

# Text Visualization

## Creating a Bag of Words model

- Break the sentence into words
- Get Freq distribution of words
- Sort the words based on freq
- Pick Top few words
- Create a bag of words Table

import nltk

import re

dataset = nltk.sent\_tokenize(paragraph)

for i in range(len(dataset)):

dataset[i] = dataset[i].lower()

dataset[i] = re.sub(r'\W', ' ', dataset[i]) // remove special characters

dataset[i] = re.sub(r'\s+', ' ', dataset[i])

// to replace multiple spaces with single space



word2count = {}

for data in dataset:

words = nltk.word\_tokenize(data)

for word in words:

If word not in word2count.keys():

word2count[word] = 1

else:

word2count[word] += 1

(9)

# → when we are building a Bow on large corpus of data (50k movie reviews). we will end up in a large set of words.  
so, we will select most frequent "n" words.

import heapq

freq\_words = heapq.nlargest(100, word2count, <sup>Dictionary</sup>key = word2count.get)

X = [ ] // this will be our bow of words model.

for data in dataset:

vector = [ ] // going to contain complete vector of sentence/doc

for word in freq\_words:

if word in nltk.word\_tokenize(data):

vector.append(1)

else:

vector.append(0)

X.append(vector) // list of lists

import numpy as np

X = np.asarray(X)

---

→ word frequency can also build using a function as below

wordlist = nltk.FreqDist(wordlist)

features = wordlist.keys()

values = wordlist.values()

Returntype = ( [ ] )

↓

dict.

## Bag of words - Problems

→ All words have the same importance

→ No semantic information preserved

↳ which word is important & which word is not

To overcome this we have the below solution

TF-IDF (Term Frequency - Inverse document Frequency)

\* Some semantic information is preserved as ~~uncommon~~ ~~as~~  
uncommon words are given more importance than common words.

Ex: 'she is beautiful'. Here beautiful will have more importance than 'she' or 'is'

when we change 'beautiful' to 'ugly' then the meaning of the sentence changes

→ Converting the sentence to lower

→ Converting sentence to a words

TF = Term Frequency

= Term Frequency of word in that document

= TF of the word for each document is not same.

IDF = Inverse document Frequency

= ~~Inverse~~ IDF is not calculated for the document, but for the whole corpus

= IDF for each word will be the same. Through of the document

Formula: 
$$\frac{\text{Number of occurrences of a word in a document}}{\text{Number of words in that document}}$$

"to be or not to be"

$$TF \text{ of 'to'} = \frac{1+1}{6} = 0.33$$

$$TF \text{ of 'be'} = \frac{1+1}{6} = 0.33$$

$$TF \text{ of 'or'} = \frac{1}{6} = 0.16$$

### TF matrix

words/Document	Document 1	Document 2	Document 3
going	0.16	0.16	0.12
to	0.16	0	0.12
today	0.16	0.16	0
i	0	0.16	0.12
am	0	0.16	0.12
it	0.16	0	0
is	0.16	0	0
again	0.16	0	0

### IDF Formula

$$\log_e \left( \frac{\text{number of documents}}{\text{number of document containing word}} \right)$$

"to be or not to be"

"i have to be"

"you got to be"

$$to = \log_e(3/3) = \log_e(1) = 0$$

$$be = \log_e(3/3) = 0$$

$$have = \log_e(3/1) =$$

words	IDF	
going	$\log(3/3)$	0
to	$\log(3/2)$	<del>0.41</del> 0.41
today	$\log(3/2)$	<del>0.41</del> 0.41
I	$\log(3/2)$	0.41
am	$\log(3/2)$	0.41
it	$\log(3/1)$	1.09
is	$\log(3/1)$	1.09
again	$\log(3/1)$	1.09



going to today? am it is rain

(12)

word

Document 0.16 0.09 0.07 0 0 0.17 0.12 0.12

Doc1 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.1

Doc2 0 0 0.07 0.07 0.07 0 0

Doc3 0 0.05 0 0.05 0.05 0 0 0

$$TFIDF(word) = TF(Document, word) * IDF(word)$$

→ This method is extensively used in Text classification, opinion mining etc.,

Why log in IDF

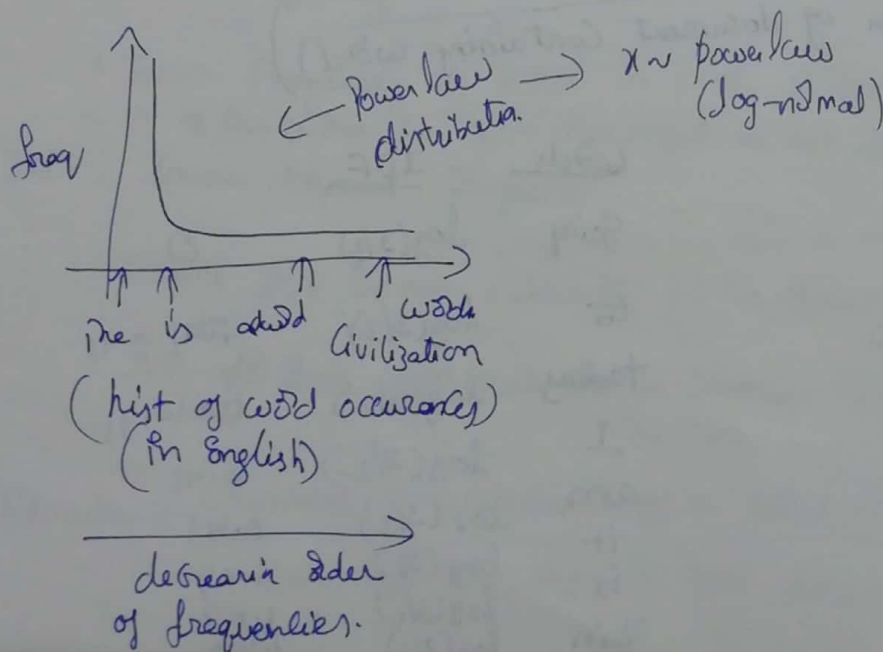
$$IDF(w_i, D_c) = \log\left(\frac{n}{n_i}\right)$$

$$\frac{n}{n_i} = \frac{\text{Total \# of Documents}}{\text{\# of doc which contain } w_i}$$

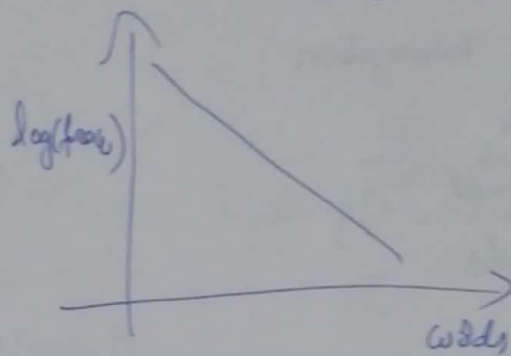
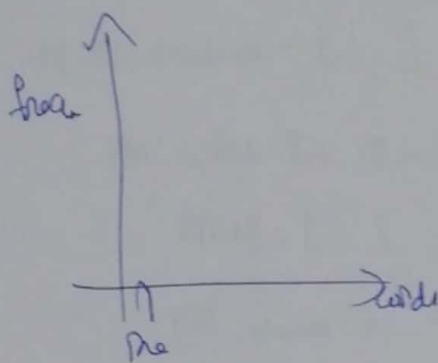
→ 1972 research paper where IDF is introduced.

→ not very strong base on theory

→ Referring to Zipf's law



According to box-Cox transformation take log of the Variable.



### Practical Reason

$$\frac{N}{n_i}$$

$$\log_e \left( \frac{N}{n_i} \right)$$

$$n_{w_i} \rightarrow 1$$

$$is \quad 1$$

⋮

$$Civilization. \quad \frac{1000}{1}$$

$$\approx 6.9$$

(occurs 1 in Every word and document)

TF \* IDF

↳ If word does not occur, the IDF dominates

## TF-IDF Model

(14)

- Sentence Tokenization
- Preprocess
- word count
- Freq-words

### # IDF-matrix

dataset = nltk.sent\_tokenize(Paragraph)

words\_ids = {}

for word in Freq-words

doc\_count = 0

For data in dataset:

If word in nltk.word\_tokenize(data):

doc\_count += 1

word\_ids[word] =  $\text{np.log}(\frac{\text{len}(\text{dataset})}{\text{doc\_count} + 1})$

↑  
bias

### # TF-matrix

tf\_matrix = {}

for word in Freq-words:

doc\_tf = [ ]

for data in Sentence:

frequency = 0

for w in nltk.word\_tokenize(data):

If w == word:

frequency += 1

tf\_word = frequency / len(nltk.word\_tokenize(data))

doc\_tf.append(tf\_word)

tf\_matrix[word] = doc\_tf

q/p  
doc1 doc2 doc3 doc4  
[after]: [0.02, 0.0, 0.0, 0.02]  
⋮  
}

# tfidf Matrix (Calculation)

98 word  
4 reviews

tfidf\_matrix = [ ]

For word in tf\_matrix.keys();

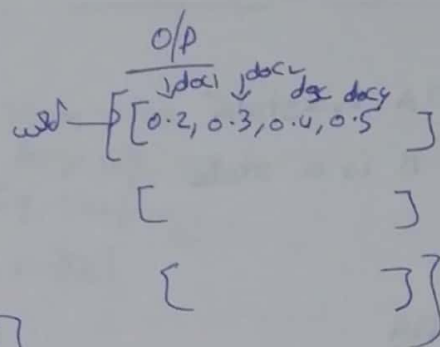
tfidf = [ ]

For value in tf\_matrix[word]

Score = Value \* word\_idf[word]

tfidf.append(Score)

tfidf\_matrix.append(tfidf)



## Conversion tfidf to a matrix nd array

X = np.asarray(x)

4 rows & 98 columns