

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.

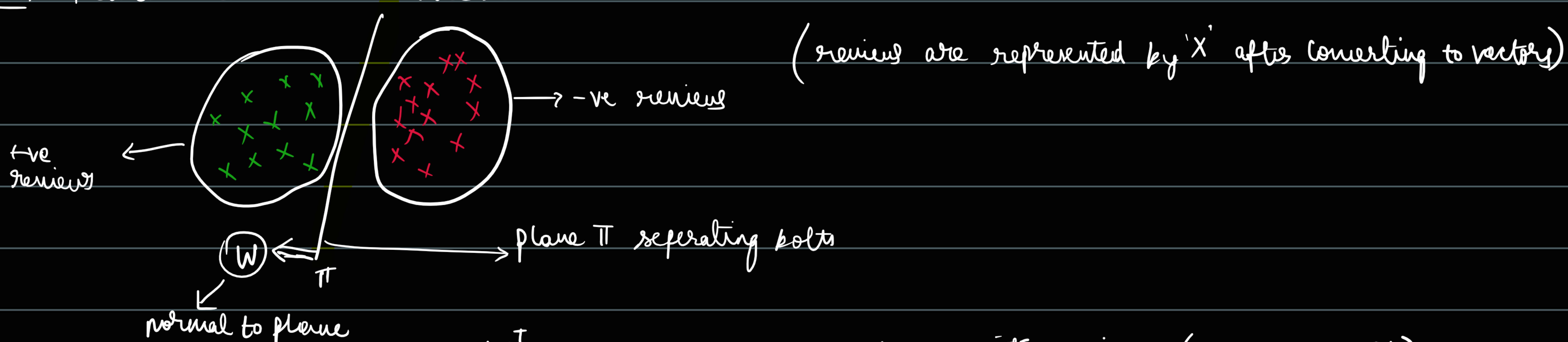
$\rightarrow df.drop_duplicates(subset=[\text{''}, \dots], keep='first')$ (Pandas)

↳ 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 n documents/reviews. [Collection of documents]

For each review

d (d dimensional vector)

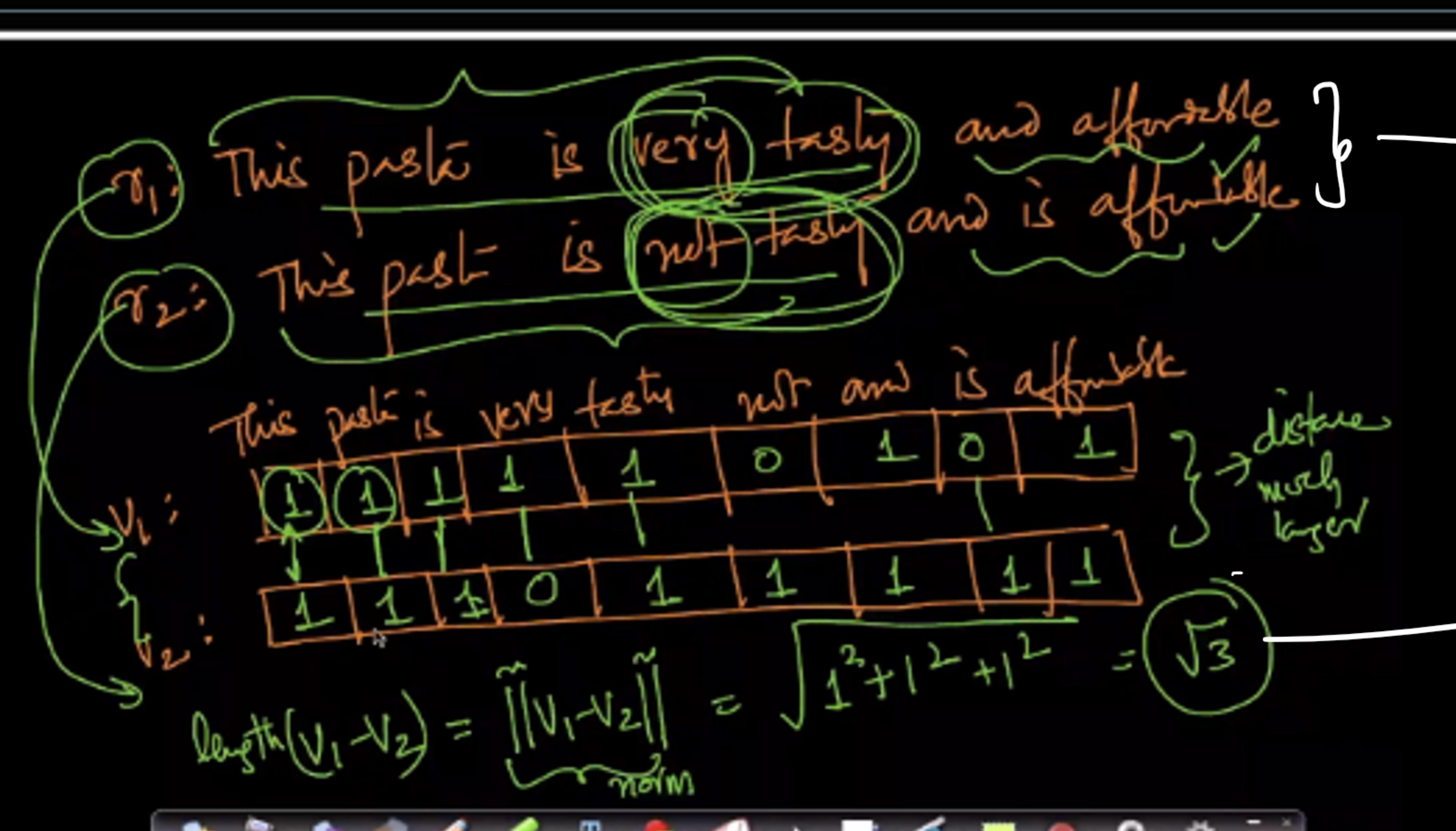
is called a corpus.]

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

$\pi_j :=$ *j*

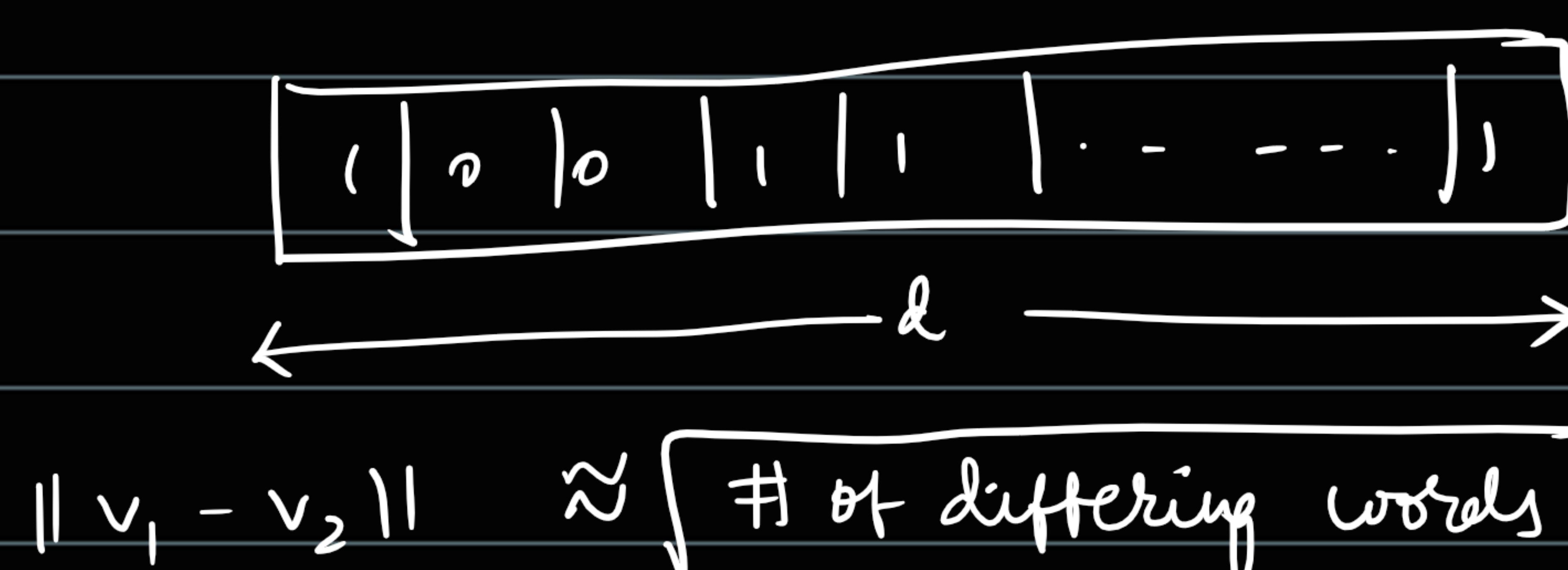
The number in bold represents the number of times
that word occurred in the review.

Each word is a different dimension



Boolean bag of words / binary bag of words :-

Instead of counting number of 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.

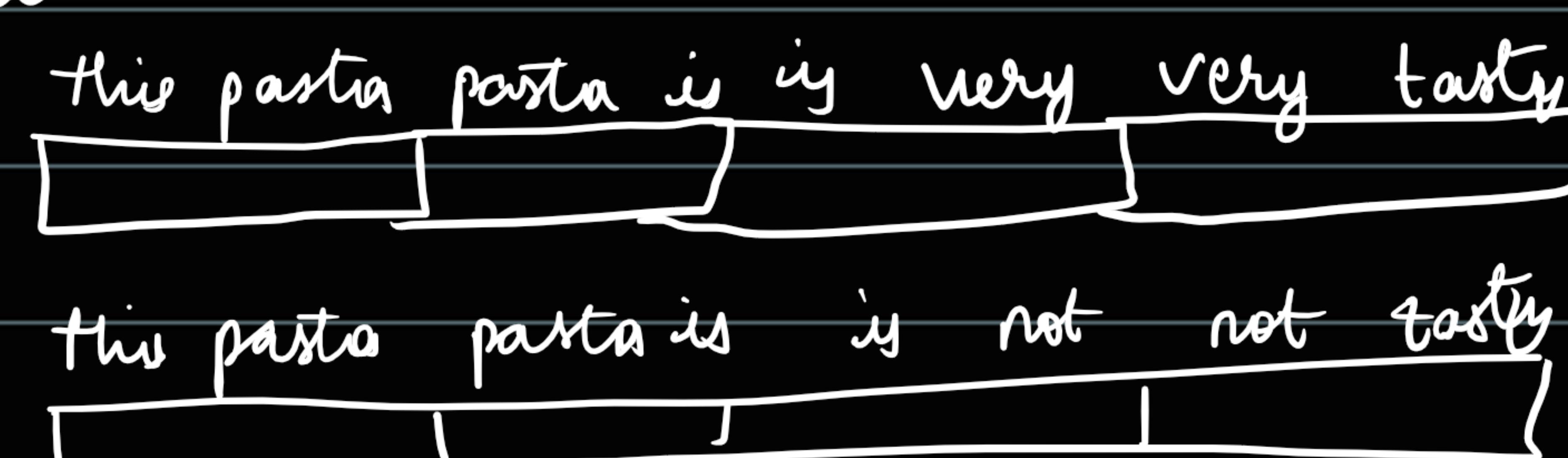
→ Stopwords are not removed in stemming. All words are reduced to base words and the tails are removed.

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 $n(\text{bigrams}) > n(\text{unigrams})$

\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)

$$r_1: w_1 \quad w_3 \quad w_2 \quad w_5 \quad w_1$$

w_1	w_2	w_3	w_4	w_5
2	1	1	0	1

\vdots

$r_2: w_3 \quad w_1 \quad w_4 \quad w_1 \quad w_6 \quad w_2$

\vdots

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

→ TF-IDF was originally meant for information retrieval

→ JF gives how often word occurs. (Probability)

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

IDF :-

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

$$\text{Document Corpus} \left\{ \begin{array}{l} D_c = \{ r_1: w_1, r_2: w_2, \dots, r_j: w_j, \dots, r_N: w_N \} \\ \text{DF}(w_i, D_c) = \log \left(\frac{N}{n_i} \right) \end{array} \right.$$

→ 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 very 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

$\gamma_1 :=$ _____

$\gamma_2 :=$ _____

.

.

.

$\gamma_n :=$ _____

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

→ Drawback:- Still no semantic importance.