

S&L Italic

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

```
: pd.read_sql_query("""SELECT name FROM sqlite_master WHERE type='table' ORDER BY name; """, con)
```

→ Prints list of tables.

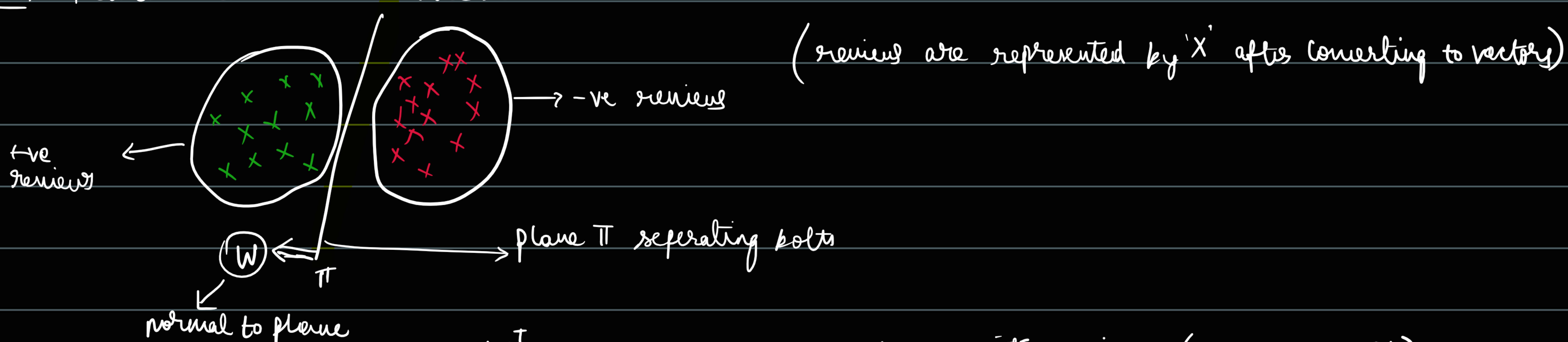
$\rightarrow df.drop_duplicates(subset=[\text{''}, \dots], keep='first')$ (Pandas)

↳ drops duplicate values with these column names & keeps the first occurrence

→ Use common sense to clean data: No general technique

→ 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 i)

$U^T x_2 < 0 \Rightarrow$ review x_2 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) < 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:

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

Step 2 - Assume there are n unique works across all documents. [Collection of documents]

For each review

d (d dimensional vector)

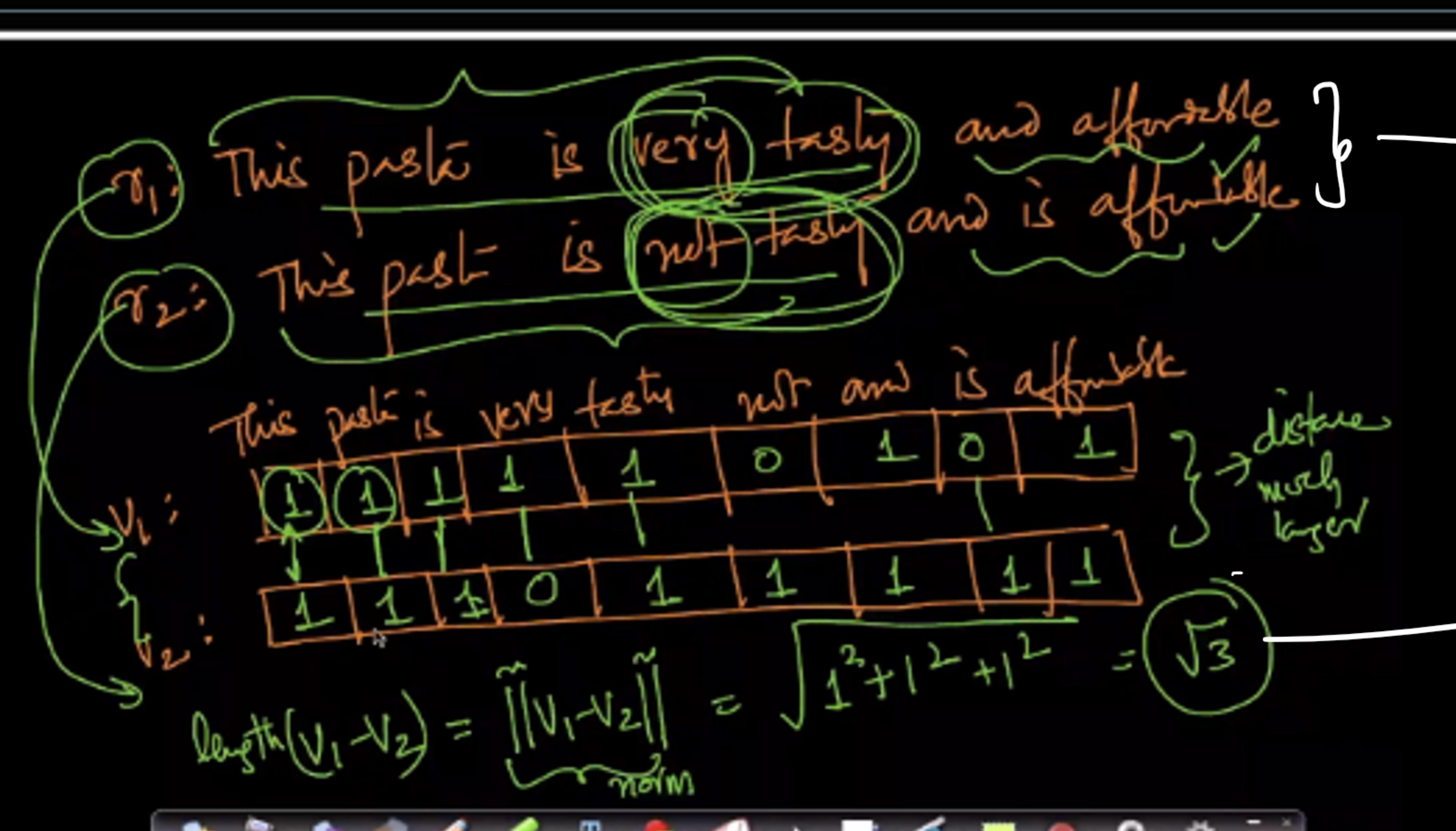
is called a wren.]

$\gamma_1 := [0\backslash 1\backslash 1\backslash 2\backslash 1\backslash \dots \dots \dots]_0$

$\pi_j :=$ *j*

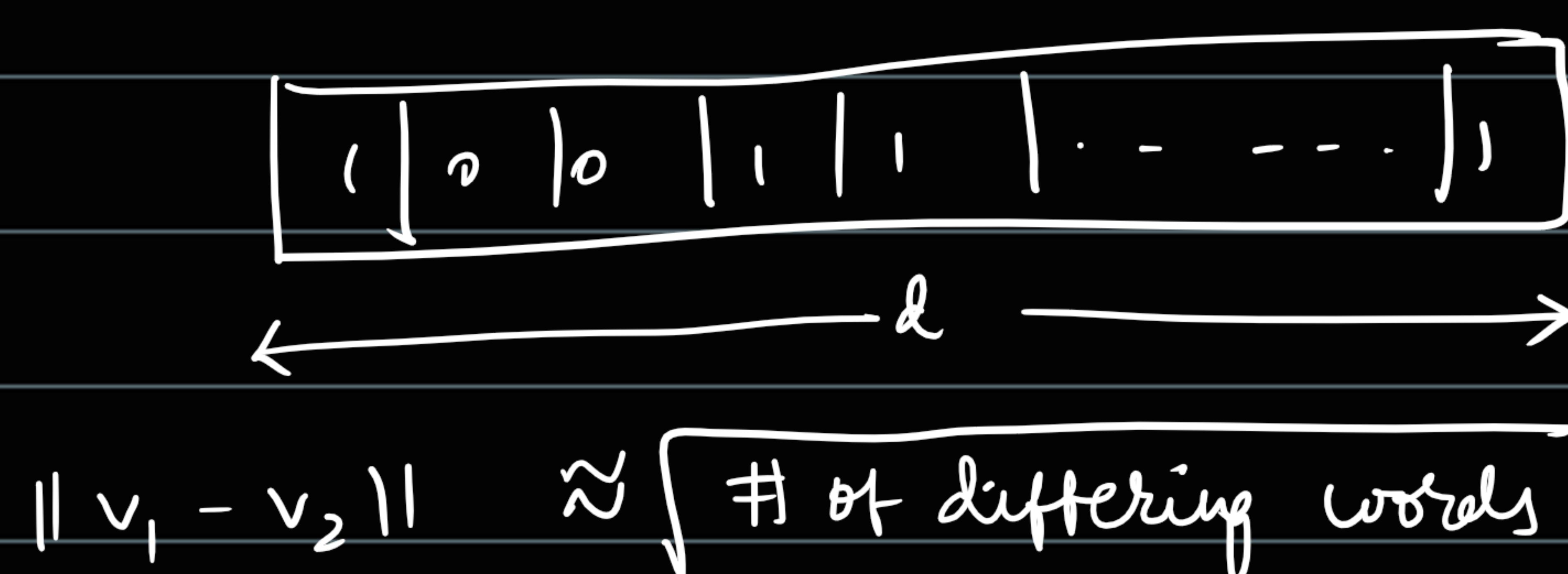
The number in bold represents the number of times
that word occurred in the review.

Each word is a different dimension



Boolean bag of words / binary bag of words :-

Instead of counting number of occurrences, 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 words = 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.

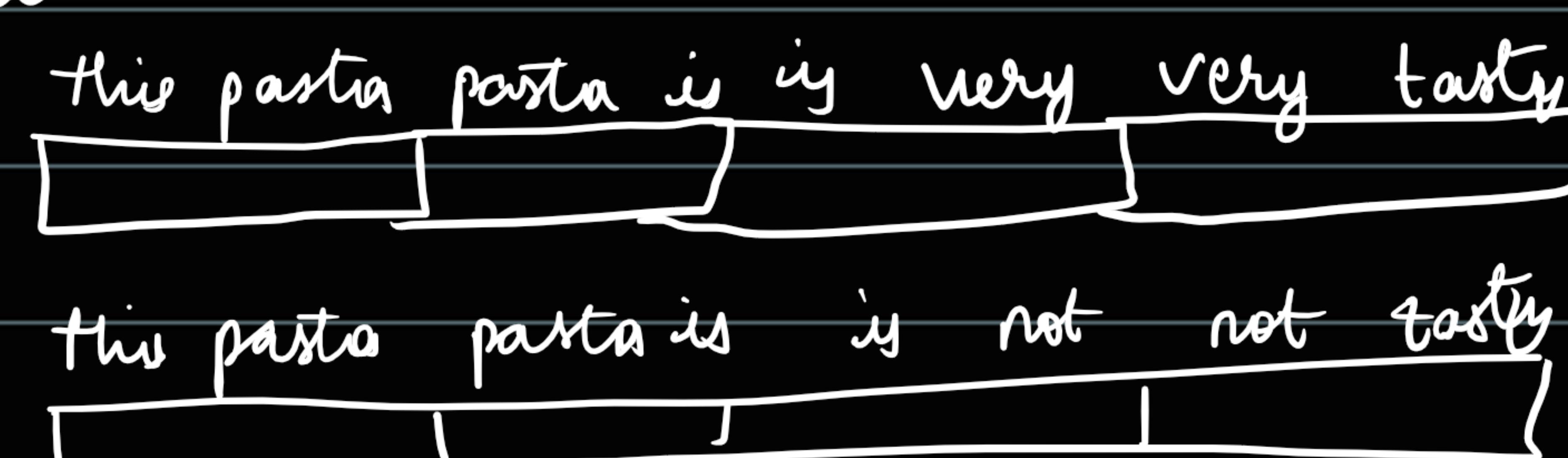
→ Stopwords are not removed in stemming. All words are reduced to base words and the tails are removed.

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.

→ but the catch is $n(\text{bigrams}) > n(\text{unigrams})$

→ As n increases, dimensionality increases drastically.

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

TF-IDF:- (term frequency × document frequency) (industrialized version of BoW)

	w_1	w_2	w_3	w_4	w_5	w_6	w_7
r_1 :	W ₁	W ₃	W ₂	W ₅	W ₁		
r_2 :	W ₂	W ₁	W ₇	W ₁	W ₆	W ₂	
.							:

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

→ TF-IDF was originally meant for information retrieval.

→ TF gives how often word occurs. (Probability)

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

IDF:-

$$D_c = \{r_1, r_2, \dots, r_n\}$$

$$\begin{aligned} D_c = \{ & r_1: \underline{\quad w_1 \quad}, \dots \\ & r_2: \underline{\quad \quad \quad}, \dots \\ & \vdots \quad \underline{\quad w_i \quad}, \dots \\ & \vdots \quad \underline{\quad \quad \quad} \} \end{aligned} \rightarrow IDF(w_i, D_c) = \log \left(\frac{N}{n_i} \right)$$

\circlearrowleft Number of docs that contain w_i

$$n_i \leq N \Rightarrow \frac{N}{n_i} \geq 1$$
$$\Rightarrow \log \left(\frac{N}{n_i} \right) \geq 0$$

⇒ As n_i increases, N/n_i decreases. As N/n_i decrease, $\log(N/n_i)$ also decreases as log is 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

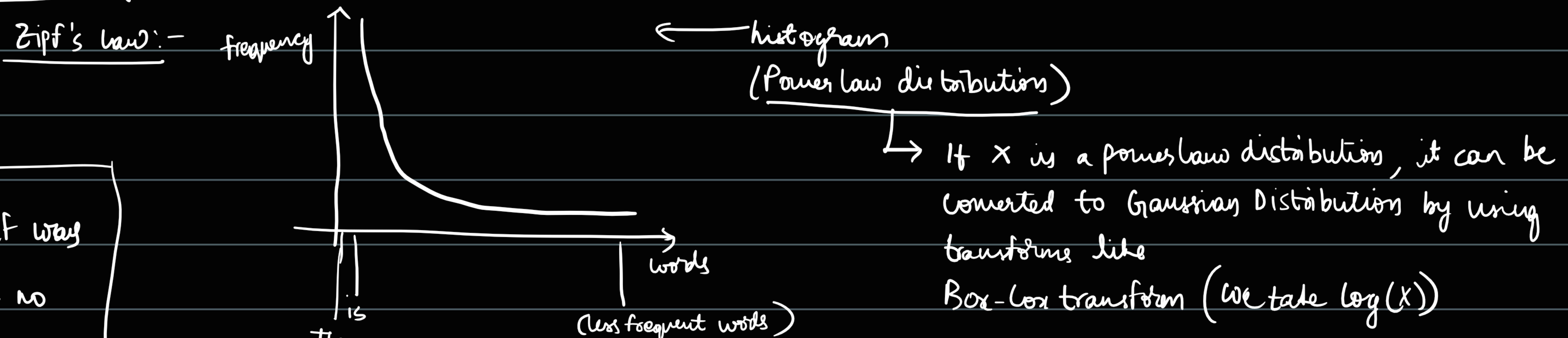
$$w_i = TF(w_i, r_i) * IDF(w_i, D_c)$$

$$D_c = \{ r_1: \underline{\quad \quad \quad}, \dots, r_n: \underline{\quad \quad \quad} \}$$

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

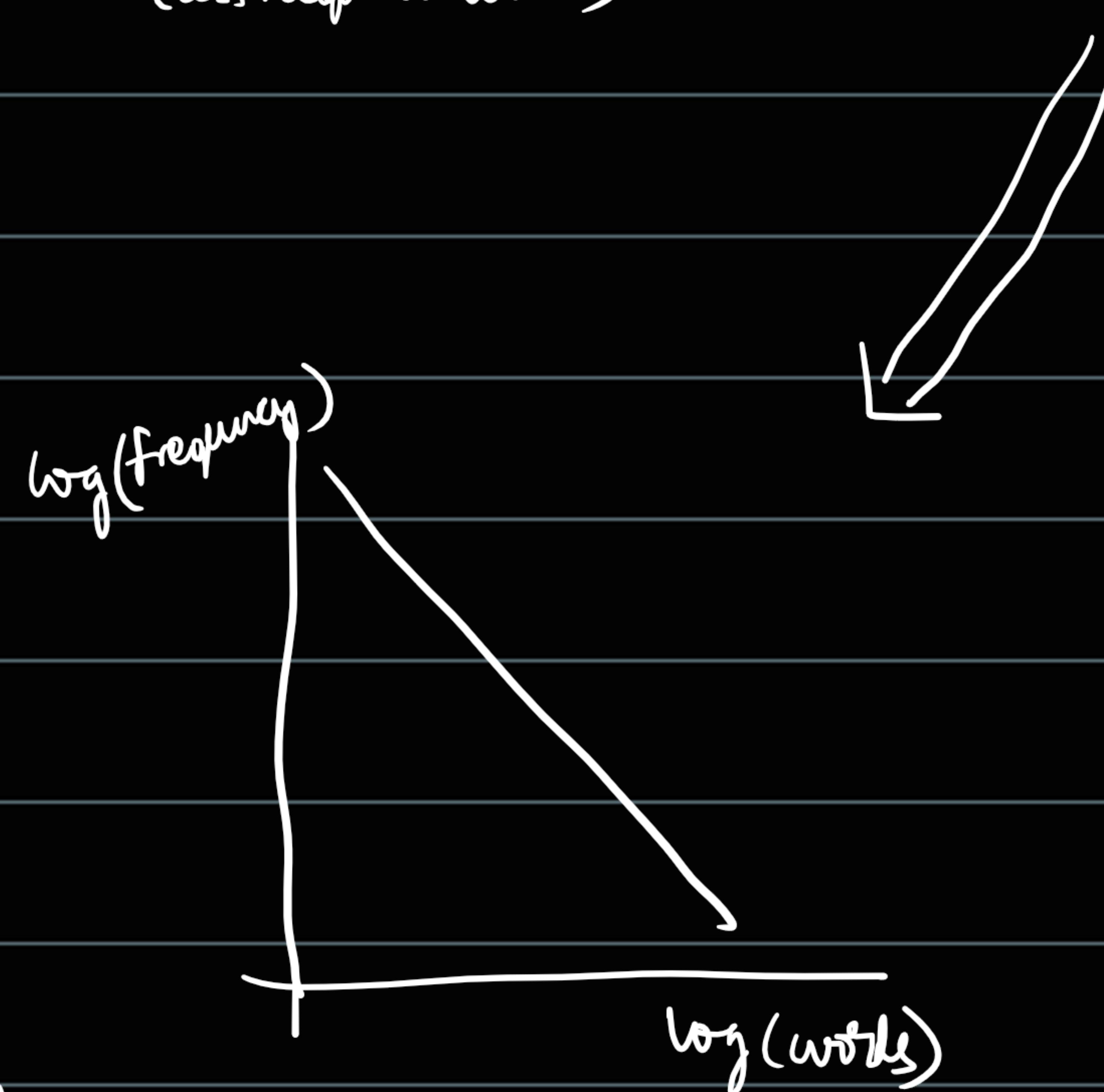
→ Drawback:- still no semantic importance.

Why do we use log in IDF? :-

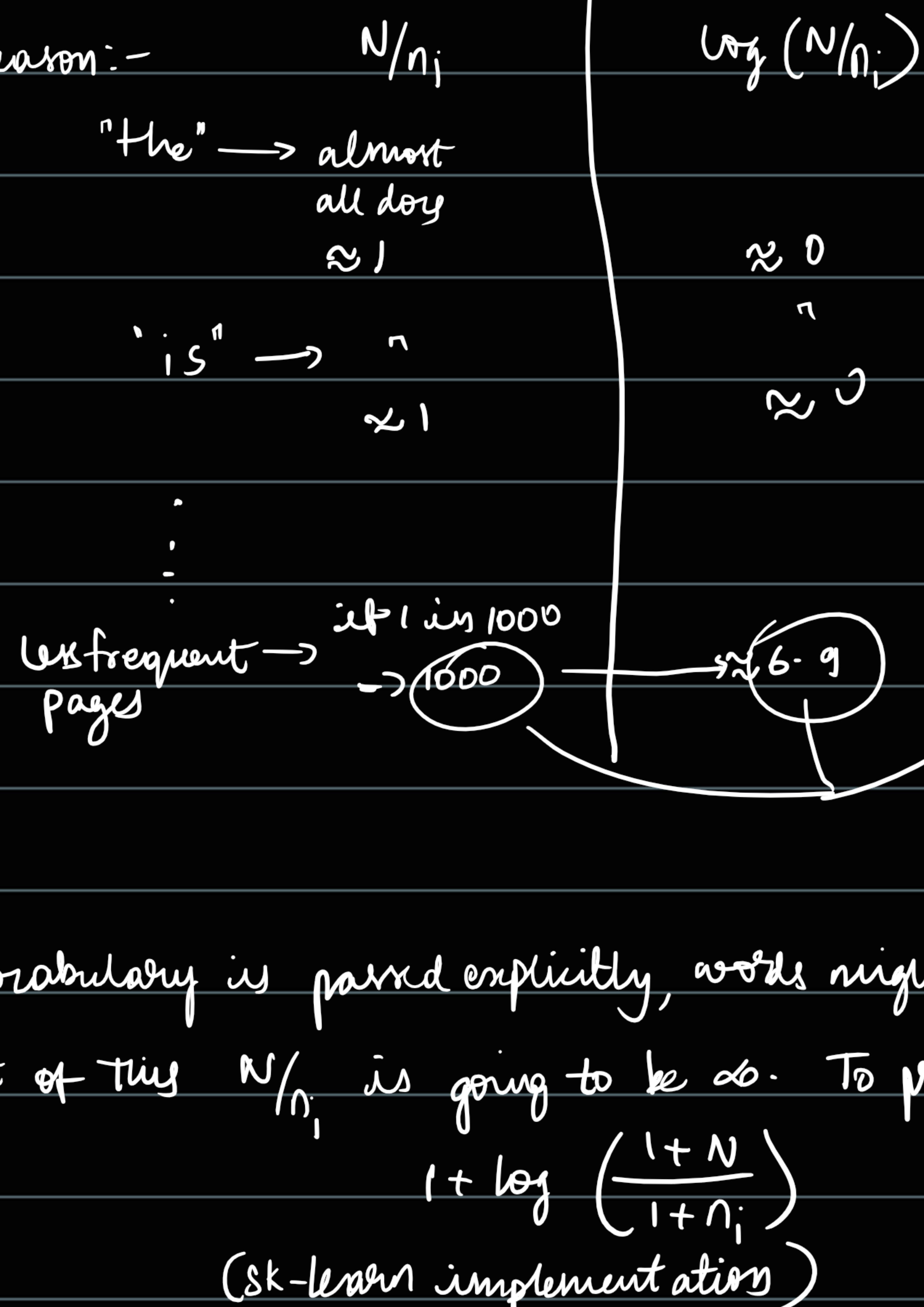


When tf-idf was introduced, no mathematical explanation was

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

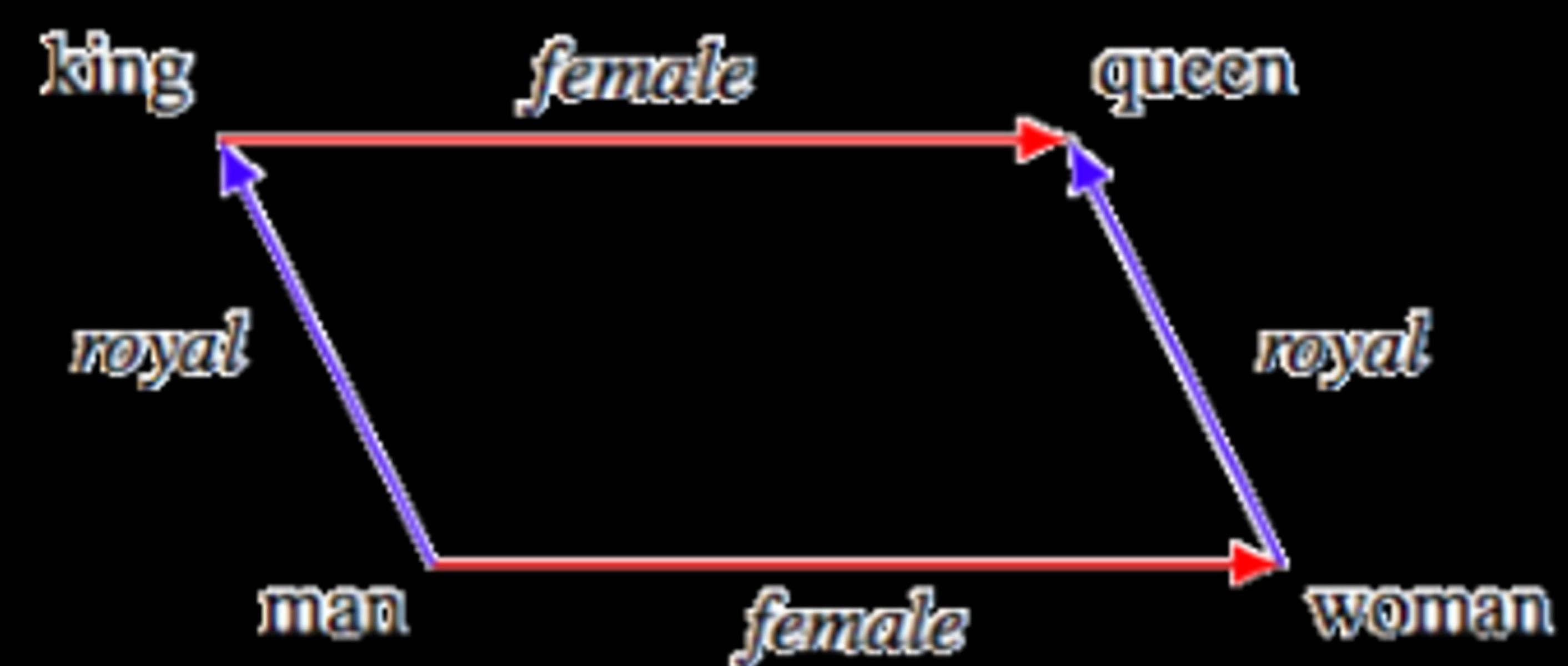
tasty

delicious

sweet

sour

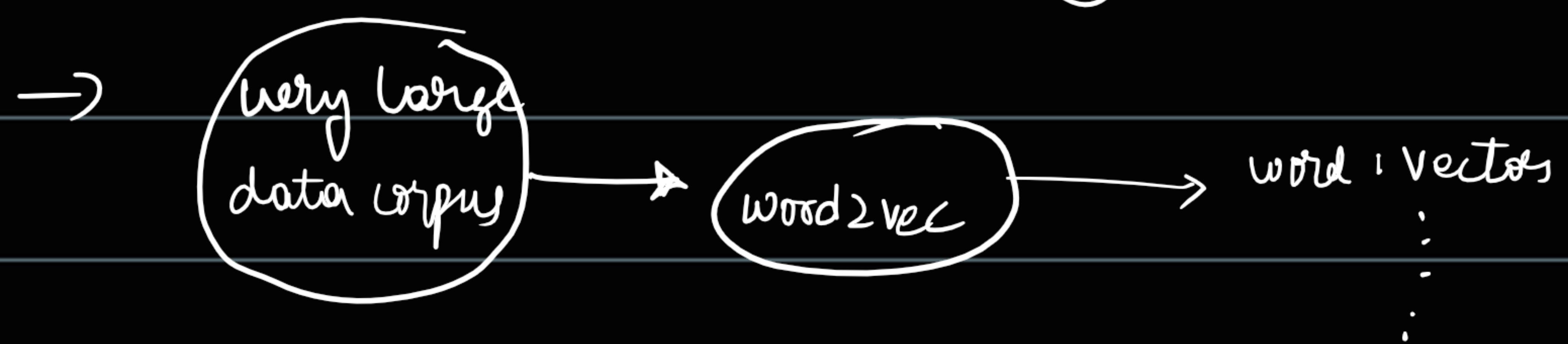
② It also tries to keep relationships



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

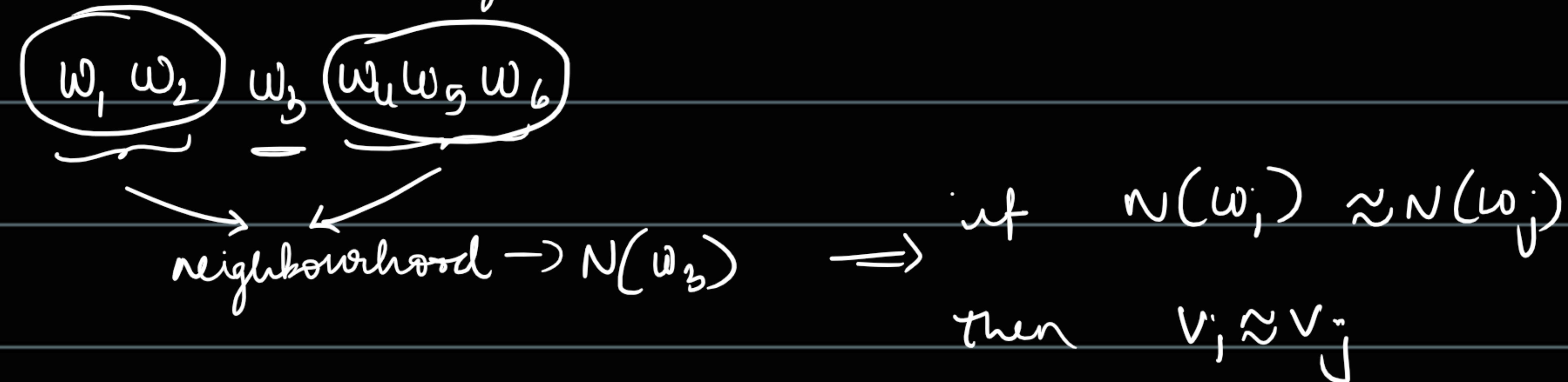
→ Even for country capitals etc,

→ 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 words and looks for words with similar neighborhood.



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