

* Between all the attributes given, text (text of the review) is the most important attribute to know if a review is positive or negative

* Data cleaning \Rightarrow 20-30% time of machine learning is applied on data cleaning & pre-processing. Otherwise without cleaning, garbage in $\xrightarrow[\text{algo}]{\text{ML}}$ garbage out

\Rightarrow To remove duplicate entries.

Steps of Data cleaning \Rightarrow

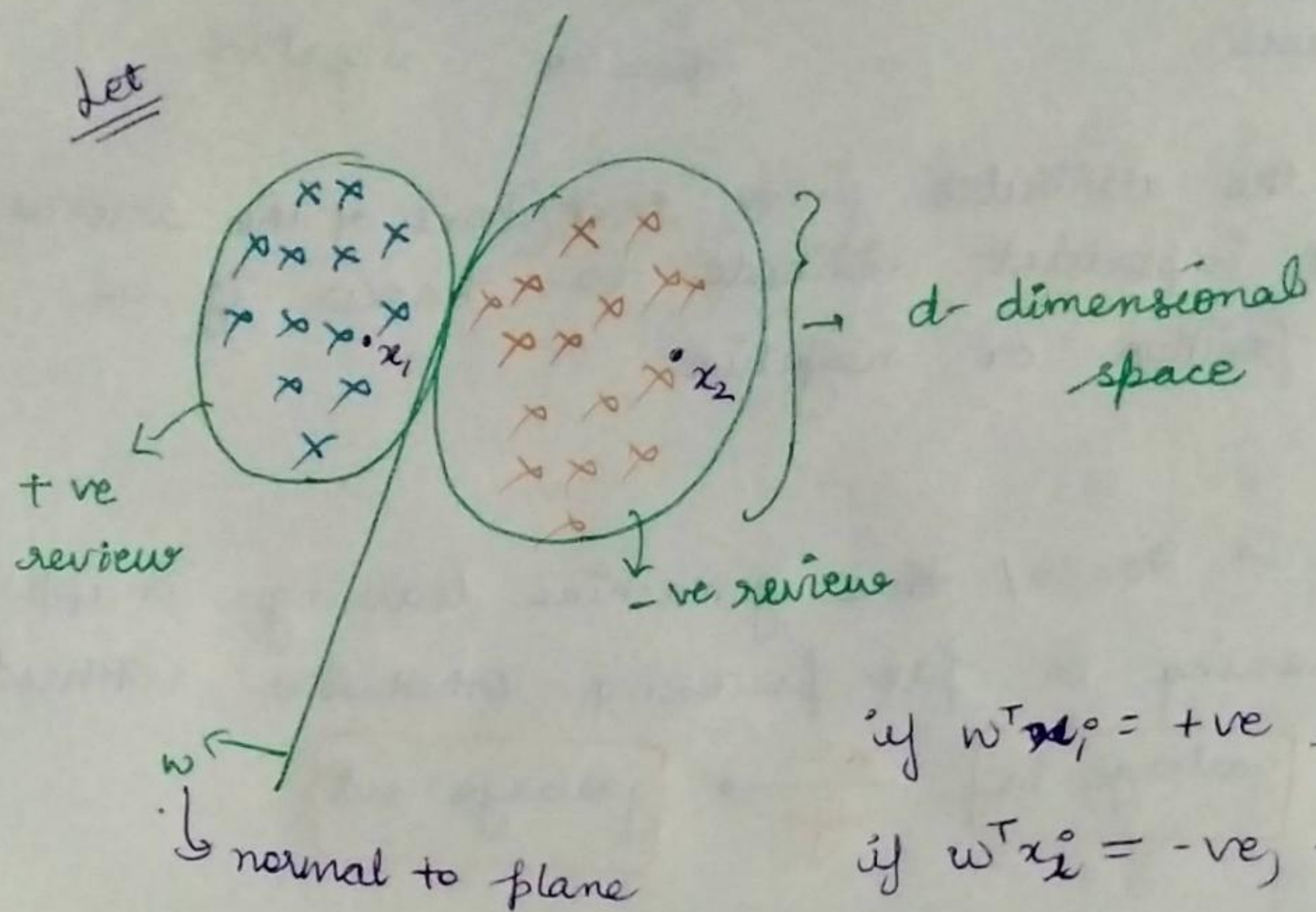
1. Sort the data (sort-values)
2. Drop duplicates values (drop-duplicates), 'keep' \Rightarrow imp
keep = 'first' \rightarrow it will keep the first one & remove other duplicates.
3. Check how much % of data still remains.

\Rightarrow To remove invalid entries.

* As the most important attribute is text here, so as we know by converting data points into vectors, we can use whole thing of linear algebra.

So, Question is How to convert text to vectors?

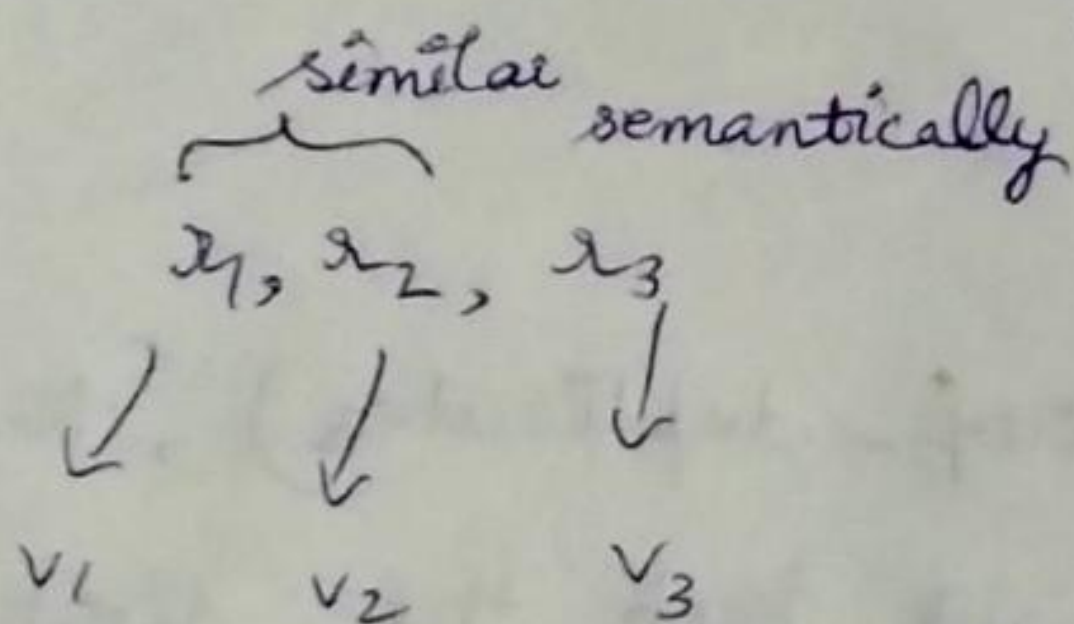
d-dimension representation of a vector



if $w^T x_1 = +ve$, +ve review

if $w^T x_2 = -ve$, -ve review

Step 1 Convert review-text $\xrightarrow{\text{rules}}$ d-dim vector



If English

$\text{Sim}(x_1, x_2) >$

English

$\text{Sim}(x_1, x_3)$

then,

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

* x_1 & x_2 are more similar then v_1 & v_2 must be close.

Now Question is How to convert text to d-dim vector such that similar text must be closer geometrical.

Techniques to convert text to d-dim vector

1. Bag of words (BOW)

Steps \rightarrow (i) Constructing a dictionary :- Set of all the unique words in your reviews.

(ii) for each review r_i make a vector v_i

for ex $\rightarrow r_1$: This pasta is very tasty and affordable

v_1 :

1	2	3	4	5			9		d-1	d
			1		1		1		1	
a	an	the	pasta	This	tasty	is				

no. of times
a word occurs
in r_i .

each word is a different dimension.

$v_i \rightarrow$ sparse vector

Similar text must form closer vectors.

for ex \rightarrow

r_1 : This pasta is very tasty and affordable.

r_2 : This pasta is not tasty and is affordable.

v_1 :

This	pasta	is	very	tasty	not	and	is	affordable
1	1	1	1	1	0	1	0	1

v_2 :

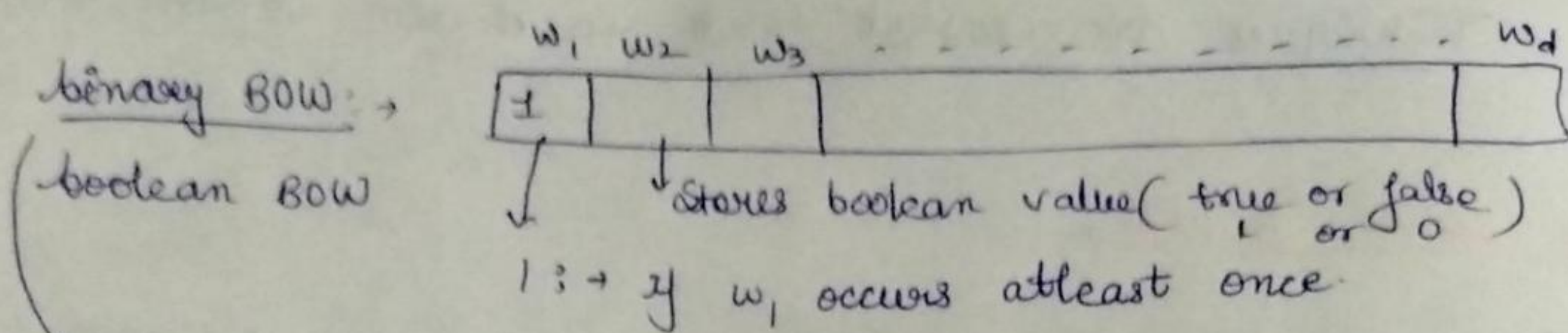
This	pasta	is	very	tasty	not	and	is	affordable
1	1	1	0	1	1	1	1	1

$$\text{length}(v_1 - v_2) = \underbrace{\|v_1 - v_2\|}_{\text{norm}} = \sqrt{1^2 + 1^2 + 1^2} = \sqrt{3}$$

BOW \rightarrow

pasta								
2								

no. of occurrences



$$\|v_1 - v_2\| \approx \# \text{ no. of differing words.}$$

Now if we notice carefully, words like this, is, and etc are not important. So, we will have we can improve BOW.

These words (this, is, ...) are called stop words.

⇒ After removing stop words our BOW will be smaller & more meaningful.

This is called text pre-processing. (✓)

⇒ Convert all words into lowercase so that they are considered to be different.

⇒ Stemming is to convert similar meaning words into a base word.

for ex: \rightarrow tastes, tasty, tasteful \rightarrow taste.

beautiful, beauty \rightarrow beauty.

Porter Stemmer & Snowball stemmers are types of Stemmers.

✓ lemmitization \Rightarrow breaking a sentence into words.

for ex \Rightarrow This pasta is very tasty. This is the best in
New York \rightarrow single word.

✓ Semantic meaning of words

\downarrow {tasty \longleftrightarrow delicious \rightarrow very similar in meaning

BOW do not take care of this.

* Uni-gram / Bi-gram / n-gram \Rightarrow

Unigram \Rightarrow single words \rightarrow in 1-D

$x_1 \Rightarrow$ This pasta is very tasty and affordable.

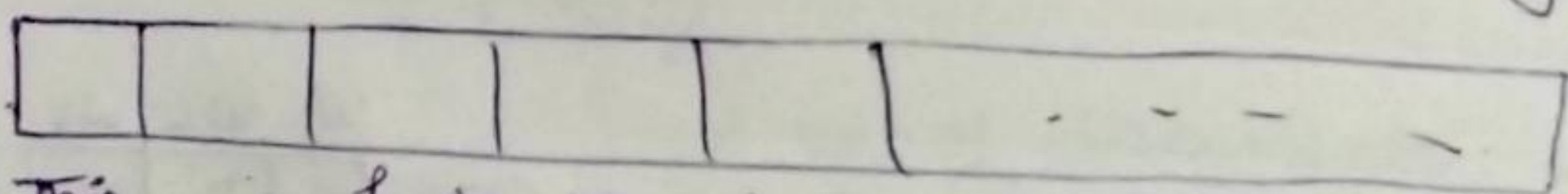
$x_2 \Rightarrow$ This pasta is not tasty and is affordable.

\Rightarrow stop words.

removing stopwords $\rightarrow x_1 \neq x_2$ are exactly same
 \Rightarrow conclude $x_1 \approx x_2$ are very similar

but this conclusion is not right at all.

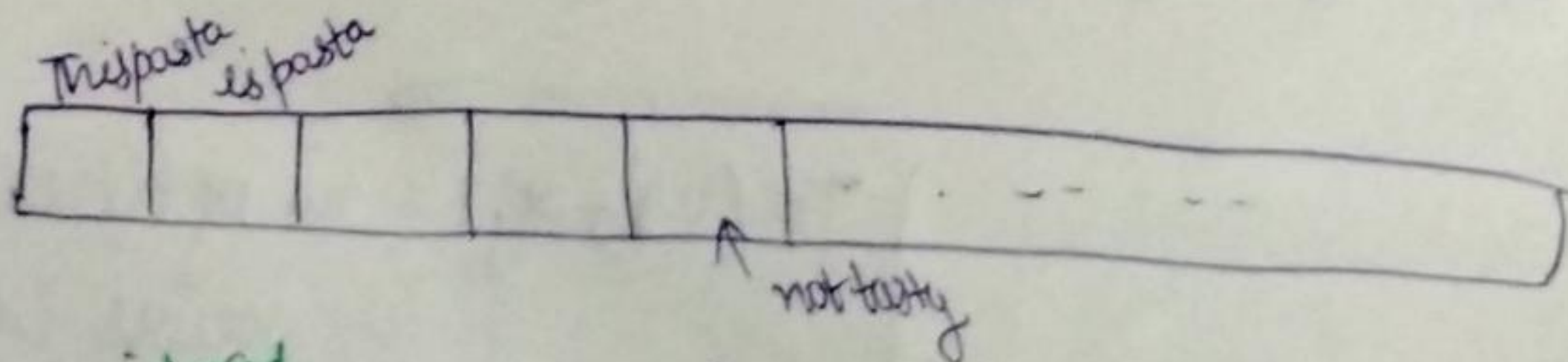
Unigram \Rightarrow



This is pasta very tasty affordable

each word is considered a dimension.

Bigram \Rightarrow



pair of words is considered a dimension = 2-2 words

trigrams \rightarrow 3 consecutive words \rightarrow in a dimension

n-grams \rightarrow $n=1 \rightarrow$ unigram

$n=2 \rightarrow$ bigram

$n=3 \rightarrow$ trigrams.

\Rightarrow Unigram in BOW \rightarrow discards the sequence info
like \rightarrow "very tasty", "not tasty"

\Rightarrow Bigrams / tri-grams \rightarrow retains some of the
sequence info.

no. of bigrams \geq no. of unigrams

! so on

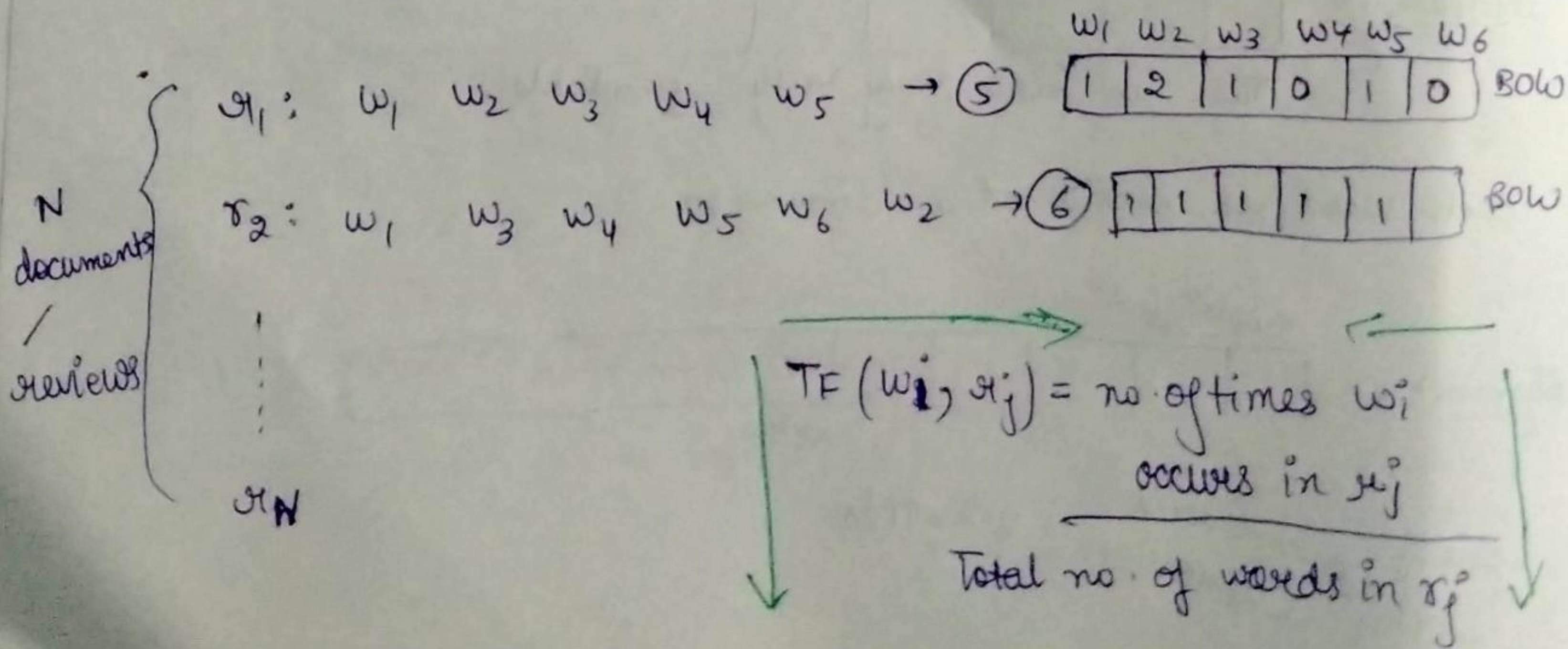
so forth.

n-grams $(n > 1)$ \uparrow \rightarrow dimensionality
'd' \uparrow

* TF-IDF \rightarrow

Term frequency

Inverse document frequency
 \downarrow
(reviews)



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

$$0 \leq TF(w_i, r_j) \leq 1 \rightarrow \text{probability}$$

TF/IDF \rightarrow Information Retrieval techniques (NLP)
 \downarrow
 natural language processing

TF says how often w_i occurs in r_j

$$TF \propto \{\text{more } w_i \text{ occurs in } r_j\}$$

* IDF \Rightarrow

$D = \begin{cases} r_1; \\ r_2; \\ r_3; \\ \vdots \\ r_N; \end{cases}$

\downarrow
 Data corpus
 no. of docs/reviews

$$IDF(w_i, D_c)$$

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

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

$\xrightarrow{\text{no. of documents}}$
 $\xrightarrow{\text{no. of documents which contain } w_i}$

$$n_i \leq N \rightarrow \text{Always}$$

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

$$\log \left(\frac{N}{n_i} \right) \geq 0 \rightarrow \textcircled{1} \textcircled{IDF \geq 0} \approx$$

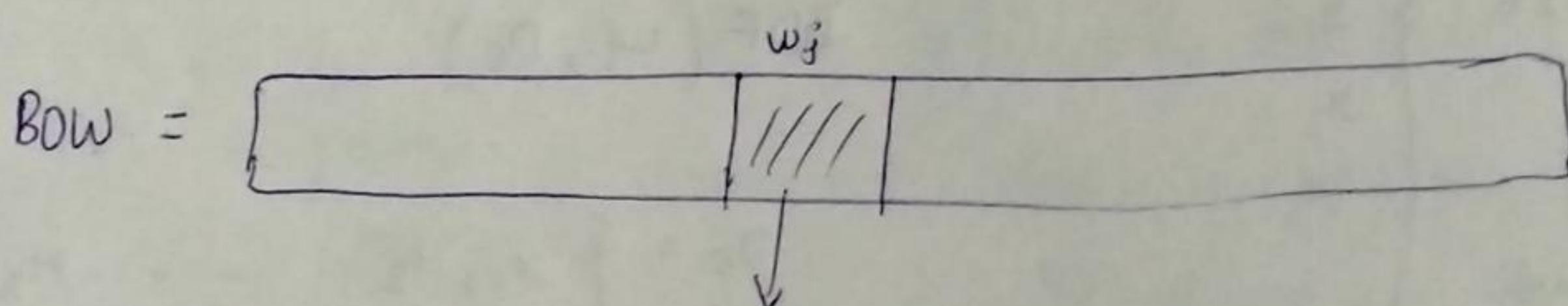
② if $n_i \uparrow$, $\frac{N}{n_i} \downarrow$, $\log\left(\frac{N}{n_i}\right) \downarrow$

as log is a monotonic functions

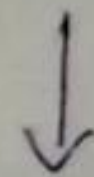
$f(x)$ behaviour is same as x .

So, if a word w_i is more frequent in a document corpus then its IDF will be low.

Now, using TF & IDF we will fill TF * IDF in BOW table.



$$TF(w_j, r_i) * IDF(w_j, D_c)$$



(w_j is frequent in r_i) (w_j is rare in D_c)

TF - IDF \rightarrow

↳ more importance to rare words in D_c .

↳ more importance if a word is frequent in a document / review.

Still - Semantic meaning problem is not solved.

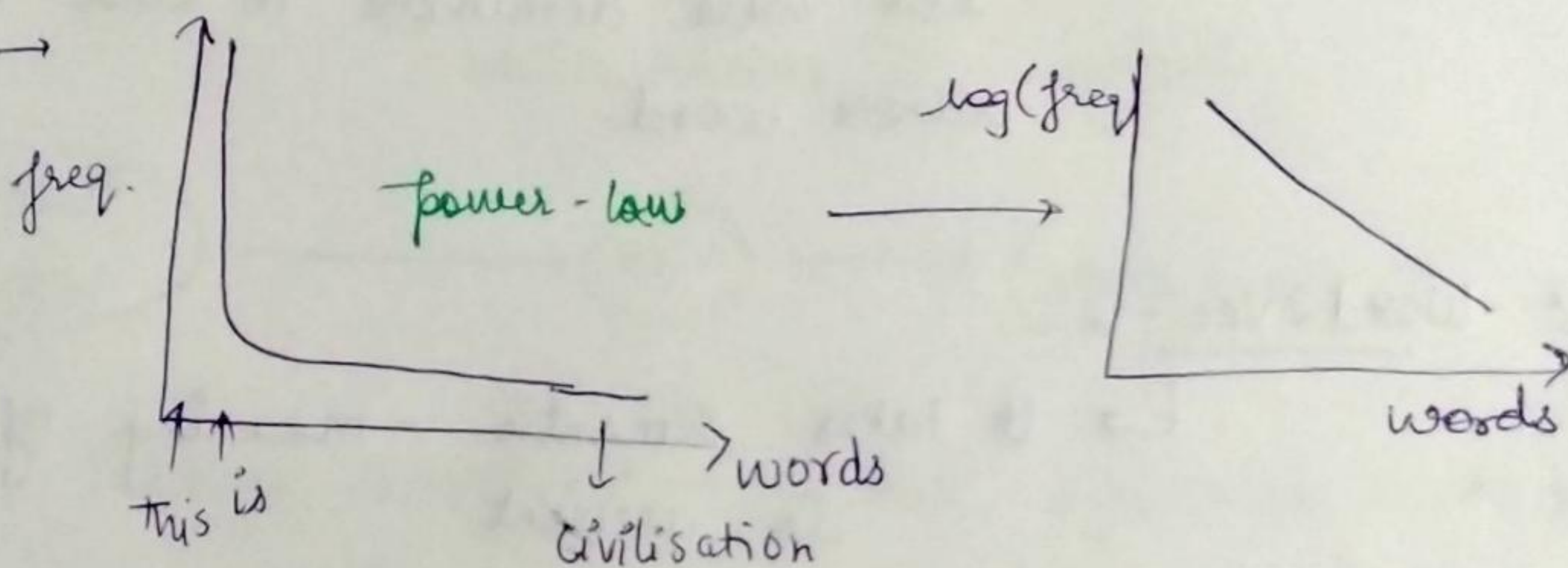
(tasty \rightarrow delicious)

(cheap \rightarrow affordable)

These are similar words but in TF-IDF representation also tasty will be 1 dimension & delicious will be other dimension.

Q: Why do we use $\log\left(\frac{N}{n_i}\right)$ for IDF?

\Rightarrow Zipf's law \rightarrow



\rightarrow decreasing order of frequency

So, Zipf's law says frequency of word & words graph has power law distribution.

and whenever there is a power law distribution, taking log is a good idea.

as managing straight line graph is easier.

→ Another practical reason is →

	$\frac{N}{n_i}$		$\log\left(\frac{N}{n_i}\right)$
The	→ 1	} → diff is 1000 here	0
is	→ 1		0
⋮			⋮
civilisation	→ 1000		≈ 6.9

} difference is 7 here

↓
1 in 1000

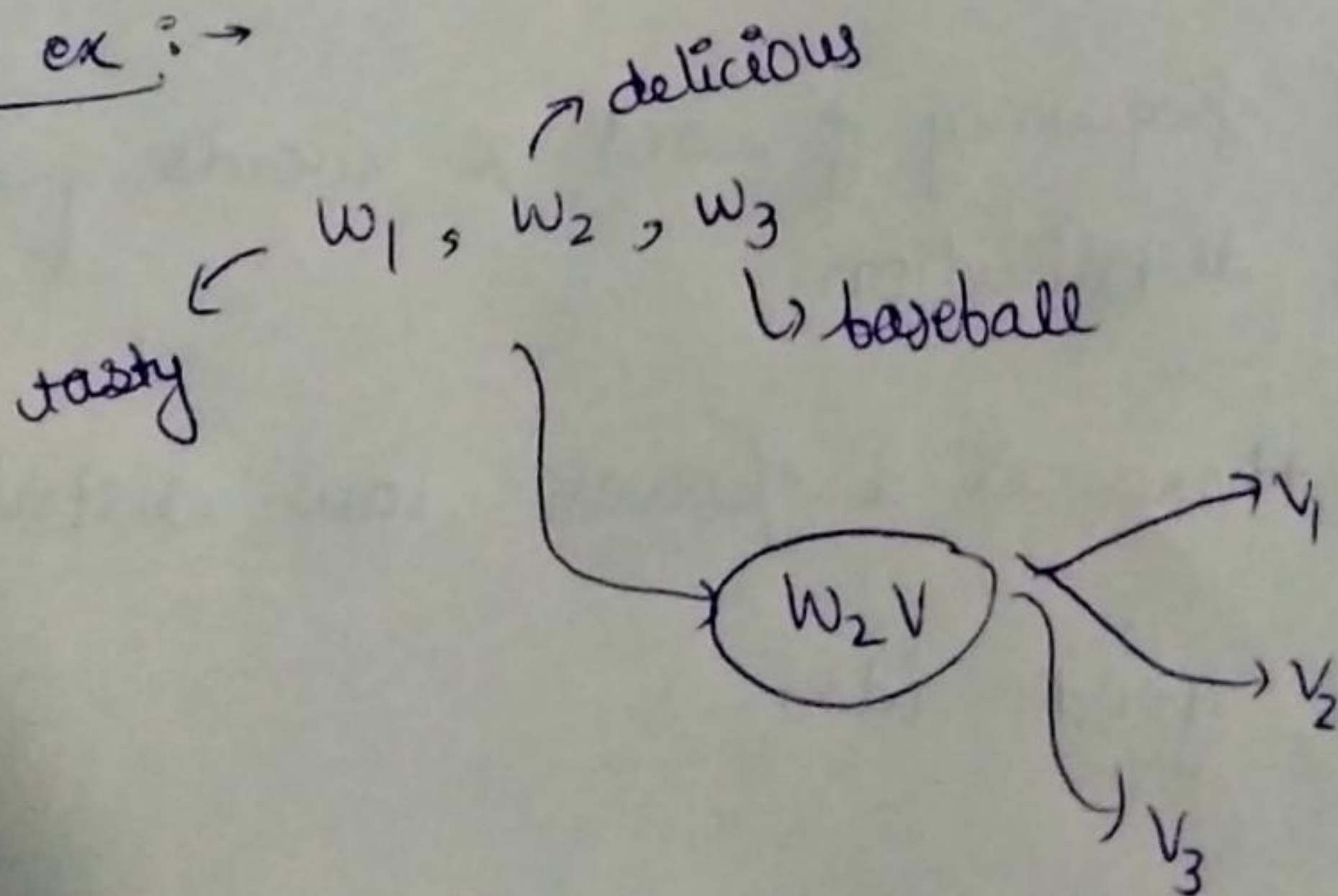
So, on calculation $TF * IDF$,
IDF will dominate in case of
rarer word.

★ Word2Vec : →

↳ It takes semantic - meaning of words
in account.

which others were not doing.

For ex : →

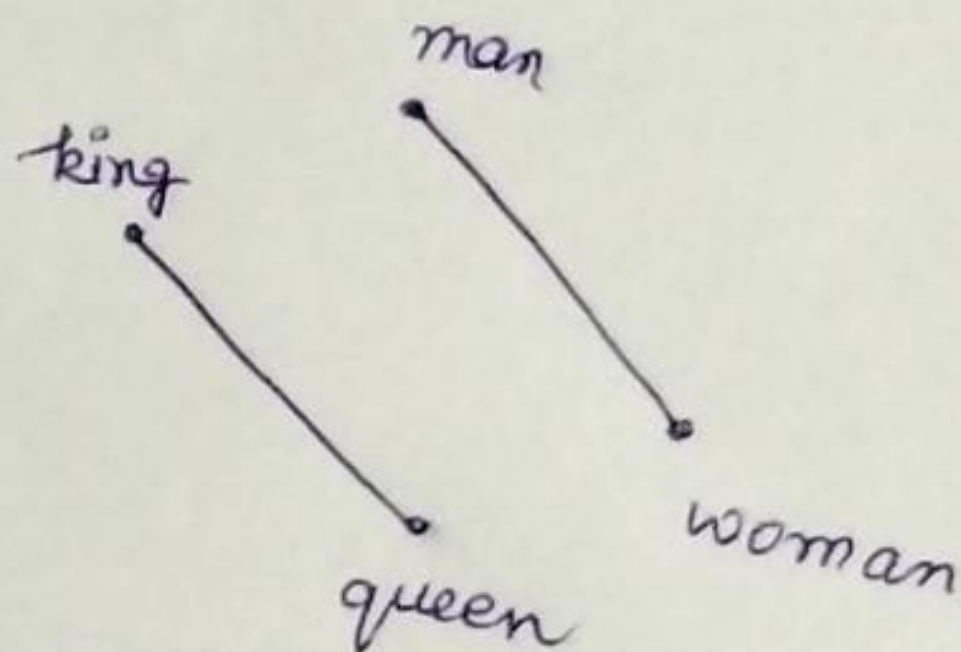


Things W2V tries to achieve \Rightarrow

① $\{ \text{if } w_1 \neq w_2 \text{ are semantically similar then } v_1 \& v_2 \text{ close} \}$

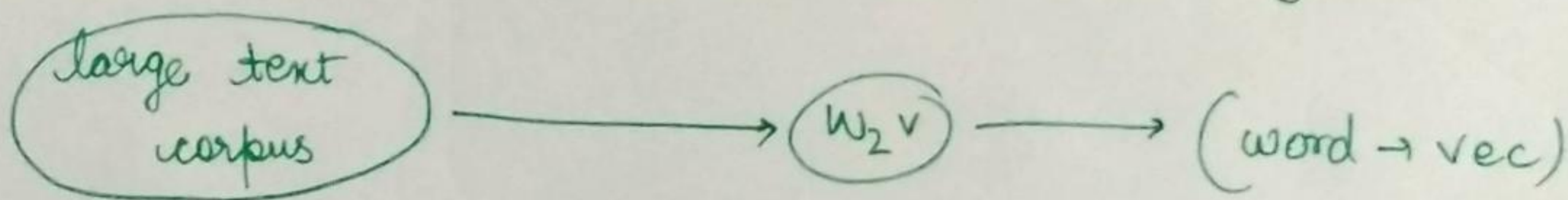
② It satisfy relationships

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



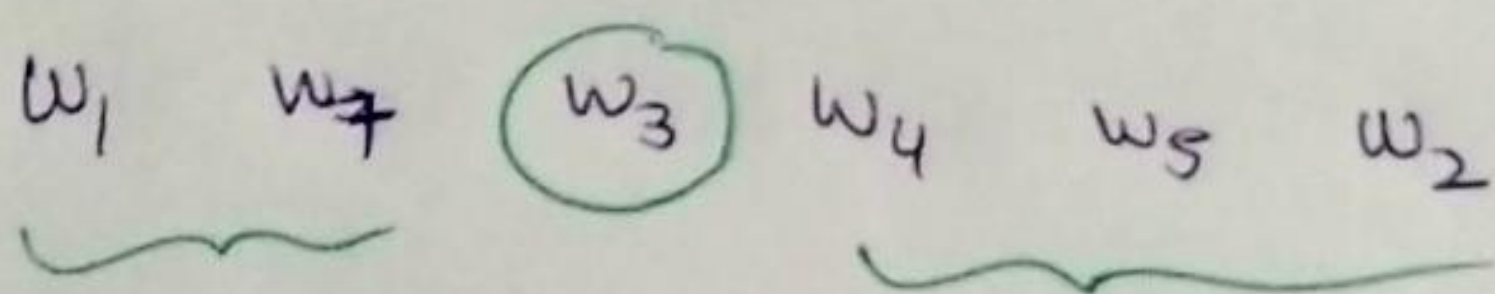
\Rightarrow ~~imp~~ W2V \rightarrow learns relationships automatically from raw-text

\downarrow which makes it magical



larger dimensions \rightarrow more info such vectors

Core of word2vec \Rightarrow



neighborhood of w_3

$$w2v(w_3) = \left\{ \begin{array}{l} \text{if } Neigh(w_i) \approx N(w_j) \\ v_i \approx v_j \end{array} \right\}$$