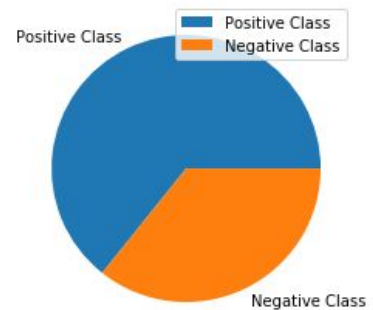


Naive Bayes using Sentiment Analysis

1. Data Cleaning, Data Analysis and Visualization

The dataset consists of Phrases that are labeled positive or negative based on their sentiment.

Total Phrases	2800
Positive Class Examples	1800
Negative Class Examples	1000



Stopping characters like ‘,’ ‘.’ ‘--’,etc were filtered out from the phrases as part of data cleaning process.

2. Feature Extraction

The dataset was splitted into train and test set in 60:40 ratio. As part of feature extraction step, total of 1842 vocabularies were extracted from whole word bag of 15136 taken by splitting words from all phrases present in training dataset(1680 examples). Labels were encoded as 0 and 1.

<i>Label Name</i>	<i>Encoded Label</i>
Positive	1
Negative	0

For feature part, occurence of word and its probability for respective class i.e. $P(W | C_i)$ was calculated with laplacian smoothing given by:

$$P(W | C_i) = (n + \alpha) / (N + \alpha |V|)$$

Where, n is occurence of W in C_i

N is total number of words in C_i

Alpha is pseudocount and $|V|$ is vocab length

A dictionary of word probabilities w.r.t. Class was created. For eg.:

```
word_probs={  
    0:{ "Hello":0.0045,...},  
    1:{ "Hello":0.00034,...}  
}
```

This is used later on to calculate posterior probability of a class.

3. Model Training and Evaluation

For prediction of posterior probability of a class, following formula was used:

$$P(\text{Class} | \text{Text}) = P(\text{Text} | \text{Class}) * P(\text{Class}) \\ = P(W_1|\text{Class})P(W_2|\text{Class})P(W_3|\text{Class})...P(W_n|\text{Class}) * P(\text{Class})$$

If $P(C=0|\text{Text}) > P(C=1 | \text{Text})$, the model will predict the given phrase as negative, otherwise positive.

Following evaluation metrics show the performance of Naive Bayes model:

Confusion Matrix:

	<i>T</i>	<i>F</i>
<i>T</i>	233	173
<i>F</i>	29	685

Model Evaluation Parameters:

Accuracy	0.82
Recall	0.57
Precision	0.88
F1-Score	0.69

For plotting ROC, alpha from Laplacian smoothing was treated as hyperparameter and varied from 1 to 100 and following graph was obtained:

Fig.: ROC Curve

Code is available [here](#)

