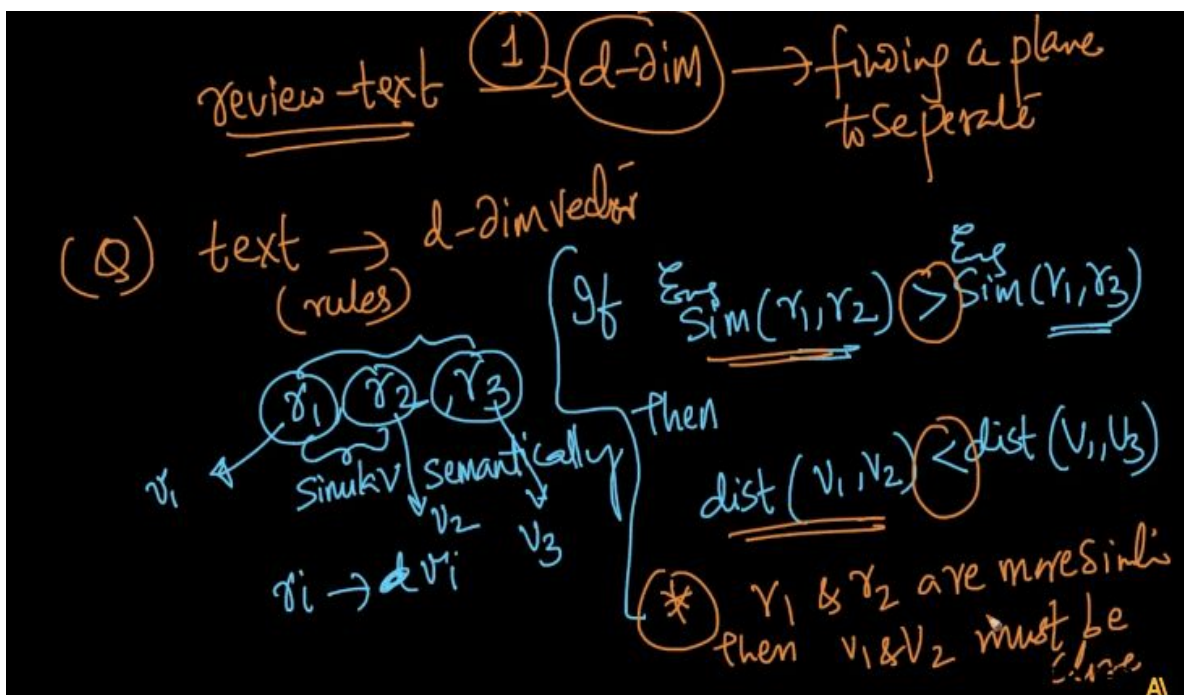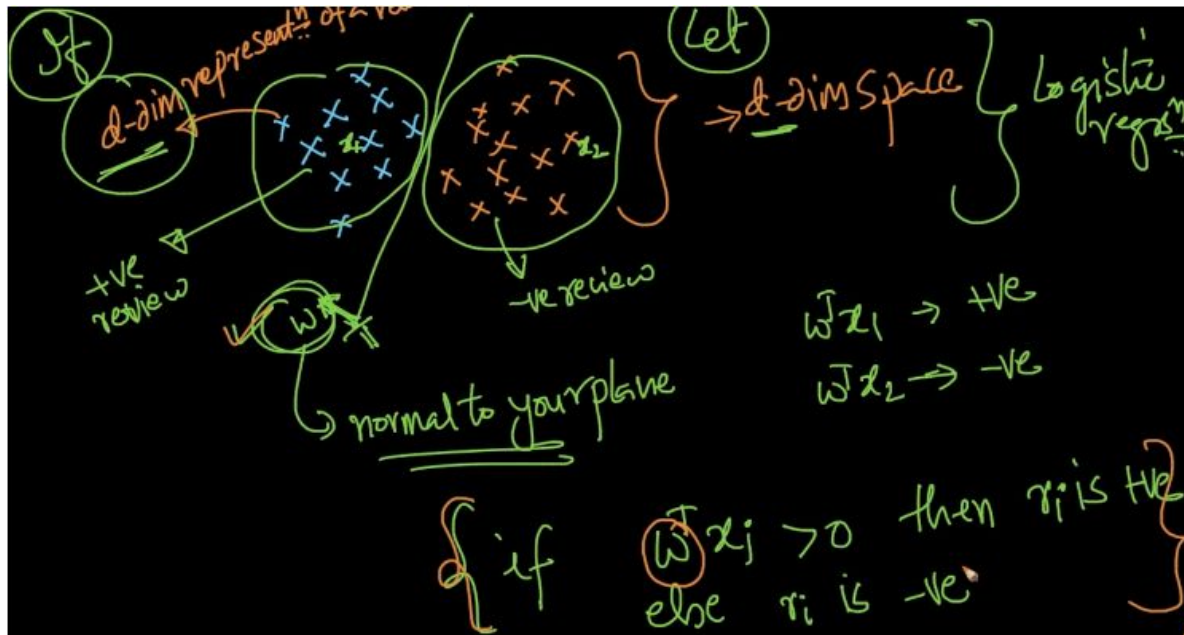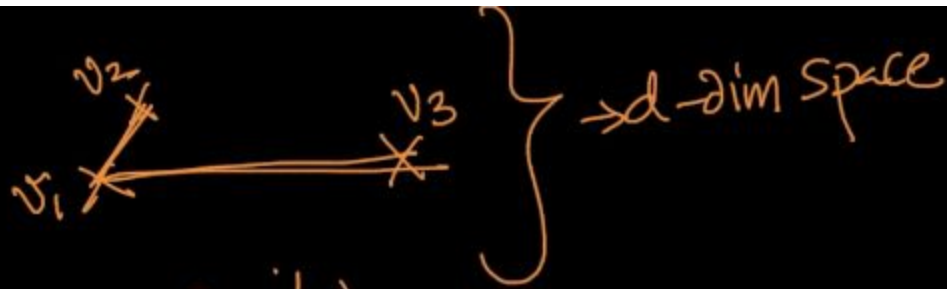# Converting Text (Words and Sentences) into numerical Vector ?

If we can convert any text to a vector, we can leverage the power of linear algebra.

For our case study of Amazon fine food reviews,

If we convert Review text into d dimensional vector

$v_2$

$v_3$

$v_1$

$\Rightarrow$ d-dim space

Similar

$EngSim(v_1, v_2) > EngSim(v_1, v_3)$

$||$

$length(v_1 - v_2) < lgt(v_1 - v_3)$

closer

find $\{ \text{text} \rightarrow \text{d-dim vector} \}$

s.t. similar text must be closer (geometrically)

BoW
tf-idf
w2V
avg w2V
tfidf w2V

## Bag of Words (BoW)

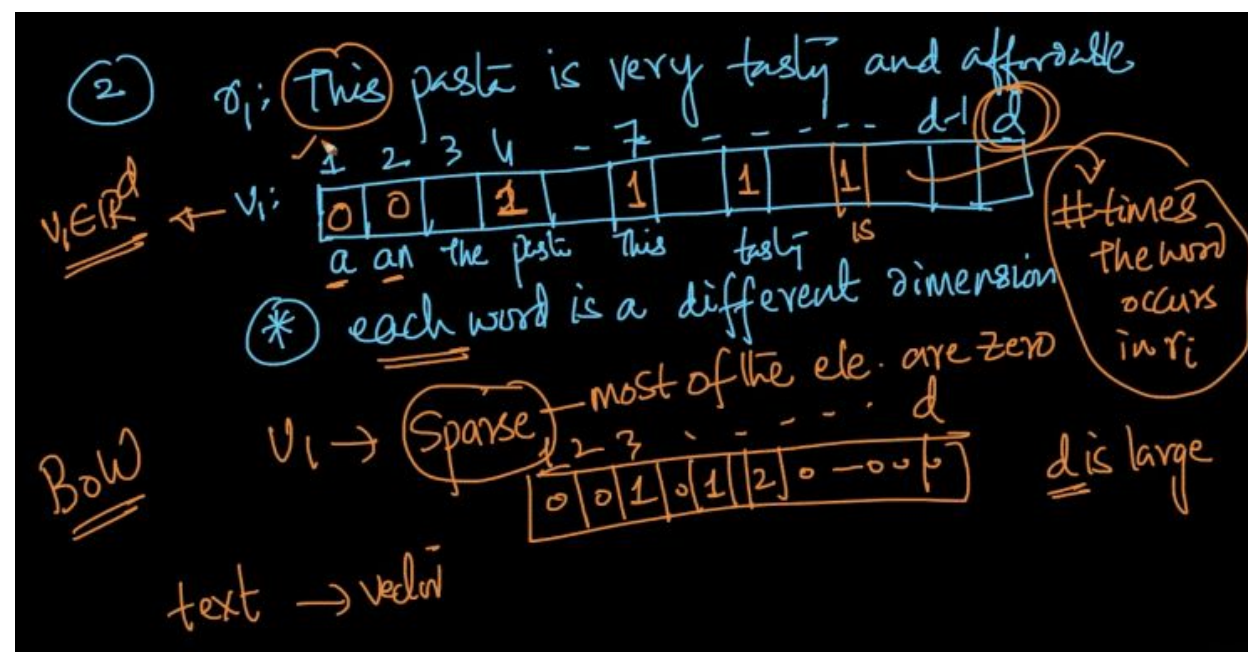NLP — Bag of Words (BoW)    Text → vec    (500K+)

(Toy)

(document)

$r_1$: This pasta is very tasty and affordable.

$r_2$: This pasta is not tasty and is affordable.

$r_3$: This pasta is delicious and cheap.    ($r_1, r_2, \ldots r_n$)

$r_4$: Pasta is tasty and pasta tastes good.

(corpus)

BoW

① constructing a (dictionary) — Set of all the unique words in your reviews

(d-unique words) → {This, pasta, is, very, . . . . } ← S

---

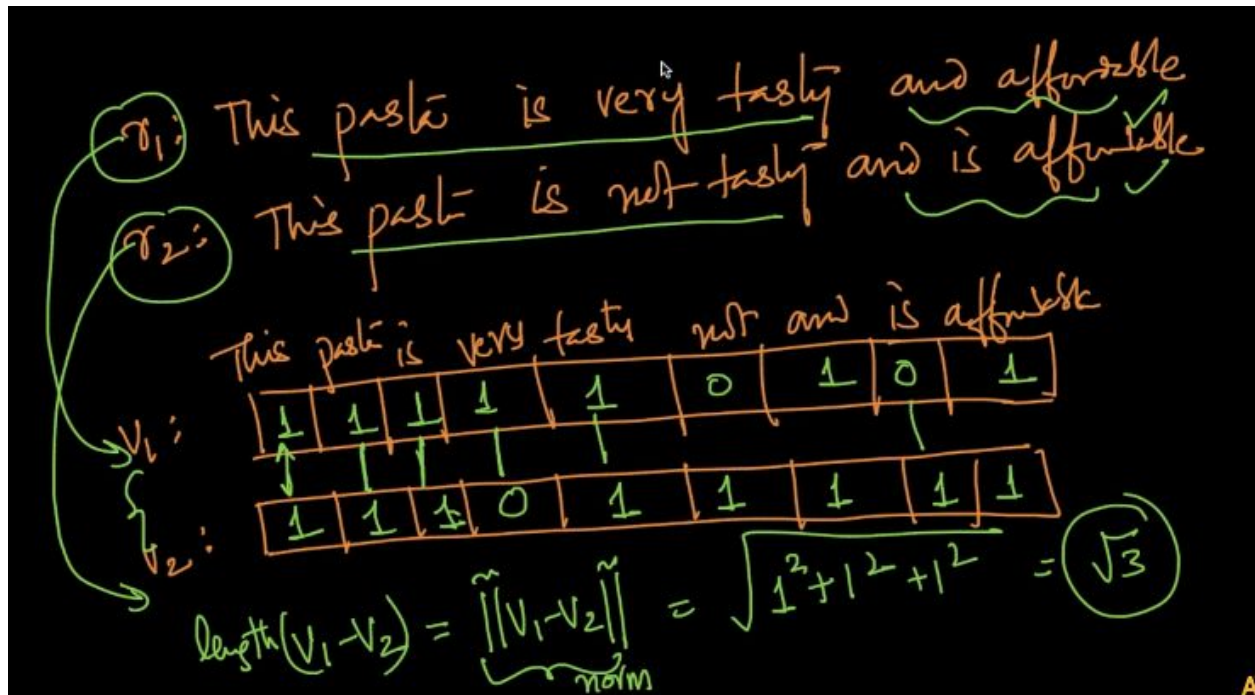②    $r_1$: (This) pasta is very tasty and affordable

$V_i \in R^d$   ← $V_i$:   

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | - | 7 | | d-1  (d) |
| 0 | 0 | 1 | | 1 | 1 | 1 | |

a  an  the  pasta  This  tasty  is

(#times the word occurs in $r_i$)

(*) each word is a different dimension

$V_1 \to$ (Sparse) — most of the ele. are zero

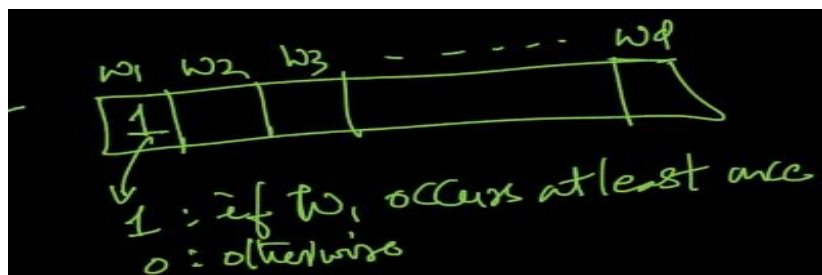| 1 | 2 | 3 | | | | | | d |
|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 1 | 0 | 1 | 2 | 0 | -0...0 | |

d is large

BoW

text → vector

**Bag of words can be thought of counting the common words when all the values exist once.**

*BoW doesn't work well when there are small changes. BOW depends on the count of words in each document corpus.This discards the semantic meaning of the documents such that documents that are completely opposite in meaning can lie closer to each other like above example.*

### Binary BoW or boolean BoW

It doesn't reflect the no of occurrence or count but summarizes if the words occurs or not.



$||v1-v2|| = $ sqrt(no of differing words) or sqrt(unique words between 2 documents)

# Improving BoW (Text Preprocessing)

D1: **This** pasta **is** very tasty **and** affordable.
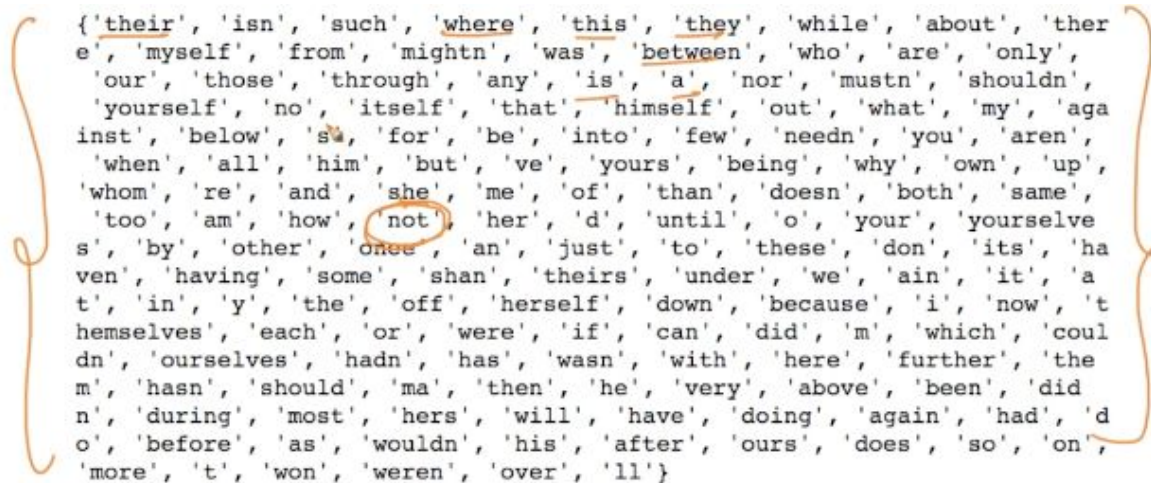D2: **This** pasta **is** not tasty **and is** affordable.
D3: **This** pasta **is** delicious **and** cheap.
D4: Pasta **is** tasty **and** pasta tastes good;

*The **Bold** and **Underlined** words are **stopwords**.*

***Stopwords*** *are words which are filtered out before or after processing of natural language data (text).Though "stop words" usually refers to the most common words in a language, there is no single universal list of stop words used by all natural language processing tools, and indeed not all tools even use such a list.*

*They also depend on domain knowledge but are usually meaningless and have no value to the semantic meaning of the document.*
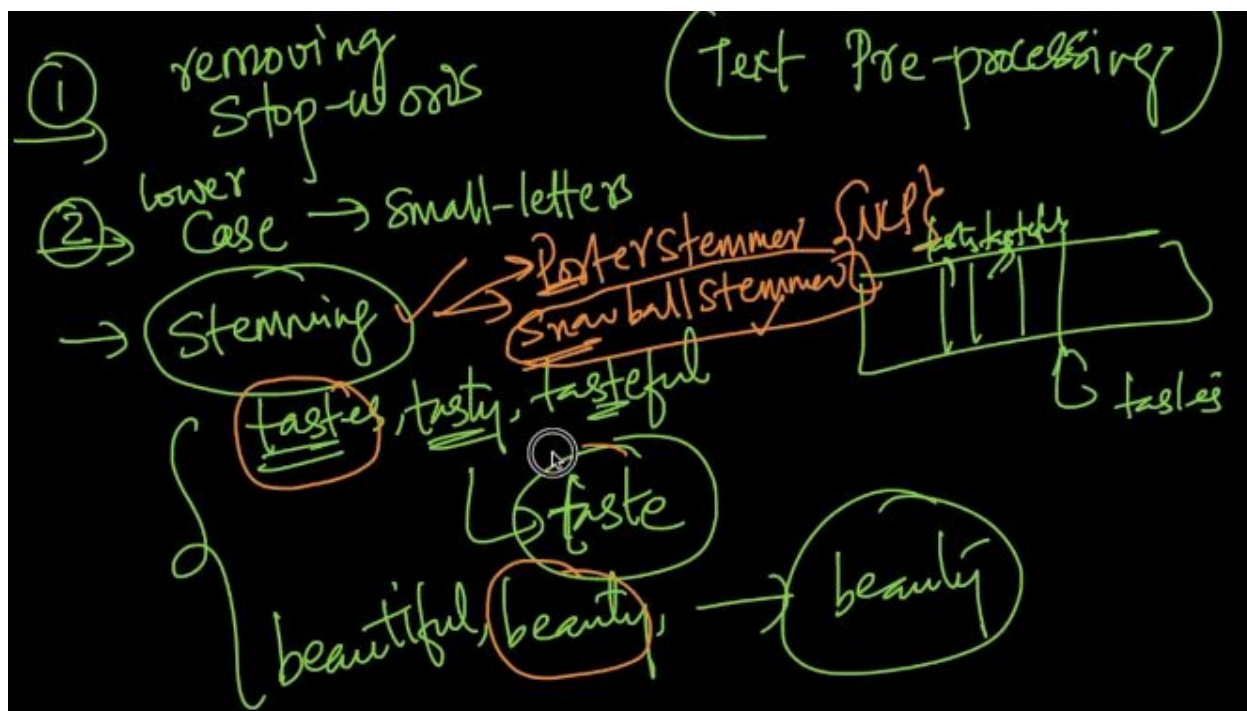
```
{'their', 'isn', 'such', 'where', 'this', 'they', 'while', 'about', 'ther
e', 'myself', 'from', 'mightn', 'was', 'between', 'who', 'are', 'only',
 'our', 'those', 'through', 'any', 'is', 'a', 'nor', 'mustn', 'shouldn',
 'yourself', 'no', 'itself', 'that', 'himself', 'out', 'what', 'my', 'aga
inst', 'below', 's', 'for', 'be', 'into', 'few', 'needn', 'you', 'aren',
 'when', 'all', 'him', 'but', 've', 'yours', 'being', 'why', 'own', 'up',
'whom', 're', 'and', 'she', 'me', 'of', 'than', 'doesn', 'both', 'same',
 'too', 'am', 'how', 'not', 'her', 'd', 'until', 'o', 'your', 'yourselve
s', 'by', 'other', 'once', 'an', 'just', 'to', 'these', 'don', 'its', 'ha
ven', 'having', 'some', 'shan', 'theirs', 'under', 'we', 'ain', 'it', 'a
t', 'in', 'y', 'the', 'off', 'herself', 'down', 'because', 'i', 'now', 't
hemselves', 'each', 'or', 'were', 'if', 'can', 'did', 'm', 'which', 'coul
dn', 'ourselves', 'hadn', 'has', 'wasn', 'with', 'here', 'further', 'the
m', 'hasn', 'should', 'ma', 'then', 'he', 'very', 'above', 'been', 'did
n', 'during', 'most', 'hers', 'will', 'have', 'doing', 'again', 'had', 'd
o', 'before', 'as', 'wouldn', 'his', 'after', 'ours', 'does', 'so', 'on',
'more', 't', 'won', 'weren', 'over', 'll'}
```

Removing the stop word is not always the best choice. For example :
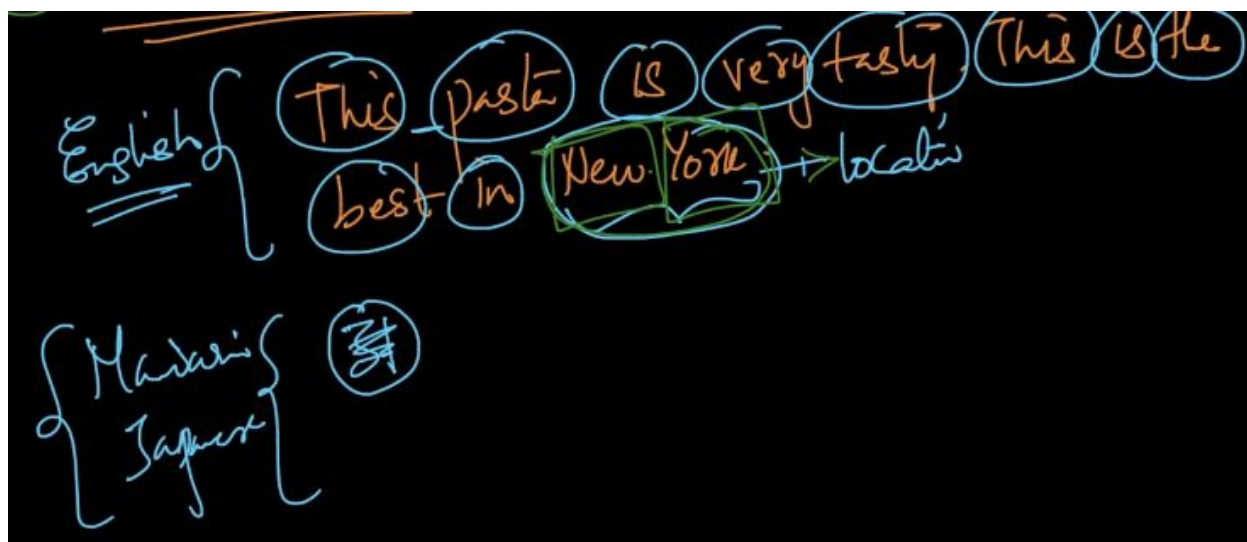D2: **This** pasta **is** not tasty **and is** affordable.
*If we remove not from D2, it completely changes the meaning of the document.*
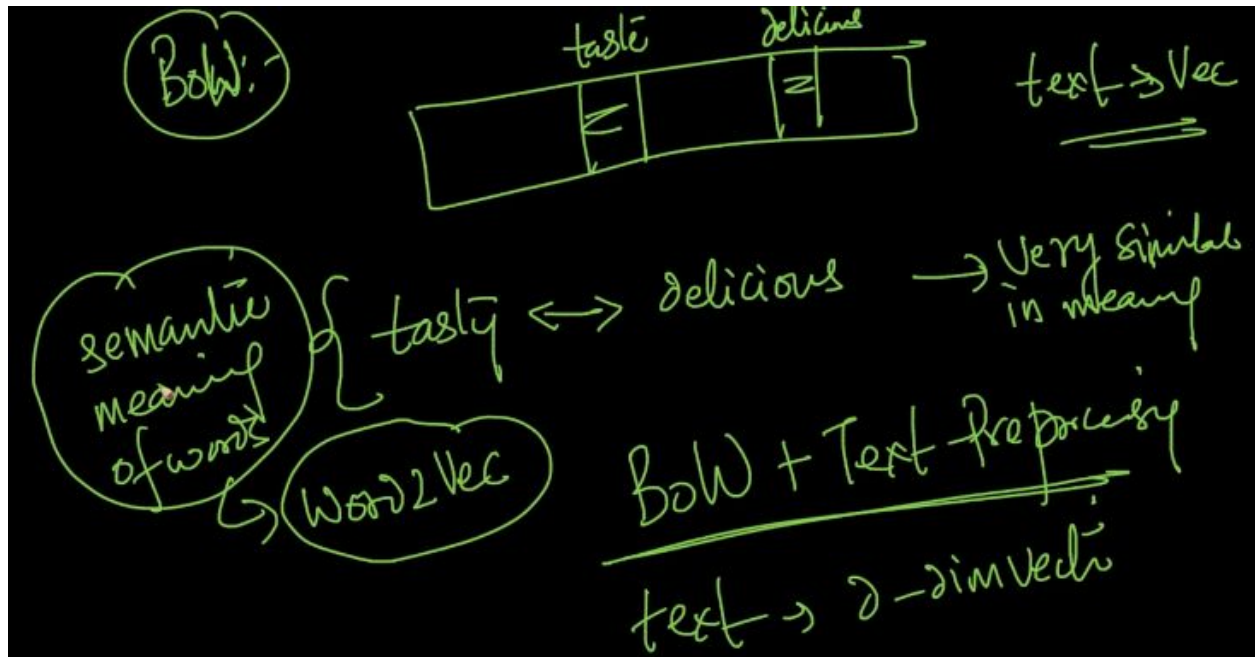But, for most of context, if we remove stop words we could have a smaller and meaningful.

Lemmatization 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.

Tokenization : Breaking a sentence into words. Language and content dependent.

BoW doesn't take semantic meaning in consideration so we use Word2vec.

Stemming just removes or stems the last few characters of a word, often leading to incorrect meanings and spelling. Lemmatization considers the context and converts the word to its meaningful base form, which is called Lemma. Sometimes, the same word can have multiple different Lemmas. We should identify the Part of Speech (POS) tag for the word in that specific context. Here are the examples to illustrate all the differences and use cases:

If you lemmatize the word 'Caring', it would return 'Care'. If you stem, it would return 'Car' and this is erroneous.

If you lemmatize the word 'Stripes' in verb context, it would return 'Strip'. If you lemmatize it in noun context, it would return 'Stripe'. If you just stem it, it would just return 'Strip'.

You would get the same results whether you lemmatize or stem words such as walking, running, swimming... to walk, run, swim etc.

Lemmatization is computationally expensive since it involves look-up tables and what not. If you have a large dataset and performance is an issue, go with Stemming. Remember you can also add your own rules to Stemming. If accuracy is paramount and the dataset isn't humongous, go with Lemmatization.

# uni-gram, bi-gram, n-grams.



Uni-gram / Bi-gram / n-gram.

✓ $r_1$: (This) pasta (is) (very) tasty (and) affordable.

✓ $r_2$: (This) pasta (is) (not) tasty (and) (is) affordable.

removing stopwords;

$V_1$ & $V_2$ are exactly same

⟹ conclude $r_1$ & $r_2$ are very similar



$r_1$: This pasta is (very tasty) and affordable

$r_2$: This pasta is (not tasty) and is affordable

This | is | pasta | very | tasty | afffable | . . . —

✓ Uni-gram:-
↳ each word is consider a dim.

This pasta | pasta is

very tasty

not tasty

is affordable

✓ bi-grams:-
↳ pair of words



⟹ # bi-grams ≥ # Uni-grams

# tri-gms
≥
4-grams

n-grams (n>1) → dimensionaly 'd' increases ↑

# tf-ldf (term frequency- Inverse document frequency)

$$N \text{ docs}/\alpha \text{ reviews} \begin{cases} r_1: w_1 \ w_2 \ w_3 \ w_2 \ w_5 \rightarrow \boxed{5} \\ r_2: w_1 \ w_3 \ w_4 \ w_5 \ w_6 \ w_2 \rightarrow \boxed{6} \\ r_3: \\ \vdots \\ r_N: \end{cases}$$

| $w_1$ | $w_2$ | $w_3$ | $w_4$ | $w_5$ | $w_6$ | |
|---|---|---|---|---|---|---|
| 1 | 2 | 1 | 0 | 1 | 0 | BoW |
| 1 | 1 | 1 | 1 | 1 | 1 | BoW |

$\rightarrow$ how often does $w_i$ occur in $r_j$

$$TF(w_2, r_1) = 2/5$$

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

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

IDF: Inverse document freq:-

$$D = \begin{cases} r_1: & w_1 \\ r_2: & \\ r_3: & \underline{\quad w_1 \quad} \\ \vdots \\ r_6 & \quad w_1 \quad \\ \vdots \\ \boxed{N}: \end{cases}$$

N docs/ $\alpha$ reviews

$$IDF(w_i, D_c)$$

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

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

$\rightarrow \#$ docs

$\# $ docs which contain $w_i$

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

$\rightarrow \#docs$ containing $w_i$

$n_i \leq N \Rightarrow \dfrac{N}{n_i} \geq 1$

$\log\left(\dfrac{N}{N_i}\right) \geq 0$

$\log(1) = 0$

---

$\log\left(\dfrac{N}{n_i}\right)^{\rightarrow low}$

①

$IDF \geq 0$

if $\boxed{n_i \uparrow}$ ; $\dfrac{N}{n_i} \downarrow$ ; $\log\left(\dfrac{N}{n_i}\right) \downarrow$ ②

if $w_i$ is more freq in my corpus then $IDF \downarrow$

$\dfrac{1000}{10} > \dfrac{1000}{20}$

monotonic fn ③

$\cancel{IDF \propto} $

$IDF \downarrow \quad n_i \uparrow$

$n_i \downarrow \quad IDF \uparrow$

---

$f(x) \uparrow$

$f(x)$

monotonically inc

$g(x) \downarrow$

$x$

$g(x)$ : monotonically decreasing fn $\downarrow$

When $W_i$ is more frequent, IDF will be low
When $W_i$ is a rare word, IDF will be higher.

$$r_1: W_1 W_2 W_3 W_2 W_4$$
$$r_2: W_1 W_3 W_5 W_6 W_2$$
$$r_3:$$
$$\vdots$$
$$r_N:$$

$D_c$

$$TF(W_2, r_1) * IDF(W_2, D_c)$$

$$V_1$$

BoW:- #occurences of $W_2$ in $r_1$

$$\left( r_1, r_2, \cdots r_N \right) = D_c$$

$$r_i \rightarrow V_i$$

$W_j$

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

$$\left( W_j \text{ is frequent in } r_i \right) \quad \left( W_j \text{ is rare in } D_c \right)$$

AI

TF-IDF

↳ more importance to rarer words in my Dc

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

→ Semantic - meaning $\begin{pmatrix} tasty \to delicious \\ cheap \leftrightarrow affordable \end{pmatrix}$

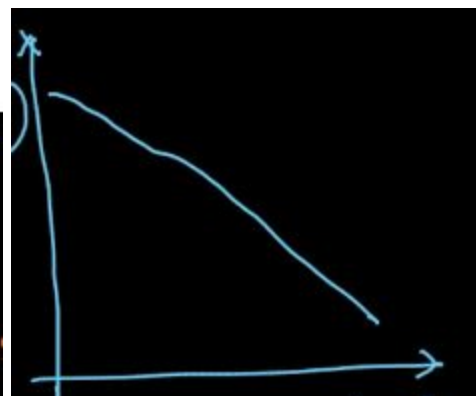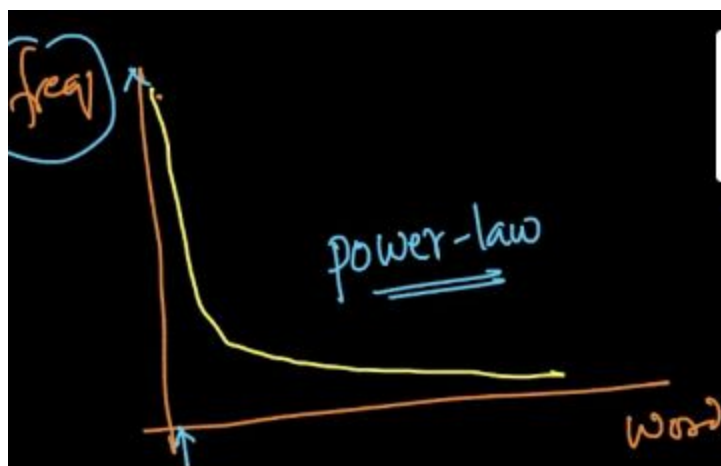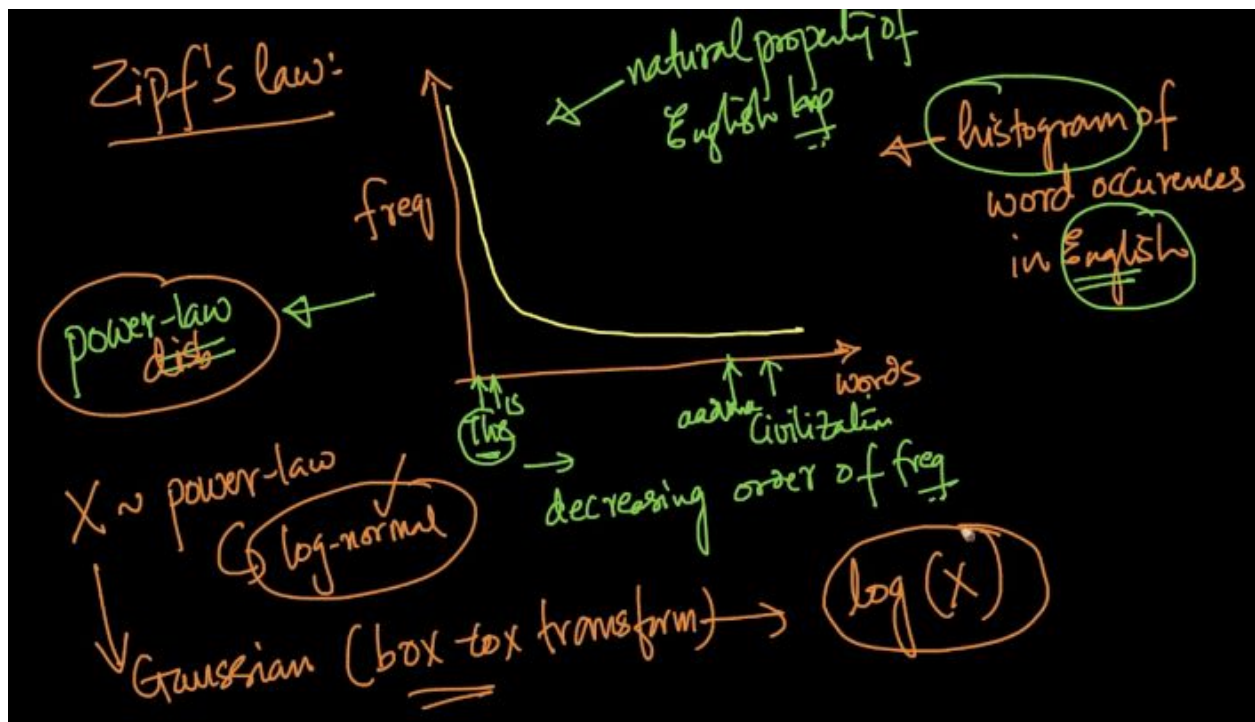## Why use log in the IDF?



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

$$IDF(w_i, Dc) = \log\left(\frac{N}{n_i}\right) \quad \begin{cases} \frac{N}{n_i} & N \to \#docs \\ & n_i \to \#docs \text{ which contain } w_i \end{cases}$$

→ 1972 research paper
→ heuristic (or) hack → not very strongly on theory

→ zipf's law
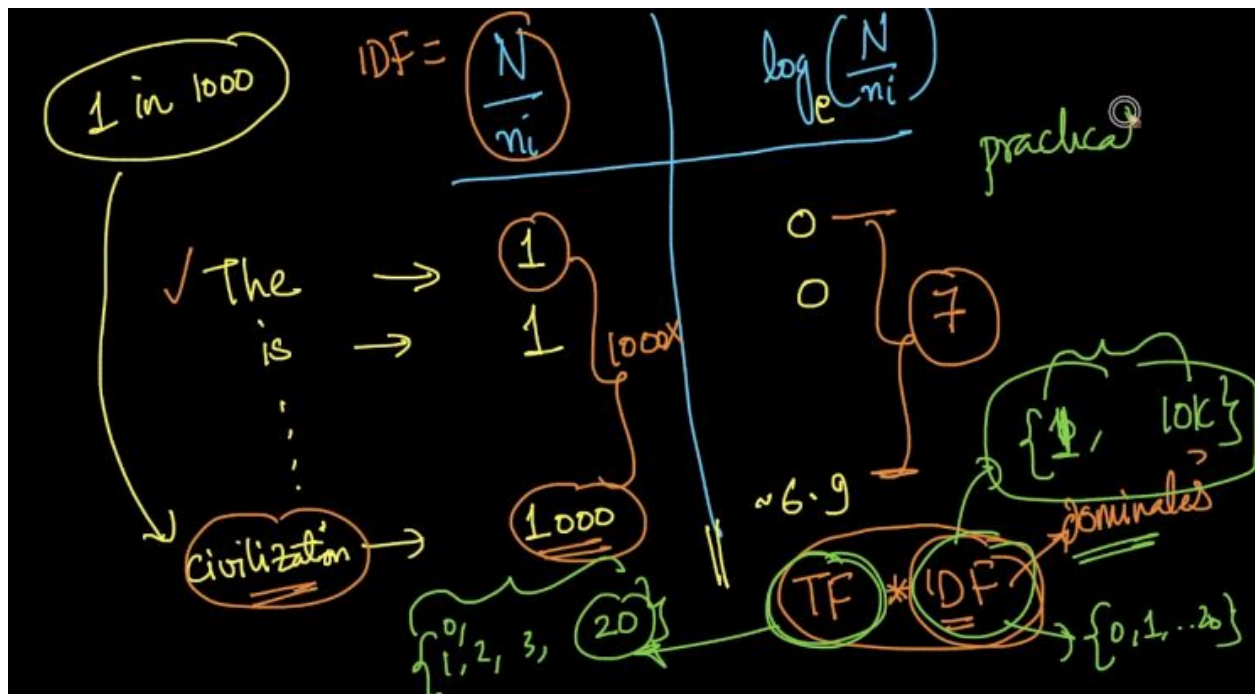
log(freq) vs log(word) gives straight line

Note:
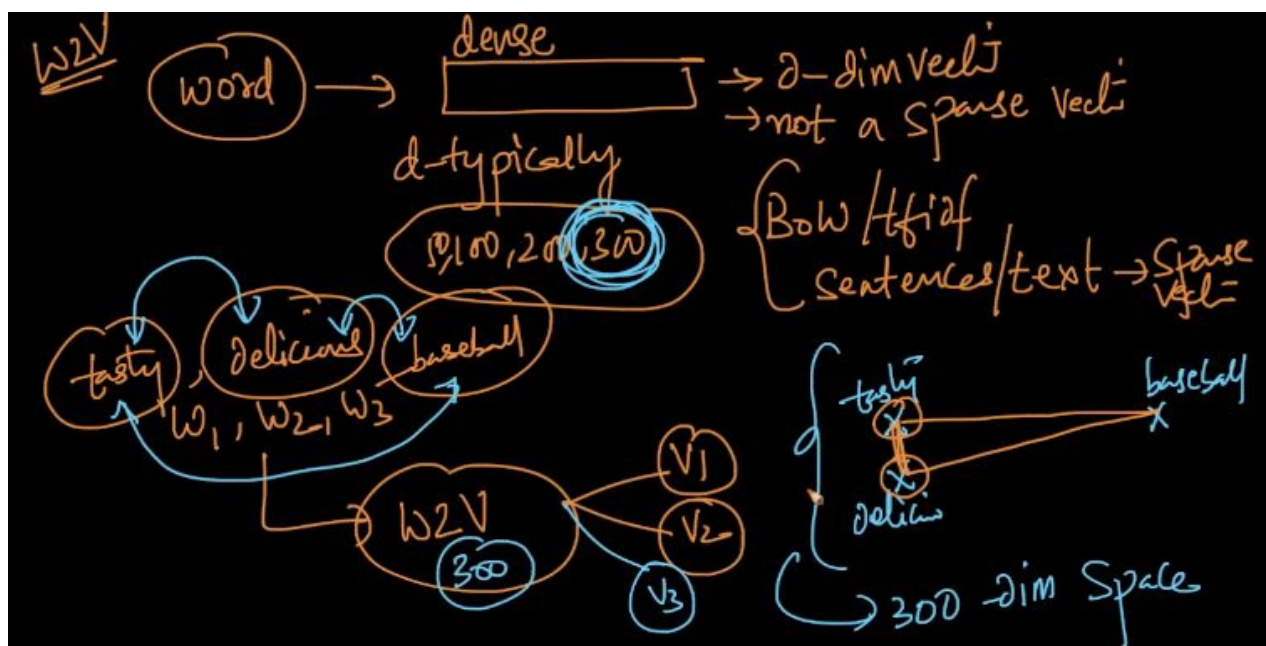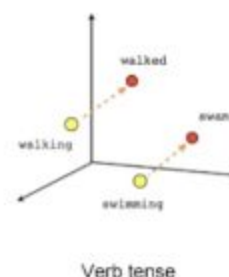IDF without a log in the formula will have large numbers that will dominate in the ML model
Taking log will bring down dominance of IDF in TF- IDF values;

## Word2Vec

**State of art methodology to convert word to vec. Based on semantic meaning .**

W2V:

300-dim

semantically

① ✓ { W₁ & W₂ are similar then
      V₁ & V₂ closer

② relationships

$(V_{man} - V_{women}) || (V_{king} - V_{queen})$

man
  ↘ • women
king •
  ↘
  ◉
  queen



300-dim

Male-Female

Country-Capital

Verb tense

Word2Vec learns all the relationships from raw-test without the involvement of programing.



※ (W2V) → learning relationships
             automatically from raw text

NLP (20B)

large text corpus → W2V → Word : Vec
                     300
                     200 -
                     100 -
                     50 -

larger dimensions → more info rich the vector is

data corpus size $\uparrow$ $\Rightarrow$ $d\uparrow$
$\hookrightarrow$ dimensionality

Bilking

Google - News
Massive $\longrightarrow$ W2V $\rightarrow$ 300-dim

review
text $\rightarrow$ W2V $\rightarrow$ word. vec

book

Core: W2V                    (intuition)

$W_1$ $W_2$ $W_3$ $W_4$ $W_5$ $W_2$

Neighborhood of $W_3$

$W2V(W_3) =$   $\begin{cases} N(w_i) \approx N(w_j) \\ V_i \approx V_j \end{cases}$

# Avg-Word2Vec, tf-idf weighted Word2Vec

$n_1 \leftarrow \gamma_1: \quad \underline{W_1} \ \underline{W_2} \ W_1 \ W_3 \ W_4 \ \overline{W_5}$   (work)

$\boxed{Avg \ W2V}$

$\gamma_1 \rightarrow V_1$

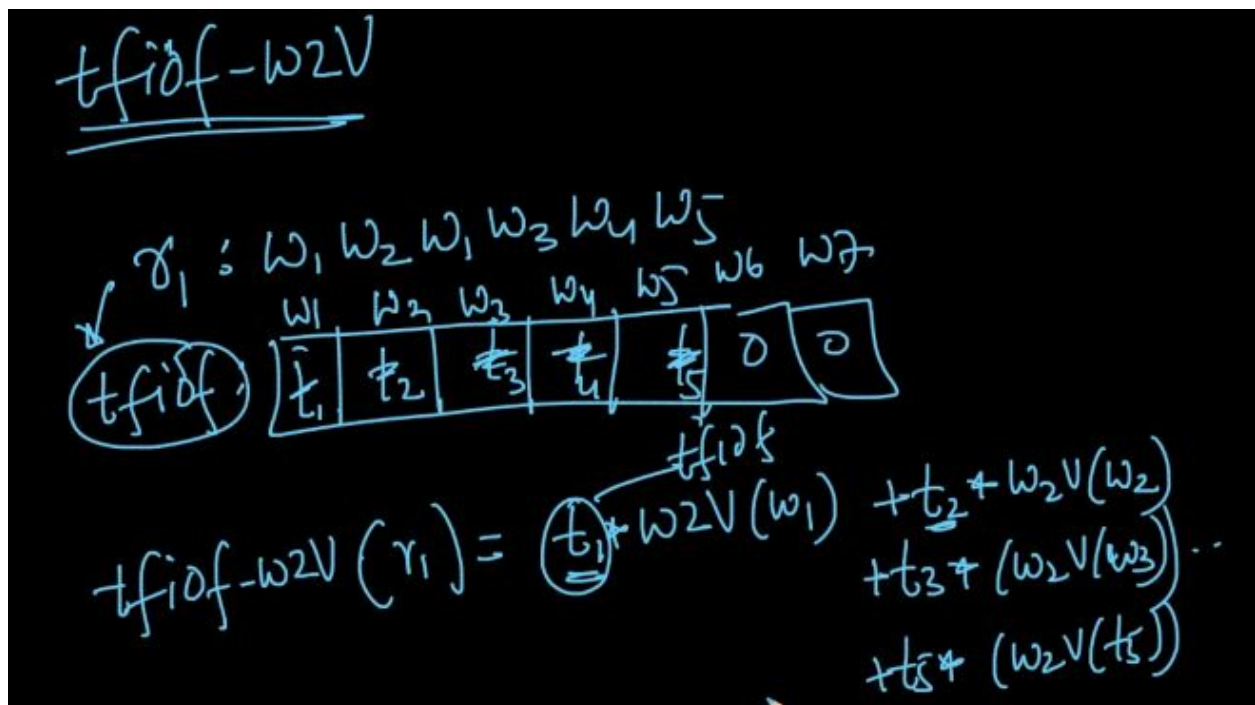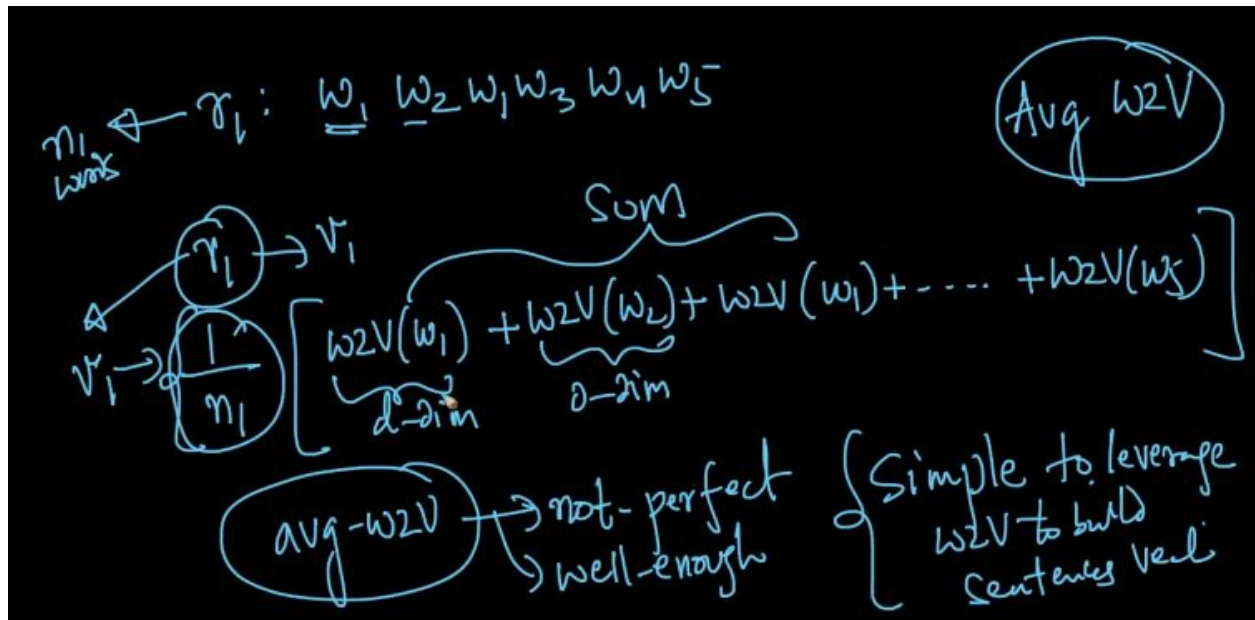$V_1 \rightarrow \phi \left[ \dfrac{1}{n_1} \right] \Big[ \underbrace{W2V(W_1)}_{d-dim} + \underbrace{W2V(W_2) + W2V(W_1)}_{0-dim} + \cdots + W2V(W_5) \Big]$

$\overbrace{\qquad\qquad\qquad\qquad}^{Sum}$

$\boxed{avg-W2V} \rightarrow$ not-perfect, well-enough

$\left\{ \begin{array}{l} Simple \ to \ leverage \\ W2V \ to \ build \\ sentences \ vec \end{array} \right.$

---

# $\underline{tfidf - W2V}$

$\gamma_1 : W_1 \ W_2 \ W_1 \ W_3 \ W_4 \ \overline{W_5}$

$\boxed{tfidf}$

| $W_1$ | $W_2$ | $W_3$ | $W_4$ | $W_5$ | $W_6$ | $W_7$ |
|---|---|---|---|---|---|---|
| $t_1$ | $t_2$ | $t_3$ | $t_4$ | $t_5$ | $0$ | $0$ |

$tfidf$

$tfidf-W2V(\gamma_1) = \underset{tfidf}{t_1} * W2V(W_1) + t_2 * W_2V(W_2) + t_3 * (W_2V(W_3)) \cdots + t_5 * (W_2V(t_5))$

$$\text{tfidf-w2v}\left(\overset{\circ}{r_i}\right) = \frac{\sum\limits_{i:\text{wri's}}\left(t_i * \text{w2v}(w_i)\right)}{\sum\limits_{i:\text{wri's}} t_i}$$

tfidf(wi, ri)

$$\text{tfidf-w2v}\left(\overset{\frown}{r_i}\right) = \frac{\sum\limits_{i:\text{wri's}}\left(t_i * \text{w2v}(w_i)\right)}{\sum\limits_{i:\text{wri's}} t_i}$$

tfidf(wi, ri)

If all of $t_i = 1$

tfidf-w2v → avg w2v

sparsity of the matrix = no.of zero elements / total elements

We can store sparse matrices with a less order of space by efficient method using by storing the row and column indices with corresponding values;

CountVectorizer() does this automatically . It has an inbuilt system to remove the super sparsity of the dataset.

## CountVectorizer()

Documentation:

https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVectorizer.html