

feature Engineering in NLP

Common Terms

1) Corpus (c) = Sabda Sagar (प्राची दातारत्मय)

Vayeko words lai

Joddha aavne

(Combination ki corpus vaninxo) - Total combination of words in a dataset

e.g -

my name is miraj	pos	→ (my name is miraj)	}
I am happy	neg	I a happy	

corpus = $\{ \}$

2) Vocabulary (v) = Tati pani unique words

Corpus banyo teslaip vocabulary vanta.

(No of unique words in the corpus)

3) Document (d) = Each individual row of

the dataset is document

...	...	document 1
...	...	document 2

4) Word (w) = Each individual word in a document is word

One Hot Encoding

D1 people watch
campusx

D2 campusx watch
campusx

D3 people write
comment

D4 campusx write
comment

Corpus =

people watch campusx campusx
watch campusx people write
comment campusx write
comment

vocabulary = people watch campusx write
comment

hami vocabulary data ~~as~~ vector form gartam

People	watch	campusx	write	comment
0	0	0	0	0

→ One hot encoding mo hami harek word i^{th} ^{in vdim} vector ko represent ma convert gartam by above vocabulary table:

let's take example of document D_1

$D_1 = \text{people watch campusx}$

representation of people =

people	watch	campusx	write	comment
1	0	0	0	0

representation of watch =

people	watch	campusx	write	comment
0	1	0	0	0

representation of campusx =

people	watch	campusx	write	comment
0	0	1	0	0

∴ The vectorical representation of D_1 will be

$$D_1 = [[1, 0, 0, 0, 0], [0, 1, 0, 0, 0], [0, 0, 1, 0, 0]]$$

Similarly for D_2 from vocabulary table

$$D_2 = [[0, 0, 1, 0, 0], [0, 1, 0, 0, 0], [0, 0, 1, 0, 0]]$$

And similar for other also #

① People watch campus - $D_1 = \begin{bmatrix} 1, 0, 0, 0 \\ 0, 0, 1, 0, 0 \end{bmatrix} = (3, 5)$

② People watch campus

watch = $D_2 = \begin{bmatrix} 1, 0, 0, 0 \\ 0, 0, 1, 0, 0 \\ 1, 0, 0, 0 \\ 0, 0, 1, 0, 0 \end{bmatrix} = (4, 5)$

Date
Page

[1, 0, 1, 0, 0]

[1, 0, 0, 0, 0]

[0, 0, 1, 0, 0]

[0, 0, 0, 1, 0]

[0, 0, 0, 0, 1]

[1, 0, 0, 0, 0]

docs

flows

can't be used
in ml algorithm

→ Intuitive

→ sparsity [1, 0, 0, 0, 0, 0, 0, 0, 0]

→ easy to implement

→ And if it needs same dimension input i.e words in document must be same

→ COV(It can't deal with word which is not in vocabulary since word vector representation is not possible)

→ NO capturing of semantic meaning

e.g.

In our vocab we have

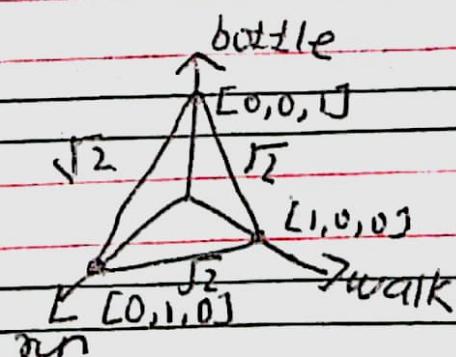
walk	run	bottle
------	-----	--------

walk vector form = [1, 0, 0]

run vector form = [0, 1, 0]

bottle vector form = [0, 0, 1]

If we plot in 3D space



All the distances are same i.e equal
distance since the relation of walk, bottle
and run bottle is not present in real world
Scenario but our algorithm shows ~~so take it as~~
a relation so, it can't capture semantic meaning

Bag of words:

D ₁	people watch miraj	1
D ₂	miraj watch miraj	1
D ₃	people convert miraj	0
D ₄	miraj comment miraj	0

Corpus:

people watch miraj miraj
watch miraj people document
miraj miraj document miraj

Vocabulary = people watch miraj write comment
V=5

people	watch	miraj	write	comment

Bag of words = yesma hami document ma
vayeko harek word kati choti ayo (repeat)
document is
teslai herera^ vector ma convert garna vdim
vector ma.

Eg - Taking document D₁

People watch miraj

Vector representation of D₁ =

people	watch	miraj	write	comment
1	1	1	0	0

(miraj watch miraj)

Taking document D₂

vector representation of D₂ =

people	watch	miraj	write	comment
0	1	2	0	0

Similary

D ₃ =	people	watch	miraj	write	comment
	1	0	1	0	1

D ₄ =	people	watch	miraj	write	comment
	0	0	2	0	1

$$\therefore D_1 = [1, 1, 1, 0, 0]$$

$$D_2 = [0, 1, 2, 0, 0]$$

$$D_3 = [1, 0, 1, 0, 1]$$

$$D_4 = [0, 0, 2, 0, 1]$$

→ Bag of words ma dui ota same class

denote garni document ko word frequency
similar huncha.

e.g.

B D₃ people write comment | 0
by campus write comment | 0

Here, write = 1, comment = 1 times in document
and this document denote 0
both

To use Bag of words we have
count vectorizer in Sci-kit learn

It has disadvantage i.e

I am miraj | 1 : vocab = [I am | miraj | not]

I am not miraj | 0

$D_1 = [1, 1, 1, 0]$ ota dimension na
 $D_2 = [1, 1, 1, 1]$ word ko frequency same
xa so, ml model will be
confused and they both are same

The three dimension of vector are same

→ $[1, 1, 1, 0]$ → $[1, 1, 1, 1]$ → $[1, 1, 1, 1]$
culling → distance is not
both document
are a same denoting
same class are to
this class

2) OOV - kunai word vocab vanda out batq

aayo rane testaj vector ma we naigarden q

like - \rightarrow we have vocab \rightarrow [won | game | wow]

Now, we have sentence lottery! you won iphone
Now, In vector this sentence will be $[1, 0, 0]$

\therefore lottery!, iphone is not used since it can be useful
information like in spam mail.

3) Sparse - It also creates the sparse matrix which may lead to overfit

P.T.O

N-Grams: (Bag of n-grams)

D ₁	People watch campusx	1
D ₂	Campusx watch campusx	0
D ₃	People write comment	0
D ₄	campusx write comment	0

Vocabulary:- (Uni-gram default) = Bag of words

People | watch | campusx | write | comment |

N-grams ma hami multiple words ko combination kar vocabulary ma takam to prevent word sentence meaning.

∴ Bi-gram: we use two words to make vocabulary:

b ₁	People watch campusx	1
b ₂	campusx watch campusx	0
b ₃	people write comment	0
b ₄	campusx write comment	0

vocab candidate: (word) must be serial

people watch, watch campusx, campusx watch,
watch campusx, people write, write comment,
campusx write, write comment

vocab:

people	watch	watch campusx	campusx watch	people write
write comment	campusx write	write campusx	campusx	

Now, converting documents into vectors using Bi-gram.

	A	B	C	D	E	F	G
$D_1 =$	1, 0, 0, 0, 0, 0, 0						
$D_2 =$	0, 1, 1, 0, 0, 0, 0						
$D_3 =$	0, 0, 0, 1, 1, 0, 0						
$D_4 =$	0, 0, 0, 0, 1, 1, 0						
$\therefore D_1 = [1, 0, 0, 0, 0, 0, 0]$							
$D_2 = [0, 1, 1, 0, 0, 0, 0]$							
$D_3 = [0, 0, 0, 1, 1, 0, 0]$							
$D_4 = [0, 0, 0, 0, 1, 1, 0]$							

∴ Tri-grams: we use three words to make vocab

People watch campus Hello | 1

mirut People watch campus | 0

Vocab Candidate:

people watch campusx, watch campus Hello,

mirut people watch, people watch campusx

Vocab:

People watch campusx|watch campus Hello/mirut people watch

$D_1 = 1, 1, 0, 0, 1, 1, 0$

$D_2 = 1, 1, 0, 0, 1, 1, 0$

$$\therefore D_1 = [1, 1, 0, 0, 1, 1, 0]$$

$$D_2 = [1, 1, 0, 0, 1, 1, 0]$$

TF-IDF: (Term frequency - Inverse document frequency)

D ₁	People watch campusx	1
D ₂	campusx watch campus	1
D ₃	people write comment	0
D ₄	campusx write comment	0

vocab:

people	watch	campusx	write	comment

∴ In TF-IDF each token is vectorized by using TF-IDF formula.

$$\Rightarrow \text{TF}(t, d) = \frac{(\text{Number of occurrences of term } t \text{ in document } d)}{(\text{Total no of terms in document } d)}$$

∴ It's like probability. TF value lies between 0 and 1 i.e. 0 ≤ TF ≤ 1.

$$\Rightarrow \text{IDF}(t) = \log_e \left(\frac{(\text{Total no of documents} - \text{No of documents with term } t \text{ in them})}{\text{No of documents with term } t} \right) + 1$$

Tf-IDF ma kunki? particular word out a document ma dherai oako xa an? or document ma ne frequent oako xa vane tyo
 Particular word ko value high assign gaoxa.

(i.e means the word must be valuable to the document as well as whole corpus)

now, let's vectorize the documents for that we need to calculate IDF value of vocab words:

~~Per~~

	IDF	$IDF = \log \left(\frac{\text{Total no of documents}}{\text{no of document where term } t \text{ is present}} \right) + 1$
people	$\log(\frac{4}{2}) + 1$	$\Rightarrow 1.69$
watch	$\log(\frac{4}{2}) + 1$	$\Rightarrow 1.69$
computer	$\log(\frac{4}{3}) + 1$	$\Rightarrow 1.28$
wrote	$\log(\frac{4}{2}) + 1$	$\Rightarrow 1.69$
comment	$\log(\frac{4}{2}) + 1$	$\Rightarrow 1.69$

Now, let's give value to each vocab words on the basis of $(Tf \times IDF)$.

Let's convert all the documents into vectors

P T D
 →

	people	watch	campus	write	comment
D1	$\frac{1}{3} \times 1.69$	$\frac{1}{3} \times 1.69$	$\frac{1}{3} \times 1.28$	$\frac{2}{3} \times 1.69$	$\frac{2}{3} \times 1.69$
D2	$\frac{1}{3} \times 1.69$	$\frac{1}{3} \times 1.69$	$\frac{2}{3} \times 1.28$	$\frac{2}{3} \times 1.69$	$\frac{2}{3} \times 1.69$
D3	$\frac{1}{3} \times 1.69$	$\frac{0}{3} \times 1.69$	$\frac{0}{3} \times 1.28$	$\frac{1}{3} \times 1.69$	$\frac{1}{3} \times 1.69$
D4	$\frac{0}{3} \times 1.69$	$\frac{0}{3} \times 1.69$	$\frac{1}{3} \times 1.28$	$\frac{1}{3} \times 1.69$	$\frac{1}{3} \times 1.69$

D1	0.56	0.56	0.42	0	0	D $\frac{3\text{dim}}{\text{mt}}$
D2	0	0.56	0.85	0	0	D $\frac{3\text{dimmt}}{\text{mt}}$
D3	0.56	0	0	0.56	0.56	D $\frac{3\text{dim}}{\text{mt}}$
D4	0	0	0.42	0.56	0.56	D $\frac{3\text{dimmt}}{\text{mt}}$

Same class

Izoxin

Chinta

Yebata

No

Ternia

Miltag

$$D_1 = [0.56, 0.56, 0.42, 0, 0]$$

$$D_2 = [0, 0.56, 0.85, 0, 0]$$

$$D_3 = [0.56, 0, 0, 0.56, 0.56]$$

$$D_4 = [0, 0, 0.42, 0.56, 0.56]$$