

ch1 Text data cleaning, preprocessing and vectorization

Video 2: Data cleaning: Deduplication

Video 3: Why convert text to a vector

We had studied vectors in linear algebra
n-dimensional vectors ~~at~~ time.

How do you convert text (English words and sentences) into numerical vectors.

Review text in the Amazon Fine Food Reviews is an n-Dim vector.

If all vectors are like :-

x_1, x_2, \dots, x_d d-dimensional space

+ve reviews / -ve reviews

hyper plane

d-dimensional representation of vectors.

w - normal

$$w^T x_1 \rightarrow +ve$$

to the plane

$$w^T x_2 \rightarrow -ve$$

If I somehow could find my w
then if $w^T x_i > 0$ then r_i is +ve
else r_i is -ve.

r_i = review.

Steps:

- 1) Converting text to d-dim - is the first step then we find the plane to separate.
- 2) Converting text to d-dim has some rules.

eg r_1, r_2, r_3 are vectors in d-dim space
similar semantically
english language

If $\text{Eng Sim}(r_1, r_2) > \text{Eng Sim}(r_1, r_3)$

then

distance $(v_1, v_2) < \text{distance}(v_1, v_3)$

If these reviews are similar, distance between the reviews in d-dim space should be less, ie, they should be close.

length of $v_1 - v_2 < \text{length of } v_1 - v_3$
minus

It will be easier to find plane when distance b/w the text/reviews is closer.

Conclusion:-

So, find text in d-dim vector such that similar text must be closer geometrically (computing geometric distance)

Video 4: Bag of Words (Bow)

Converts text to numerical vector

Steps in Bag of Words:-

1) Constructing a dictionary (not a Python dict)

Dictionary → Set of all the unique words across the document corpus (Collection of documents)

So, computing d unique words.

2) Creating d-sized vector for each unique word in the vector.

Each word is a different dimension

So wherever the word exists in the vector, we put 1 or where it occurs 2 times, we put 2 or 3 if 3 times.

$$v_i \in \mathbb{R}^d$$

v_i is a sparse vector.

Note

semantically
ie, English.

→ Similar text must result in closer vectors.

$$\text{length of } (v_1 - v_2) = \|v_1 - v_2\|$$

Toy example :-

This pasta is very tasty and normal

$$v_1 : \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 0 & 1 \end{bmatrix}$$

Subtract

$$v_2 : \begin{bmatrix} 1 & 1 & 2 & 0 & 1 & 1 & 1 \end{bmatrix}$$
$$\begin{bmatrix} 0 & 0 & 1 & 1 & 0 & 1 & 0 \end{bmatrix}$$

v_1 : This pasta is very tasty and affordable.

v_2 : This pasta is not tasty and is affordable

$$\text{So length } (v_1 - v_2) = \|v_1 - v_2\|$$
$$= \sqrt{1^2 + 1^2 + 1^2}$$

norm

$$= \sqrt{3}$$

$$\text{length of } v_1 = 8$$

Bag of words doesn't work well with small difference in words.

~~like~~

Variations of Bag of words:-

1) Binary Bag of words / Boolean BOW

affordable

Instead of count put
1 if it occurs (Otherwise
at least one)

1

1

0

$\|v_1 - v_2\| \sim \# \text{ of differing words.}$

$$\sqrt{(1)^2 + (1)^2} = \sqrt{2} \text{ distance b/w } v_1 \text{ & } v_2 \text{ in Binary BOW}$$

Words like This is and are not very useful. Also called as Stop words. But words which are trivial matter more.

Doubt

Video 5 : Text preprocessing : Stemming,

Stop word removal, Tokenization, Lemmatization

→ Stop words → and, is, this, etc.
→ words which are not important.

Very is ~~not~~ a stop word.

In English not is a stop word but for some sentences we don't get to know that it is a negative sentence, so be careful.

- 1) Removing Stop words carefully is called Text preprocessing. not is a stop word in eng which after removal might change the meaning of the sentence.
- 2) Make all your words lowercase.

- 3) Stemming : Eg - Tasty and Tasteful
Tasty base word as Tast
So instead of separating these words of same base lets make them one.

Multiple algorithms for stemming

a) Porter Stemmer (very old)

b) Snowball Stemmer. (more powerful)

5) Lemmatization: Algorithmic process of determining the lemma of a word based on its intended meaning. Unlike stemming, lemmatization depends on correctly identifying the intended part of speech and meaning of a word in a sentence as well as within the larger context.

4) Tokenization:

~~Tokenization~~ Given a sentence surrounding that sentence such as neighbouring sentences or even an entire document.

~~Tokenization~~ is all about breaking a sentence into words and groups them so they form a meaning (e.g - It is done by linguists New York)

Note:- New Tasty and Delicious are synonyms but we create 2 different dimensions in Bag of words. But in W2V considers semantic meaning of words unlike Bag of words.

Video 6: unigram, bigram and n-gram

This comes under lemmatization.

Unigram - Each word is considered a dimension.

Bigrams - Pair of words are considered as dimension.

Trigrams - 3 consecutive words make a dimension.

n-grams. $n=1$ Unigrams

$n=2$ Bigrams

$n=3$ Trigrams

unigram BoW - discards the sequence info.

But if you use bigrams, trigrams or n-grams maintain the sequence information in BoW it will be good.

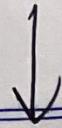
~~# bigrams \geq # unigrams.~~

and

$\# \text{bigrams} \geq \# \text{trigrams} \geq \# \text{unigrams}$.

So, n-grams ($n > 1$) dimensionalities increases drastically.

used in Information Retrieval (BOW, W2V, TFIDF)



What is tf-idf (Term frequency-inverse document frequency)

$$TF(w_i, \sigma_j) = \frac{\text{# of times } w_i \text{ occurs in word document } \sigma_j}{\text{Total # of words in doc } \sigma_j}$$

Defn

$$\downarrow \text{So } 0 \leq TF(w_i, \sigma_j) \leq 1$$

Term Frequency

is the probability of finding a word w_i in a document σ_j .

D(data corpus) = Collection of documents.

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

where $N = \# \text{ of documents in corpus}$
 $n_i^c = \# \text{ of docs in corpus which contain } w_i$

TFIDF → Doesn't take the semantic meaning of words like Tasty & Delicious.

Note $n_i \leq N$

$$\frac{N}{n_i} \geq 1$$

① So $\log\left(\frac{N}{n_i}\right) \geq 0$ $\because \log(1) = 0$.

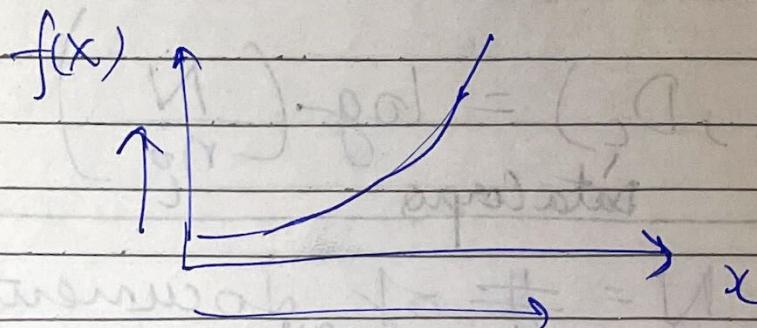
So $IDF \geq 0$.

②

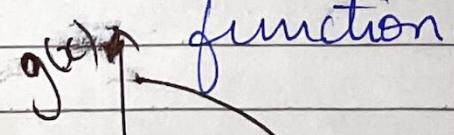
If $n_i \uparrow$, $\frac{N}{n_i} \downarrow$

Then $\log\left(\frac{N}{n_i}\right)$ also \downarrow

\log is a monotonic function.



∴ $f(x)$ is a monotonically increasing function.



As $x \uparrow$, $g(x) \uparrow$, it is monotonically increasing function.

$\log(x)$ is valid only if x is (+ve).

→ If w_i is more frequent in my corpus then its IDF will be lower.

→ As w_i is a rare word then its IDF will be higher.

If we combine TF & IDF

$$TF-IDF = \underbrace{TF(w_i, \text{doc})}_{\text{doc}} * \underbrace{IDF(w_i, \text{Corpus})}_{\text{Corpus}}$$

→ TF is ↑ when word occurs freq. in document
freq. = FTF ↑

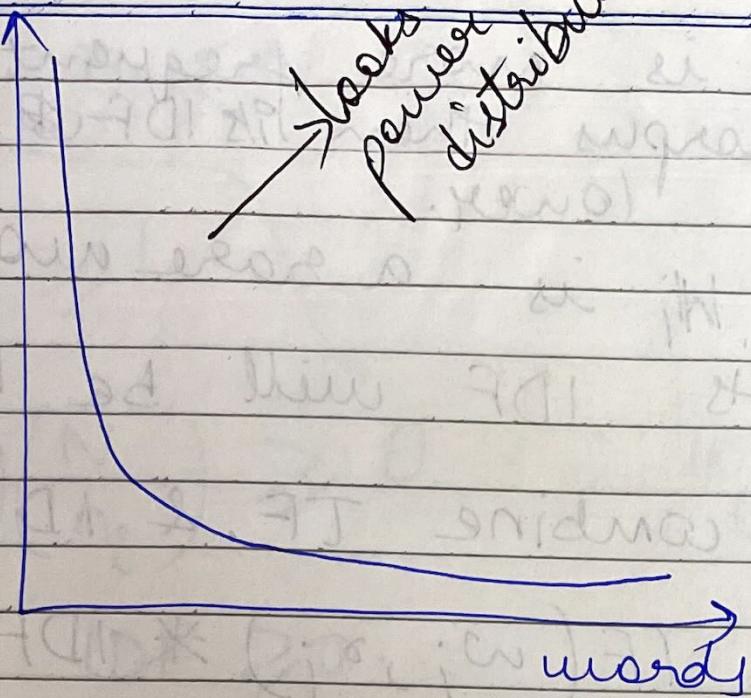
→ IDF is ↑ when word is rare in corpus
rare = RIDF ↑

Video 8: Why use log in IDF?

log was just used as a heuristic or a hack.

It's explanation given in Zipf's law. In Zipf's law if frequency of words is in y-axis and words occurs on x-axis. The histogram falls very sharply

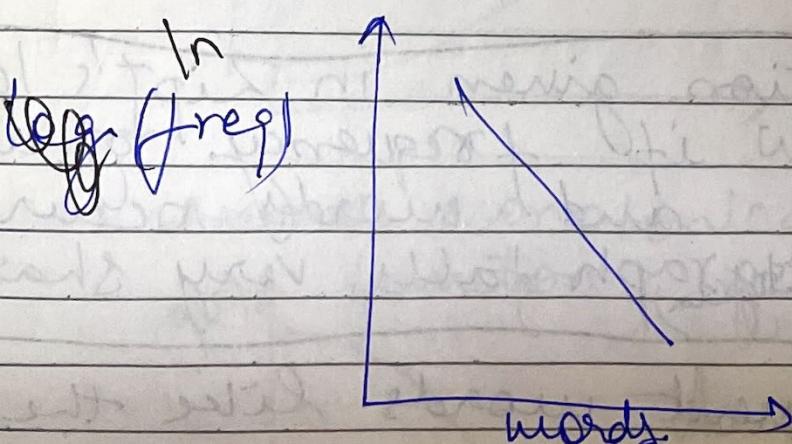
More frequent words like the, is, etc are more higher and many more words which are trivial occur very frequently



If we plot \log of freq. it becomes a straight line in graph.

It is just like we convert power law to Gaussian distribution by taking \log of it (Box-Cox transform).

We are taking ~~\log_{10}~~ (natural \log).



$\log_e(x) = \ln(x)$ same
natural log

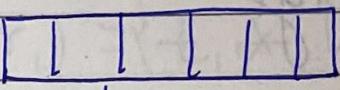
In this way, IDF in $TF \times IDF$
has less dominance

So TF-IDF is in range 0-20
not in 1000 (as it may happen
if we don't take log).

When this TF and IDF both
are in same range then their
multiplication makes sense and
IDF over TF is equal
in almost not too dominant
and TF not too less dominant.

Video 9 : Word2Vec

It takes semantic meaning of the
words/sentences unlike BoW and TFIDF

word \rightarrow 
d-dim vector
but not a sparse vector
and as a dense d-dim vector.

BoW } Sparse
TFIDF }
W2V - Dense

Word2Vec takes only 1 word
into consideration unlike the first 2
which take entire sentence into consideration.

similar

① If w_1 and w_2 are semantically closer then these points will be closer.

② It also satisfies relationships between words, like Male-Female gender of certain words. and it can also consider Country Capital relationships.

It learns all these relationships almost on its own.

$$(V_{man} - V_{woman}) \parallel (V_{king} - V_{queen})$$

Mikolov published this for W2V and Google uses it essentially in almost every NLP.

You need to give a very large ^{text} _{corpus} eventually the dimension will also become large.

Essentially W2V looks at the neighbours of each word

W2V from Google News is a 300 dim vector
Neighbour $N(w_i) \approx N(w_j)$
vector $- v_i \approx v_j$

$$tfidf(w_1) * W2V(w_1) + \dots$$

$$tfidf(w_1 \text{ in doc}) + tfidf(w_2 \text{ in doc}) + \dots$$

The below strategies work for converting sentences to vectors.

Video 10: Avg Word2Vec, tf-idf weighted Word2Vec

[Avg W2V]

doc γ_1 : $w_1, w_2, w_3, w_4, w_5 \rightarrow n_1$, words

vector representation of word w_i

$$v_1 = \frac{1}{n_1} [w_2v(w_1) + w_2v(w_2) + w_2v(w_3) + w_2v(w_4) + w_2v(w_5)]$$

each $w_2v(w_i)$ is a d-dim vector so we take avg. of its sum.

works fairly well, not always.

[TFIDF weighted W2V]

doc γ_1 : w_1, w_2, w_3, w_4, w_5

$w_1, w_2, w_3, w_4, w_5, w_6, w_7$

tfidf: $[t_1 | t_2 | t_3 | t_4 | t_5 | t_6 | t_7]$

Vector
for γ_1 doc

$$\begin{aligned} \text{tfidf } w_2v(\gamma_1) &= \{t_1 * w_2v(w_1) + t_2 * w_2v(w_2) \\ &+ t_3 * w_2v(w_3) + t_4 * w_2v(w_4) + t_5 * w_2v(w_5)\} \\ &\div \{t_1 + t_2 + t_3 + t_4 + t_5 + t_6 + t_7\} \end{aligned}$$

$$\begin{aligned} \text{tfidf-w2v}(\gamma_1) &= \sum_{i=\text{words}} (t_i * w_2v(w_i)) \\ &\div \sum_{i=\text{words}} t_i \end{aligned}$$