

# Naïve Bayes Classifier for Spam Detection

Made By –

Anusha Agrawal (2018A3PS0032H)

Kriti Jethlia (2018A7PS0223H)

Tanay Gupta (2018AAPS0343H)

## Model Description –

- Naïve Bayes classifier assumes that words in a sentence are independent of each other. Probability of a sentence to be spam or not spam is denoted by -

$$P(\text{spam}|\text{sentence}) = P(\text{spam}|w_1, w_2, w_3, \dots, w_n) = \prod_1^n P(\text{spam}|w_i)]$$

$$P(\text{spam}|w_i) = P(w_i|\text{spam}) * P(\text{spam}) / P(w_i)$$

But if a word hasn't occurred in a spam sentence than the above probability will be zero, so we use **Laplacian Smoothing**.

$$P(\text{spam}|w_i) = (\text{Frequency of word in spam sentence} + 1) / (\text{Number of words in vocabulary} + \sum \text{Frequency of each word in spam sentence})$$

- Similar computation is done for non-spam sentences.
- Also, we use  $\log(P(\text{spam}|w_i))$  because multiplying probabilities directly may cause numerical underflow.
- A sentence is considered spam if  $P(\text{spam}|\text{sentence}) > P(\text{not spam}|\text{sentence})$ , this implies  $\log(P(\text{spam}|\text{sentence})) - \log(P(\text{not spam}|\text{sentence})) > 0$ .
- To mitigate the issue of uneven distribution of spam and not spam sentences in the training set, we set a prior  $\log(P(\text{spam})/P(\text{not spam}))$ .
- So the final formula becomes –

$$\log(P(\text{spam}|\text{sentence})/P(\text{not spam}|\text{sentence})) = \log(P(\text{spam})/P(\text{not spam})) + \sum \log(P(\text{spam}|w_i)/P(\text{not\_spam}|w_i))$$

If the above equation results in a number  $> 0$  then the given sentence is spam else not spam.

## Implementation –

Sentences are first preprocessed to remove punctuations, stopwords and digits as these characters don't provide any useful information about whether a sentence is spam or not.

The sentences are then shuffled and divided into 7 cross folds and cross fold validation is used to train and test the model.

Probabilities required by the model and all other required information is generated using the train set. While testing a sentence if a word in it is not present in the vocabulary it is ignored.

Finally the model is tested on the test set.

The probabilities for each word is saved in a dictionary for fast retrieval.

## Accuracy of Model –

Mean Accuracy = 82.6%

Minimum Accuracy = 78.87%

Maximum Accuracy = 86.61%

Accuracy over each crossfold = 82.39%, 81.69%, 78.87%, 79.57%, 86.61%, 84.50%, 80.28%

## Limitations –

1. Naïve Bayes assumes that all predictors (or features) are independent which rarely happens in real life.
2. Naïve Bayes implementation without any kind of smoothing suffers from zero frequency problem, meaning that if a word was not encountered in the spam class of train dataset then the probability of that word occurring in the spam class is zero,  $P(\text{word}|\text{spam}) = 0$ .
3. Naïve Bayes cannot deal with continuous features and binning is used to convert them into discrete features, but if this is not done properly then a lot of information can be lost in the process.