# MACHINE LEARNING (CO 327) PROJECT REPORT

# SPAM DETECTION USING NAIVE BAYES

## Under the Supervision of Mrs Juhi Jain

**Submitted by:**
Tejas Harish Borkar (2K18/CO/373)
Tushar Ahuja (2K18/CO/374)
(Group : 9)

## Department of Computer Science , DTU

# OBJECTIVE:

- Our goal was to code a **spam filter** from scratch that classifies messages as spam and non-spam with an accuracy greater than 80%.
- We have used multinomial **Naive Bayes** Algorithm and conditional probability to achieve it.

# EXPERIMENTAL DESIGN:

## DATASET:

- The datasets are a set of SMS tagged messages that have been collected for spam research.
- It contains a set of SMS messages in English tagged according to ham (legitimate) or spam.
- It is an example of **supervised learning**.

The datasets have been downloaded from three sources :

- https://archive.ics.uci.edu/ml/datasets/sms+spam+collection
- https://www.kaggle.com/ozlerhakan/spam-or-not-spam-dataset
- https://www.kaggle.com/venky73/spam-mails-dataset

**Oversampling** had to be done as the datasets were imbalanced.

## OVERSAMPLING

- Oversampling is used to adjust the class distribution in a data set. In other words, it helps to adjust the ratio between the different classes/categories represented by the dataset.
- We performed over-sampling of the spam messages by adding spam entries to the dataset.
- The ratio of spam:ham before over-sampling: 14:86
- The ratio of spam:ham after over-sampling: 36:64

## FORMAT OF THE DATASET

The files contain one message per line. Each line is composed of two columns: one with label (ham or spam) and other with the raw text. Here are some **examples:**

ham   What are you doing?how are you?

ham   Ok lar… Joking with u oni…

ham   dun say so early hor… U c already then say…

spam   Sunshine Quiz! Win a super Sony DVD recorder if you can…
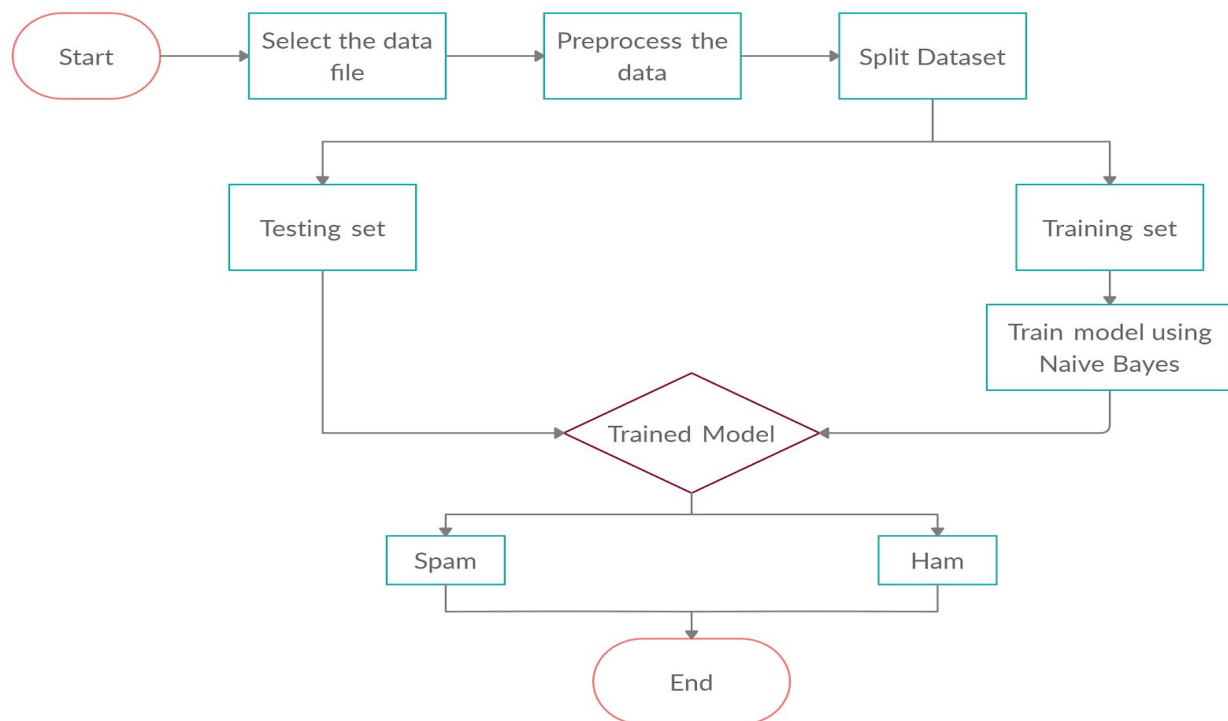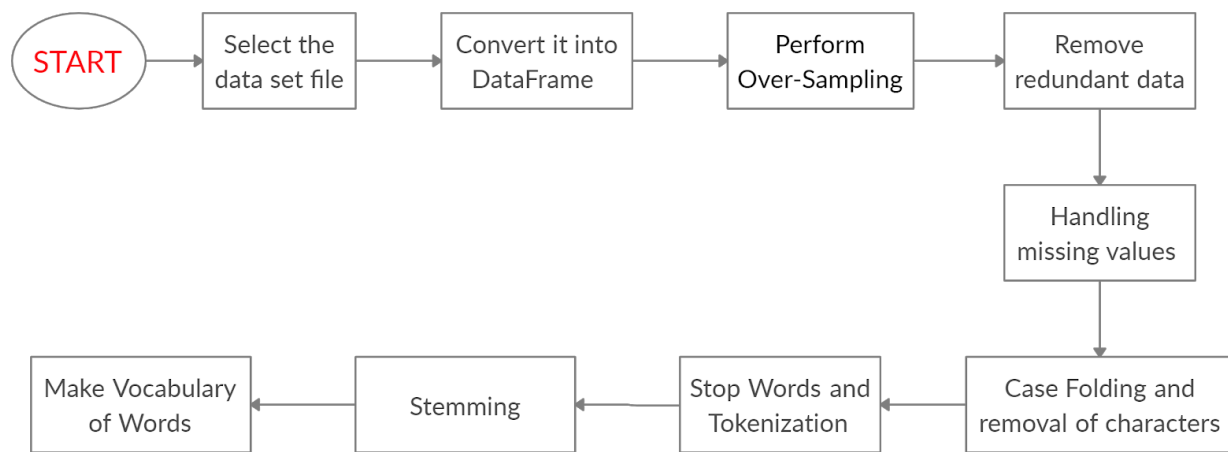
## RESEARCH VARIABLES:

### Dependent variable

- The dependent variable is the label that is whether the message is classified as spam/ham.
- It is a nominal variable.

### Independent variable

- These are the indices of the words stored in vocabulary (0 to total no. of words in vocabulary-1).
- It is a metric ratio variable as real 0 has a meaning that is the absence of a word.

## DIAGRAM: Spam detection using Naive Bayes

**Preprocessing of Data**

## PREPROCESSING TECHNIQUES:

- Converting the dataset to a dataframe.
- Handling missing values
- Remove redundant data
- Case Folding and removal of characters
- Tokenization
- Vocabulary creation.
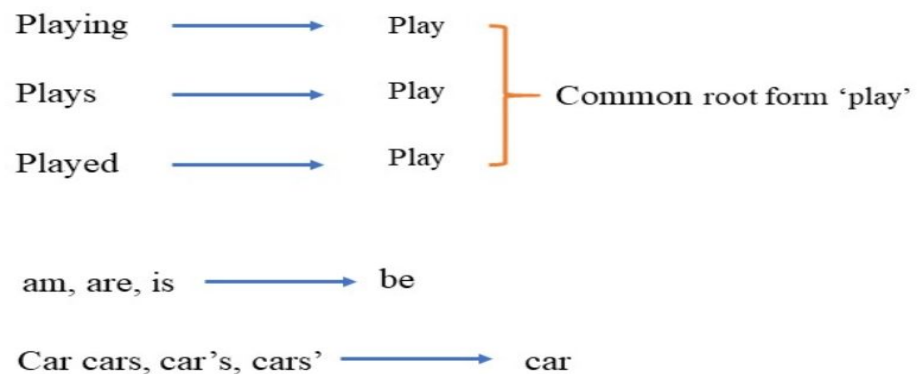- Stemming

# RESEARCH METHODOLOGY:

## DATA PREPROCESSING:

The Preprocessing of the training data and testing is separately performed to simplify the process. The phases being taken in the preprocessing include :

- **Handling missing values:** We handle missing values by :
  - Deleting the row if Email is missing.
  - Replacing missing values with ham if the label is missing so that there is no risk of data loss.
- **Handling redundant values:** We delete the duplicate rows that is we keep only a single instance of every message.
- **Case Folding and removal of characters:** All the text is converted to

lowercase and characters other than letters, numbers and punctuation are deleted.

- **Removing Stop Words and Tokenization:** In this process, we remove stopwords and tokenization is carried out to break the string into tokens or individual words.
- **Stemming:** Stemming is the process of converting different variants of a word to the root/base form of the word. For example, the words "playing", "plays" and "played" are converted to the root form, i.e. "play".Stemming is an essential part of the pipelining process in Natural language processing.

Playing ⟶ Play

Plays ⟶ Play ⟩ Common root form 'play'

Played ⟶ Play

am, are, is ⟶ be

Car cars, car's, cars' ⟶ car

Using above mapping a sentence could be normalized as follows:

the boy's cars are different colors ⟶ the boy car be differ color

**Stemming**

- **Vocabulary creation.**
  - The vocabulary is created, which means a list of all the unique words in our training set.
  - We will now use the vocabulary we created to make a new DataFrame.
  - Therefore, data access is simple, and this DataFrame is used to make the required data transformation.

| | Label | SMS |
|---|---|---|
| 0 | spam | SECRET PRIZE! CLAIM SECRET PRIZE NOW!! |
| 1 | ham | Coming to my secret party? |
| 2 | spam | Winner! Claim secret prize now! |

| | Label | secret | prize | claim | now | coming | to | my | party | winner |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | spam | 2 | 2 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 1 | ham | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 0 |
| 2 | spam | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 |

**Note:** This image is for representational purposes only. In the actual model , the indices of words in vocabulary are used as the independent variables.

## DATA ANALYSIS TECHNIQUE:

We have used **Naive Bayes** classifier as our data analysis technique to build the spam detection model.
- The Naive Bayes algorithm is primarily based on Bayes' theorem with the assumption of independence between predictors.
- A Naive Bayesian model is relatively easy to build, as there is no complicated iterative parameter estimation. This makes it useful for massive datasets.
- Despite its simplicity, it does surprisingly well and is extensively used because it performs better than some sophisticated classification methods.

**Algorithm**
The Bayes theorem enables us to calculate the posterior probability P(A|B) from P(A), P(B), and P(B|A). Naive Bayes classifiers are based on the assumption that the effect of the value of the predictor (B) on a given class (A) is independent of the values of other predictors. This is known as **class conditional independence**.

**LIKELIHOOD**
The probability of "B" being
True, given "A" is True

**PRIOR**
The probability "A" being
True. This is the knowledge.

$$P(A|B) = \frac{P(B|A).P(A)}{P(B)}$$

**POSTERIOR**
The probability of "A" being
True, given "B" is True

**MARGINALIZATION**
The probability "B" being True.

- P(A|B): the posterior probability of class (target) given predictor (attribute).
- P(A): the prior probability of class.
- P(B|A): likelihood which is the probability of the given predictor class.
- P(B): the prior probability of the predictor.

## ADVANTAGES OF NAIVE BAYES

- When the assumption of independent predictors holds, a Naive Bayes classifier performs better as compared to other models.
- A small amount of training data is used to estimate the test data. So, the training period is less.
- It is easy to implement.

## DISADVANTAGES OF NAIVE BAYES

- One of the main limitations is the assumption of independent predictors. An essential premise of the algorithm is conditional independence. But in many situations, it is impossible to get an entirely independent set of predictors.

- Suppose a categorical variable has a category in the testing data set, which was not previously observed in the training data set. In that case, the model assigns a 0 probability and will be unable to make a prediction which is often known as Zero Frequency. We can use smoothing technique.s to solve this issue. One of the most straightforward smoothing techniques is called **Laplace estimation**.

## Calculating Constants and parameters(Naive Bayes Algorithm)

When a new message comes in, the model makes the classification based on the results it gets on the basis of these equations, where "w1" denotes the first word, and w1,w2, ..., wn is the entire message:

$$P(\text{Spam}|w_1, w_2, ..., w_n) \propto P(\text{Spam}) \cdot \prod_{i=1}^{n} P(w_i|\text{Spam})$$

$$P(\text{Ham}|w_1, w_2, ..., w_n) \propto P(\text{Ham}) \cdot \prod_{i=1}^{n} P(w_i|\text{Ham})$$

If $P(\text{Spam} | w_1, w_2, ..., w_n)$ is greater than $P(\text{Ham} | w_1, w_2, ..., w_n)$, then the message is spam.

To calculate $P(w_i|\text{Spam})$ and $P(w_i|\text{Ham})$, we need to use separate equations:

$$P(w_i|\text{Spam}) = \frac{N_{w_i|\text{Spam}} + \alpha}{N_{\text{Spam}} + \alpha \cdot N_{\text{Vocabulary}}}$$

$$P(w_i|\text{Ham}) = \frac{N_{w_i|\text{Ham}} + \alpha}{N_{\text{Ham}} + \alpha \cdot N_{\text{Vocabulary}}}$$

$N_{w_i|\text{Spam}}$ = the number of times the word $w_i$ occurs in spam messages

$N_{w_i|\text{Ham}}$ = the number of times the word $w_i$ occurs in ham messages

$N_{\text{Spam}}$ = total number of words in spam messages

$N_{\text{Ham}}$ = total number of words in ham messages

$N_{\text{Vocabulary}}$ = total number of words in the vocabulary

$\alpha = 1$   ($\alpha$ is a smoothing parameter)

Now the constants are calculated

- P(Spam) and P(Ham)
- $N_{\text{Spam}}$, $N_{\text{Ham}}$, $N_{\text{Vocabulary}}$

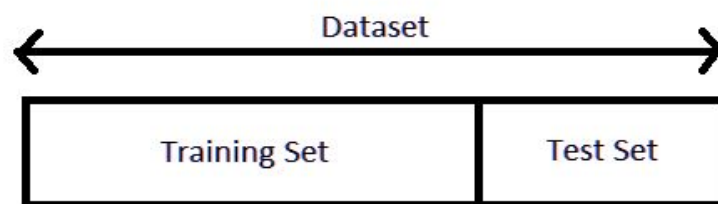We also use Laplace smoothing and set $\alpha = 1$.

**Need for Laplace Smoothing**

$$P(\text{Spam}|w_1, w_2, ..., w_n) \propto P(\text{Spam}) \cdot \prod_{i=1}^{n} P(w_i|\text{Spam})$$

- It is performed to overcome the disadvantage of Naive Bayes.
- If any term in the Right-Hand side is 0, that is if a given word doesn't occur in spam messages, then the entire probability becomes zero, that is not a desired trait.
- Therefore Laplace smoothing is needed so that the probability remains non zero.
- The same is true in the case of ham messages as well.

## VALIDATION TECHNIQUES USED:

### 1) Hold-out Validation Technique



- The dataset is split into two sets: a training set and a testing set.
- 80% of the data is used for training, and the remaining 20% for testing.
- The dataset is randomized before the split is done to ensure that spam and ham messages are uniformly spread throughout the dataset.
- The percentage of ham and spam messages in the training and testing sets are analyzed.
- The percentages are expected to be close to what we have in the complete dataset.

### 2) Stratified K-Folds cross-validation
- In Stratified Sampling, let the population of a state be 55.7% male and 44.3% female. If we wish to choose 1000 people from that state and if 557 male ( 55.7% of 1000 ) and 443 female ( 44.3% for 1000 ), i.e. 557 male + 443 female (Total=1000 people) are picked for their opinion, then this chosen group of people represent the entire state.

- Stratified k-fold cross-validation is the same as just k-fold cross-validation, but in Stratified k-fold cross-validation, stratified sampling is done instead of random sampling.

## Classifying A New Message

Take as input a new message (w1, w2, ..., wn).
The values of P(Spam|w1, w2, ..., wn)=x and P(Ham|w1, w2, ..., wn)=y are calculated and compared.

- ❖ If y > x, then the message is classified as ham.
- ❖ If x < y, then the message is classified as spam.
- ❖ If x = y, then the algorithm may request human help.

Some messages may contain words that are not part of the vocabulary. Such words are ignored while calculating the probabilities.

## PERFORMANCE MEASURES:

- **Accuracy:** The ratio of the number of correct predictions to the total predictions made.

  Accuracy = (TP+TN)/total

- **False Positive (FP):** The number of misclassified non-spam emails.
- **False Negative (FN):** The number of misclassified spam emails
- **True Positive (TP):** The number of spam messages correctly classified as spam.
- **True Negative (TN):** The number of non-spam emails correctly classified as non-spam.
- **Recall:** Percentage of spam messages managed to block.

  r = TP/(TP+FN)

- **Precision:** Percentage of correct message for spam emails.

  p = TP/(TP+FP)

- **F Beta-measure:** Weighted average of precision and recall.

$$\text{F-Beta-measure} = (1+beta2)p*r \,/\, (beta2)(p+r)$$

## CHOOSING THE APPROPRIATE VALUE OF Beta

- In case if FP and FN both are equally important then we use beta = 1
- In case if FP is more important than FN then we reduce beta value (between 0 to 1)
- In case if FN is more important than FP than we increase beta value (greater than 1)

In our case if a ham message arrives and if the model reports this message as spam then this is a big trouble for the user/client i.e number of misclassified non spam emails is important (False Positive) .

So we will use beta = 0.5 (also known as F0.5 score)

$$\textbf{F0.5 score} = 5*p*r/(p+r)$$

# LITERATURE REVIEW

There are various supervised machine learning algorithms & Statistical Spam Filtering Techniques for filtering spam and ham messages like naïve Bayes Algorithm, support vector machines algorithm, decision trees, artificial neural networks, and many more.

### A. Naive Bayes

The Naive Bayes algorithm is based on Bayes Theorem. This algorithm classifies each object by looking at all of its features individually. Bayes Rule helps us to calculate the posterior probability for just one feature. The posterior probability of every item is calculated, and these probabilities are multiplied together to get the final probability. This procedure is repeated for the other class as well. The class with the greater probability determines what class the object should belong.

## B. Support Vector Machine

SVMs have out-performed other learning algorithms with good generalization, the number of tuning parameters, and its firm theoretical background. The state of the art of Support Vector Machines evolved, mapping the learning data from input space into higher dimensional feature space where the classification performance is increased SVM's perform best when using binary features. SVM's provide good accuracy and speed. They also have significantly less training time as compared to other techniques.

## C. Decision Trees

It is a commonly used tool for decision-making. A new instance is passed along the tree, and as per the values of the nodes, it is classified accordingly. C4.5, ID3, CART are some of the most popular decision tree algorithms. Spam detection using decision trees is done using recursive binary splitting.

## D. Artificial Neural Networks (ANN)

As per the Neural Network theory, for static pattern classification, the best performance shows the layered feed-forward networks, called Multilayer Perceptrons, typically trained with static backpropagation. Their advantage is that they are reasonably easy to use, and they can approximate any input/output map. The main disadvantage is that they have considerable training time and require large amounts of training data. The method for spam detection employs attributes composed from characteristics of the elusive patterns that spammers use rather than the frequency or context of keywords in the messages. They find out which ANN configuration will have the best performance and least error to find the desired output.

## PERFORMANCE

**HOLD OUT VALIDATION:**

```
ACCURACY    :  0.8335694050991501
PRECISION   :  0.8431876606683805
RECALL      :  0.6533864541832669
F0.5_SCORE  :  1.840628507295174
```

**STRATIFIED K FOLD CROSS VALIDATION:**

```
ACCURACY    :  0.8223796033994334
PRECISION   :  0.8523880813379252
RECALL      :  0.610994287478771
F0.5_SCORE  :  1.7750591106239704
```

As it can be observed from the performance measures, hold out validation has a slightly better accuracy whereas as stratified k folds cross-validation has slightly better precision. If we consider recall and F0.5_Score then hold out validation outperforms stratified k folds cross-validation.

## CONTRIBUTION

**Tejas :**

- Data Cleaning
- Naive Bayes algorithm
- Performance Measures

**Tushar :**

- Oversampling
- Vectorisation
- Validation Techniques

## REFERENCES

- Arifin, D. D., & Bijaksana, M. A. (2016, September). Enhancing spam detection on mobile phone Short Message Service (SMS) performance using FP-growth and Naive Bayes Classifier. Google.com. In 2016 IEEE Asia Pacific Conference on Wireless and Mobile (APWiMob) (pp.80-84). IEEE.
- https://app.creately.com/ for creating flow charts.
- https://www.kaggle.com/ for datasets.