**PROJECT 2 REPORT**
**SHREYAS MOHAN  1001669806**


## Description


- The whole newsgroup dataset is organized into 20 newsgroups-

```
['alt.atheism',
 'comp.graphics',
 'comp.os.ms-windows.misc',
 'comp.sys.ibm.pc.hardware',
 'comp.sys.mac.hardware',
 'comp.windows.x',
 'misc.forsale',
 'rec.autos',
 'rec.motorcycles',
 'rec.sport.baseball',
 'rec.sport.hockey';
 'sci.crypt',
 'sci.electronics',
 'sci.med',
 'sci.space',
 'soc.religion.christian',
 'talk.politics.guns',
 'talk.politics.mideast',
 'talk.politics.misc',
 'talk.religion.misc']
```

- 
  We use  **Naive Bayes classifier** which is suitable for discrete classification.
- Then the model starts with training of first 500 files in each category where the words and its frequency are stored in a dictionary for easy access to test it in the future .
- Since the dataset are series of words, before running any machine learning algorithms, we need to convert text representations into numerical representations. One way to do it is using the "Bag of Words" method, first split each file into single words, then count the number of times each word appears in each file, so every unique word will have a unique value in our dictionary.
- At the end of training model, we will have dictionary of trained files under each category, dictionary of words under each category and a master dictionary of total words.
- Later, the model is tested with the remaining 500 files.


## Code Implementation


- It contains a function for data handling  named data_handler . First, the input folders are retrieved as a list from the parent directory.  This function filters the data by removing the special characters and stop words.

- Stop words are the most frequent words occurring in the file.
- Now, using the trained model we test the remaining 500 files in each category and classify the words by computing log probabilities(because it improves accuracy of model) over the testing data using the formula,

$$P(i \mid j) = \frac{word_{ij} + \propto}{word_j + |N| + 1}$$

Where N is the total number of words in vocabulary , $word_{ij}$ counts i in category j, $word_j$ is count of words in category j and $\propto = 0.0001$(Laplace smoothing)

$$P(j) = \log \pi_j + \sum_{n=1}^{N} \log(1 + f_i)\log(P(i|j))$$

Which is the optimal model.

- The number of collisions of maximum probabilities of a particular word in a category is computed and is divided by the total population in order to calculate accuracy of the model. **Accuracy of this model comes out to be 65.00 % .**

## How to improve the performance?

**Under-sampling:** Remove samples from over-represented classes , use this when have huge dataset like the Newsgroup dataset.

# Result Screenshot