

* Module : 3

→ How classification ~~works~~ :-

→ Amazon Fine Food Reviews Analysis:

No. of review : 568,454

No. of users : 256,059

No. of products : 74,258

Time span : Oct 1999 - Oct 2012

No. of Attributes / columns in data : 10

Attribute Information:

1. Id

2. product Id

3. User Id

4. profile Name

5. Helpfulness Numerator

6. Helpfulness Denominator

7. Score (rating)

8. Time

9. Text

Objective : Given a review, determine whether the review is positive (4, 5 rating) or negative (1 or 2)

If we ~~don't~~ have rating we can predict it using simple If-else but we need to predict ~~the~~ the review without score (rating)

* Data cleaning : Deduplication

It is observed that the reviews data had many duplicate entries. Hence it was necessary to remove duplicates in order to get unbiased results for the analysis.

#code:-
sorted_data = filtered_data.sort_values("ProductId", axis=0);
final = sorted_data.drop_duplicates(subset=["userId", "profileName", "text", "Time"],
keep="first", inplace=False);

★ Text to vector

The problem we had is given a feature (without rating) we need to determine polarity of review. The most useful feature is text feature and summary. Both of them are basically text. If we can somehow convert text to vector we can use the power of linear algebra.

Q) How do you convert ~~str~~ text into vector.

→ Bag of words:-

R₁ = This pasta is very tasty and affordable

R2 = This pasta is not tasty and is affordable.

R3: This pasta is delicious and cheap.

R4 = Pauri is faulty and pauri fault is good.

In BOW:

step ① : Constructing a Dictionary :- set of all the words in your review:

Let we have ← } This, pasta, is, very, tasty, and, affordable, not, delicious, cheap, good.
d-unique words.

② ex: This pasta is very tasty and affordable.

[illegible]

vector of size d

↳ each index corresponds to a word in ~~new~~ dictionary.

each entry in new dictionary.
keeps the freq of that word occurred in review.

It will be a sparse approach because there may be huge words in a dictionary, but a sentence can have 10, 15, 20, 100 words.

Each word in Bow is a different dimension.

We can improve BOW using some technique.

look at the R₁, we can see that there are various under words, these words are called stop word.

like in 2, $\{$ is, and, this $\}$ these one-stop words.

If we remove them from dictionary then Bow distance will be smaller and we have only important words in vector.

Note: Some time removing stop words can cause problem.
Hence stop word removal is not a golden rule.

→ Another thing is, we should make all the words in lower case letter. because we can have 2 diff words in dictioⁿ i.e. Parla, parla are same but they can be counted 2.

→ Stemming:- there are some words which signify similar meaning like tasty, tasteless (both are talking about the parent word taste) so, we can use only one representation.

There are various stemmer: ① Porter stemmer
② Snowball stemmer **

→ Lemmatization:- Lemmatization is all about breaking up the sentences into words.

Limit of Bow:- we are not considering semantic ^{meaning} relation of words. like tasty and delicious are treated as 2 diff words and hence diff dimension in Bow but tasty and delicious are very much similar.

→ Uni-gram | Bi-gram | n-gram:-

When we are learning stop words:- we saw various stop words.

x_1 : (This) parla (is) (very) tasty (and) affordable

x_2 : (This) parla (is) (not) tasty (and) (is) offer

↓ removing stop words

x_1 : parla tasty afford

x_2 : parla tasty afford

After removal of stop words $x_1 = x_2$ but in original ~~was~~ Review they are opposite. How to tackle this problem.

we use uni-gram / bi-gram to handle this problem.

Uni-gram:

this	is	but	very	tasty	offer
------	----	-----	------	-------	-------

In uni-gram we give have one dimension for each word.

Bi-gram: we have dimension for 2 consecutive words in sentence

e.g. This pasta is very tasty and affordable

bi-gram:

This pasta	pasta is	is very	very tasty	tasty and	and afford
------------	----------	---------	------------	-----------	------------

Similarly we can have n-gram (n consecutive words for each dimension)

So, Bi-gram or higher gram retains the sequential information

But the no. of bigrams > no. of uni-grams: hence it increases the vector size.

Hence n-gram \rightarrow dimensionality 'd' increases.
($n > 1$)

★ TF-IDF :- It is a variation of Bow

TF-IDF :- Term freq - Inverse document freq.

N docs / reviews

r_1	:	w_1, w_2, w_3, w_2, w_5	$\rightarrow 5$ words
r_2	:	$w_1, w_3, w_4, w_6, w_2, w_5$	$\rightarrow 6$ words
r_3	:		
\vdots			
r_N	:		

$$TF(w_i, r_j) = \frac{\# \text{ of time } w_i \text{ occurs in } r_j}{\text{Total \# of words in } r_j}$$

ex:- $TF(w_2, r_1) = \frac{2}{5}$

$$0 \leq TF(w_i, r_j) \leq 1$$

Basically TF says what is the probability of finding w_i in document r_i

IDF \Rightarrow Inverse document frequency.

So we have N documents / reviews.

$D = \{d_1, d_2, d_3, \dots, d_N\}$ D is collection of all documents / reviews.

$$IDF(w_i, D) = \log \left(\frac{N}{n_i} \right)$$

N \rightarrow Total No. of documents
 n_i \rightarrow No. of documents which contain word w_i .

we know that N is always greater than n_i hence

$$N/n_i > 1$$

$N/n_i = 1$ if word w_i occurs in every document?
The IDF of that word $= 0$.

So, if a word w_i is more frequent its IDF will be lower.

$r_1: w_1, w_2, w_3, w_2, w_4$

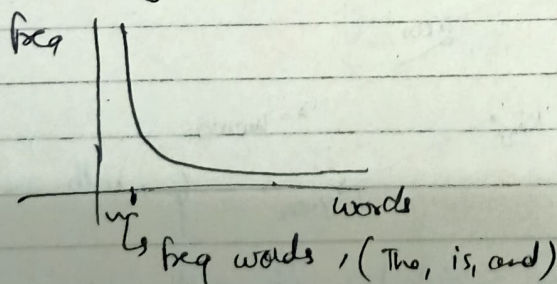
w_1	w_2	w_3	w_4	w_5

$\rightarrow TF(w_2, r_1) * IDF(w_2, D)$

each word is represented as its $TF * IDF$ value.

$TF-IDF \rightarrow$ give more imp to word rare in whole Docum but frequent in current document.

The reason behind the usage of \log in IDF is zipf's law



zipf's law states that most part of sentence is covered by small set of words:-

we can see that above plot is following power-law. In power-law we saw that to convert a $y \propto x^{-\alpha}$ power-law to gaussian distribution we took \log of x .

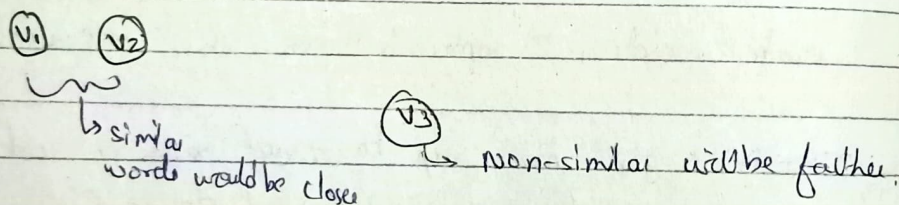
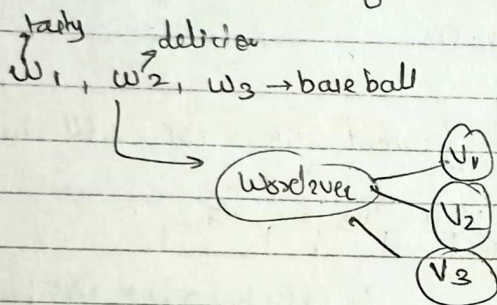
Also applying \log make the range of IDF smaller.

If IDF is very large then IDF will dominate over TF. To balance them we apply log to IDF to make it smaller.

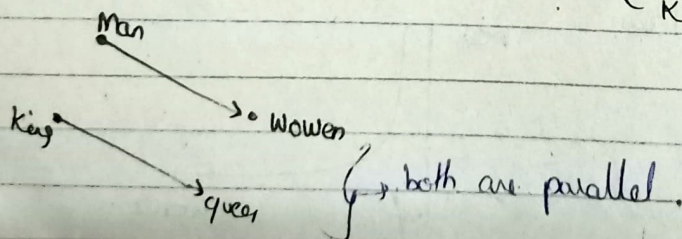
* Word2Vec :- This takes semantic meaning into consideration. We will see full mathematics of word2vec in class later. But now we only see it as black box.

word \rightarrow \rightarrow d-dim. vector
 \rightarrow not a sparse vector

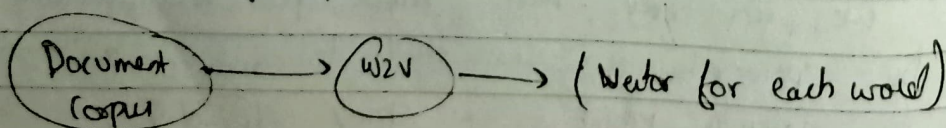
Note: Bow, TF-IDF takes sentence at a time, word2vec gives vector for each word.



- ① * w2v: * Take semantic meaning into consideration.
- ② It takes relationship into consideration. (Ex: Man, woman, King, Queen).



Similarly it captures various other types of relationship like country-capital.



→ Avg w_{2v} , tf-idf w_{2v} :-

w_{2v} gives vector for each word, but our sentences are collection of words. so how do I convert a sentence into vector using w_{2v} .

let $s_1: w_1 w_2 w_1 w_3 w_4 w_5$

$$\frac{w_{2v}(w_1) + w_{2v}(w_2) + w_{2v}(w_1) + w_{2v}(w_3) + w_{2v}(w_4) + w_{2v}(w_5)}{\# \text{ of words.}}$$

Avg- $w_{2v} \rightarrow$ works well, but not perfect.

→ tf-idf w_{2v} :-

$s_1: w_1 w_2 w_1 w_3 w_4 w_5$

↓ calculate tf-idf vector for s_1

$$\text{TF-idf} \begin{bmatrix} t_1 & t_2 & t_3 & t_4 & t_5 & t_6 \end{bmatrix}$$

$$\text{TF-idf-}w_{2v}(s_1) = \begin{bmatrix} t_1 w_{2v}(w_1) & t_2 w_{2v}(w_2) & t_3 w_{2v}(w_1) & t_4 w_{2v}(w_3) & t_5 w_{2v}(w_4) & t_6 w_{2v}(w_5) \end{bmatrix}$$

$$\text{TF-idf-}w_{2v}(s_1) = \frac{(t_1 * w_{2v}(w_1) + t_2 * w_{2v}(w_2) + t_3 * w_{2v}(w_1) + t_4 * w_{2v}(w_3) + t_5 * w_{2v}(w_4) + t_6 * w_{2v}(w_5))}{(t_1 + t_2 + t_3 + t_4 + t_5 + t_6)}$$

$$\text{TF-idf-}w_{2v}(s_1) = \frac{\sum_{i=1}^n (TF(w_i) * IDF(w_i) * w_{2v}(w_i))}{\sum_{i=1}^n (TF(w_i) * IDF(w_i))}$$

$$\sum_{i=1}^n TF(w_i) * IDF(w_i)$$