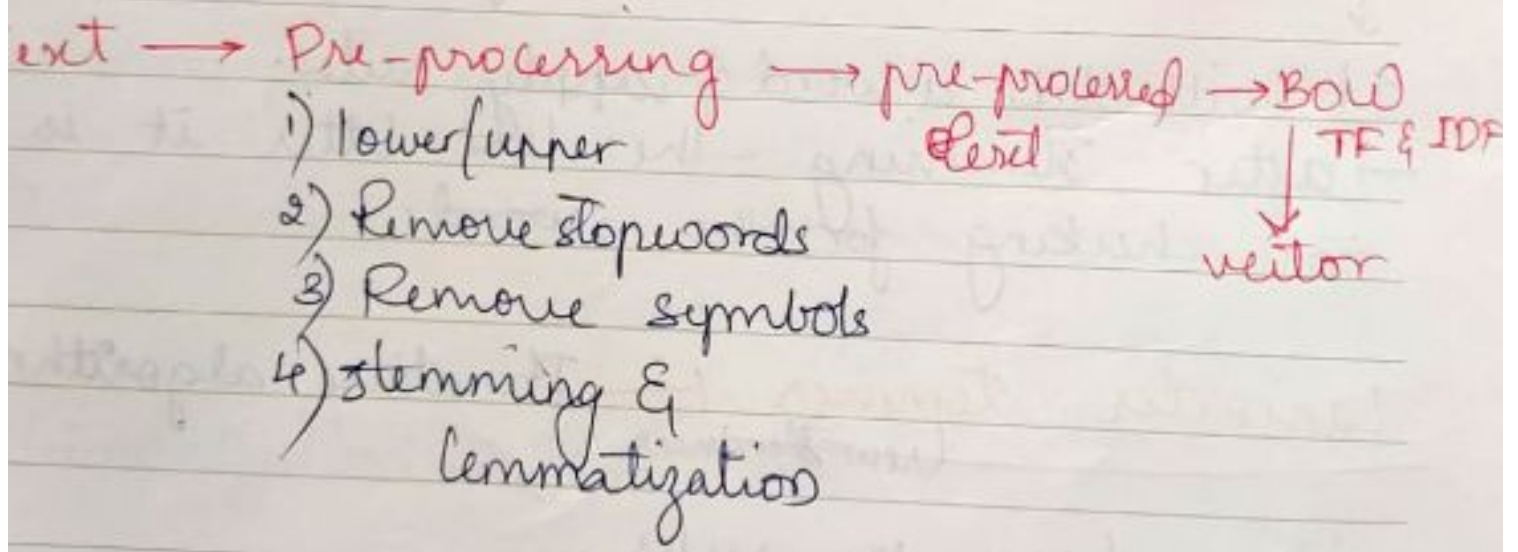
because of high stemming the root word will not be the correct english word.

Snowball Stemmer — use on non-english words.

Lemmatization —

= algorithms → word net lemmatization

after stemming the last letter it checking with the Database & confirms if it is a root word present in the corpus (i.e., contains all english words) ~~and~~ then gives output, if it is ~~not~~ present in corpus, ~~it then~~ then it redirects to last letter & stops there.



2) TF-IDF :-

Term Frequency — Inverse Document Frequency

Steps :-

1) Create a Dictionary (this step → to know the dimensions)



we have set of unique words.

2) Term Frequency (TF) :-

ex :-

τ_1	—	w_1	w_2	w_4	w_3	w_5
τ_2	—	w_5	w_6	w_2	w_1	w_2
\vdots						
τ_n						

$$TF(w_i, \tau_i) = \frac{\text{No. of times } w_i \text{ occurs in } \tau_i}{\text{Total No. of Words in } \tau_i}$$

Term Frequency of w_i with respect to τ_i is equal to No. of times w_i occurs in τ_i divided by total No. of words in τ_i .

$$\rightarrow TF(w_1, \tau_1) = \frac{1}{5}$$

$$TF(w_2, \tau_2) = \frac{2}{5}$$

Always Term Frequency value will be between

$$0 < TF < 1$$

$$0 \leq 1$$

3) Inverse Document Frequency -

ex:-

Coprus

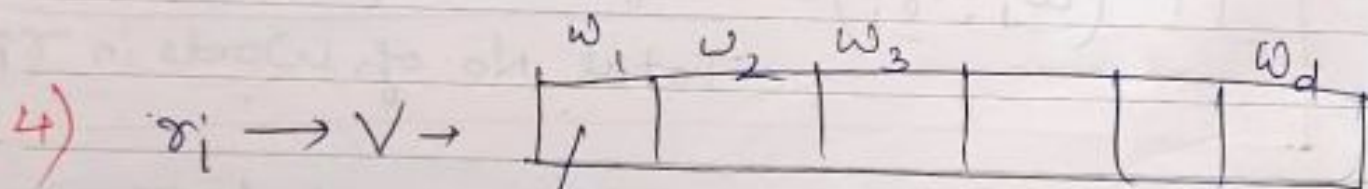
σ_1	$=$	w_1	w_2	w_5	w_6
σ_2	$=$	w_2	w_1	w_2	w_5
σ_7	$=$		w_3		
σ_{10}	$=$		w_3		
σ_n	$=$		w_3		

$$IDF(w_i, \mathcal{D}) = \log \left[\frac{N}{n} \right]$$

$\therefore N = \text{no. of documents}$
 $n = \text{Total no. of document where } w_i \text{ occurs}$

IDF of w_i with respect to Documents (entire Coprus) will be equal to $\log \frac{N}{n}$.

$$IDF(w_3, \mathcal{D}) = \frac{100}{3}$$



$$TF(w_1, \sigma_i) \times IDF(w_1, \mathcal{D})$$

ex:- $d = 12$



$$TF(w_1, \sigma_1) \times IDF(w_1, \mathcal{D}) = \frac{100}{2}$$

$$= \frac{1}{4} \times \frac{100}{2} = \boxed{w_1 \text{ value}} \rightarrow \text{keep it in } w_1 \text{ vector}$$

Ex 1-

$$\tau_1 = \omega_1 \quad \omega_2 \quad \omega_3 \quad \omega_4$$

$$\tau_2 = \omega_1 \quad \omega_3 \quad \omega_4 \quad \omega_7$$

$$\tau_3 = \omega_2 \quad \omega_3 \quad \omega_6 \quad \omega_5$$

$$\tau_4 = \omega_1, \omega_2$$

from the above we got $d=7$

Vectors for τ_1 :-

ω_1	ω_2	ω_3	ω_4	ω_5	ω_6	ω_7	ω_2

$$\omega_1 \Rightarrow TF(\omega_1, \tau_1) \times IDF(\omega_1, \mathcal{D})$$

$$= \frac{1}{4} \times \frac{4}{3}$$

$$\omega_2 \Rightarrow TF(\omega_2, \tau_1) \times IDF(\omega_2, \mathcal{D})$$

$$= \frac{1}{4} \times \frac{4}{3}$$

$$\omega_3 \Rightarrow TF(\omega_3, \tau_1) \times IDF(\omega_3, \mathcal{D})$$

$$= \frac{1}{4} \times \frac{4}{3}$$

Vectors of $\tau_4 \rightarrow$

ω_1	ω_2	ω_3	ω_4	ω_5	ω_6	ω_7	ω_2
$\frac{1}{3}$	0	0	0	0	0	0	$\frac{2}{3}$

$$\omega_1 \Rightarrow TF(\omega_1, \tau_1) \times IDF(\omega_1, \mathcal{D}) = \frac{1}{2} \times \frac{4}{3}$$

$$\omega_2 \Rightarrow TF(\omega_2, \tau_1) \times IDF(\omega_2, \mathcal{D}) = \frac{1}{2} \times \frac{4}{3}$$

This value will be higher only when
 $TF(w_i, d_i) \times IDF(w_i)$
 more less.

The word which is given more importance is words with more occurrence in a single review will have more value & word which is occurring rarely in entire corpus.

$\log \frac{N}{n}$ — less so N will be more.

$$\frac{N}{n} = 1$$

$$\log(1) = 0 \rightarrow n \leftarrow N$$

$$\text{when } n = N \Rightarrow \log(1) = \underline{0}$$

$$\text{then } \log\left(\frac{N}{n}\right) \geq 0$$

if $IDF = 0$, then, that particular word is there in all the documents so the word is given less imp.

use \log to lessen the gap between TF & IDF . (as TF will always be ≥ 1)

Disadvantages of TFE, IDF

- ✓ No sequential order
- ✓ No semantic meaning consideration

→ to give some sequential concept to the ~~base~~ ^{TF & IDF} so that the meaning is not changing.

for this we use n-grams.

1-gram → 1 word = 1 dimension

2-grams → 2 consecutive word = 1 dim

3 → grams → 3 consecutive word = 1 dim

$n \uparrow \quad d \uparrow \rightarrow$ curse of dimensionality

she is a good girl — 5

she is a good girl →