# PREDICT RATING PRODUCT REVIEWS ON AMAZON

## DATASET OVERVIEW: AMAZON FOOD REVIEWS (EDA)



Dataset understanding :

1) HelpfulnessNumerator - Who found review useful useful  (Ex - 4500)

2) HelpfulnessDenominator - They found the review useful + not useful (Ex - 4500 + 500(Not use))

3) Score - It is between 1 to 5 and we want it to be binary so 1*,2* is negative and 4*,5* positive by adding a new label.

4) Text  + Summary : Most important in our dataset since it is an **NLP** problem

# WHY CONVERT TEXT INTO VECTOR ?



In the Amazon Food review if we could convert our text into vectors we could use all the Linear Algebra things for our sentiment analysis. Let's say we convert our review text to d-dim vector



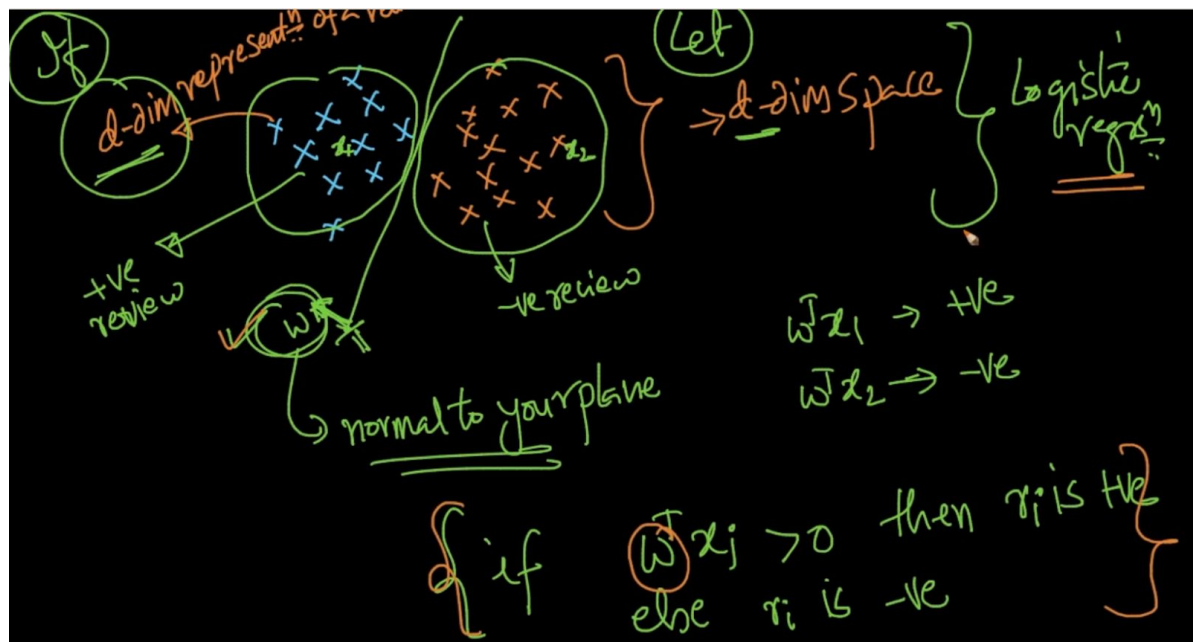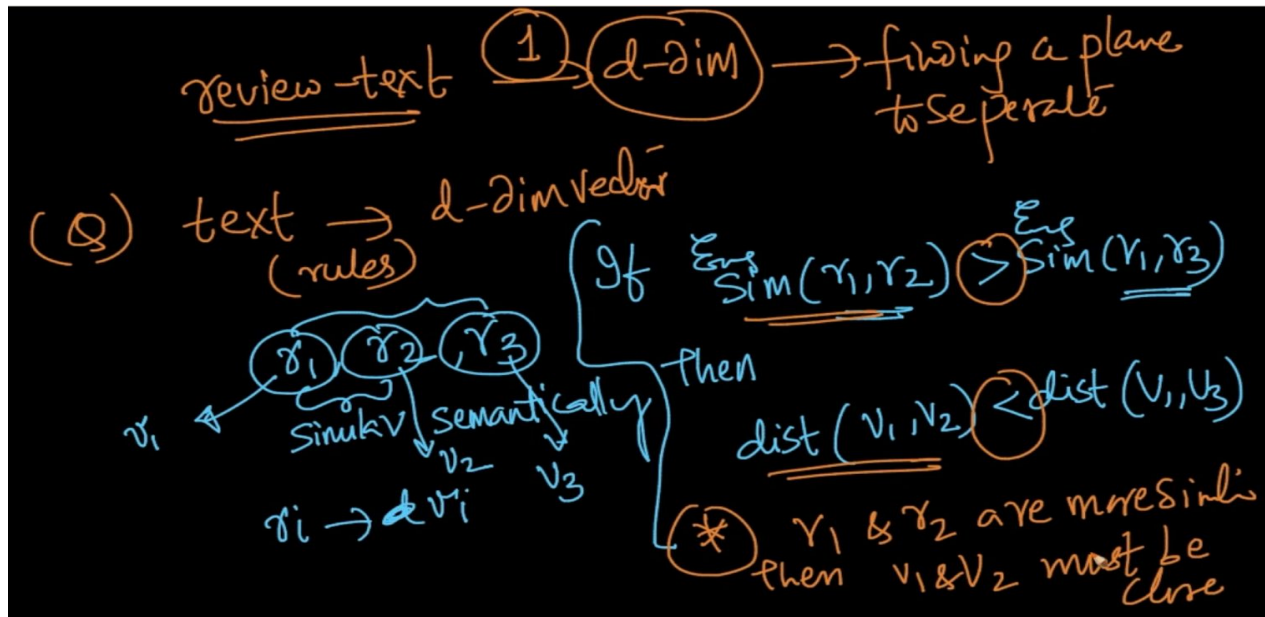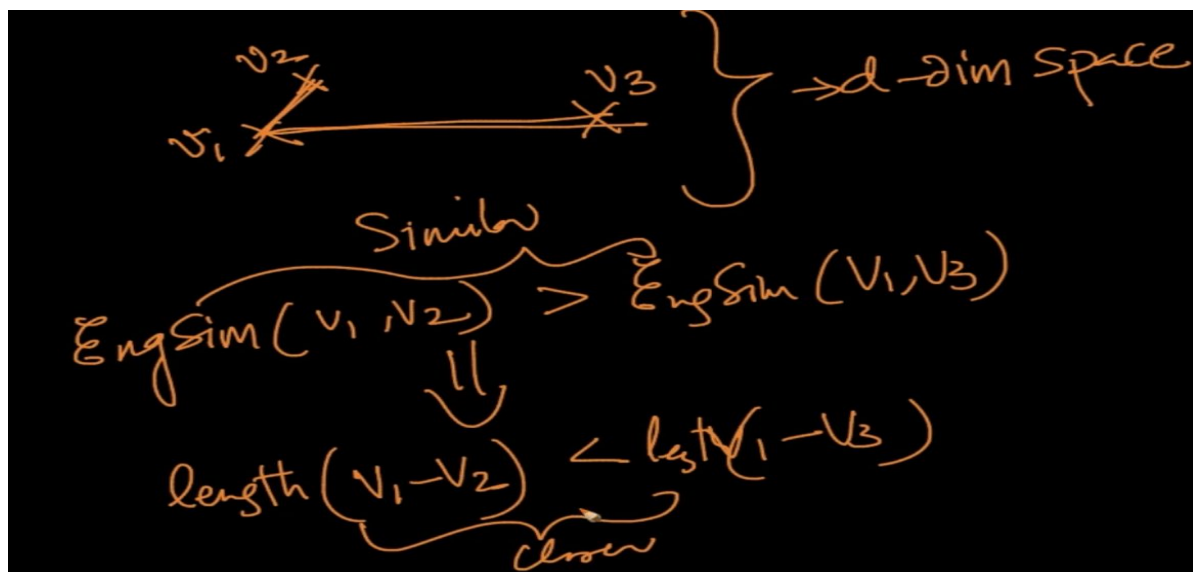In a d - dimensional space if the data-points(text reviews) are represented like this then we could separate them using a Linear Algebra by finding a normal to the plane that has separated the dataset

review text ① (d-dim) → finding a plane to seperate

(Q) text → d-dim vector
   (rules)

If $Eng\ sim(r_1, r_2) > Eng\ sim(r_1, r_3)$

$r_1$ $r_2$ $r_3$   then

$v_i$   sinukv semantically   $dist(v_1, v_2) < dist(v_1, v_3)$
   $v_2$   $v_3$

$r_i \rightarrow \& v_i$

(*) $r_1$ & $r_2$ are more similar then $v_1$ & $v_2$ must be close

So first our text should be converted into d-dimensional vectors. Then suppose if we've 3 reviews r1,r2 and r3 and r1 and r2 are more similar then distance of r1 and r2 should be less than distance between r1 and r3



$v_2$
$v_1$   $v_3$   } → d-dim space

Similar

$Eng\ sim(v_1, v_2) > Eng\ sim(v_1, v_3)$

||

$length(v_1 - v_2) < lgth(v_1 - v_3)$
   closer

Geometric interpretation. So if the similar texts i.e positive text with positive and negative with negative are aligned or geometrically closer then we can find the plane to separate them.

So converting text to vector s.t they are geometrically closer we'll learn techniques like Bag of words(BOW), word2vec, tf-idf and so on

So, these are the 4 reviews. In BoW the 1st step is constructing a dictionary: Set of all the unique words in reviews {This,pasta,is,very …..}



The dict is a d-dim vector containing all the words. Each index will correspond to the word and the value will tell us about the frequency of a word. Since there will be many words in all the reviews (documents) the vector for a single review is very sparse

So in BoW we are trying similar texts to be geometrically closer.



Here, we took 2 reviews and converted them into vectors and then we calculated the differences in vectors and calculated distances as seen but there's a problem . In the 2 reviews, *very tasty* (r1) and *not tasty* (r2) have very different meaning but the distance is less. So it's a problem in BoW
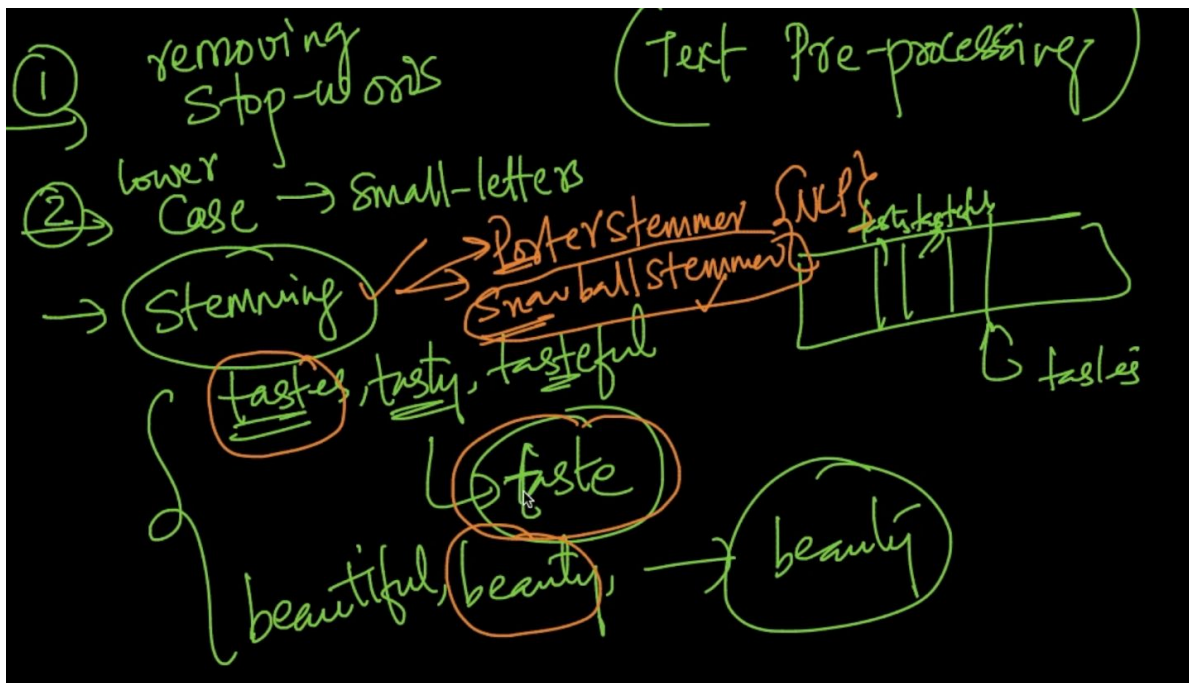
**Note:**

1) Binary BoW: The value in the vector is no. of times the word is occuring so if the word "Pasta" is occuring twice then value will be two. So in binary BoW we if the word is occuring then we'll just give value 1 because we are trying to min the distances of similar texts

2) Words like "The, is, are ,that etc" are not important in text-analysis so we try to remove them

# Stop-word removal, Tokenization, Lemmatization (Featurizations )



First step is we remove stopwords like this, is, not etc. It makes the vector smaller and more meaningful.
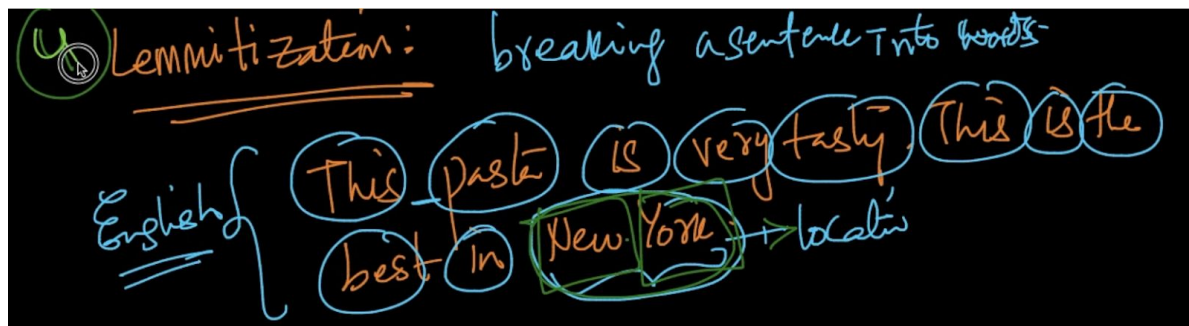
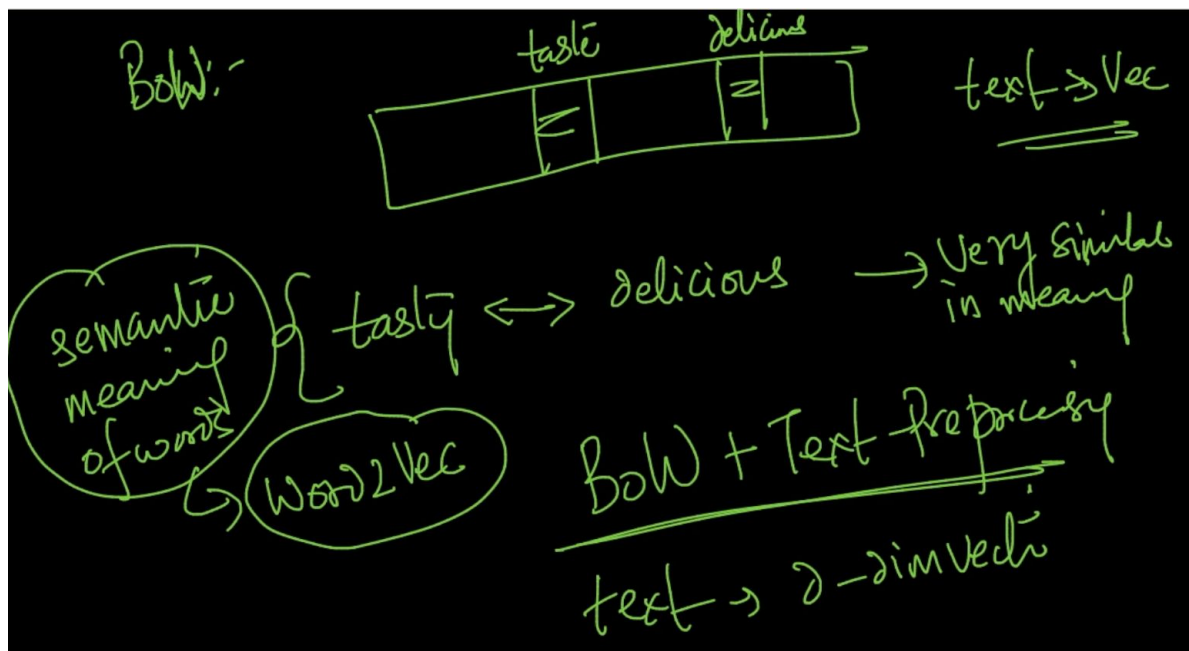

## Text Pre-processing steps:

1 ) Removing stop-words

**2 ) Lower case** - We convert the letters to lower case . Why? In the above example, *Pasta* and *pasta* are two different words when they are same . So make all of them lowercase

3) Stemming - We get the essence of the word from stemming. For ex: tastes, tasty and tasteful have just one essence taste. So instead of using 3 different vectors for them, we can use a single vector
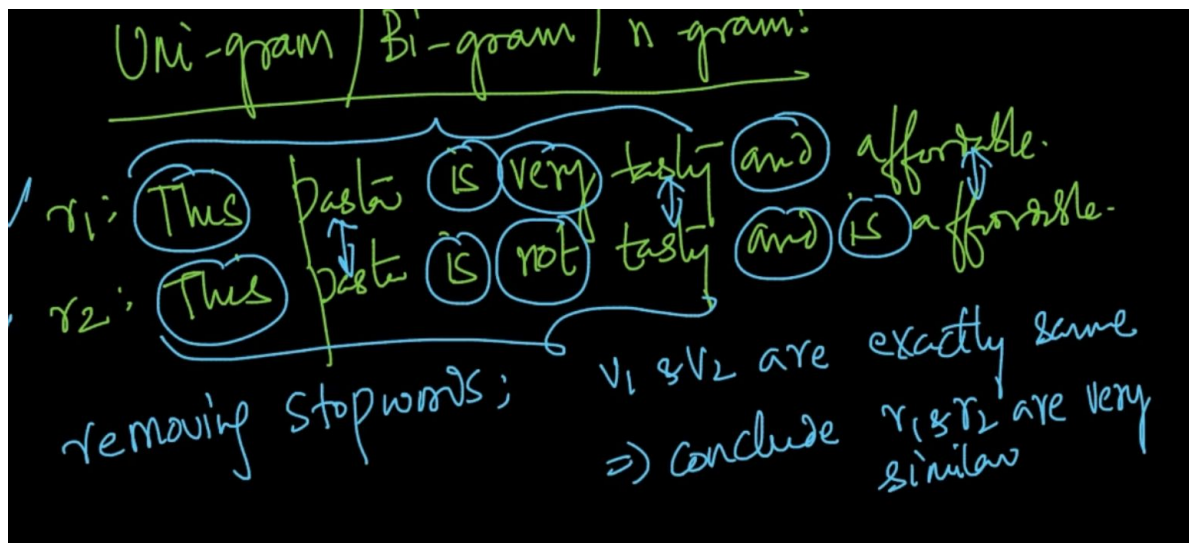
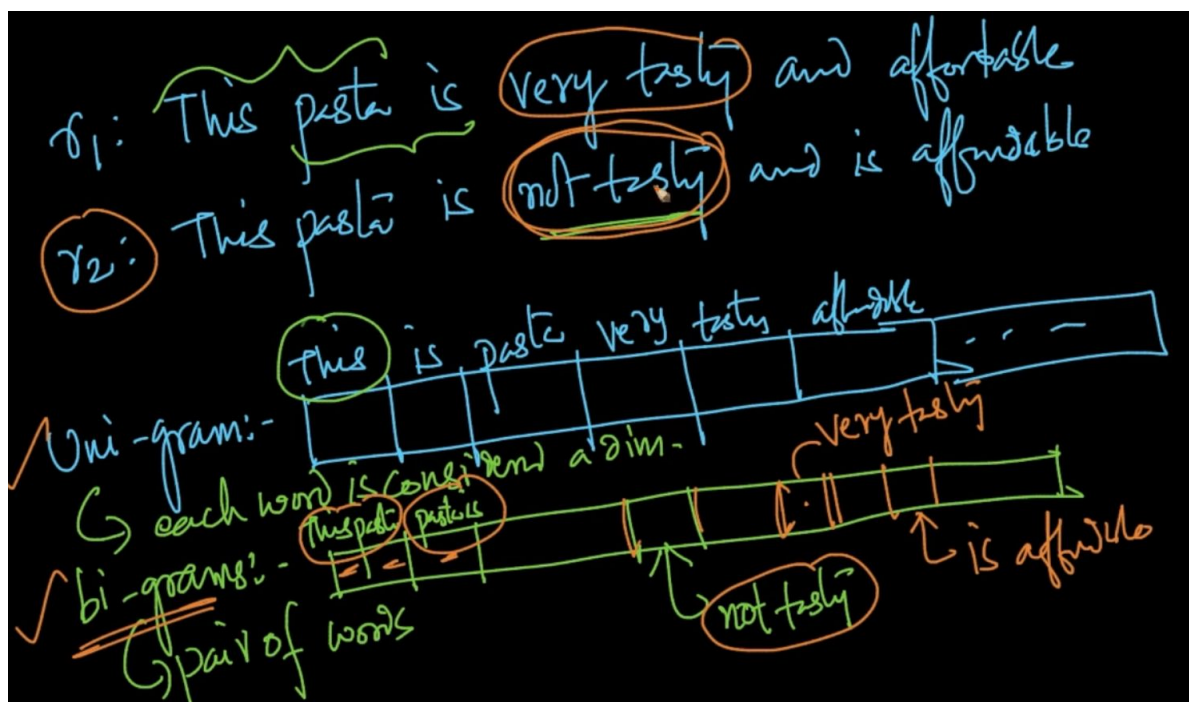**4) Lemmatization** - We break the sentence into words

There's a problem with BOW as we can see in r1 and r2  reviews above we can clearly see that they've the same meaning like tasty and delicious are synonyms but BOW will consider as different vectors. So it will not cover the semantic meaning behind the texts. Therefore, with BOW + Text Preprocessing we are converting text into vectors.
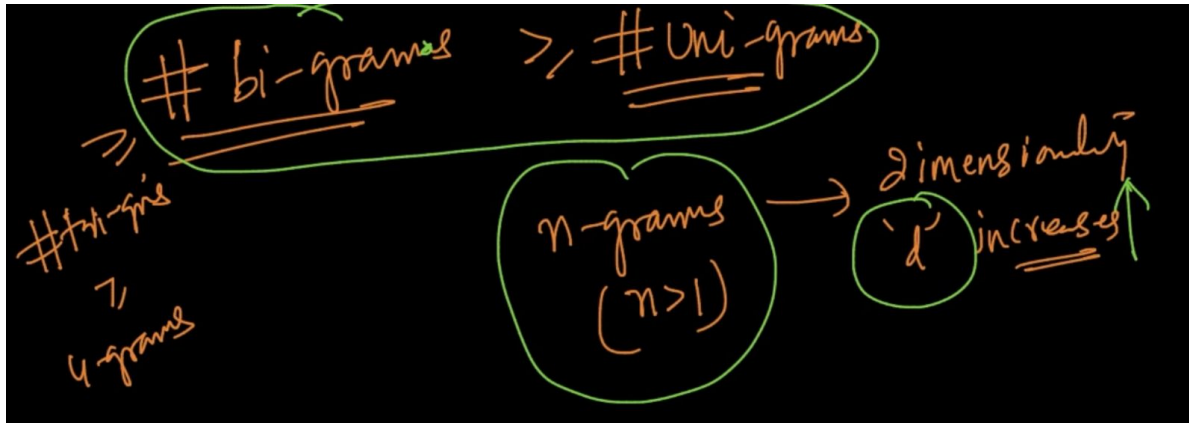
# Uni-gram/bi-gram/n-gram



If we are removing stopwords then the above 2 reviews get very similar albeit the fact that they are different. We'll see below how to solve that
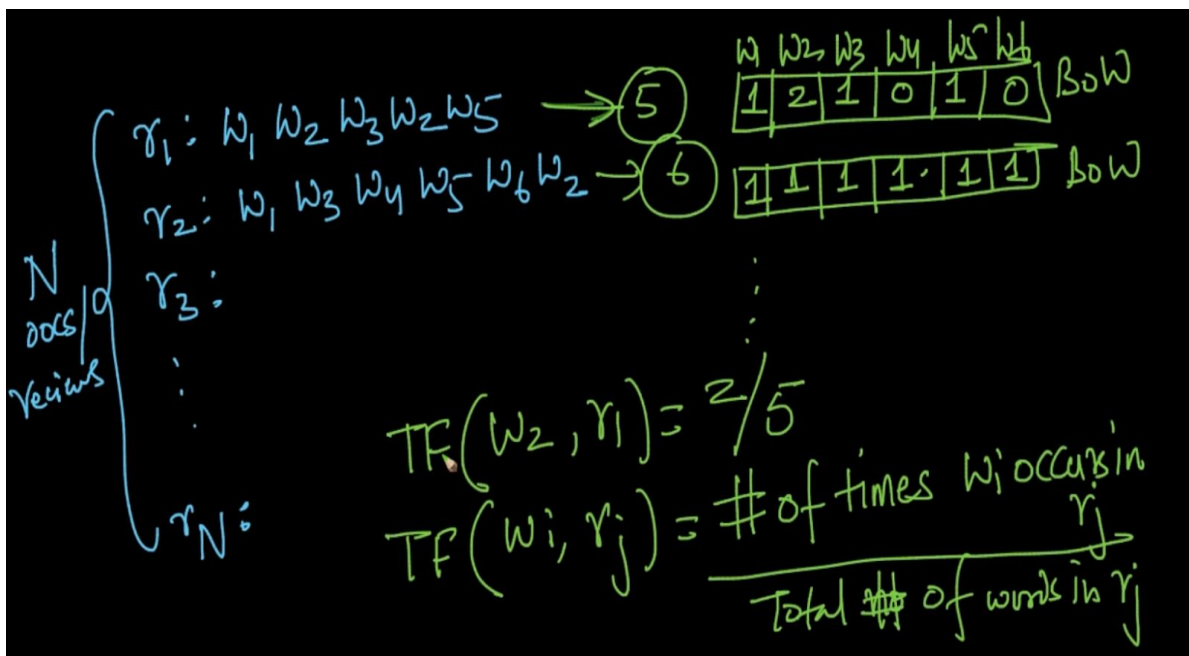


Our traditional approach was uni-gram approach . We use bi-gram in which a vector of **2 consecutive word** is created.

Note - <mark>In our unigram all the sequence info was lost as seen but in our n-grams some info is retained</mark>

No. of vectors in bi-grams > unigrams. So as n-grams -> dimensionality 'd' increases.

## Tf-idf (Term Freqency- Inverse Document Frequency)



We are making BOW for each review. TF is occurence of a word/ Total no. of words

Ex : TF (W2, r1) = ⅖ as w2 is occuring twice and total no. of words in r1 is 5.

$0 < TF(W_i, r_j) < 1$. So, TF is a probabilistic model. Therefore it can be thought as a probability of a word W occurring in text 'r'

**IDF**



$D_c$ is data corpus i.e collection of all the documents. IDF of a word is log as mentioned above.

1 ) As mentioned in the above text IDF $\geq 0$

3) IF word W is occuring more than IDF will be low . For ex: the word 'the'. So 'the' has low idf

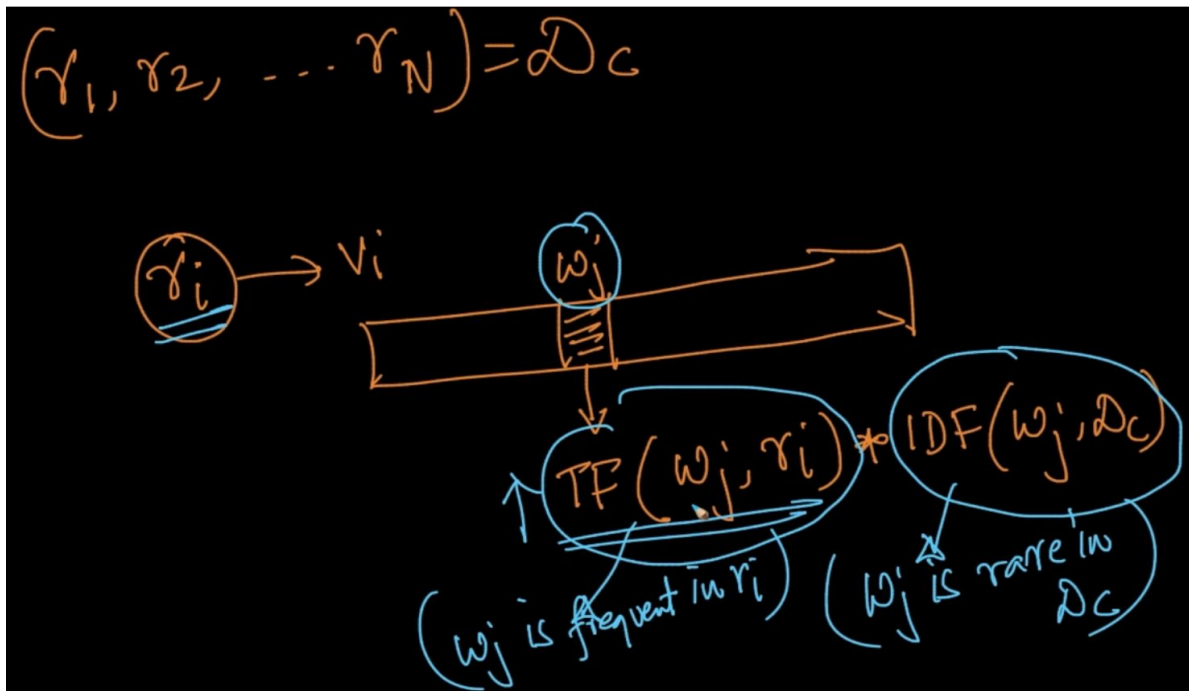4) IDF ↓ n ↑ . Basically the above 3rd point and n ↓ IDF ↑ i.e if a rare word is there then it's IDF will be more.



In BOW : We took the occurence of word in the vector but in TF-IDF we multiplty the TF and IDF of word W. Let's understand the essence of TF-IDF below

$$\left( \gamma_1, \gamma_2, \ldots \gamma_N \right) = \mathcal{D}_C$$

$$\gamma_i \rightarrow V_i$$

$$w_j' \quad$$

$$\underset{(w_j \text{ is frequent in } \gamma_i)}{TF\left(w_j, \gamma_i\right)} * \underset{\left(w_j \text{ is rare in } \mathcal{D}_C\right)}{IDF\left(w_j, \mathcal{D}_C\right)}$$

So what is happening when TF * IDF . TF says w is more frequent in review r and IDF tells us about the rarity .

TF -IDF
↳ more importance to rarer words in my $\mathcal{D}_C$
↳ more importance if a word is frequent in a document / review
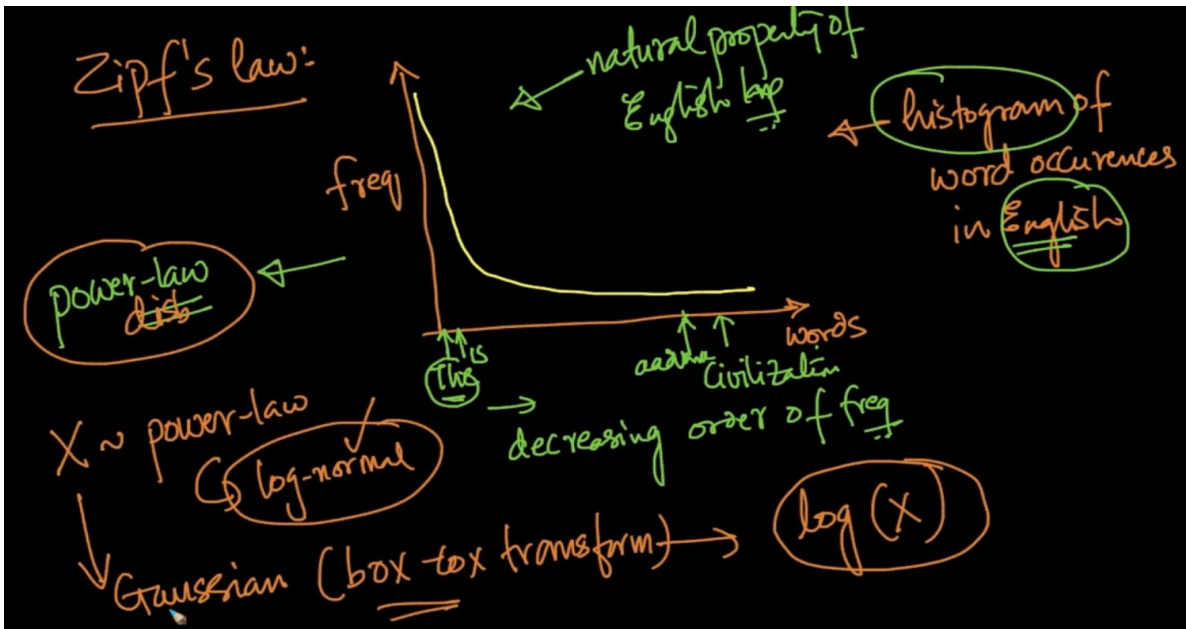↳ Semantic - meaning ( tasty → delicious )

It's importance is mentioned above but ▤ semantic meaning is still not covered .
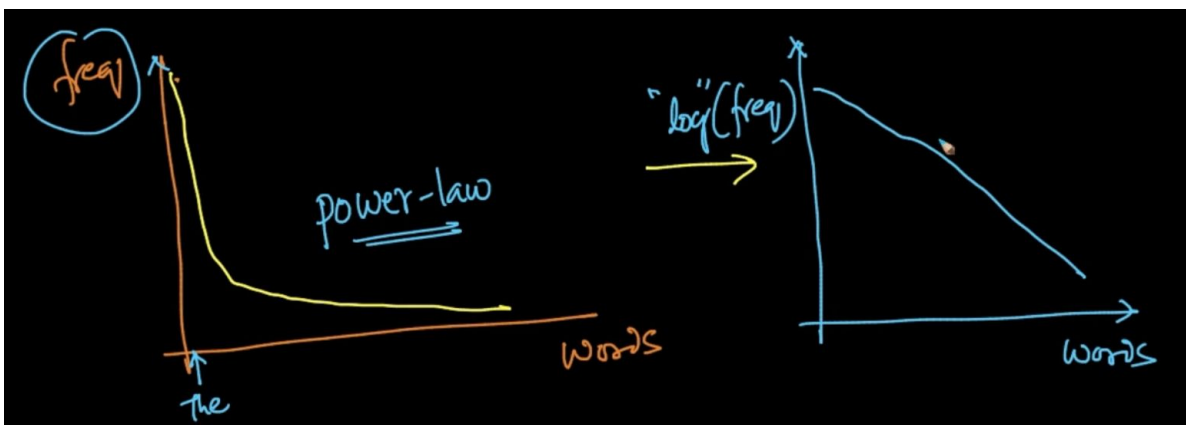
# Why use log in IDF?

In IDF we know that IDF(w,D) = log(N/n) . Here, N/n made sense because N is no. of total documents and n is no. of docs where the word 'w' occurred but why log.

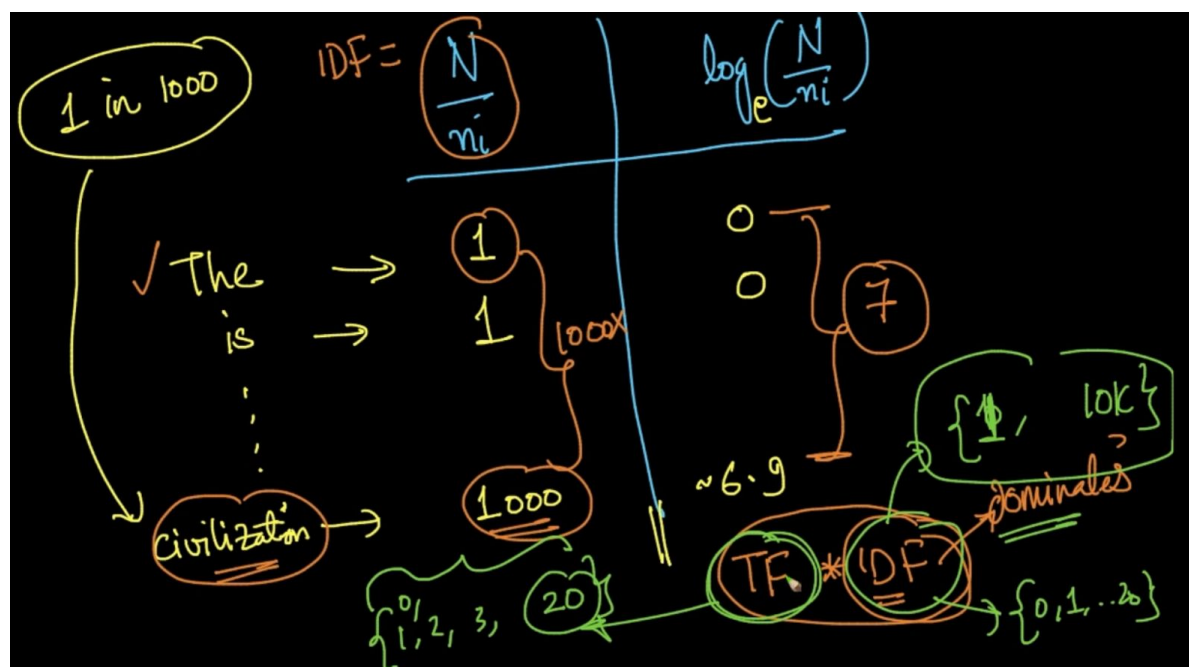Let's use Zipf's law to get an understanding of it



The above is a histogram of words in English let's say Wiki. As seen words like 'the,is, etc' occurs more but words like 'civilization, aardvark etc' has less frequency.

The graph is power-law distribution and we know that if it's power-law then we can convert it to Gaussian by converting the random variable X to log(X).



If we get the log of frequency we get a straight line which is much more manageable.
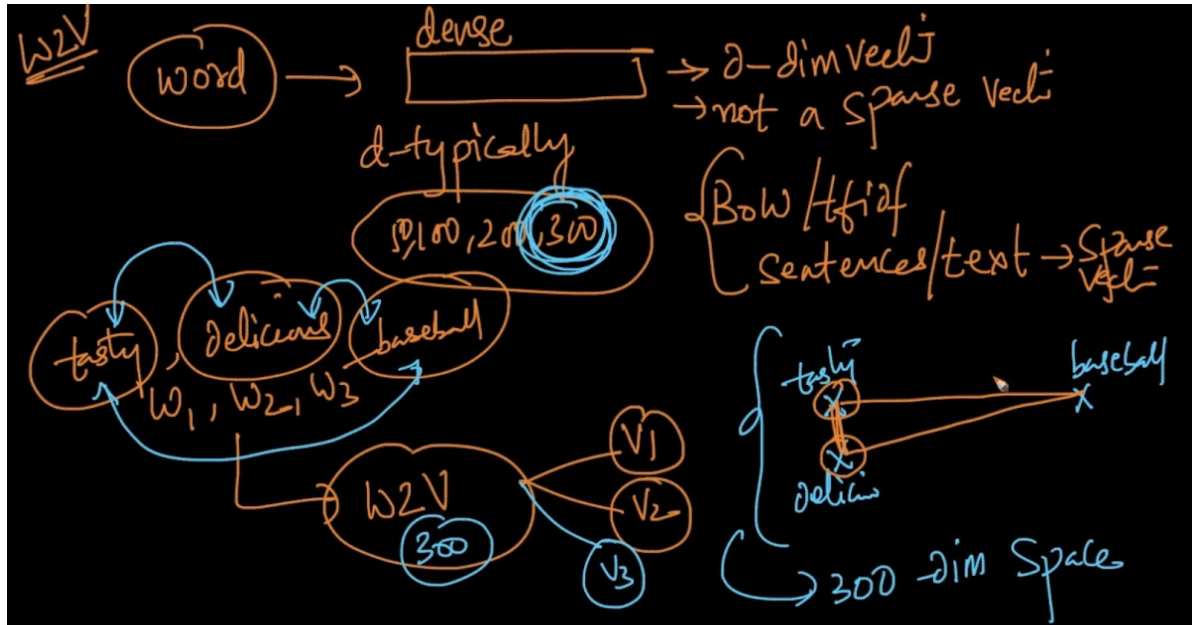
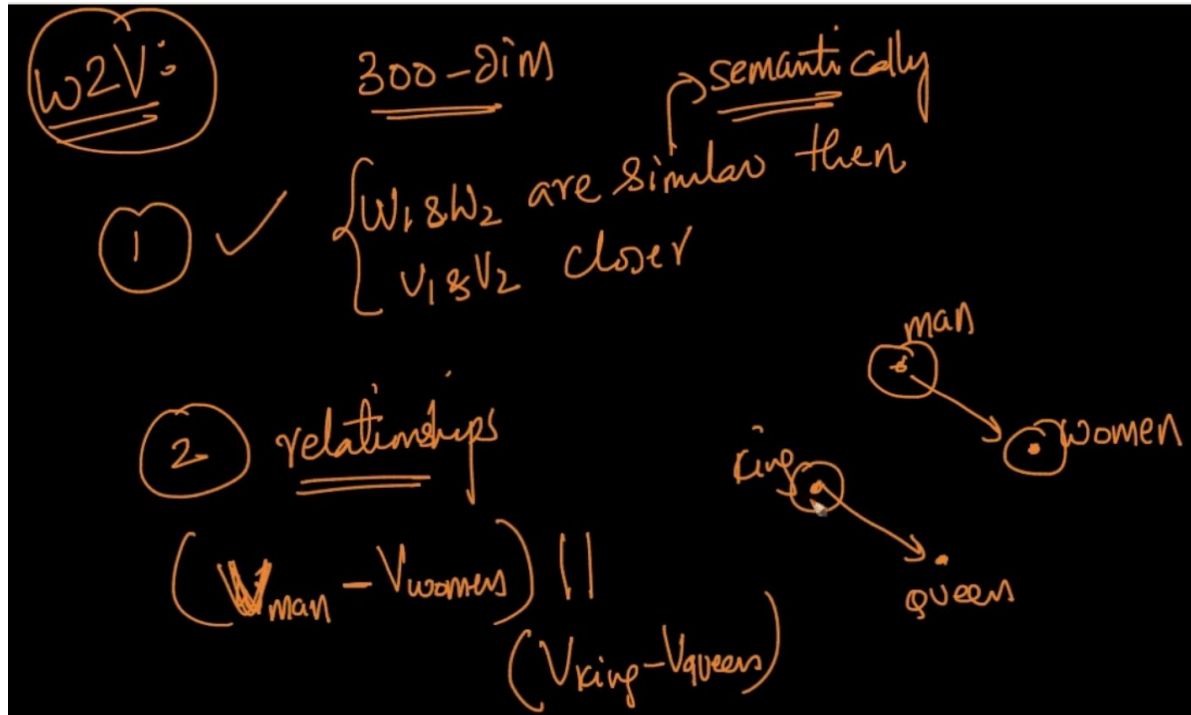Let's get a more practical solution to this

# WORD2VEC

In word2vec, we try to get the semantic meaning of the words. As the name suggests it is conversion of a word to a vector.
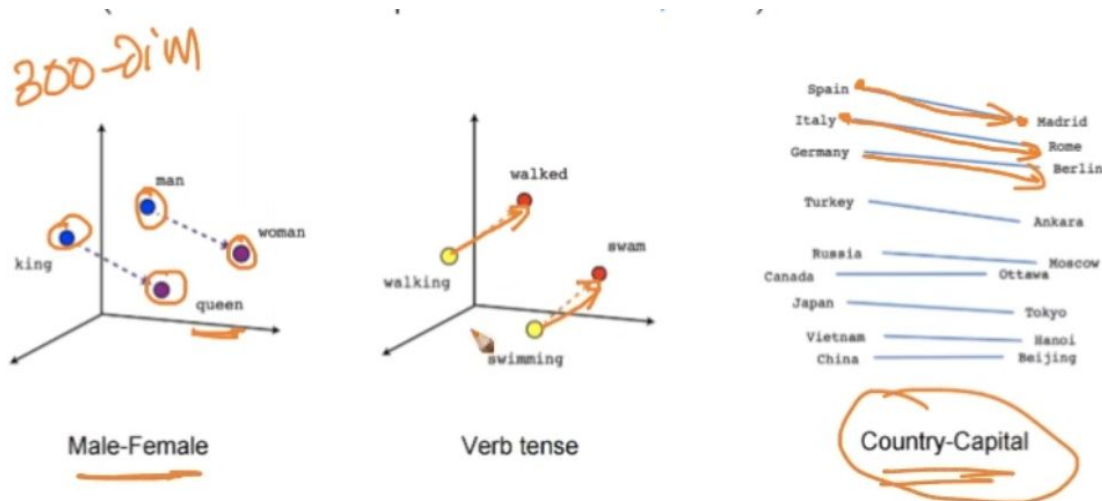


In BOW/tfidf we got a sparse vector of a text but in word2vec it is not sparse. It takes a word and tells us how the dense the word is.
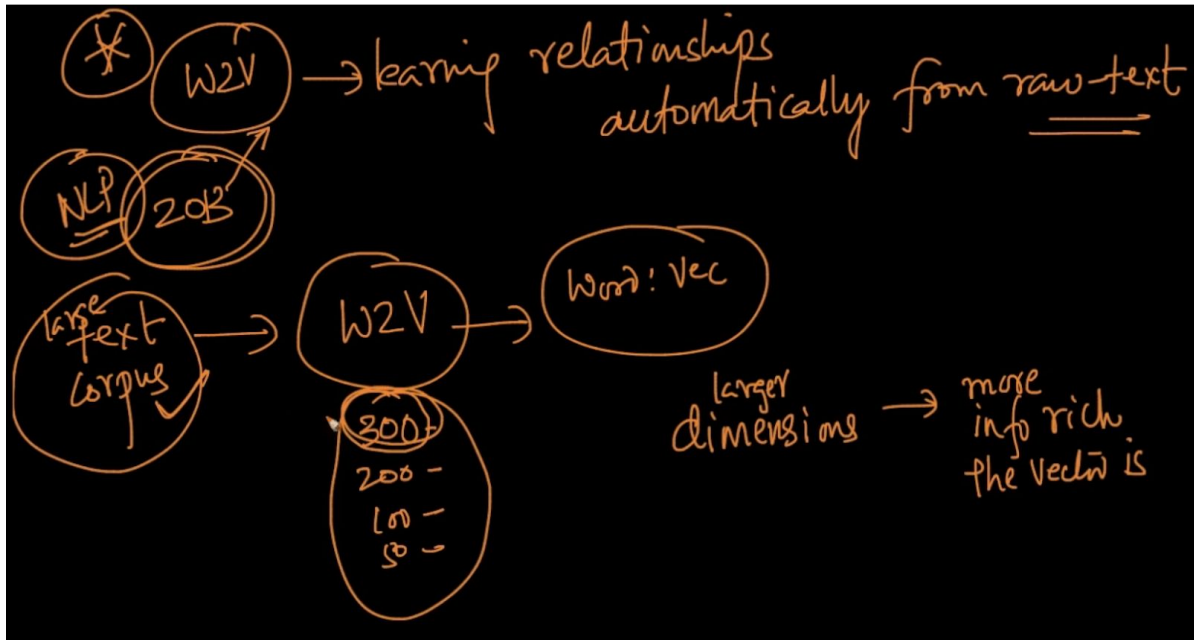
From the above example. We took 3 words *tasty, delicious , baseball* . Suppose word2vec is creating a 300-dimensional vector. So in that vector space *tasty* and *delicious* are close whereas baseball is far away .So word2vec is trying to capture the semantic essence of the words.
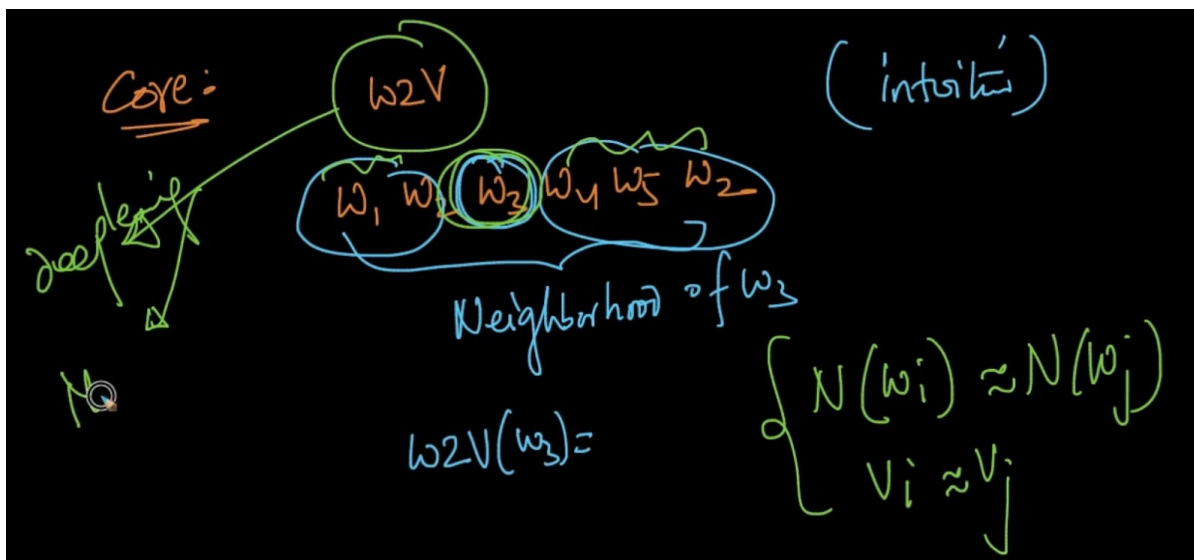
1st point is very clear and in the 2nd point it tells us about the **relationships** of the word. For ex: The words "men" and 'woman' and 'king' and 'queen' the vectors will be parallel as seen above. SO our word2vec is somehow trying to understand the difference between male and female.



It can capture other relationships as well like country-capital . As seen the vectors for country-capital are almost parallel.
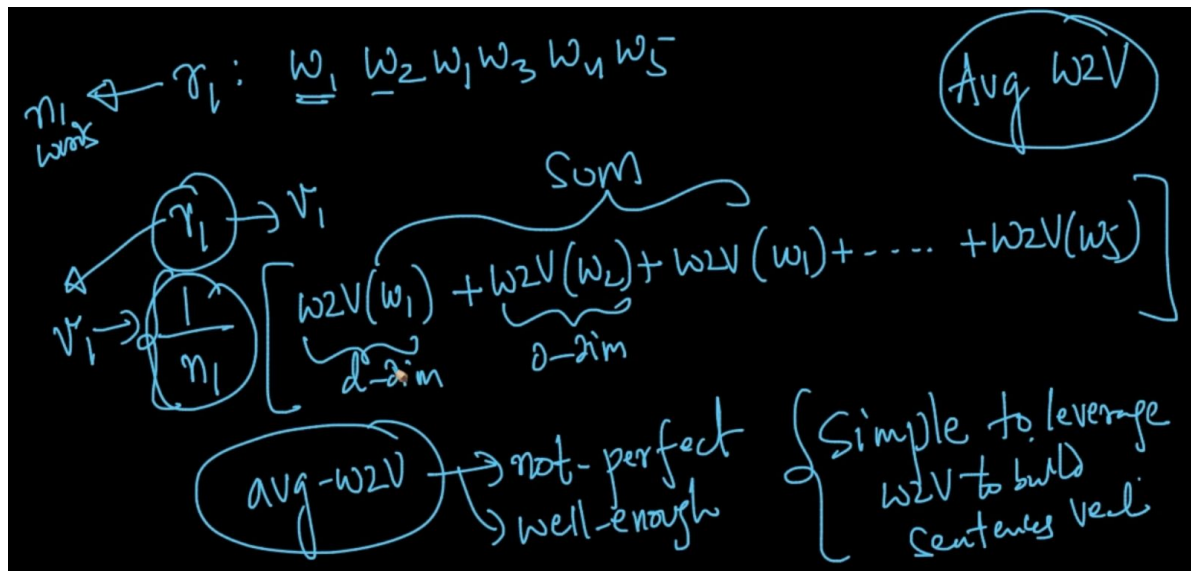
Word2vec automatically learns relationships from text. We give a large text corpus and if large text corpus is given more dimensions will be created and the larger dimensions are the more information vectors we can get .
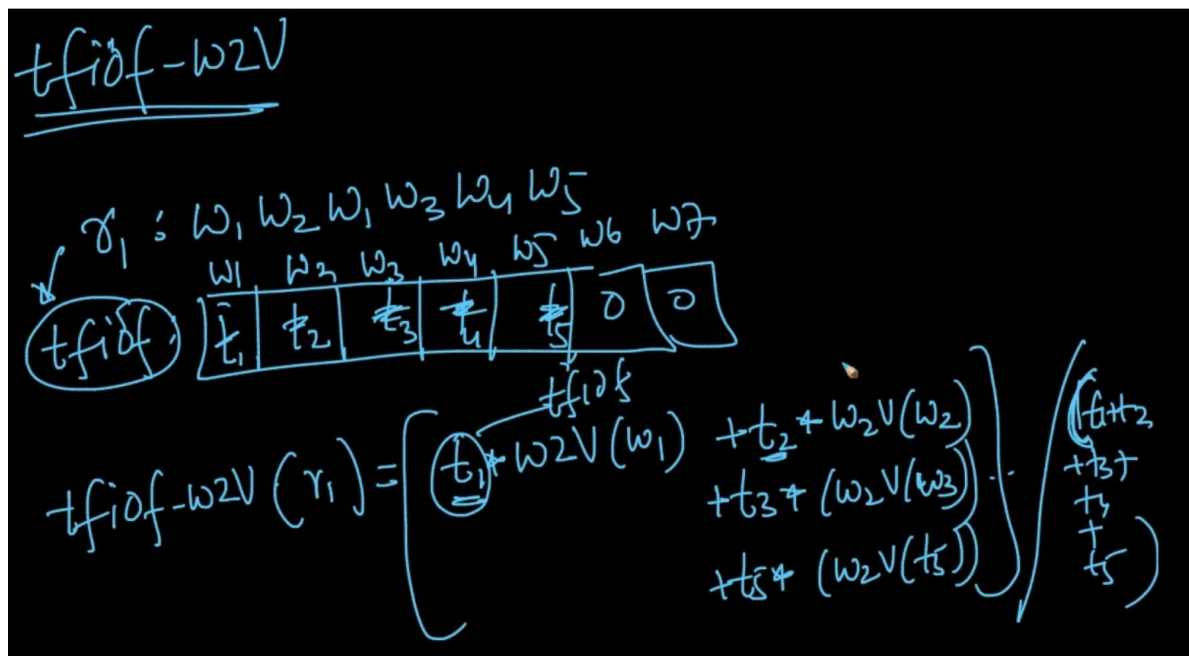


We are taking a word W3 . If we are taking word2vec(W3) it will try to get the surrounding words known as neighborhood. So if the neighborhoods are similar then vectors will be similar as well as seen

# Avg-Word2Vec, tf-IDF weighted Word2Vec



Suppose we want a vector of a review we take the w2v(W) of all the words and then take out their averages as seen above