Amazon Fine Food reviews:
Bag of Words:



Each of the review is called as document in NLP.
Constructing a vector:
A vector of length of each words that is 'd' unique words if constructed for each review. The frequency of each word in the review is counted. The vector will be highly sparse.
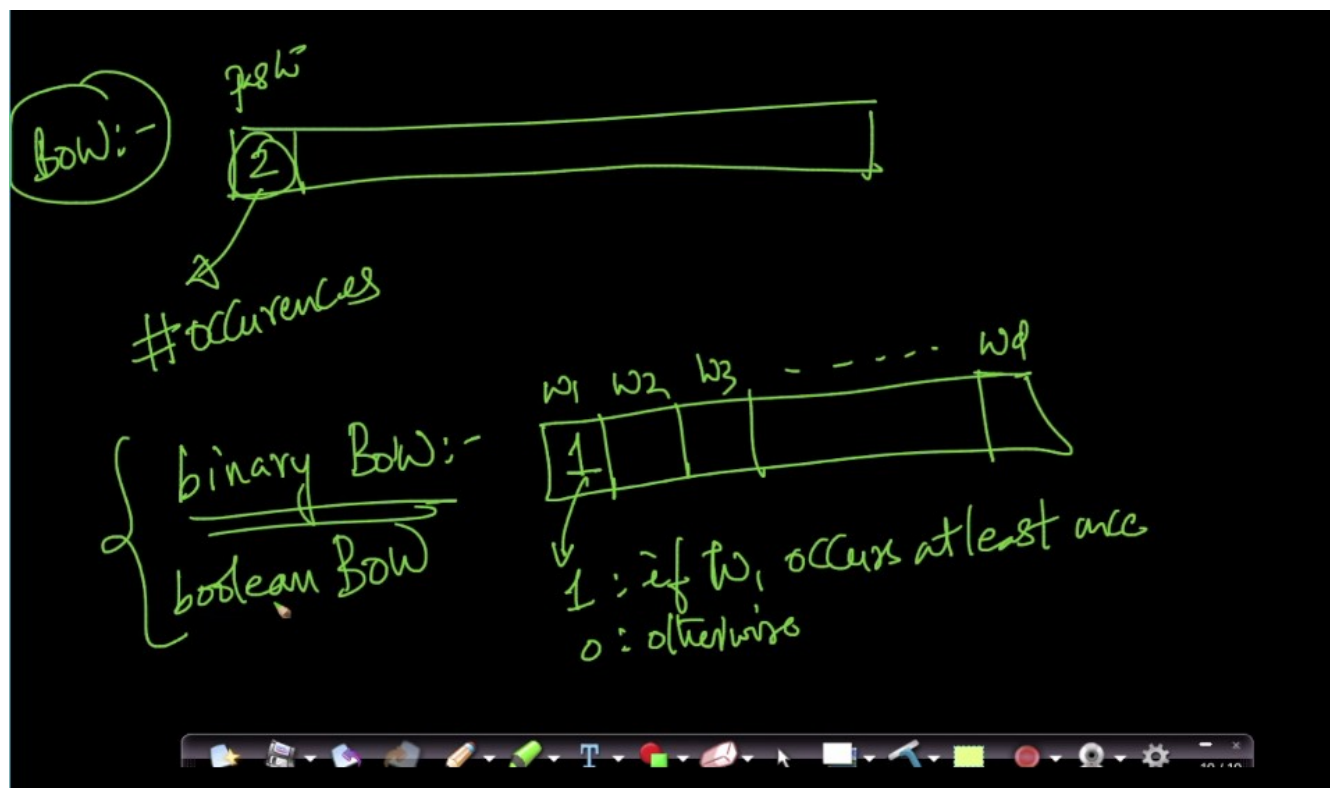
Our objective is similar text must be more similar. In BOW the opposite sentences may be at the closer distance, because it only counts the occurrence of the words in the sentence.
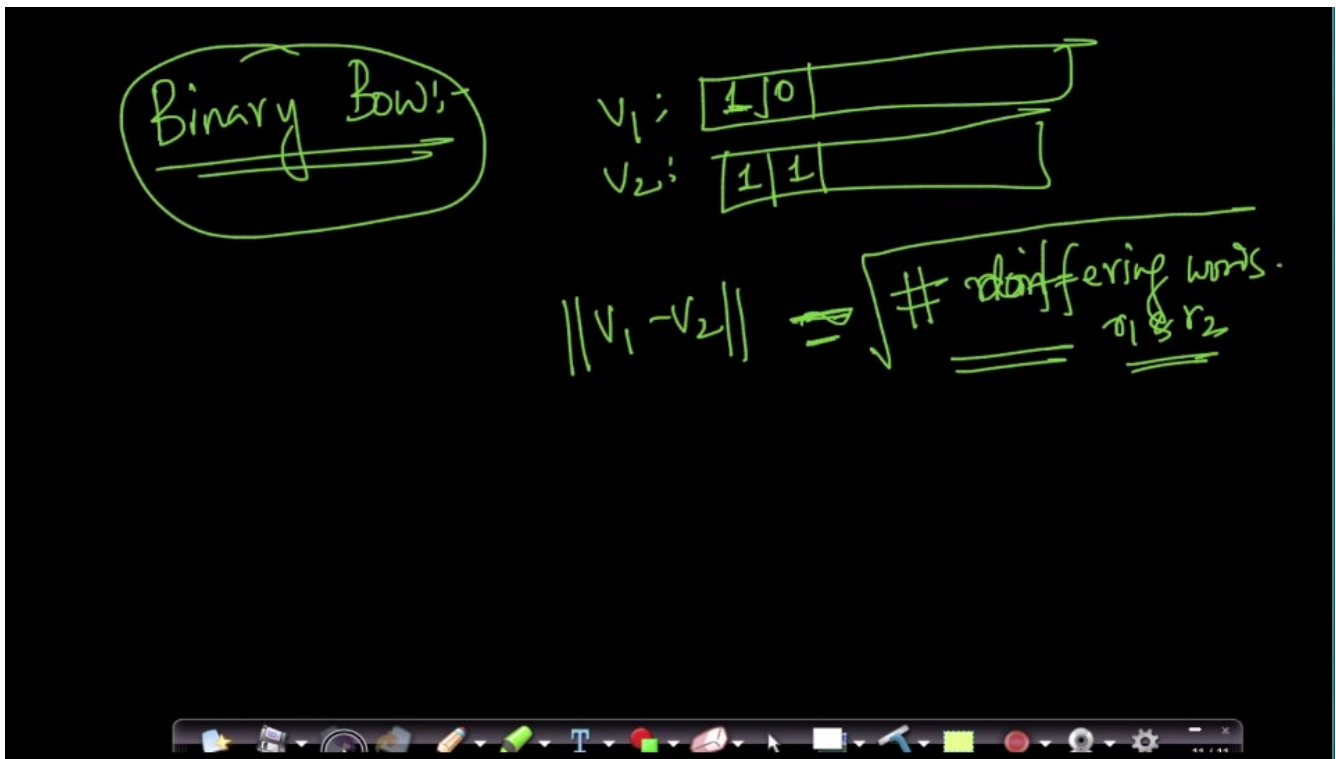
Example:

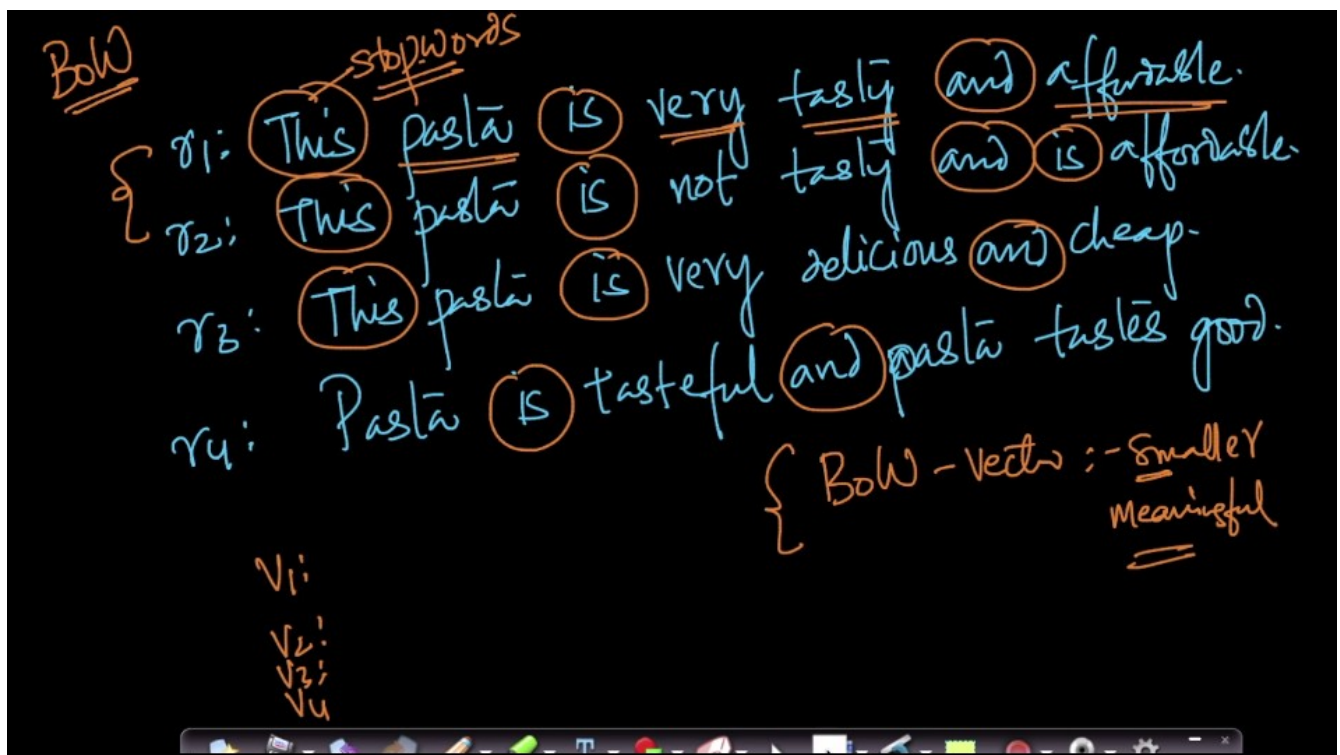

Binary bag of words: The word occurs or not in the vector.

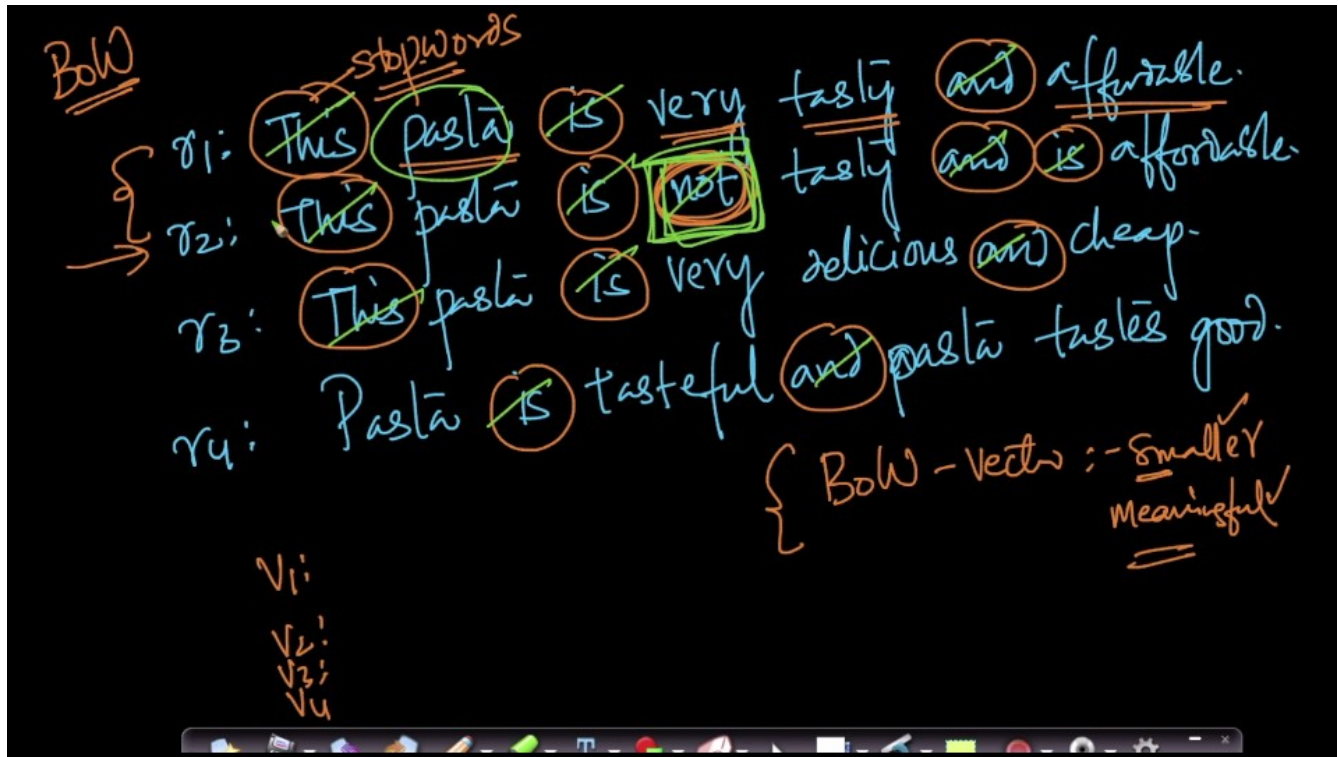In BOW the distance between the vectors v1 and v2 is almost equal to number of differing words.



The stop words are removed because they does not give any trivial meaning to the review.

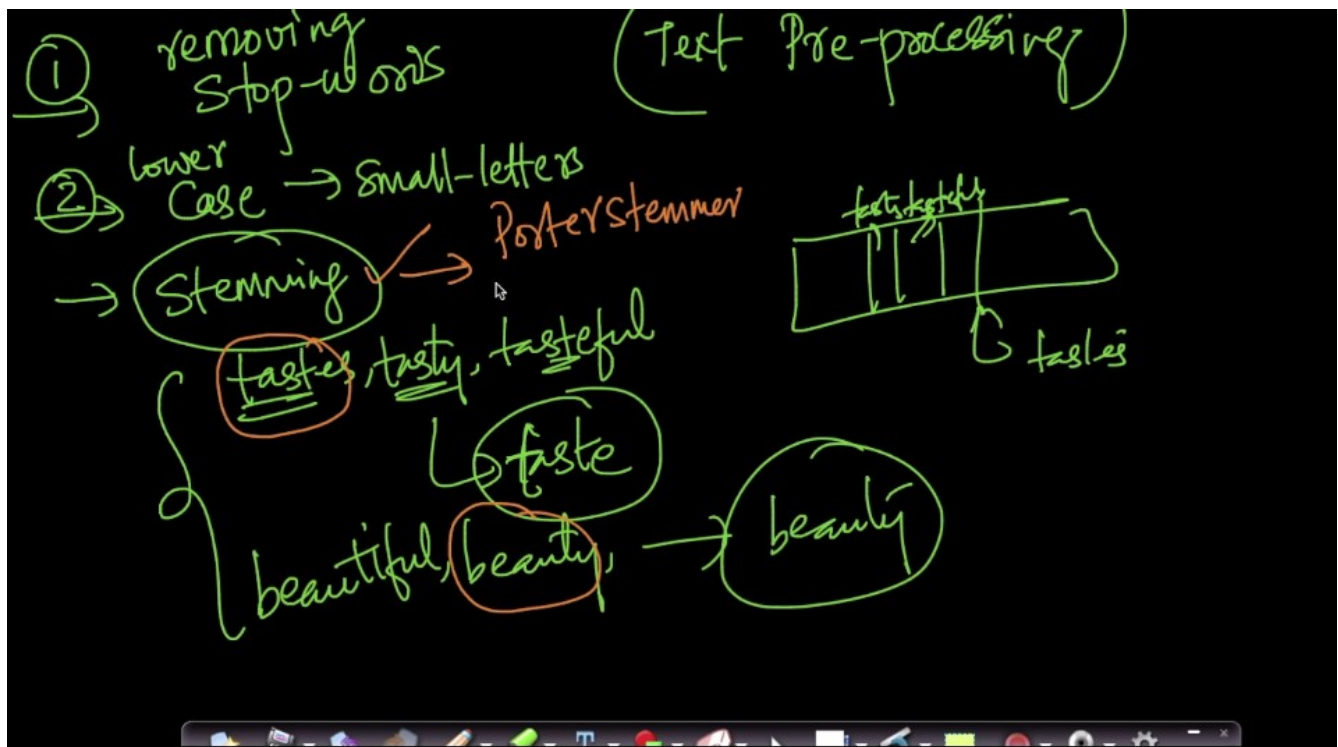Text Preprocessing- Stemming, Stop-word removal, Tokenization, Lemmatization:
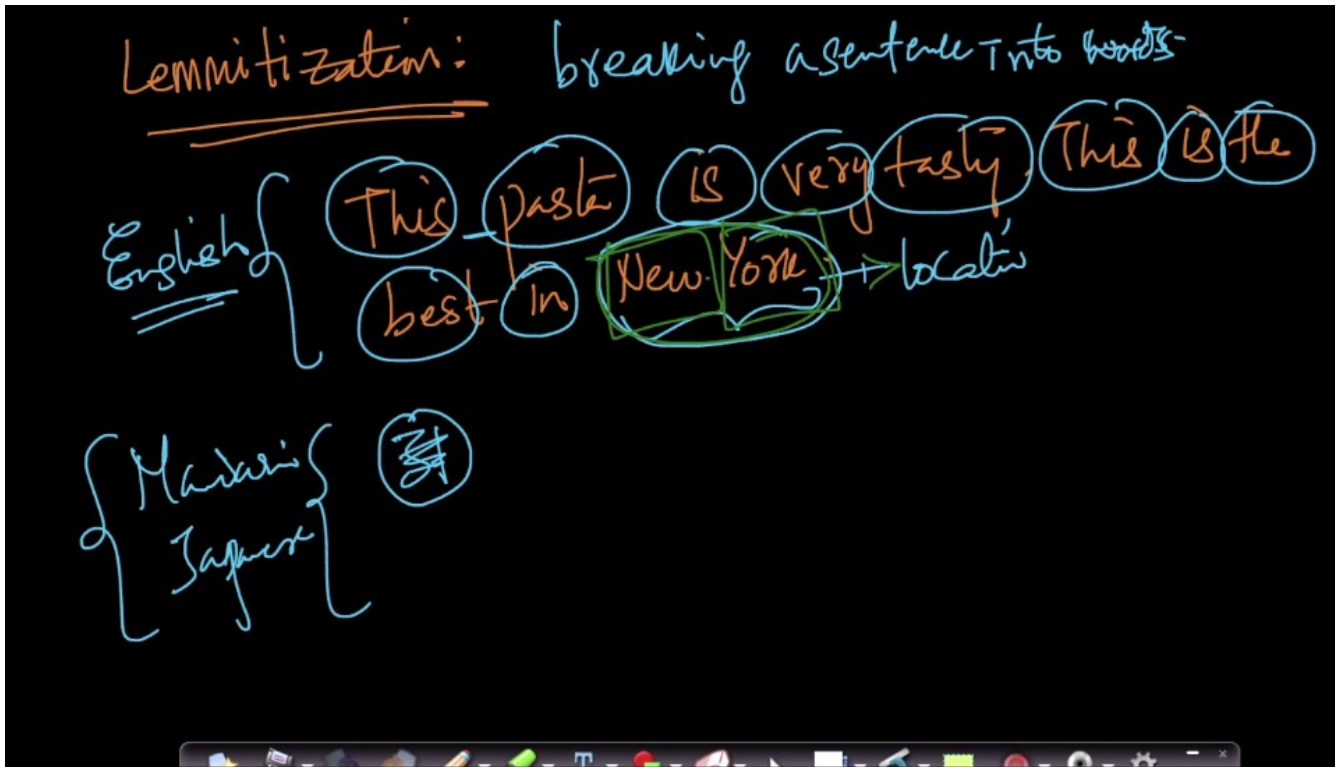
Removing stop words:



To make everything lower case: Removing the words that are in the uppercase.
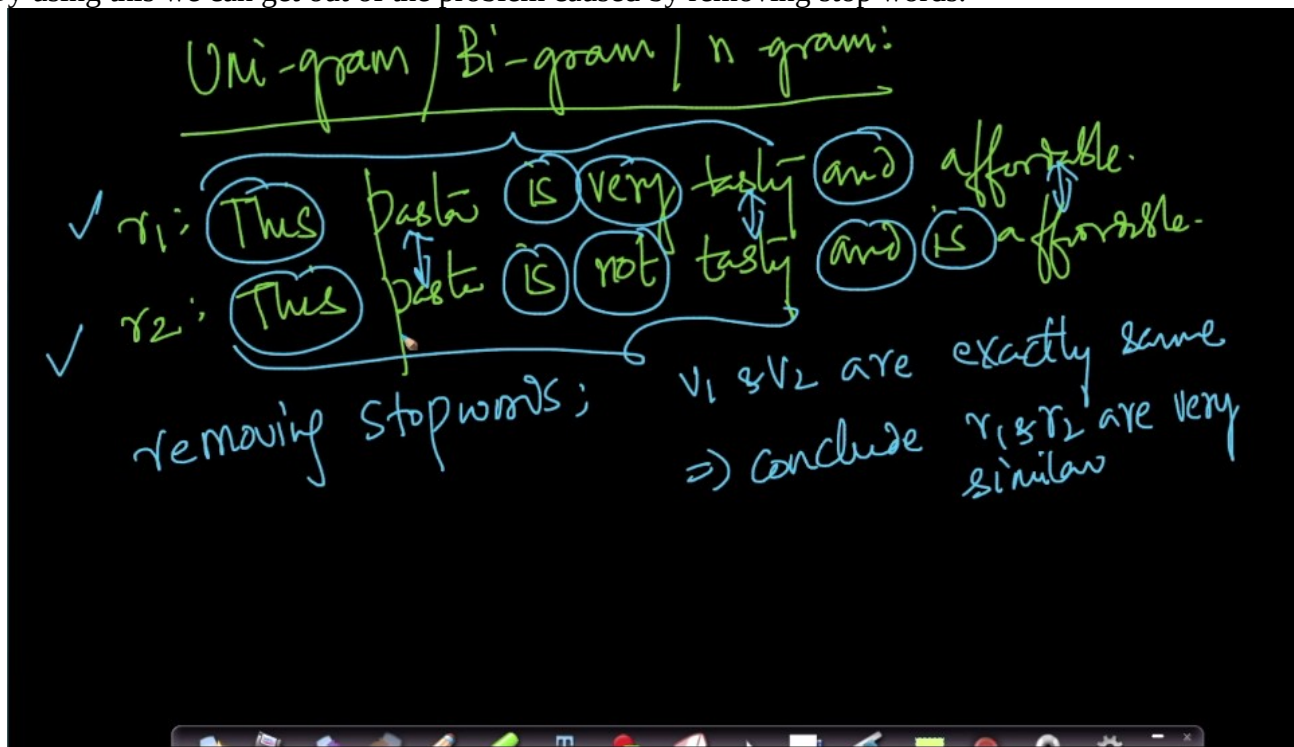
Stemming:

Lemmatization:

It is all about how to break the words, like New York.



Semantic meaning of words: We will solve this problem by using Word2Vec.
Uni – gram, Bi – gram, N – gram:
By using this we can get out of the problem caused by removing stop words.

Bi – gram: In this we use pair of words in BOW.



Tri- grams: 3 Consecutive words.

The number of bi grams are greater than equal to uni grams, the dimensions 'd' increases drastically.



Tf – Idf:

Term frequency(tf): (The number of times a word occurs)  / (total number of words in the review)

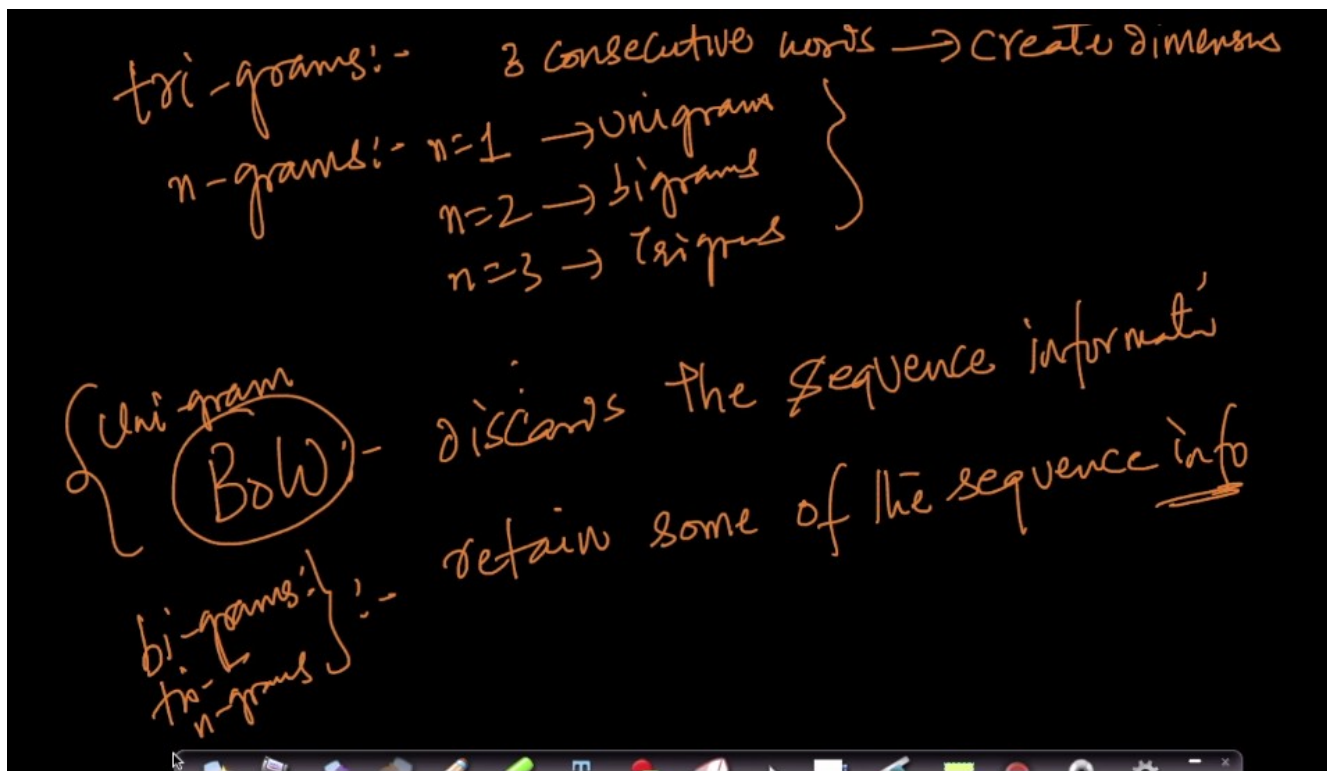The term frequency of every word is in between 0 and 1.



It can be taught of probability of occurring of the word in the review.

IDF(inverse document frequency): The IDF of the word in the document is calculated as the log(#docs / #docs word occurred).

IDF is always greater than equal to 1.

$$IDF(W_i, \mathcal{D}_c) = \log\left(\frac{N}{n_i}\right)$$

$N \to$ #docs

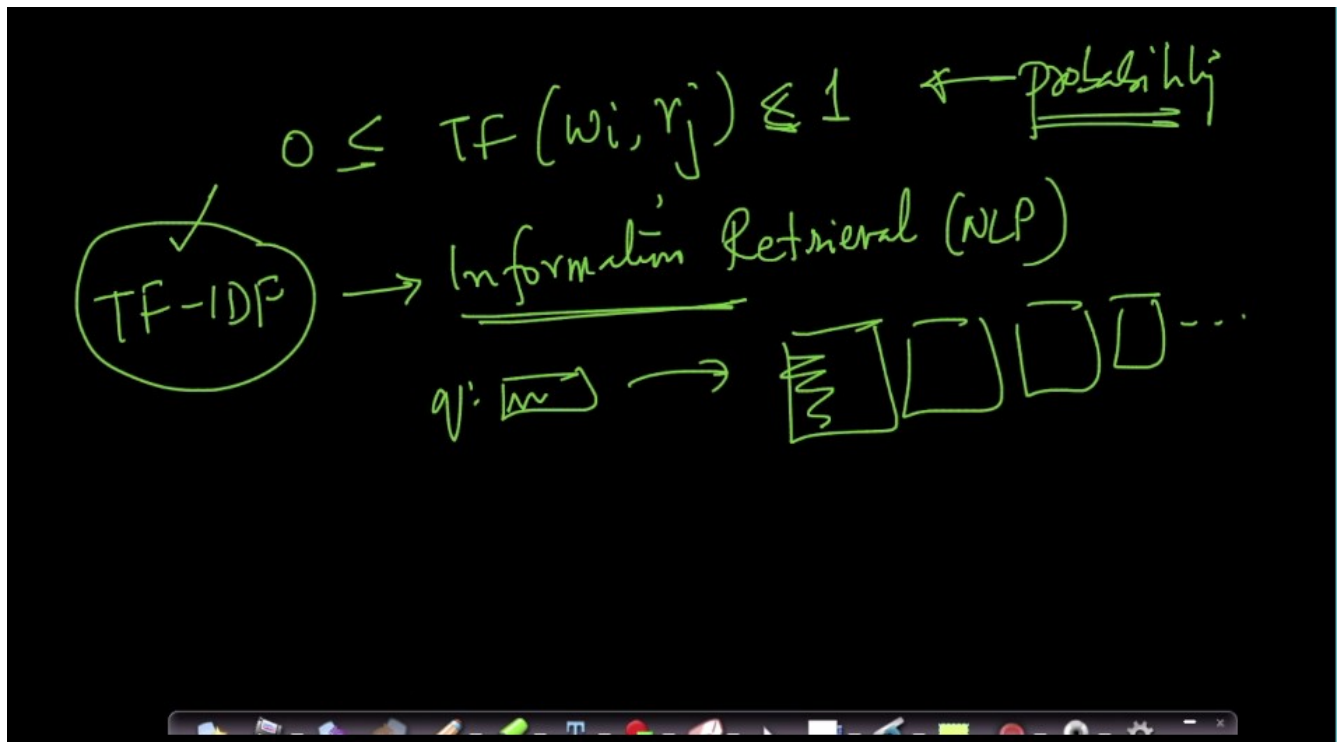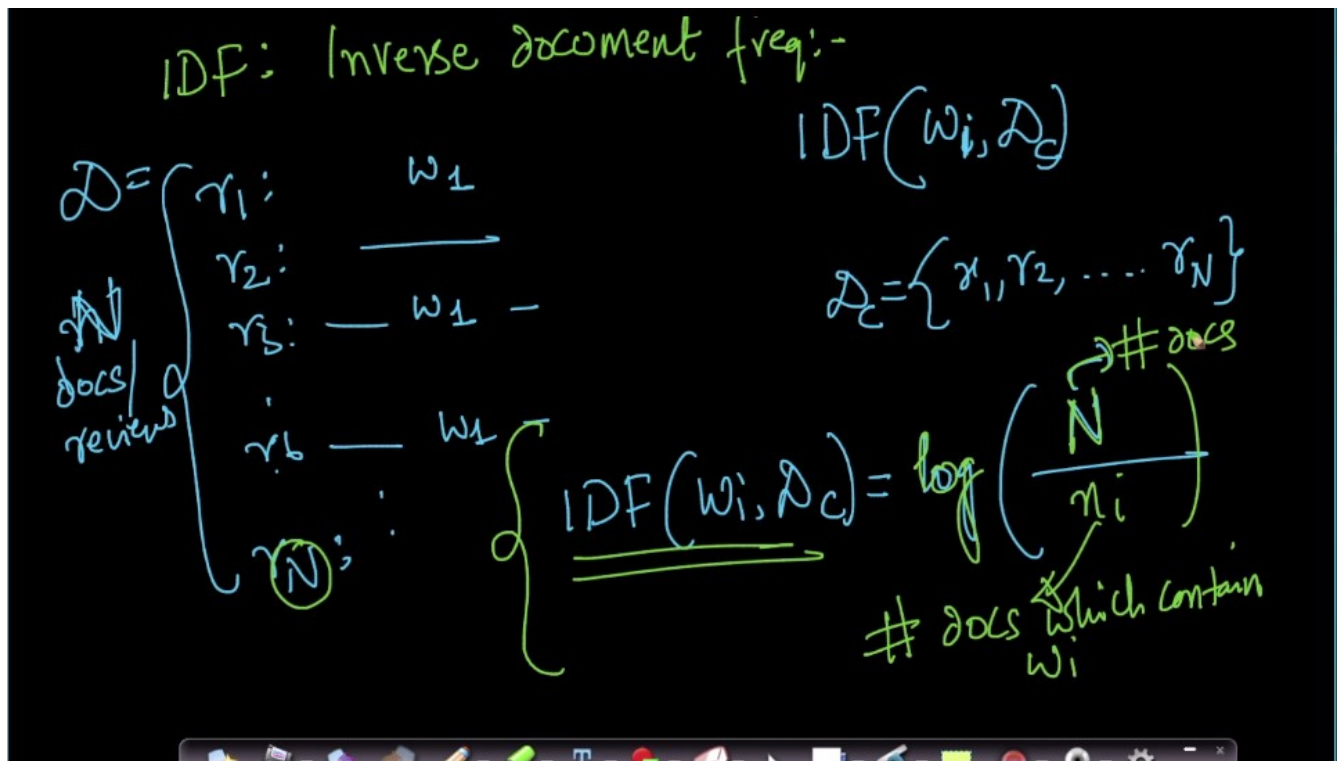$n_i \to$ #docs, containing $W_i$

$n_i \le N \Rightarrow \frac{N}{n_i} > 1$

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

$\log(1) = 0$

Properties of IDF:

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

$IDF \ge 0$

① if $n_i \uparrow$ ; $\frac{N}{n_i} \downarrow$ ; $\log\left(\frac{N}{n_i}\right) \downarrow$ ②

if $W_i$ is more freq in my Corpus then IDF $\downarrow$

$\frac{1000}{10} >$ ; $\frac{1000}{20}$ ; monotonic fn ③

IDF $\propto \frac{1}{n_i}$

IDF $\downarrow$  $n_i \uparrow$

Since, log is a monotonically increasing function.

Calculation of TF-IDF:

TF says how often the word occurs in the document / review.

$$\left(\gamma_1, \gamma_2, \ldots \gamma_N\right) = \partial c$$

$$\gamma_i \longrightarrow V_i$$

$$w'_j$$

$$TF\left(w_j, \gamma_i\right) * IDF\left(w_j, \partial c\right)$$

$$\left(w_j \text{ is frequent in } \gamma_i\right) \quad \left(w_j \text{ is rare in } \partial c\right)$$

TF – IDF:

$$\boxed{TF - IDF}$$

↳ more importance to rarer words in my $\partial c$

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

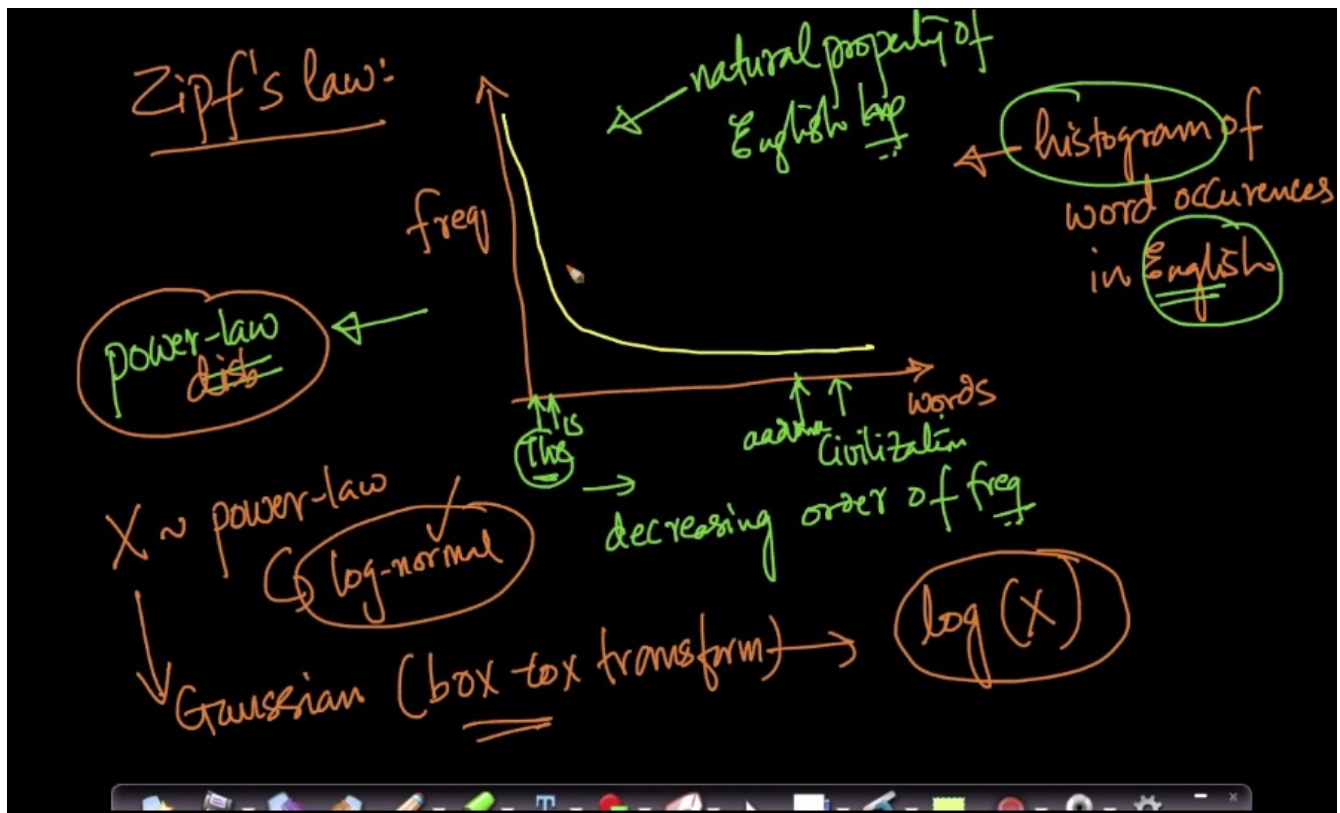→ Semantic – meaning $\left(\begin{array}{l} tasty \to delicious \\ cheap \to affordable \end{array}\right)$
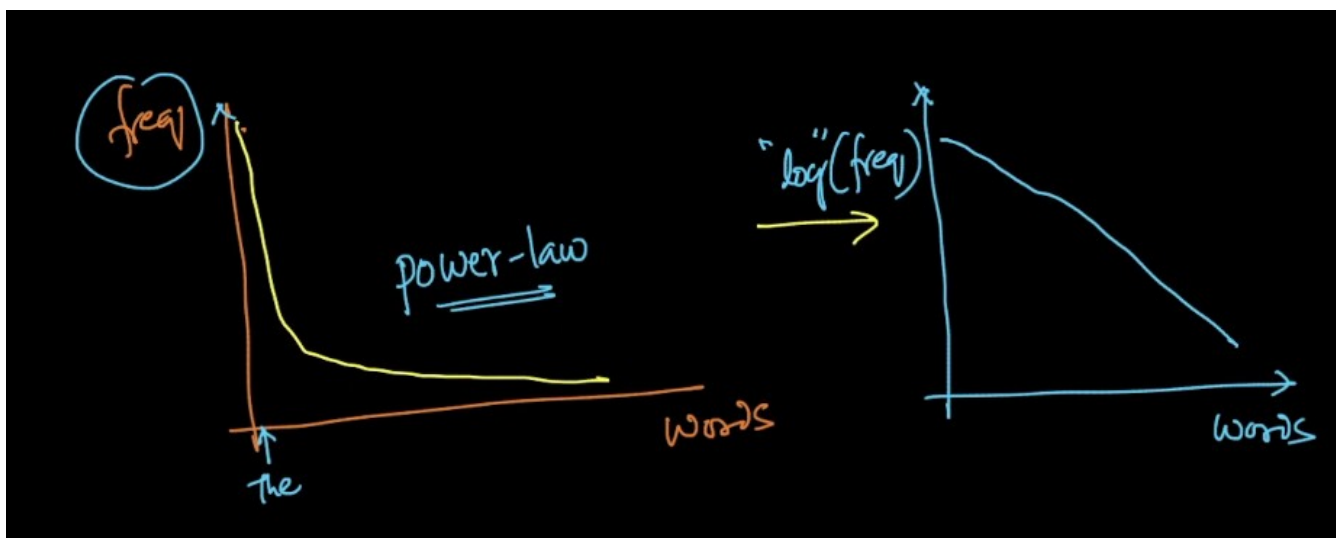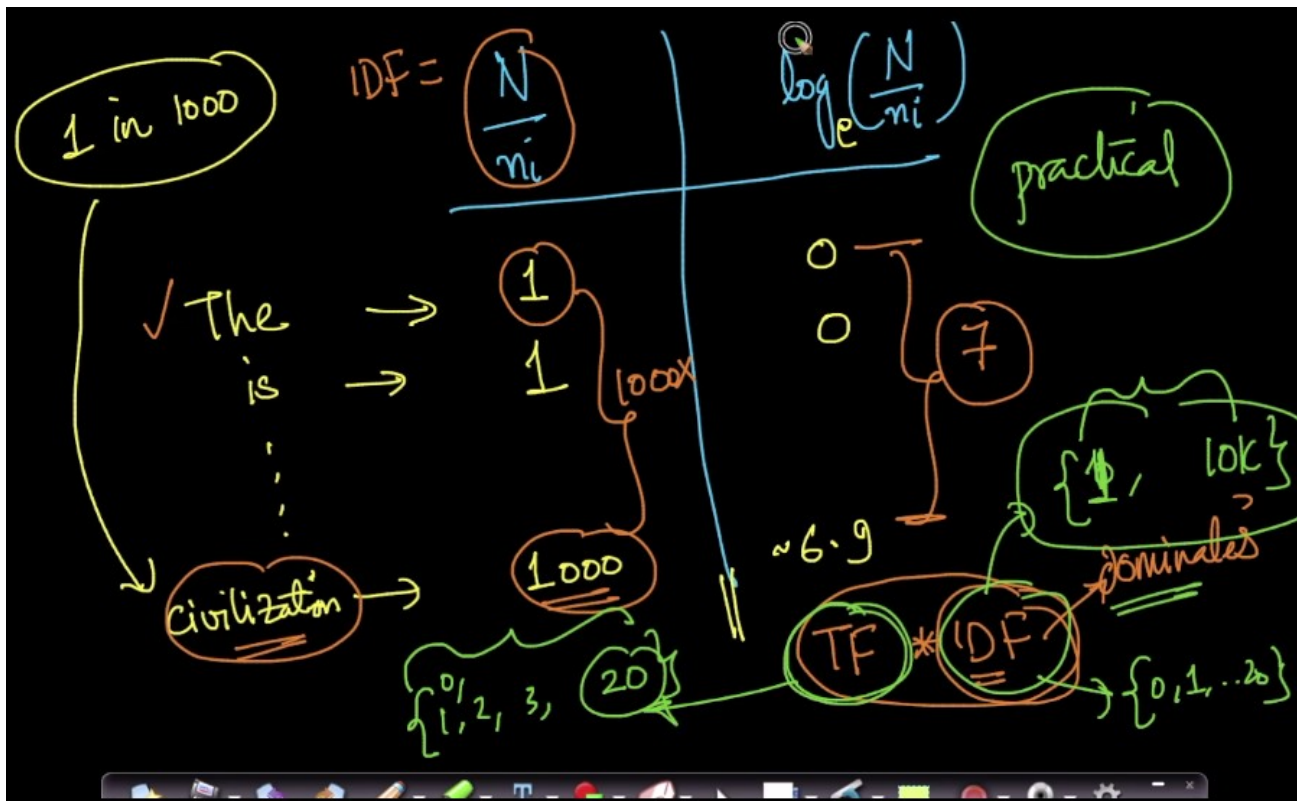
Why use log in IDF:

Log is given on the name of Zipf's Law:



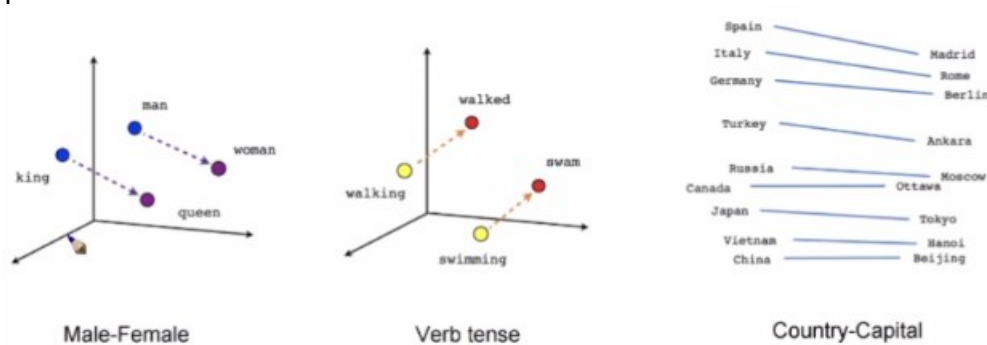Zipf's law says that the words in english follow power law distribution.

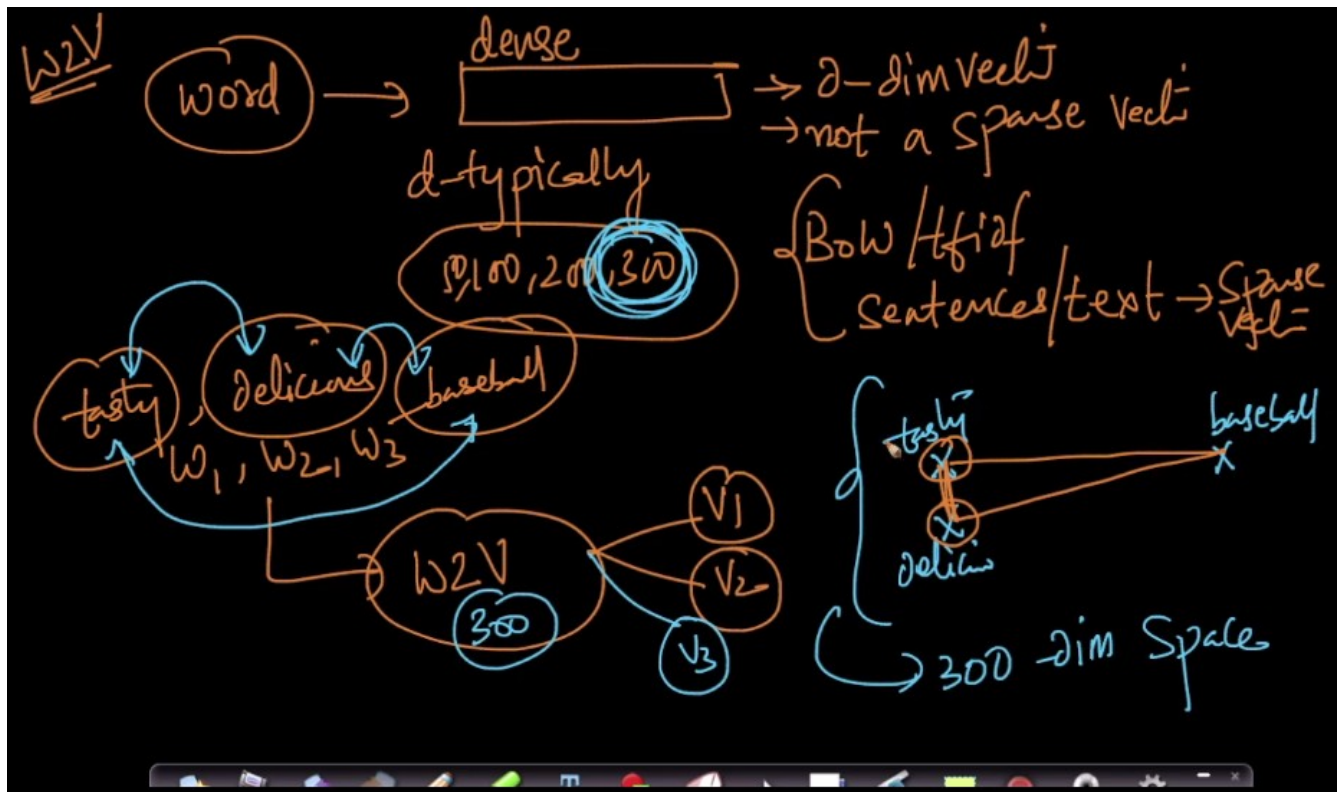By applying log the range of values differing will decrease.



Word2Vec:

      It takes the semantic meaning of the words and relationships of the data.
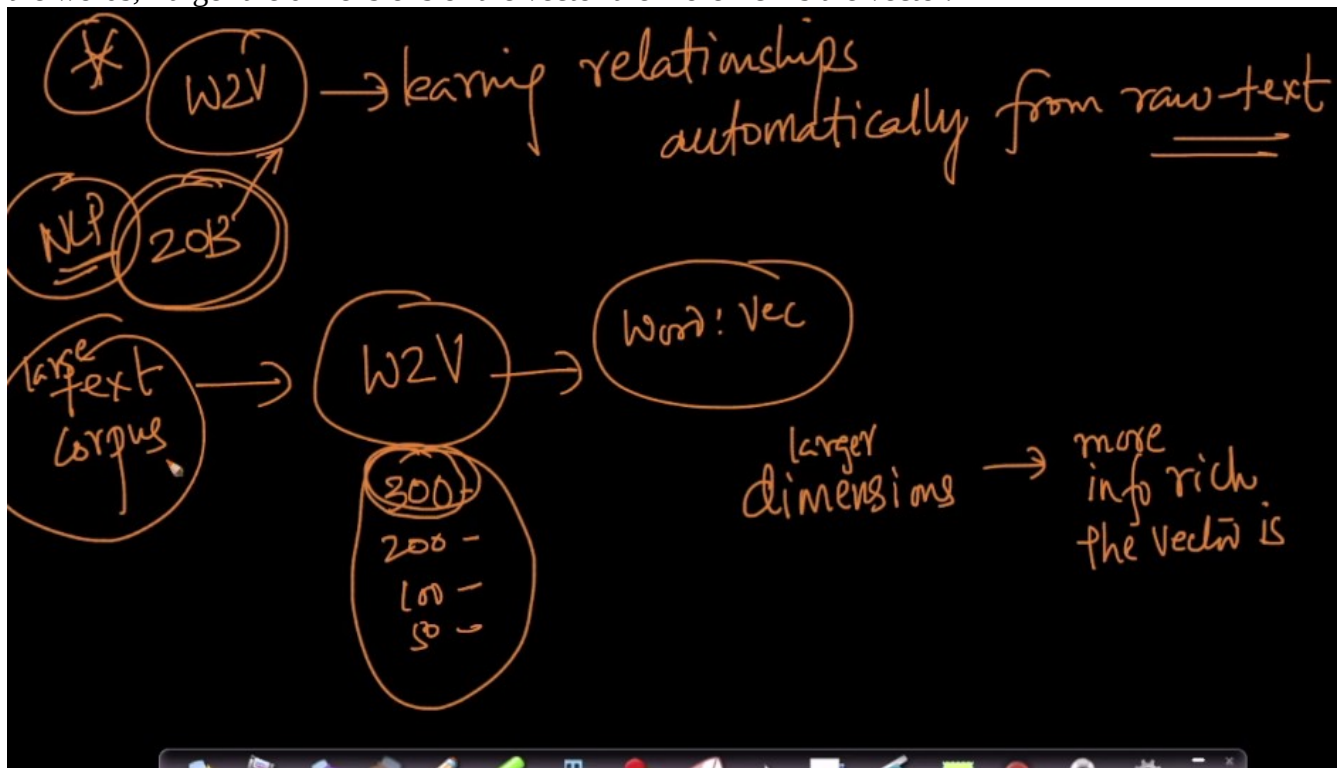
Word2Vec representation:



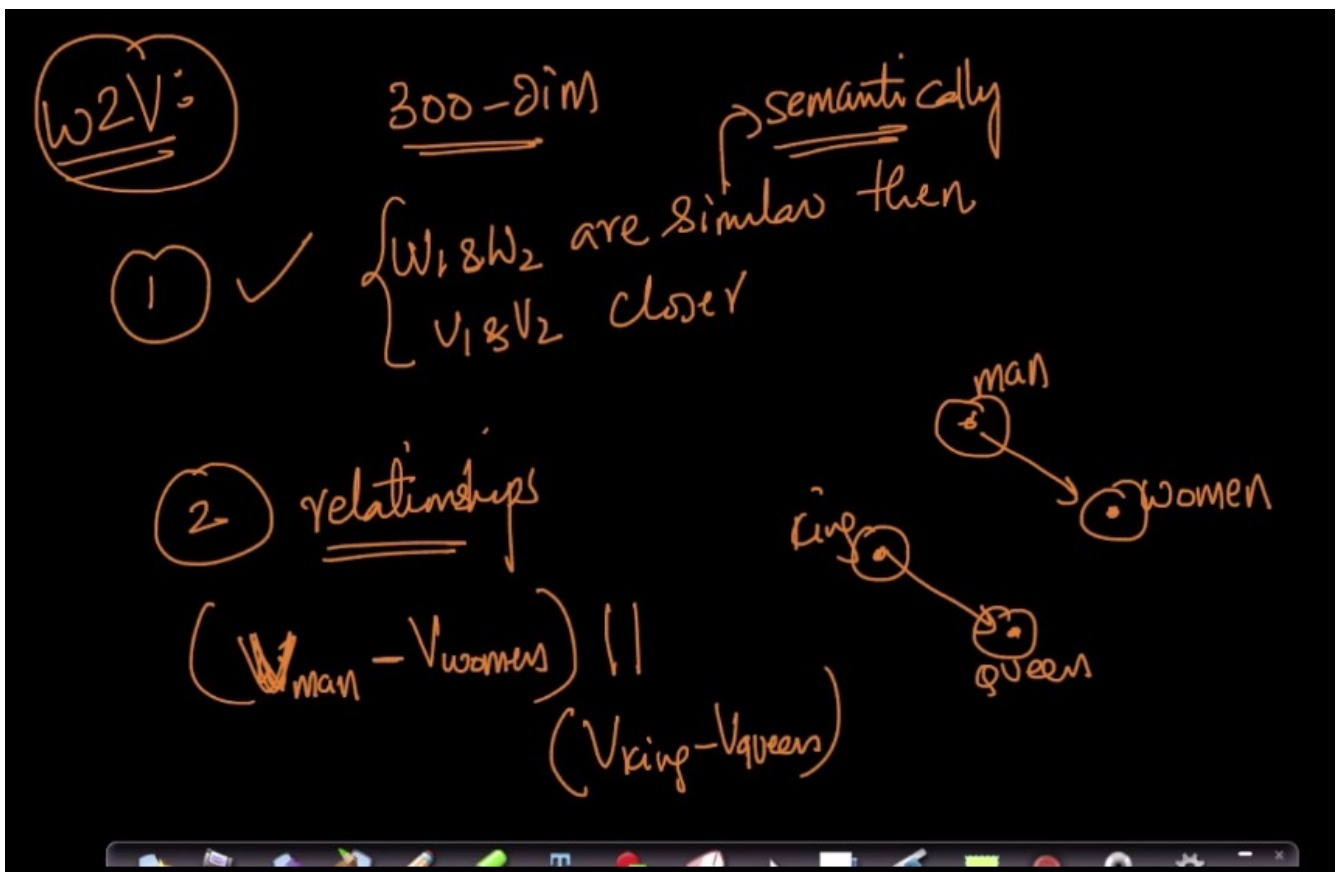Male-Female             Verb tense             Country-Capital

W2V represents the word in the form of dense values:



W2V: It also understands the relationships of the words. It learns all the internal relationship between the words, Larger the dimensions of the vector the more rich is the vector.
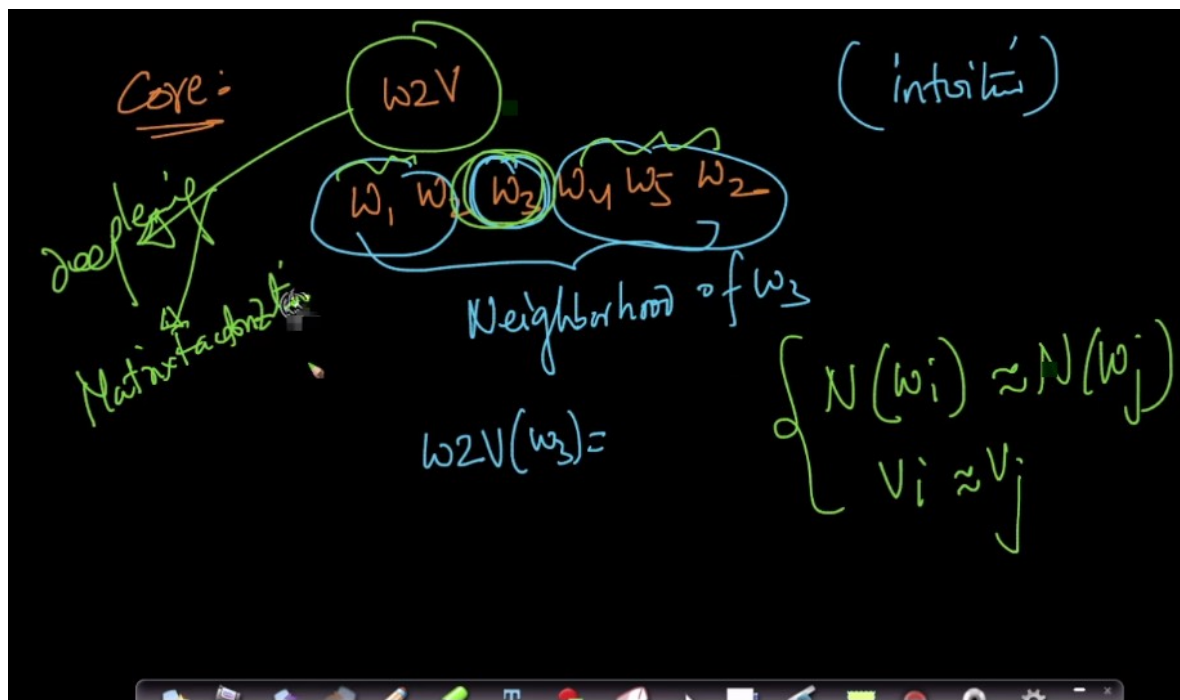
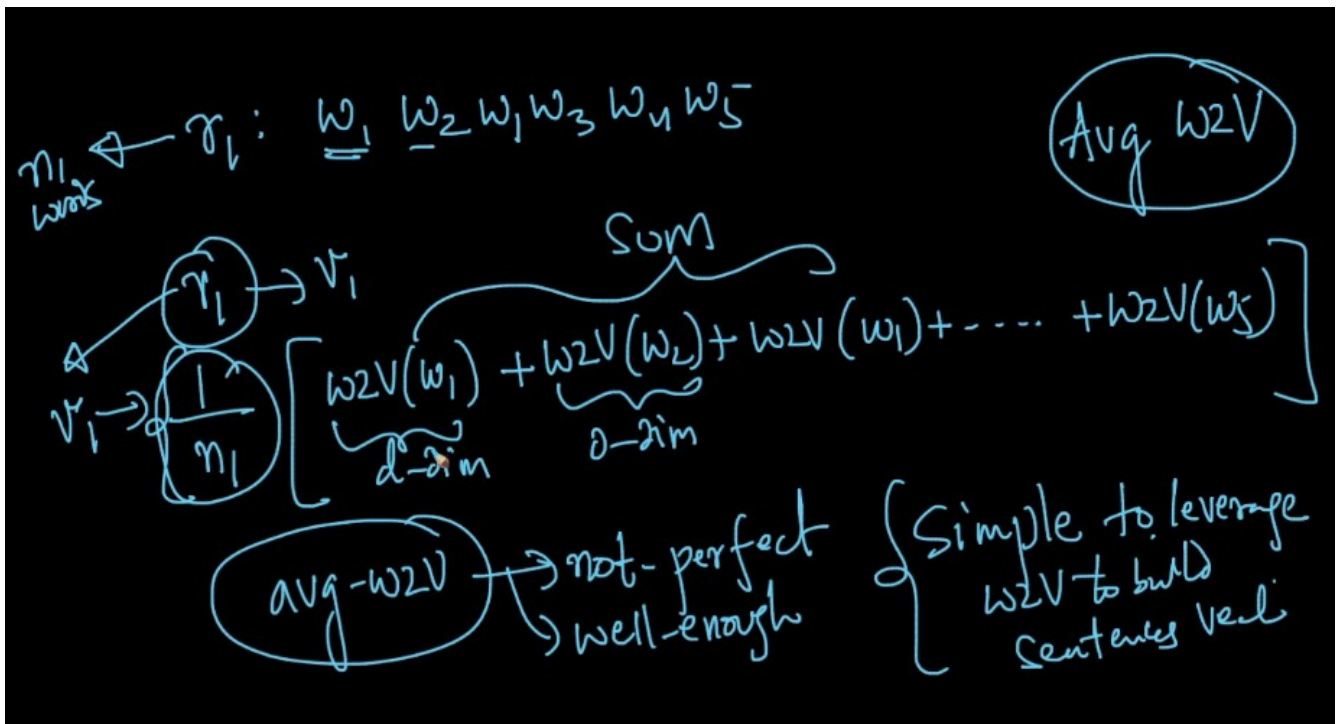If the data corpus size is large, larger the dimensions of the vector.
Google News W2V has the 300 dimensional vector.
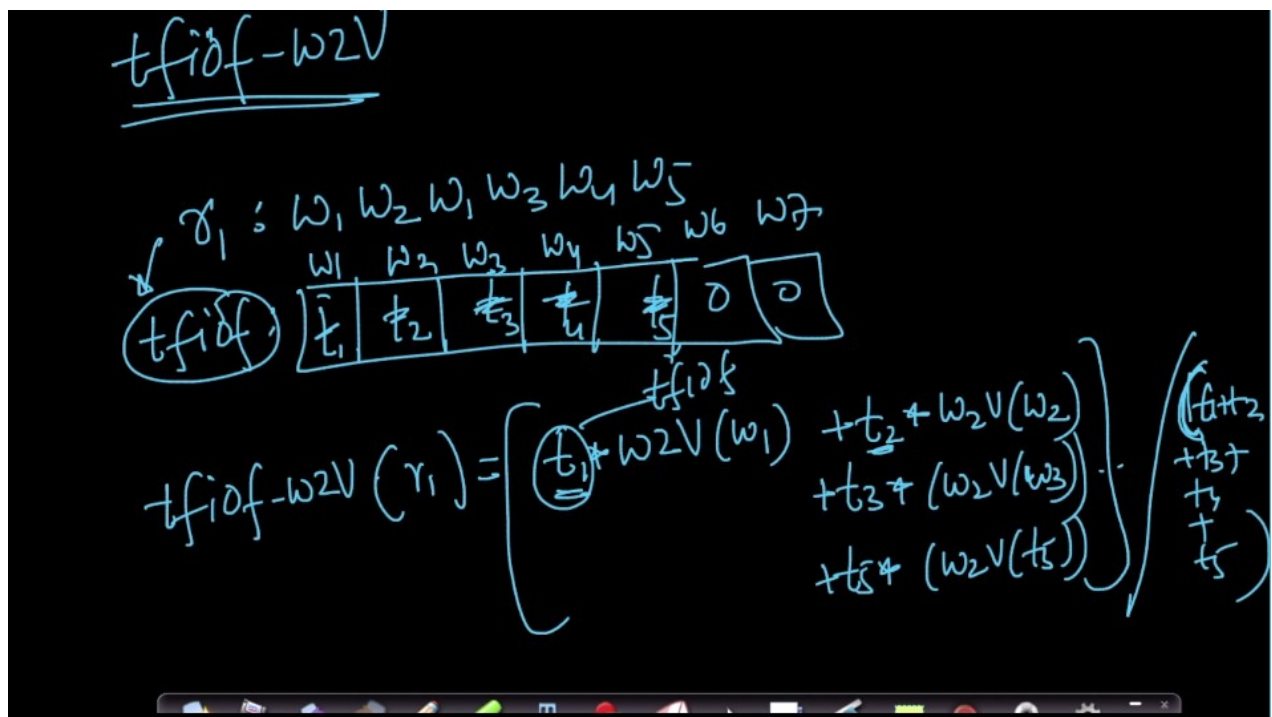W2V also looks at the neighbors of the word in the sentence.

Average – W2V, tf-idf weighted W2V:

Compute the w2v of each of the words and add the vectors of each words of the sentence and divide the vector with the number of words of the sentence.



TfIdf – w2v:
The W2V value of each of the word is multiplied by the tfidf value and summed up and then divided by the sum of the tfidf values of the sentence.

If all the tfidf values of the words in the review are equal to 1, then it is equal to avg-w2v.

$$tfidf\text{-}w2v(\hat{r_i}) = \frac{\sum\limits_{i:wrs}\left(t_i * w2v(w_i)\right)}{\sum\limits_{i:wrs} t_i} = tfidf(w_i, r_i)$$

If all of $t_i = 1$

$$tfidf\text{-}w2v \rightarrow avg\ w2v$$