

# S&Litt

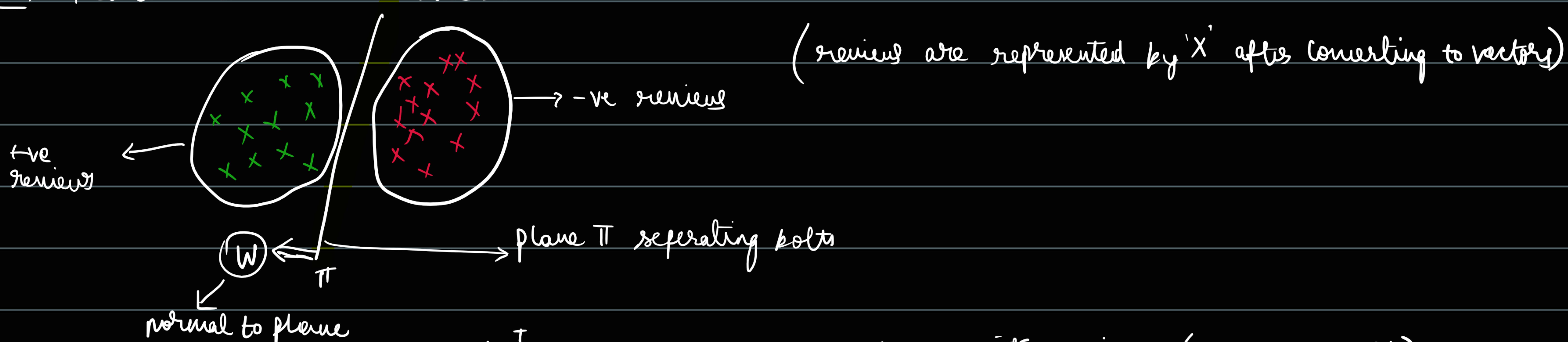
```
pd.read_sql_query("""SELECT name FROM sqlite_master WHERE type='table' ORDER BY name; """, con) → Prints list of tables
```

`→ df.drop_duplicates(subset = [", ", "..."], keep = 'first')` (Pandas)  
↳ drops duplicate values with these column names & keeps the first occurrence.

→ Use common sense to clean data: No general techniques

→ how to convert text to numerical vectors?

Ex:- Review tent  $\rightarrow$  n-dim vector



$w^T x_i > 0 \Rightarrow$  review  $x_i$  is a positive review (on movie  $m_i$ )

$u^T x_2 < 0 \Rightarrow$  review  $x_1$  is a negative review (on diff side)

→ Rules for conversion from text to vectors :-

$$\gamma_1 \quad \gamma_2 \quad \gamma_3$$

assume semantic similarity  $b(\omega(r_1, r_2) > (r_1, r_3))$

then distance b/w  $(v_1, v_2) < \text{dist b/w } (v_1, v_3)$

i.e., if they are similar, vectors must be close.

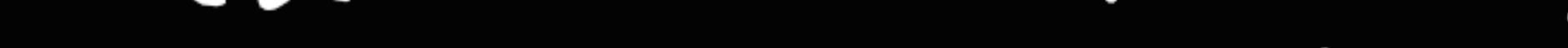
$v_1$   $x$   $v_2$   $x$   $v_3$

# Bag of words is

Step 1 :- Constructing dictionary (set of all unique words)

Step ② Assume there are n unique works across all collections of documents/reviews.

For each review

$d$  (d dimensional vector) 

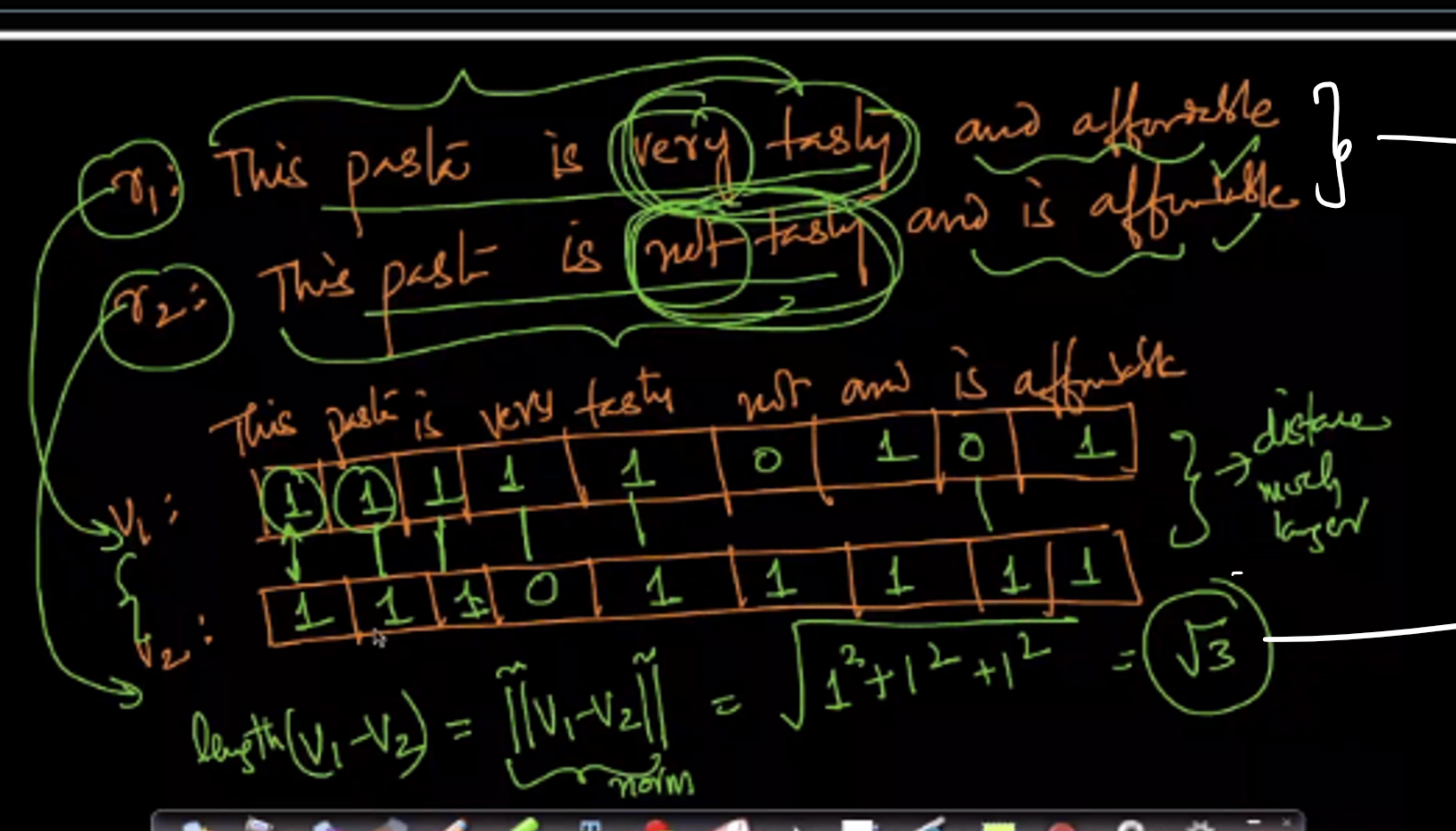
is called a corpus.]

$\gamma_1 := [0|1|1|2|1| \dots \dots \dots ]$  | 0

$\pi_j :=$  .

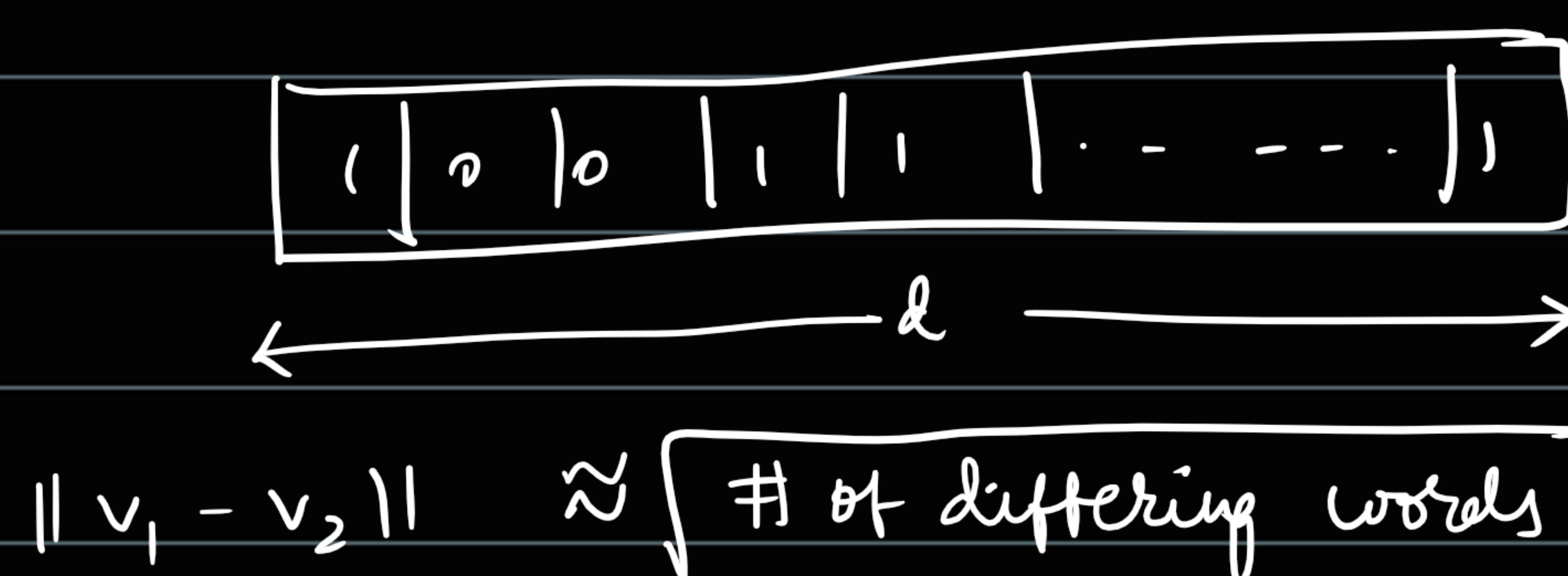
The number in bold represents the number of times  
that word occurred in the review.

Each word is a different dimension



## Boolean bag of words / binary bag of words :-

Instead of counting number of occurrence, it only counts occurrences -



→ We can notice that some of the words in BoW don't matter (i.e., and, a, the etc,)

→ Bag of Words is still used for creating baselines with Generalized Linear Models.

Stopwords

→ Removing stop word = smaller dimension data.

These are called text preprocessing steps :-

→ from nltk.corpus import stopwords

→ We also make everything lowercase.

→ Stemming :- tasty, tasteful, tasty all have same meaning.

converting all into one word 'taste'. This is called stemming.

Porter Stemmer  
Snowball Stemmer } two techniques for stemming  
more powerful.

→ Tokenization :- Breaking sentence into words. e.g. if we use space as delimiter, New York are separate words.

→ In BoW we are not taking semantic meaning.

e.g. delivering tasty convey same meaning. We use word2vec to solve this -

**Lemmatisation** is the 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 surrounding that sentence, such as neighboring sentences or even an entire document.

→ Algorithm combines inflected words into a common word called lemma.

{go, going, went} => 'go' is the lemma given by algo.

→ NLTK makes use of parts of speech. This is called tagging. Makes sense of context and helps pick the right lemma.

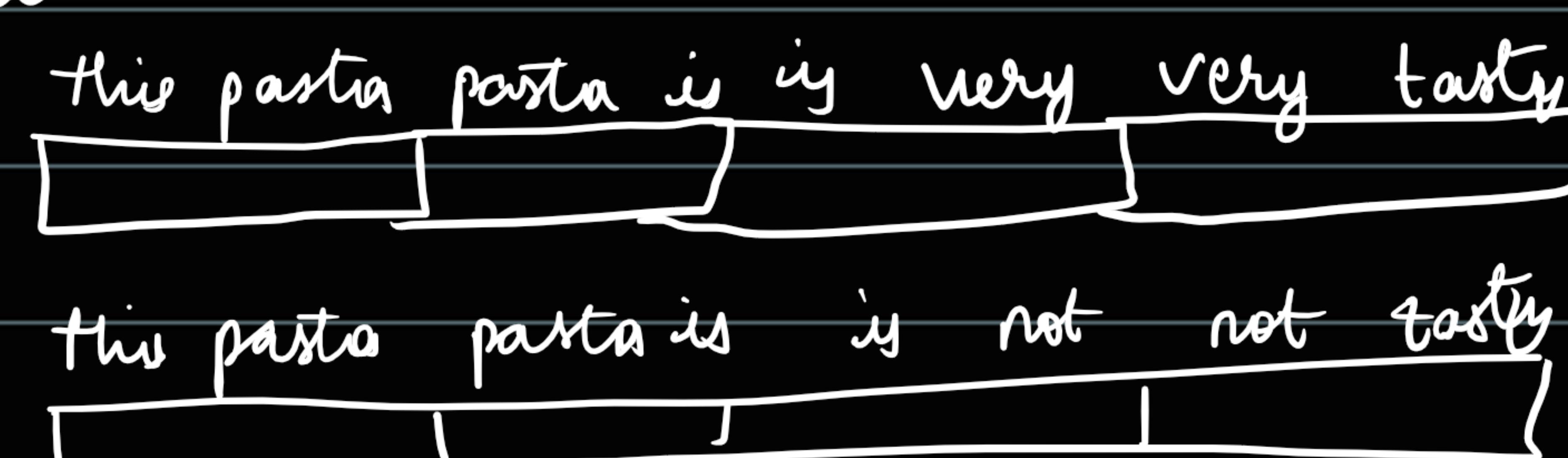
→ Stopwords are not removed in stemming. All words are reduced to base words and the tails are removed.

## unigram / Bi-gram / n-gram :-

unigrams :- each word is considered as a dimension

bigrams :- pairs of words

n = 3 → trigrams



The difference will be higher than that of unigrams.

→ unigrams BoW completely discards sequence information

Bigrams, trigrams, . . . retains partial sequence information.

$\rightarrow$  But the catch is  $n(\text{bigram}) > n(\text{unigram})$

$\rightarrow$  As  $n$  increases, dimensionality increases drastically

→ But if there are no repeated words in text, #unigrams > # bigrams

TF-IDF :- (term frequency in document frequency) (inverted index of BoW)

$$\begin{array}{l}
 \gamma_1: w_1 \quad w_3 \quad w_2 \quad w_5 \quad w_1 \\
 \vdots \\
 \gamma_2: w_3 \quad w_1 \quad w_4 \quad w_1 \quad w_6 \quad w_2 \\
 \vdots \\
 \gamma_N: w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6
 \end{array}$$

$$TF(w_i, r_j) = \frac{\# \text{ of times } w_i \text{ occurs in } r_j}{\text{total number of words in } r_j}$$

→ TF-IDF was originally meant for information retrieval

→ TF gives how often word occurs. (Probability)

$$0 \leq \text{TF}(w_i, y_j) \leq 1$$

IDF :-

$$D_L = \{r_1, r_{21}, \dots, r_n\}$$

$$\text{Document Corpus} \left\{ \begin{array}{l} D_c = \\ r_1: \overbrace{\quad w_1 \quad} \\ r_2: \overbrace{\quad \quad \quad} \\ \vdots \quad -w_1 \quad \quad \quad \\ \vdots \quad \overbrace{\quad w_i \quad} \\ \vdots \quad \quad \quad \end{array} \right. \rightarrow DF(w_i, D_c) = \log \left( \frac{N}{n_i} \right)$$

Number of dogs that contain  $w_i$

As  $n_j$  increases,  $\frac{N}{n_j}$  decreases. As  $\log(\frac{N}{n_j})$  also decreases, it is a monotonic function.

If a word is more frequent in corpus, then its IDF is low. If it's less frequent in corpus, then its IDF is high.

→ combining both

$$w_i = \text{TF}(w_i, r_i)^* \text{IDF}(w_i, D)$$

$v_i := [ \quad | \dots | \quad | \dots | ]$

$D_C$

$\gamma_1 :=$  \_\_\_\_\_

$\sigma_2 :=$  \_\_\_\_\_

.

.

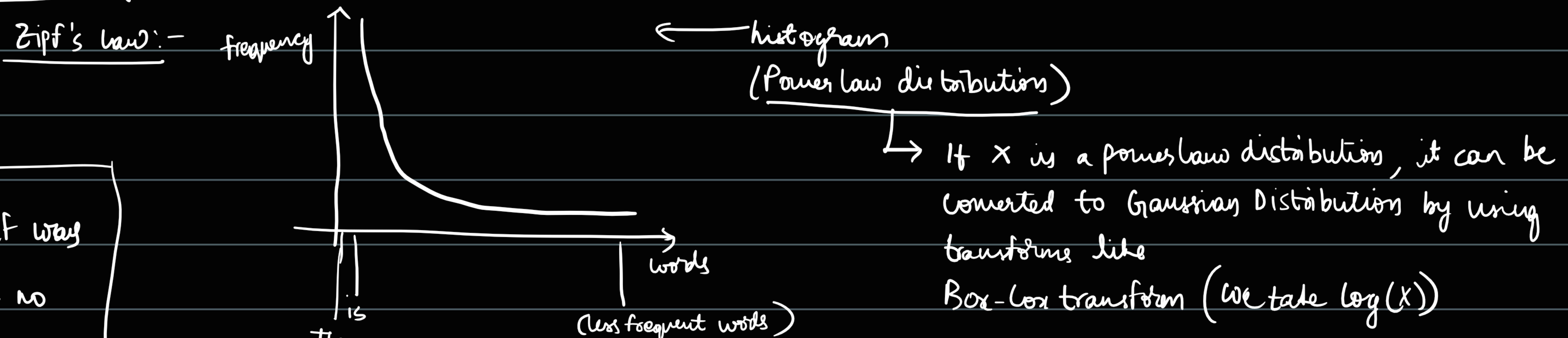
.

$\gamma_n :=$  \_\_\_\_\_

→ Advantage is that we give more importance to → word is more frequent (from TF)  
→ word is less frequent (from DF)

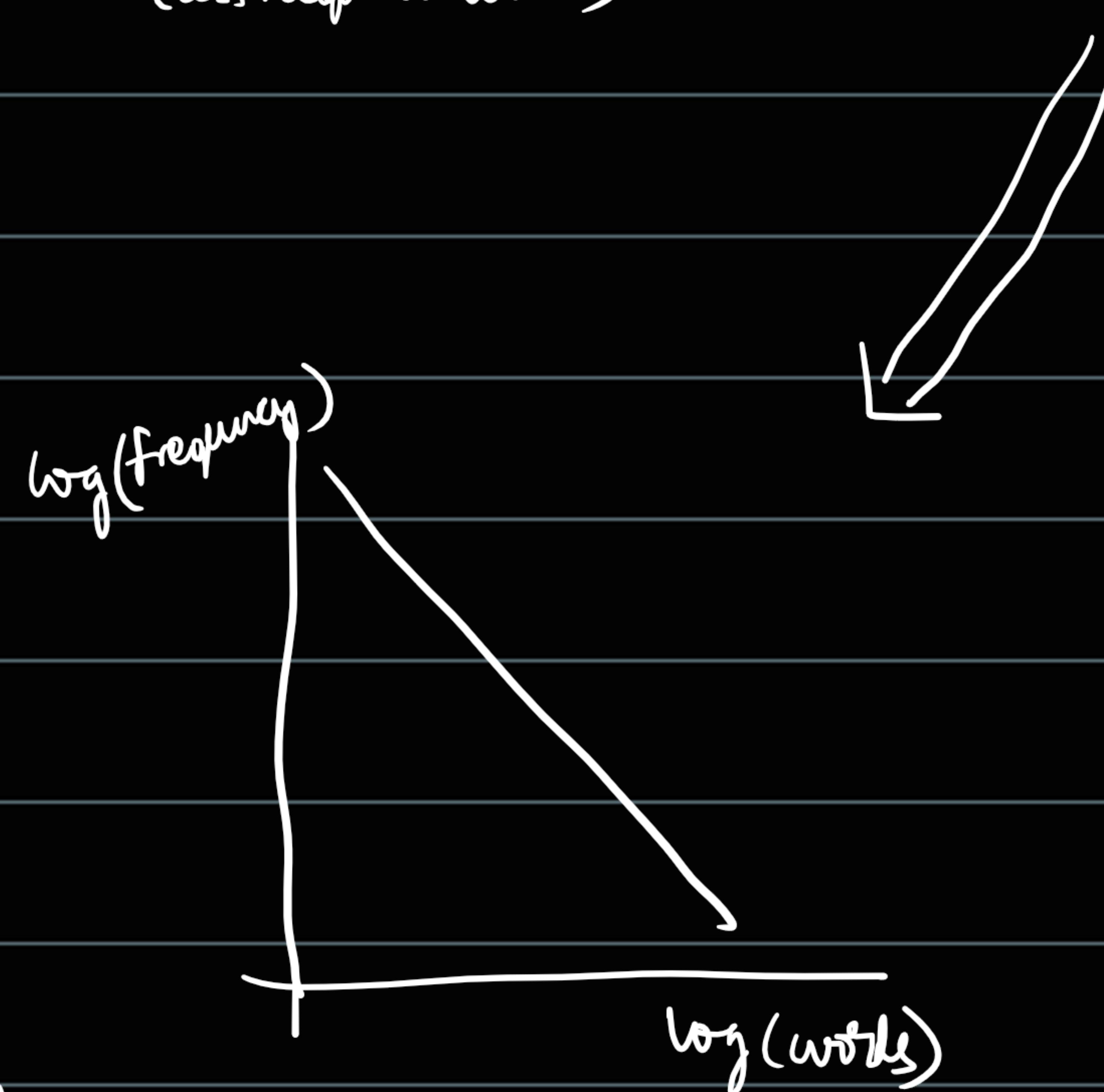
→ Drawback :- Still no semantic importance.

Why do we use log in IDF? :-

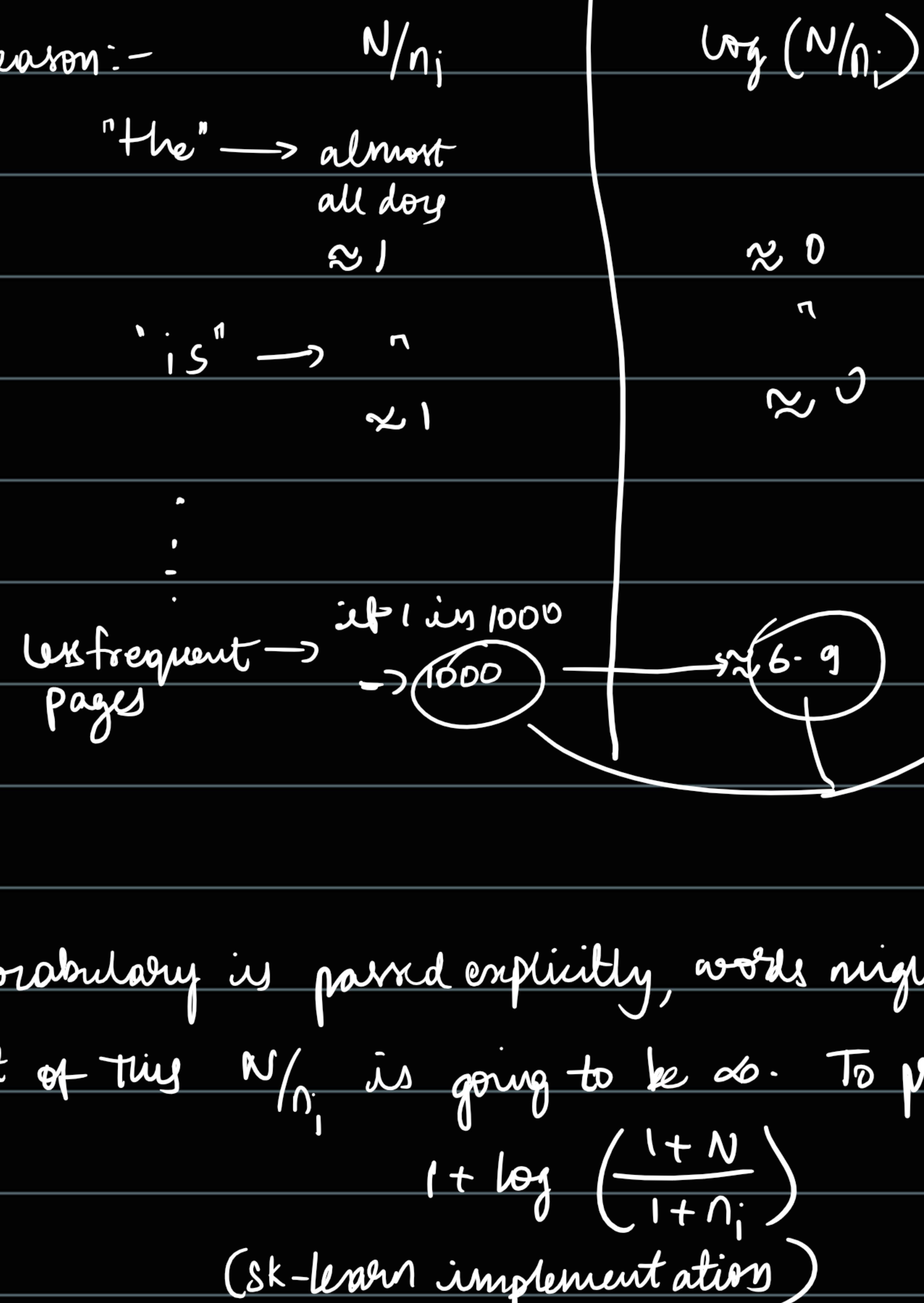


When tf-idf was introduced, no mathematical explanation was

given as to why log was used in idf



→ Another reason:-



→ lot of difference. During calculation of TF-IDF, IDF will dominate if there's no log.

→ Sometimes if vocabulary is passed explicitly, words might not occur in training phase, but might appear in testing phase. As a result of this  $N/n_i$  is going to be 0. To prevent this, we sometimes do

$$1 + \log \left( \frac{1+N}{1+n_i} \right)$$

(sk-learn implementation)

→ TF :- How frequently word occurs in a document

IDF :- How rare is the word in the whole corpus. If word is more frequent in the document corpus, it is less important.

TF-IDF :- Importance of the word in the whole corpus.

word2vec :- State of the art word to vector conversion that takes semantic meaning into consideration.

→ The dimensions are usually really high (like 300...)

→ It tries to achieve the following :-

① Semantically similar words have closer vectors

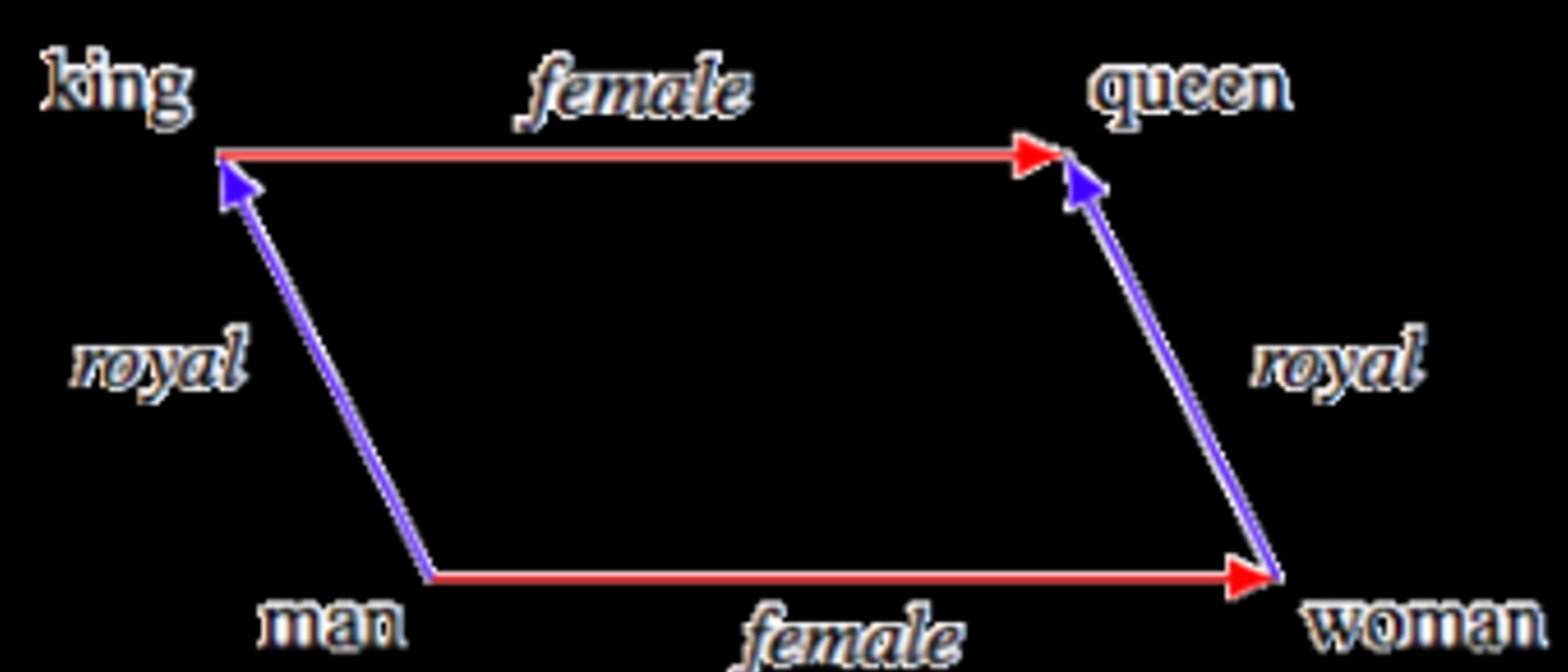
tasty

delicious

sweet

sour

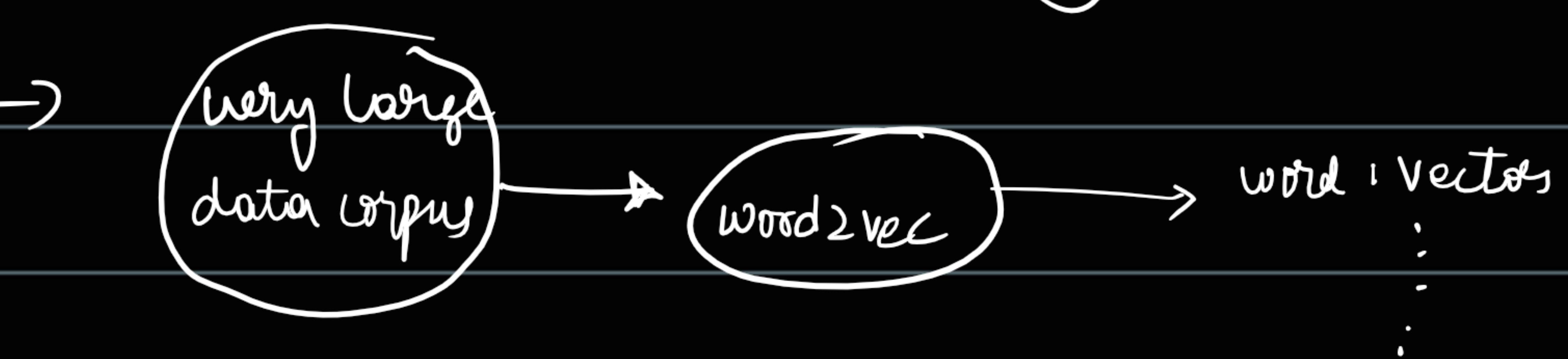
② It also tries to keep relationships



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

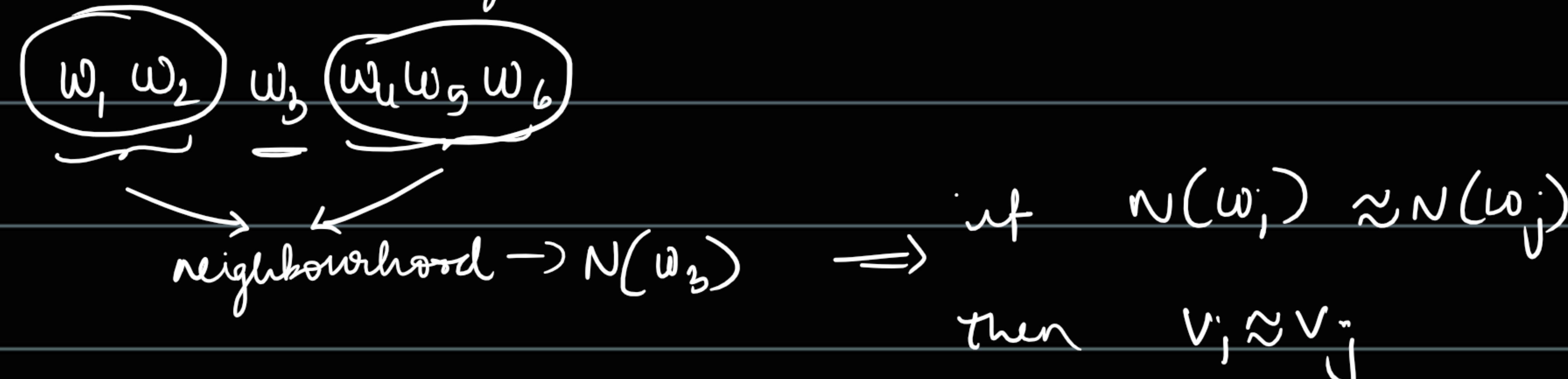
→ Even for country capitals etc,

→ It does this automatically without being programmed to do so. Came out in 2013.



The more the dimension of vector, the more information rich it is.

→ At its core, word2vec looks at the neighborhood of word and looks for words with similar neighborhood.



→ Google trained word2vec on its Google News dataset & this word2vec gives 300 dimension vectors.

Avg word2vec, tfidf weighted word2vec:-

word → word2vec → vector.

→ For sentences, sentence2vec exists but there are other methods like average word2vec.

Average word2vec :-  $r_i := w_1, w_2, w_3, w_4, w_5$

$$r_i := \frac{1}{n_i} [w_2v(w_1) + w_2v(w_2) + w_2v(w_3) + w_2v(w_4) + w_2v(w_5)]$$

$n_i$  = number of words in  $r_i$ .

→ Not perfect but simplest method.

tfidf w2v :- first tfidf vector of all words of sentence is calculated.

then each word's  $w_2v$  is multiplied with its tfidf score. and whose thing is added & divided by sum of tfidf of all words.

$$\text{tfidf-w2v}(r_i) = \frac{[t_1 * w_2v(w_1) + t_2 * w_2v(w_2) + t_3 * w_2v(w_3) + \dots + t_n * w_2v(w_n)]}{(t_1 + t_2 + t_3 + \dots + t_n)}$$

$t_n$  → tfidf of word  $w_n$

$$= \frac{\sum_{i=1}^n (t_i * w_2v(w_i))}{\sum_{i=1}^n t_i}$$

If tfidf of all words is 1, then their avg will be 1, then tfidf w2v becomes avg w2v

→ tfidf w2v & avg w2v are weighted schemes and are simple models. They are not always perfect.

→ thought vectors is a technique that is bit advanced. (Seen in deep learning chapter)

→ Each row of sparse matrix is called sparse vector.

→ If an  $n \times m$  matrix has only ' $K$ ' non-zero elements, sparsity of the matrix =  $\frac{(n \times m) - K}{(n \times m)}$

density of matrix =  $1 - \text{sparsity}$

vallakatla Charith

How to avoid removing stop word 'not' from text ?

9 Votes

Reply  

AppliedAIcourse

you can remove the 'not' from stopwords set, syntax: set.remove(element) [stopwords.remove('not')]  
will do the work

May 04, 2018 11:19 AM

May 04, 2018 11:55 AM

5 Votes

Pavan Kalyan Reddy Ramireddy

Can you explain the parameters of CountVectorizer such as 'preprocessor' and 'tokenizer' ?

Reply  

Apr 23, 2019 11:22 AM

Applied AI Tech Admin

So far you have seen preprocessing like converting the uppercase words into lower case, stopword removal, removal of punctuations, etc.

**preprocess:**

When you apply CountVectorizer, it splits the text into features directly into the words on the basis of spaces. So before in case if we need to perform any addition preprocessing on the text before splitting the text into the features we have to define those additional preprocessing steps in the form of a function(say f1) and pass this function to the parameter **preprocess**(i.e., preprocess=f1)

**tokenizer:**

If we want to go for a customized way of tokenizing the text(i.e., rather than splitting the text into words on the basis of white space, if we want to split the text in a different and a customized way), then define those new tokenizing steps in the form of a function(say f2) and pass the function name as a value to the parameter **tokenizer**. (i.e., tokenizer=f2)

Reply 

Apr 23, 2019 13:41 PM

→  $\text{min\_df} = 0.1$  → ignore words that appear in less than 10% of documents

= 10 → ignore words that appear in less than 10 documents.

$\text{max\_df}$  :- Ignore words that appear too frequently. AKA corpus specific stop words.

→ In w2v (Google's), words & their sense might not be similar

Vijay Maurya

Hello Team I have downloaded google model by this method

1. For download

!wget -c "https://s3.amazonaws.com/dl4j-distribution/GoogleNews-vectors-negative300.bin.gz"

2. For extraction

!7z e GoogleNews-vectors-negative300.bin.gz

3. Google model

```
# from gensim.models import KeyedVectors
# w2v_model_google=KeyedVectors.load_word2vec_format('GoogleNews-vectors-negative300.bin',
binary=True)
Don't need download anything run it in Google Colab
```

27 Votes

→ To load big datasets in colab w/o downloading.