# Assignment No. 4

* **Problem Statement :**
  Consider a suitable text dataset. Remove Stop words, apply Stemming & feature Selection techniques to represent documents as vector classify documents & evaluate Precision, recall.

* **Objective :** Learn how to tokenize & filter a document into its diff. words and then do word count for each word in doc. Apply stemming, feature Selection on doc (text)

* **outcome :** demonstrated text Processing using nltk, understood vectorizing & removing Stop words.

* **Theory :**
  ⊙ NLP is a Subfield of linguistic computer Science & AI Concerned with interaction b/w computer & Human language, in Particular how to Program computer to Process & Analyse large amt of Natural language data

  ⊙ In computing stop words are words that are filtered out before or after the natural language data ( text) are

Processed. while 'stop words' typically re[fer]
to the most common word, there is no
universal list of stop words.

① STEMMING ⇒ for grammitical reasons, do[c]
are going to use diff. forms of a word
(e.g. organized, organizing, etc). Additional[ly]
there are families of derivitavelly relati[on]
words with similar meaning like democ[ra]
democresy & democratization.
→ The goal of stemming (and lemmatization)
is to reduce inflectional forms & sometime[s]
derivationally linked forms of a word t[o]
its common base form
   am, are, is → be
   car, cars, cars', car's → car
→ when applied to a doc, result will be
somewhat like this
   the boy's cars are diff. color )
   the boy car be differ color

→ Stemming is a more crudde Process,
a heuristic that chops off the end of a
word in the hope of achieving the
goal correctly more of time, and ofter
includes removal of derivational affices

→ Feature selection is a Process of
selecting a subset of the terms occurning

in the remaining set and using only this subset as feature in text classification.
→ It serves 2 porpuses. first, it makes training & applying a classifion more efficient by decreasing the size of vocab. this is of Particular importance to classifion that are expensive to train.
→ Second, feature selection often increases classification accuracy by eleriating noise featre these features when added to docs representati increases the classification error on new data.

→ Vectorization is a Process of converting text into machine redable form. words are represented as 'Vectors' (numerically)

→ Count vectorizer (one - hot encoding) involves counting the number of occurrence for each word occuring in the doc.
→ Idea behind it is Simple. vector is created having as many dimentions as there are distinct words in text/doc/collection of docs being used. each unique word has a unique dimention & is represented by a 1 in the dimention with 0s everywhere else. This results in huge & spase vectors that capture no relational data.

→ If IDF vectors are related to one hot encoded vectors, but instead of just featuring a count, they feature numerical representation, where the words aren't just present or not present instead, they are represented by their in term freq. multiplied by inverse Doc. freq.

→ In simpler terms, words that occur everywhere should be given very little weight on significance, coz they dont provide a large amt. of value, However if a word appears very little or frequently but only in specific place, then they are possible of higher significance

→ the downside is there in no capture of Semantic relatedness, this is solved with the co-occurence matrix or a ~~preview~~ ~~predictive~~ model.
neutral Probabilistic

\* conclusion :
    Text docs were tokenized, filtered, vectorized, & thus Successfully Processed basics of NLP were understood.