

### Q stop words:-

stop words are words which are not contain important significance to be used in search queries, usually these words are filtered out from search queries because they return vast amount of unnecessary information.

### Removal of stop words and punctuation:-

Doc1:- After removing stop words and punctuation the output will be as follows.

o/p:- researchers focus computational phenotyping produce disease prediction models machine learning statistical tools.

Doc2:- researchers develop tools bayesian statistical information generate ~~casual~~ models large

complex phenotyping datasets

Doc3:- researchers build computational information engine uses machine learning combine gene function gene interaction information distillate genetic data sources

## N-gram:

N-gram is a contiguous sequence of  $n$  items from a given sequence or dataset of speech. The item can be phonemes, syllables, letters, words or base pairs according to the application. The  $n$ -grams typically are collected from a text or speech corpus. When the item is words,  $n$ -grams may also be called simply

An  $n$ -gram of size 1 is referred to as a 'unigram', size-2 is a bigram, size-3 is a trigram. Larger sizes are sometimes referred to by the value of  $n$  in modern language. eg: four-gram, "five-gram" and so on.

Doc of: After applying  $n$ -gram the result of doc is as follows:-

The researcher will  
researcher will focus  
will focus on  
focus on Computational  
on Computational Phonology  
Phonology and will  
and will produce  
...

The output for N=3 along removal stopword for Doc 1  
is :-

researching focus computational N=3  
focus computational phenotyping  
computational phenotyping produce  
phenotyping produce disease  
disease prediction model  
prediction model machine

DOC 2 :-

researching develop tool  
develop tool Bayesian  
tool Bayesian statistical } N=3

DOC 3 :-

build computational information  
computational information engine  
information engine uses. } N=3

2

list of itemy	from the	3	documenty	or show below
	<u>P1</u>	<u>P2</u>	<u>P3</u>	<u>Count in 3 docy</u>
<u>Researchy</u>	1	1	1	3
focus	1	0	0	1
proactive	1	0	0	1
modely	1	1	0	2
machine	1	0	1	2
learning	1	0	1	2
statistical	1	1	0	2
develop	0	1	0	1
Information	0	1	1	2
generate	0	1	0	1
casual	0	1	0	1
Complex	0	1	0	1
build	0	0	1	1
interaction	0	0	1	1
combine	0	0	1	1
function	0	0	1	1



b) TF :- It is the term frequency which measures how frequently a term occurs in a document. Since every document is different in length, it is possible that a term would appear more times in a long document.

$$TF(t) = \frac{\text{no of times term } t \text{ appears in a document}}{\text{total no of terms in the document}}$$

IDF :- which measures how important a term is, while computing TF, terms are considered equally important.

$$IDF(t) = \log_e \left( \frac{\text{total no of documents}}{\text{no of documents with term } t} \right)$$

TF-IDF :- TF-IDF stands for term frequency inverse document frequency and the TF-IDF weight is a weight often used in information retrieval and text mining.

now we are calculating TF-IDF valy for each term D:

for research

$$TF = \frac{1}{12} \quad IDF = \log_e(3/1)$$

$$TF-IDF = \frac{1}{12} \times \log_e(3/1) = 0.0146$$

for computation:

$$TF = \frac{1}{12} \quad IDF = \log_e(3/1) = 0.477$$

$$TF-IDF = \frac{1}{12} \times 0.477 = 0.039$$

for disease:

$$TF = \frac{1}{12} \quad IDF = \log_e(3/1) = 3$$

$$TF-IDF = \frac{1}{12} \times 0.477 = 0.39$$

for machine:

$$TF = \frac{1}{12} \quad IDF = \log_e(3/2) = 0.176$$

$$TF-IDF = \frac{1}{12} \times 0.176 = 0.0146$$

for learning:

$$TF = \frac{1}{12} \quad IDF = \log_e(3/2) = 0.176$$

$$TF-IDF = \frac{1}{12} \times 0.176 = 0.0146$$