* A**ssignment DA-4** *

* **Title:-** Twitter data analysis.

* **Problem statement:-** Twitter data analysis use twitter data for sentiment analysis. The dataset is 3mB in size and has 31.962 tweets. Identify the tweets which are hate tweets and which are not.

* **objective:-**
Twitter data sentiment analysis for determing hate tweets & which are not.

* **outcome:-**
Students will be able to do senti-ment analysis on twitter dataset using classification algorithm.

* **Theory:-**
Sentiment analysis is the process of determining wheather a piece of writing (Product/movie review tweet etc) is positive, negative or neutral etc. It can be used to identify the customer or followers attitude towards a brand through the use of variable such as context, tone, emotion etc. Marketers can use sentiment analysis to research public opinion of their company & products.

* steps to perform sentiment analysis:-
  * gather relevant tweets from twitter.
  * preprocessing (stopword removal)
  * feature extraction.
  * feature selection.

1) Gather relevant tweets from twitter:-
   various dataset are available of twitter data for sentiment analysis.

2) Preprocessing:-
   It makes raw text ready for mining and then it becomes easier to extract information and apply machine learning algorithm. If it is skipped we were working with noisy and inconsistent data. Now data will become less noisy and relevant for sentiment analysis. Punctuation, special characters, numbers & terms do not carry more weightage.

For preprocessing following steps are required:-

(A) Removing twitter handles (@user):-
    The tweets contain lots of twitter handles (@user) that is how a twitter user a cknowledged on twitter.

(B) Removing punctuations, Numbers and special characters punctuations, Numbers & special characters do not help much. It is better to remove them

(C) Removing short words:-
Some short words are removed for cleansing.

(D) Tokenization:-
Tokens are individual words and tokenization is the process of splitting a string of text into tokens.

(E) Stemming:-
Stemming is a rule based process of stripping the suffixes ("ing", "ly". "es" "s" etc) from a word. For example - "Play", "Player", "Played", "Playing".

3> Feature Extraction:-
Selection of useful words from the tweet is called as feature extraction. In the feature extraction method extract the aspects from preprocessed twitter data set.

4> Feature selection:-
Correct feature selection techniques are used in sentiment analysis that has important rules. They are categorized into 4 main types, Natural language processing statistical clustering based Hybrid

## Classification:-

5)

### (a) Naive Bayes Classifier

This approach uses Bayes theorem which describes how Conditional Probability of each of set of possible Causes for a given observed outcome

$$P(c/x) = \frac{P(x/c) \cdot P(c)}{P(x)}$$

← class prior probability.

↓ Posterior Probability

→ Predictor prior Probability.

### (b) SVM Classifier:-

The support vector machine (SVM) is a algorithm to perform sentiment analysis. It analyzez data, define decision boundries and uses kernels for computation.

### Conclusion:-

Twitter sentiment analysis is performed to identify which are hate tweets and which are not.