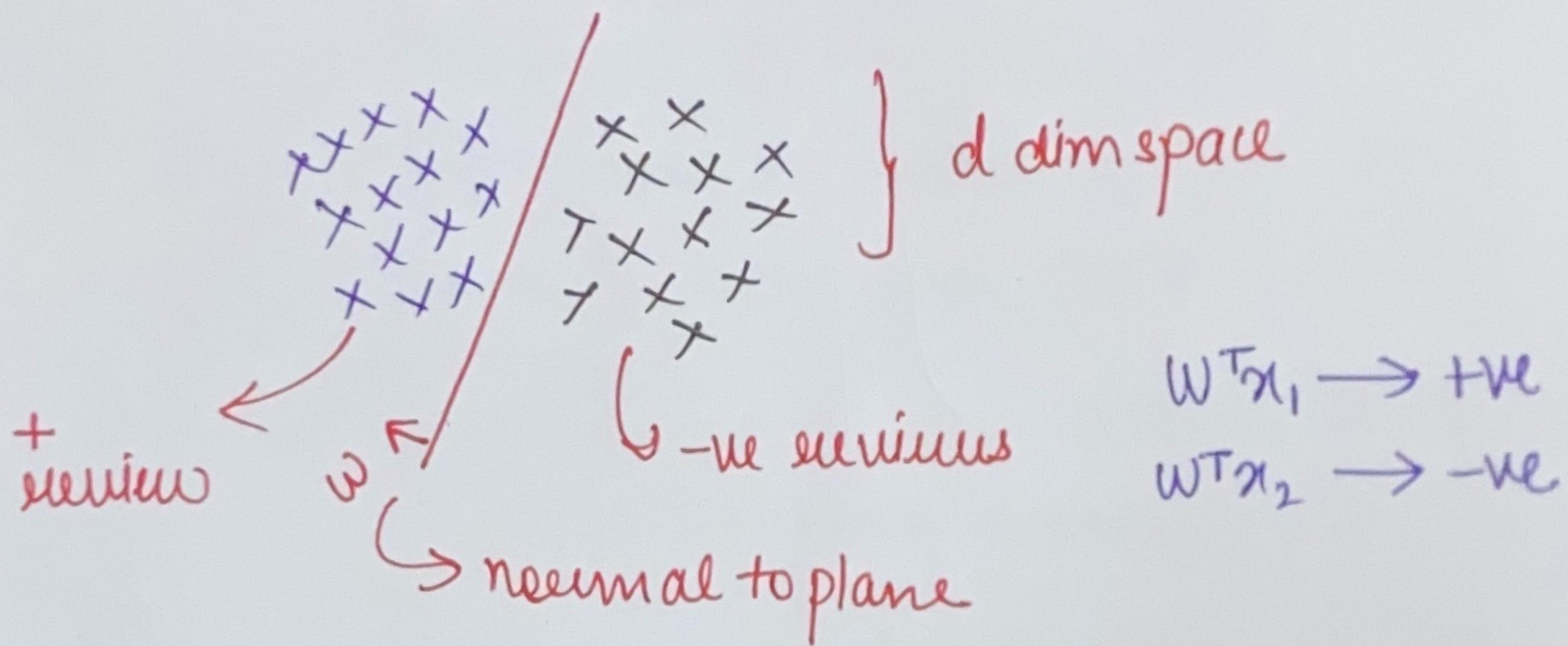


Tent to Vector

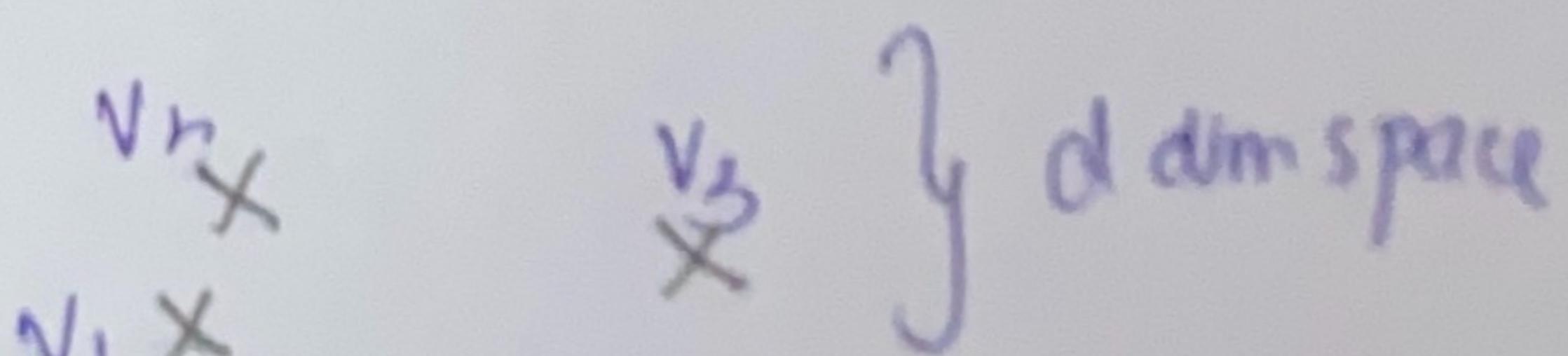
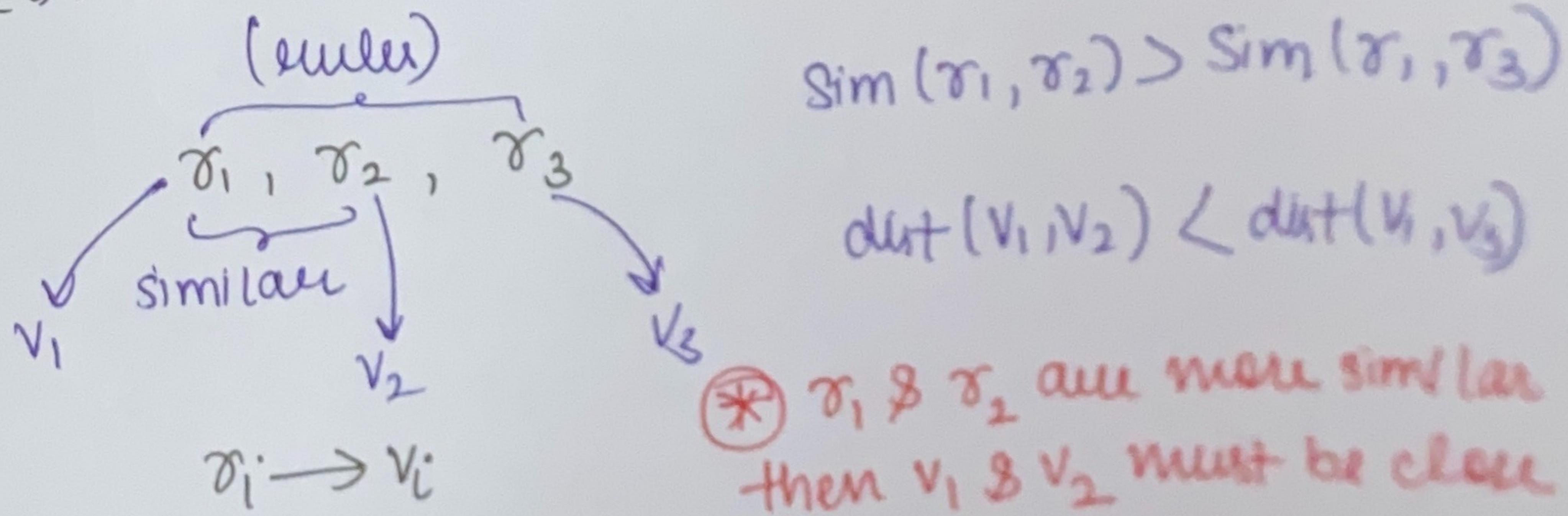
(Q) how do you convert tent into numerical vectors
 Review tent \rightarrow d dim vector



if $w^T x_i > 0$ then x_i is +ve
 else x_i is -ve

review tent \rightarrow d dim \rightarrow finding a plane to separate

(Q) tent \rightarrow d-dim vector



$$\text{Eng sim}(v_1, v_2) > \text{Eng sim}(v_1, v_3)$$

\Downarrow

$$\text{length}(v_1 - v_2) < \text{length}(v_1 - v_3)$$

★ Similar tent must be closer

Bag of Words (BoW)

- Corpus
- τ_1 : This pasta is very tasty and affordable
 - τ_2 : This pasta is not tasty and is affordable
 - τ_3 : This pasta is delicious and cheap
 - τ_4 : Pasta is tasty and tastes good.

BoW ① Constructing a dictionary: Set of all **unique words** in your reviews.

(**d unique words**) \hookrightarrow [This, pasta, is, very---]

vector v_1 \leftarrow τ_1 : This pasta is very tasty and affordable

no of times word occurs v_1 :

1	2	3	4	-	7	-	d
0	0	1			1	1	1

 a an The pasta This is tasty

* Each word is different dimension

$v_1 \rightarrow$ sparse \rightarrow most of elements are zero
 $\begin{matrix} 1 & 2 & 3 & - & \dots & d \\ 0 & 0 & 1 & & & 0 & 0 & 1 \end{matrix}$ due to many

BoW: text \rightarrow vector

Similar text must result in closer vector

This pasta is very tasty and is affordable not
 v_1 :

1	1	1	1	1	1	1	0	1	0	0
---	---	---	---	---	---	---	---	---	---	---

v_2 :

1	1	1	1	1	0	1	1	1	1	1
---	---	---	---	---	---	---	---	---	---	---

$$\text{length}(v_1 - v_2) = \sqrt{\sum_{i=1}^d (v_{1,i} - v_{2,i})^2} = \sqrt{1^2 + 1^2 + 1^2 + 1^2 + 1^2 + 0^2 + 1^2 + 1^2 + 1^2 + 1^2} = \sqrt{10}$$

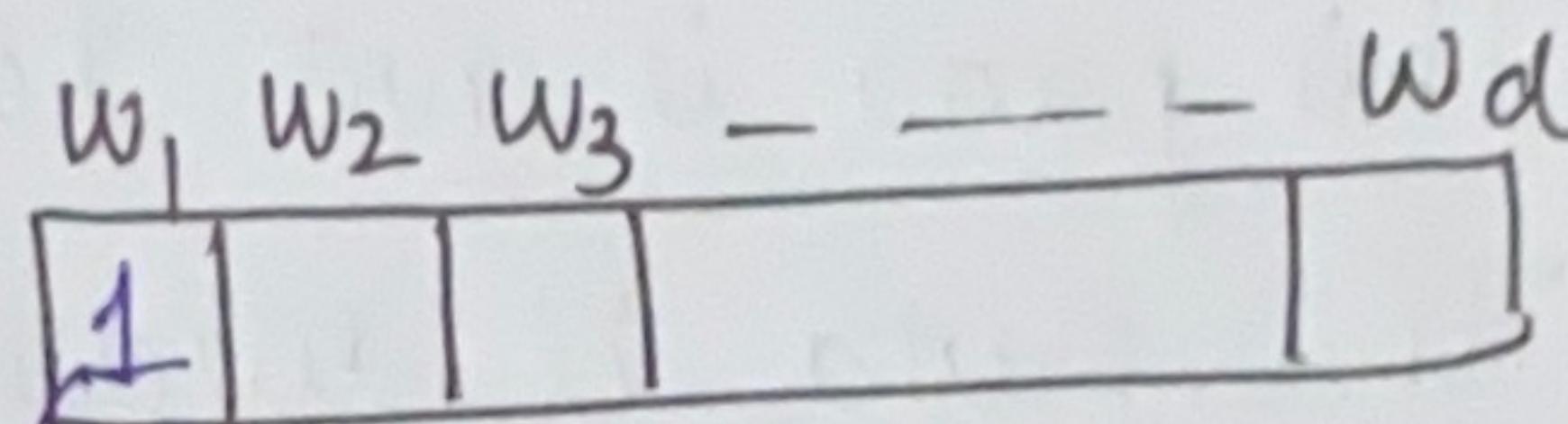
Norm

Bow:

2 |

occurrences

binary Bow :-
Boolean Bow



1: If w_1 occurs atleast once
0: otherwise

Binary Bow

v_1 : _____
 v_2 : _____

$\|v_1 - v_2\|$ ~~is~~ $\sqrt{\text{no of diff words}}$

Stop Words

- Bow
- σ_1 : This pasta is very tasty and affordable
 - σ_2 : This pasta is not tasty and is affordable
 - σ_3 : This pasta is very delicious and cheap
 - σ_4 : Pasta is tasteful and pasta tastes good.

Bow: vector: Smaller: meaningful

① \rightarrow removing Stop words \rightarrow (Text Pre-processing)

② \rightarrow Lower Case \rightarrow Small letters

③ \rightarrow Stemming \rightarrow Porter Stemmer
Snowball Stemmer

tastes, tasty, tasteful

\rightarrow taste

beautiful, beauty \rightarrow beauty

④ Lemmatization (breaking sentence into words)

This pasta is very tasty, This is but in New York.

Bow:

taste	delicious

Semantic meaning of words $\left\{ \begin{array}{l} \text{tasty} \leftrightarrow \text{delicious} \\ \dots \end{array} \right.$

↪ Word2Vec

↪ Bow + Text Processing

↪ text \rightarrow d-dim Vector

Unigram / Bi-gram / n-gram

σ_1 : This pasta is very tasty and affordable

σ_2 : This pasta is not tasty and is affordable

removing stop words; v_1 & v_2 are exactly same
 \Rightarrow conclude σ_1 & σ_2 are very similar

Unigram:

Th	is	pasta	very	tasty	-----

↪ each word considered as a dim
bigram

↪ pair of words
trigram

↪ 3 consecutive words \rightarrow multi dimension

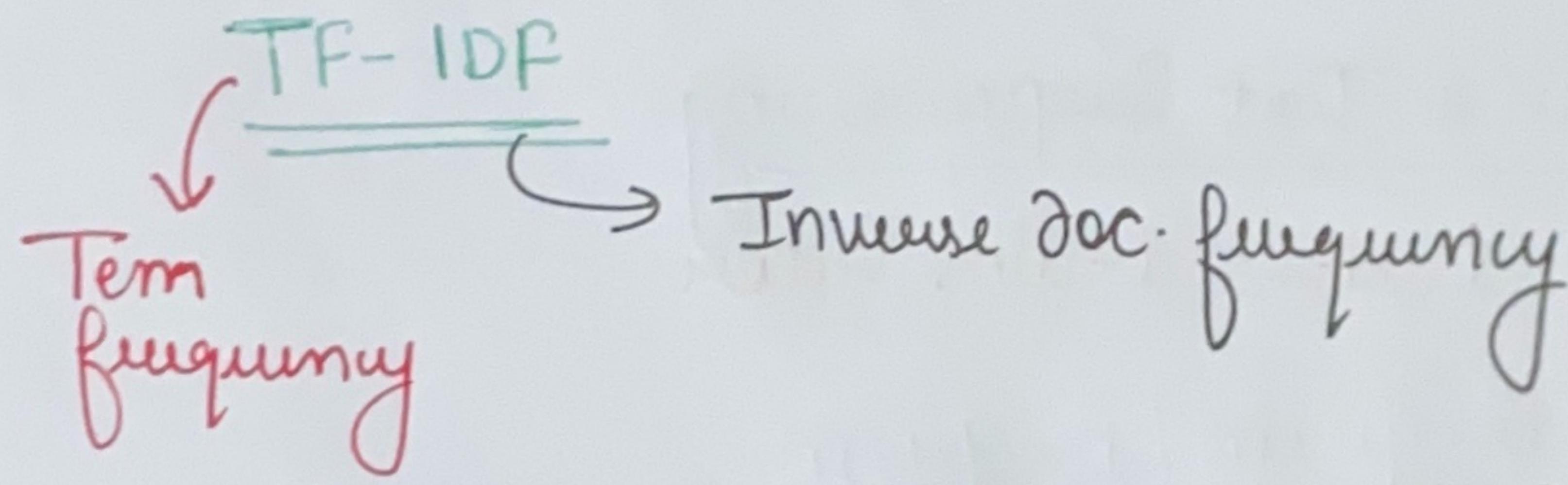
Unigram

Bow: discards sequence information

bigram
trigram }
n-gram } retain some of the sequence info

$$\# \text{ bigram} \geq \# \text{ unigrams}$$

n-grams \rightarrow dimensionality
($n \geq 1$) 'd' increases



N
docs/
reviews

	w_1, w_2, w_3, w_4, w_5	⑤	w_1, w_2, w_3, w_4, w_5
$d_1:$	w ₁ w ₂ w ₃ w ₄ w ₅	—	1 2 1 0 1
$d_2:$	w ₁ w ₃ w ₄ w ₅ w ₆ w ₂	⑥	1 1 1 1 1 1
$d_3:$:
$d_N:$:

$TF(w_i, d_j) = \frac{\# \text{ of times } w_i \text{ occurs in } d_j}{\text{total no of words in } d_j}$

$$TF(w_2, d_1) = 2/5$$

$0 \leq TF(w_i, d_j) \leq 1 \leftarrow \text{probability}$

TF-IDF \rightarrow Information Retrieval (IR)

IDF: Inverse Document freq

$$\mathcal{D} = \left\{ \begin{array}{ll} \tau_1 - w_1 & \text{IDF}(w_i, \mathcal{D}) \\ \tau_2 - & \\ \tau_3 - w_1 - & \\ \vdots & \\ \tau_N - w_1 - & \end{array} \right.$$

$$\mathcal{D}_c = \{\tau_1, \tau_2, \dots, \tau_N\}$$

$$\text{IDF}(w_i, \mathcal{D}_c) = \log \left(\frac{N}{n_i} \right)^{\text{\# of docs}}$$

$$n_i \leq N \rightarrow \frac{N}{n_i} \geq 1 \quad \begin{matrix} \text{\# of docs} \\ \text{containing } w_i \end{matrix}$$

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

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

$$\left[\begin{matrix} \text{if } n_i \text{ increases;} \\ \frac{N}{n_i} \downarrow \end{matrix} \right] \textcircled{2} \quad \log \left(\frac{N}{n_i} \right) \downarrow$$

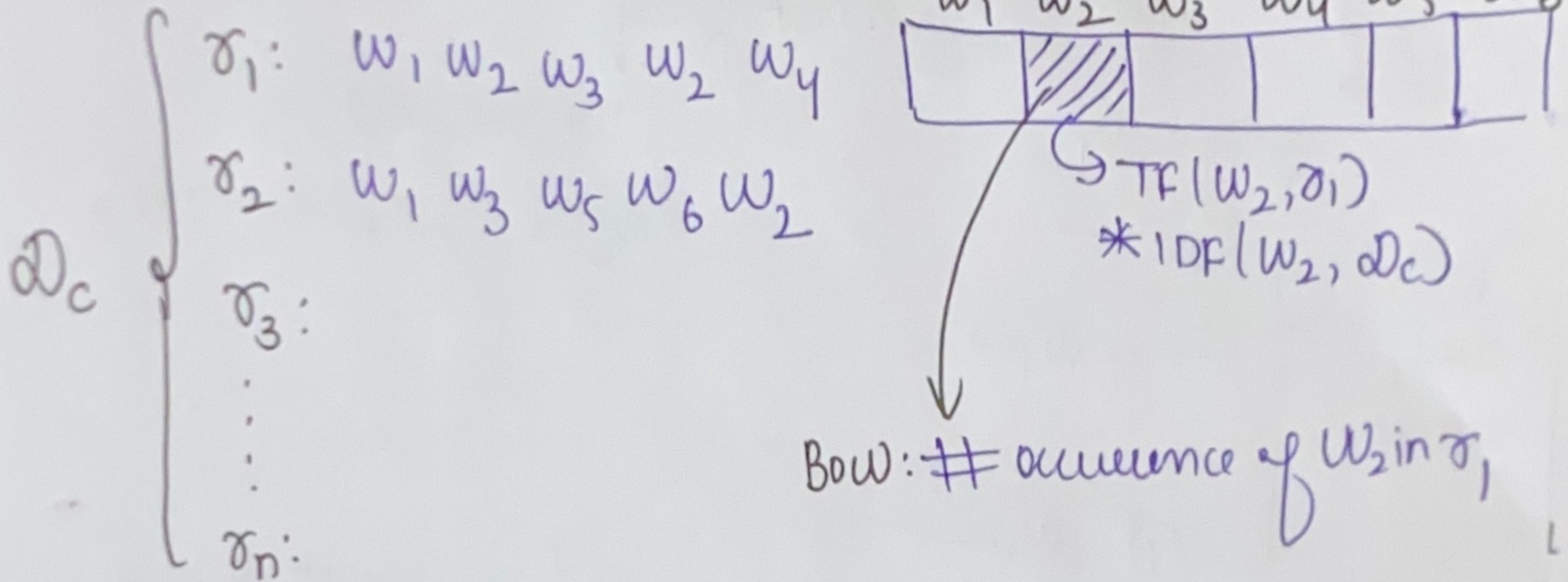
$$\boxed{\begin{matrix} \text{IDF} \downarrow n_i \uparrow \\ n_i \downarrow \text{IDF} \uparrow \end{matrix}}$$

w_i is more frequent

\hookrightarrow IDF will be low

w_i is rare word

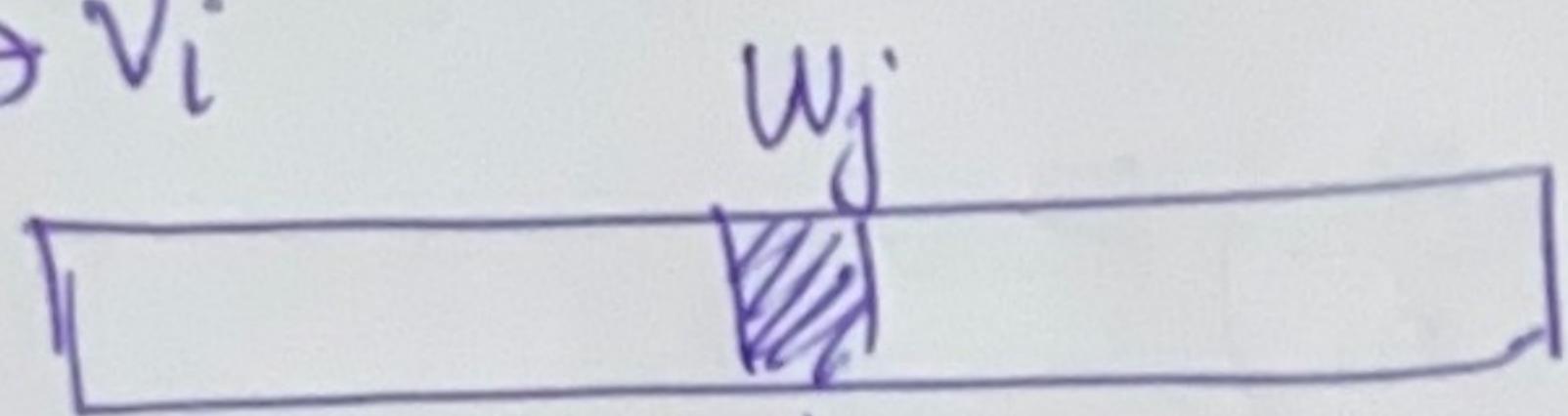
\hookrightarrow high IDF



$$\text{TF}(w_i, \sigma_j) * \text{IDF}(w_i, \mathcal{D}_c)$$

$$(\sigma_1, \sigma_2, \dots, \sigma_n) = \mathcal{D}_c$$

$\sigma_i \rightarrow v_i$



$$\text{TF}(w_j, \sigma_i) * \text{IDF}(w_j, \mathcal{D}_c)$$

w_j frequent w_j is rare in \mathcal{D}_c

TF-IDF

- ↳ more importance to rare words in \mathcal{D}_c
- ↳ more importance if a word is frequent in doc.

limitation

ignores semantic meaning

(tasty \leftrightarrow delicious)

(cheap \leftrightarrow affordable)

...

v

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

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

$\frac{N}{n_i} \rightarrow \# \text{docs}$
 $\downarrow \# \text{docs which contain } w_i$

Zipf's law

freq

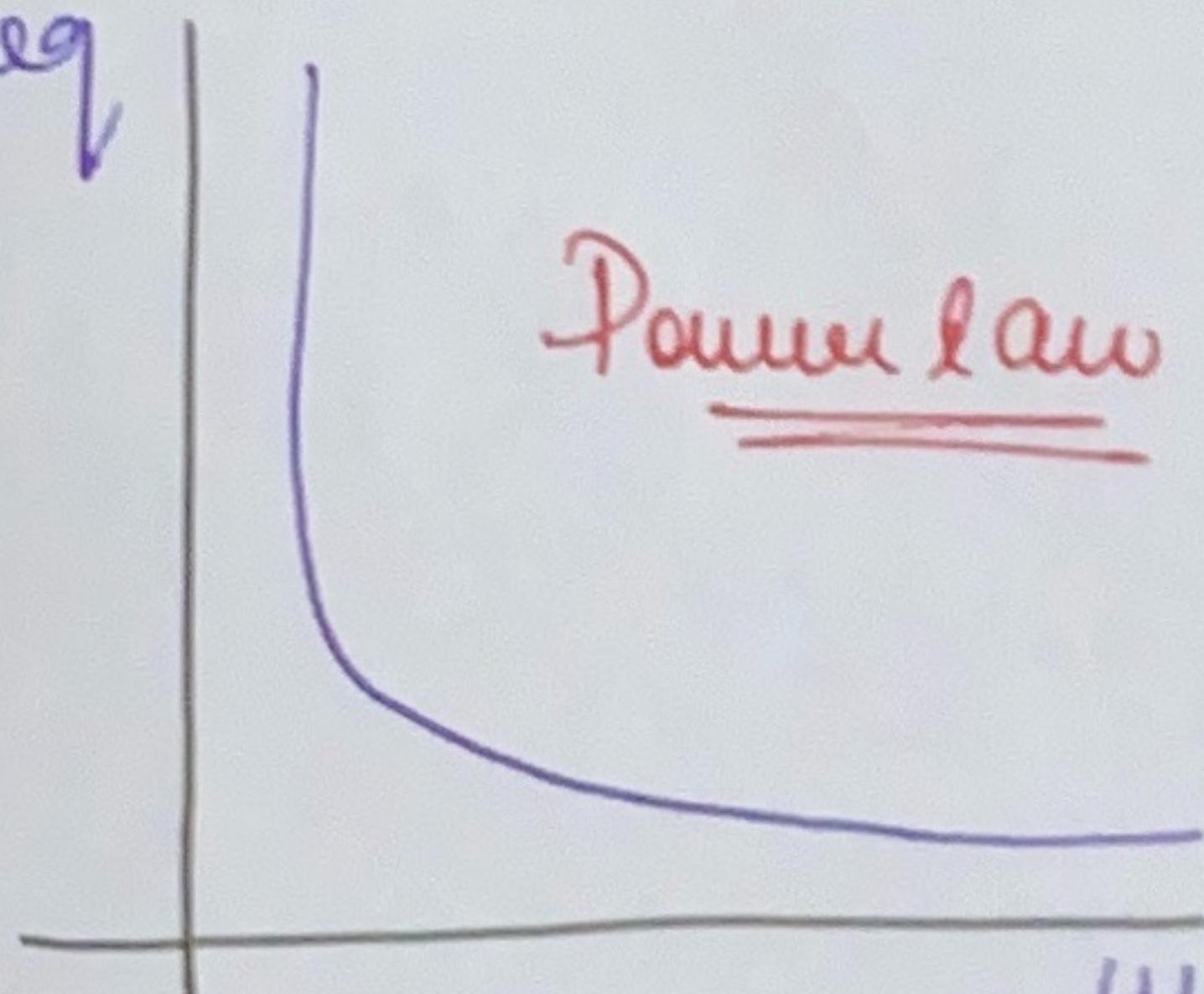
Power Law
dis

$\times \checkmark$ power law

\hookrightarrow log normal

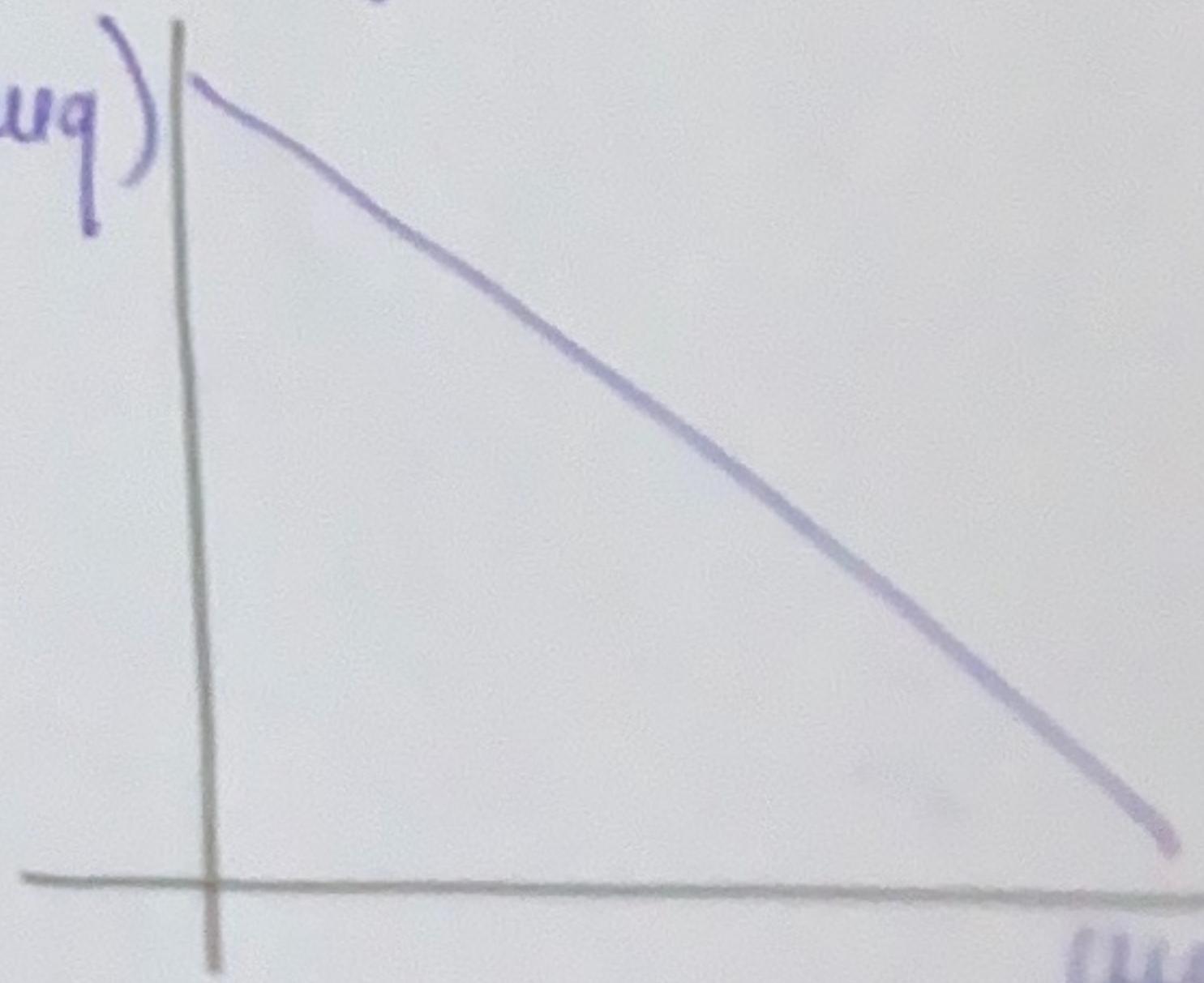
Gaussian (box cox transform) $\rightarrow \log(n)$

freq



$\log(freq)$

\rightarrow



words

$\frac{N}{n_i}$

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

The
is

①
1

1000

②
0

+

TF * IDC

civilization
(1 in 1000)

1000

$\approx 6^{-4}$

due to log domination of IDF dimensions

natural property
of eng

← histogram
of word
occurrences

The is
civilisation
distribution of freq

↑↑

↑ words

↓ civilization