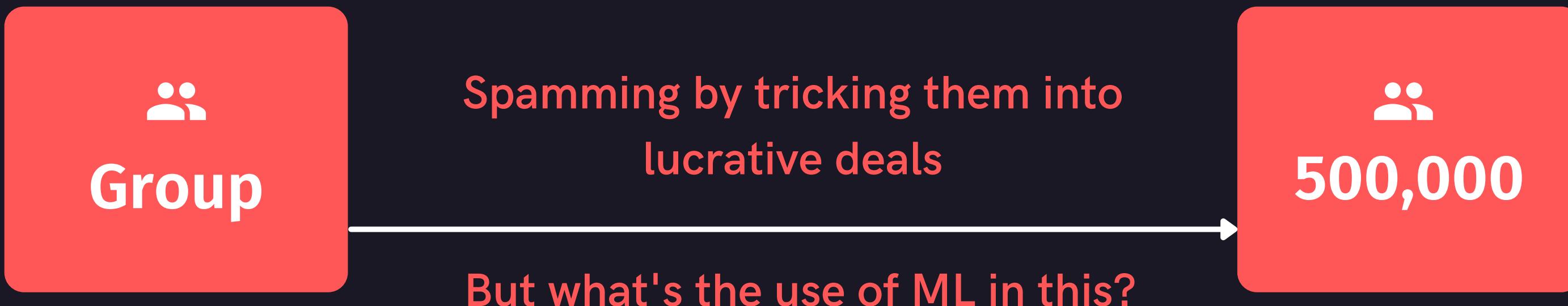


Comparative Study on best ML algorithm for SPAM Message Classification.

Yash Raj Ojha
ENG19CS0367

Understanding the topic

We first need to understand what is Spam and what it is widely used for. Spam messages are messages sent to a large group of recipients without their prior consent, typically advertising for goods and services or business opportunities. In the recent period, the percentage of scam messages amongst has increased sharply.



Diving into the dataset

With the use of ML and DL, we can build a spam message classifier/filter which will be a step towards building a tool for scam message identification and early scam detection. Now let's understand the two biggest terms.

Ham

"Ham" is e-mail that is not Spam. In other words, "non-spam", or "good mail". It should be considered a shorter, snappier synonym for "non-spam".

Spam

It is an irrelevant or unsolicited message, sent to a large number of users, for advertising, phishing, spreading malware.

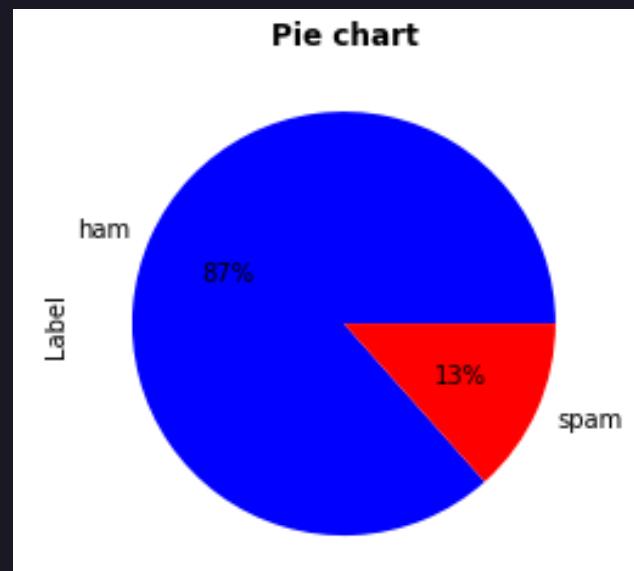
	v1	v2	Unnamed: 2	Unnamed: 3	Unnamed: 4
0	ham	Go until jurong point, crazy.. Available only ...	NaN	NaN	NaN
1	ham	Ok lar... Joking wif u oni...	NaN	NaN	NaN
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...	NaN	NaN	NaN
3	ham	U dun say so early hor... U c already then say...	NaN	NaN	NaN
4	ham	Nah I don't think he goes to usf, he lives aro...	NaN	NaN	NaN
5	spam	FreeMsg Hey there darling it's been 3 week's n...	NaN	NaN	NaN
6	ham	Even my brother is not like to speak with me. ...	NaN	NaN	NaN

A sneak peek
into the dataset



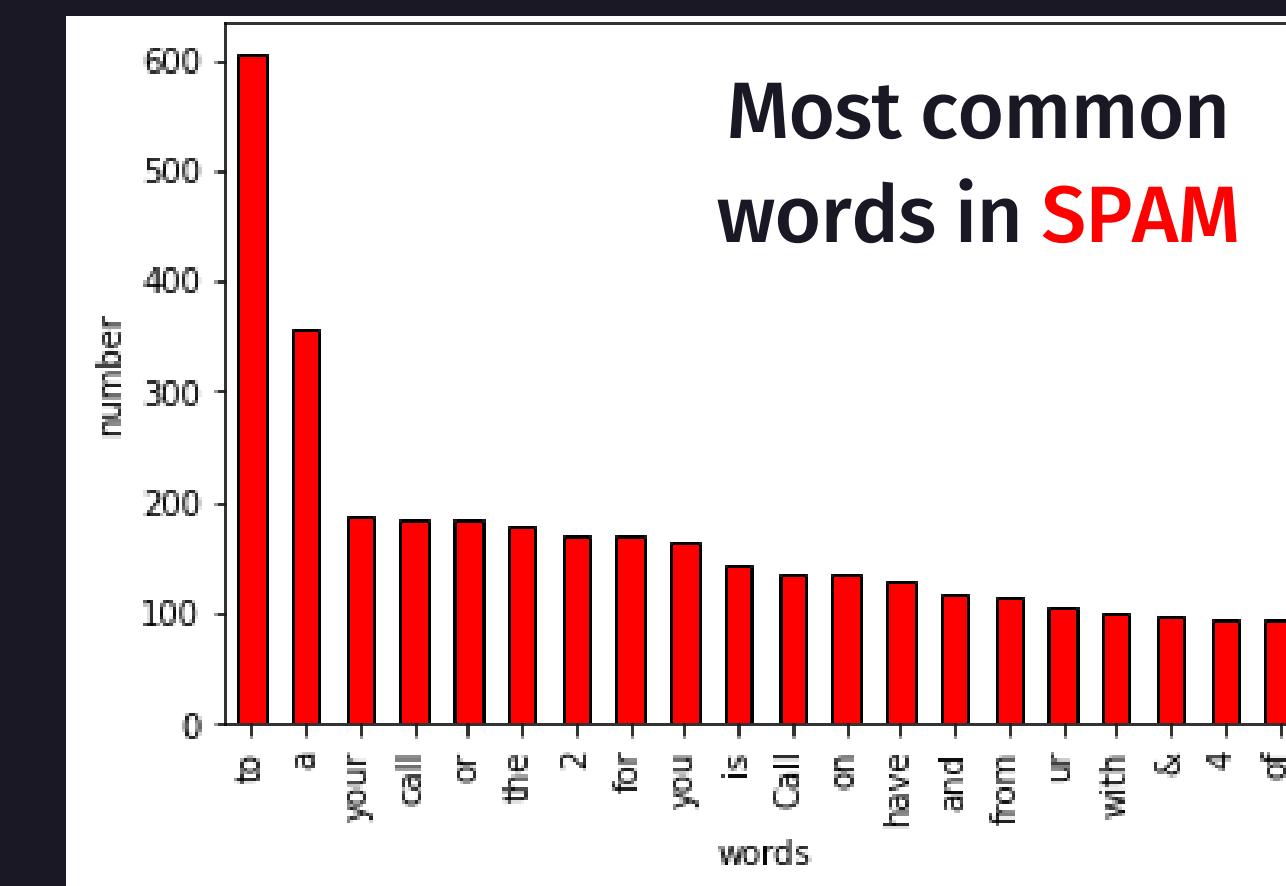
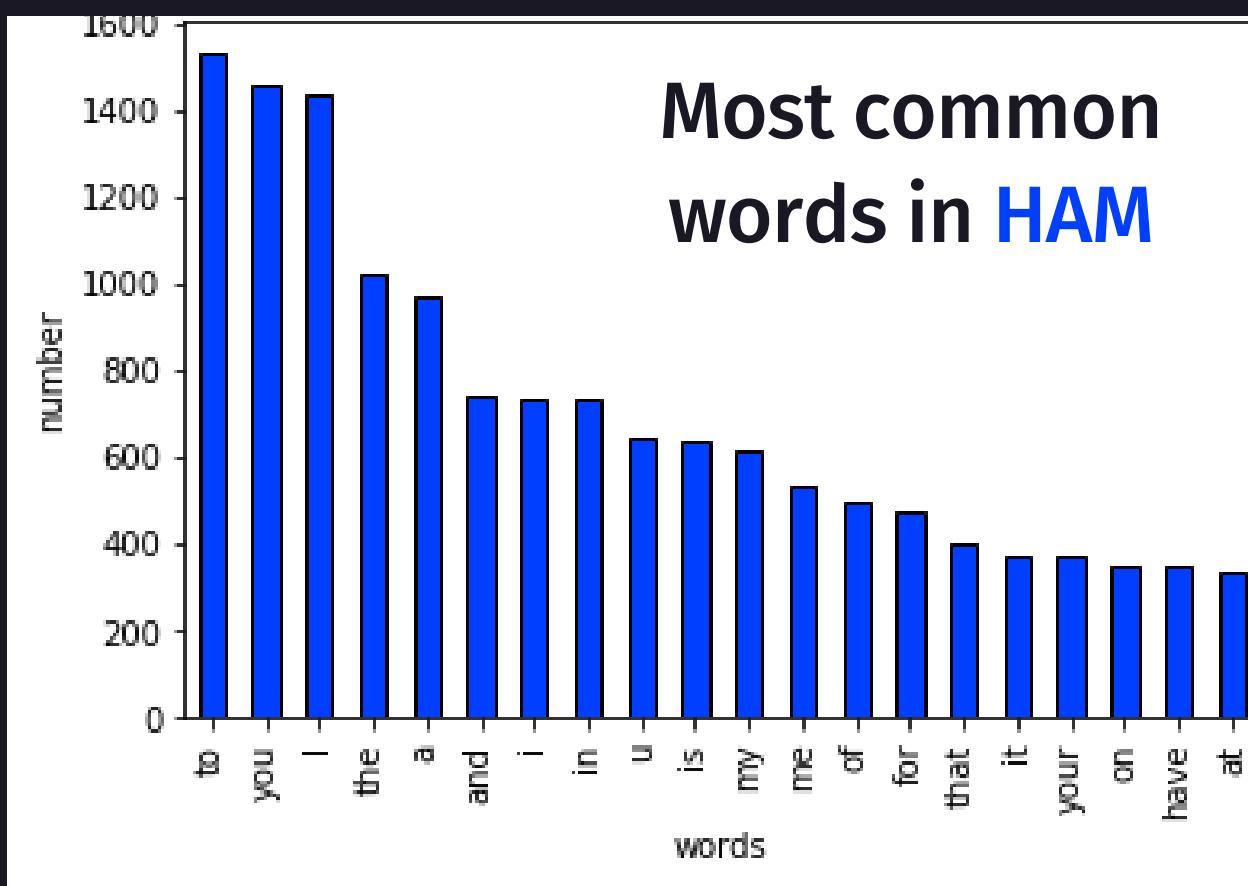
Analysing the dataset

Using different visualizations to understand the dataset



87% of the messages in the dataset is
HAM

13% of the messages in the dataset is
SPAM



Cleaning the dataset for predictive analysis

As mentioned earlier we need to ignore the 'STOP WORDS' for the algorithm to work with better efficiency. I have first imported stopwords from NLTK, then removed special characters, then converted every word to lower cases and joining in the sentence after removing the stop words in the English language.

After
removing
stopwords



```
0      Go jurong point crazy Available bugis n great ...
1                      Ok lar Joking wif u oni
2      Free entry 2 wkly comp win FA Cup final tkts 2...
3                      U dun say early hor U c already say
4      Nah dont think goes usf lives around though
...
5567    2nd time tried 2 contact u U £750 Pound prize...
5568                      I b going esplanade fr home
5569                      Pity mood Soany suggestions
5570    guy bitching acted like id interested buying s...
5571                      Rofl true name
Name: Message, Length: 5572, dtype: object
```

We cannot work with the text directly when using machine learning algorithms. Instead, we need to convert the text to numbers, this process is called Vectorization.

More on Vectorizers

Vectorization is of two types, firstly **Count Vectorizer**, it is when we were to feed the count of the words as a matrix directly to a classifier, this means very frequent terms like, a, an, the would **overshadow the frequencies more important terms.**

To avoid this we use the second method, in this method **we give a certain weight to the words, convert them into floating-point values** suitable for usage by a classifier, it is very common to use the **TF-IDF vectorizer**.

$$\text{TF - IDF} = \text{TF}(t, d) * \text{IDF}(t)$$

term frequency
no. of times term 't' appears in document 'd'

$$\log\left(\frac{1+n}{1+df(t)}\right) + 1$$

n → no. of documents
df → no. of documents that contain term 't'

The matrix contains a weight value that signifies how important a word is for an individual text message or document.

This is very common algorithm to transform text into a meaningful representation of numbers which is used to fit machine algorithm for prediction.

How will a TF-IDF Matrix look like?

Document-1: Text processing is necessary.

Document-2: Text processing is necessary and important.

Document-3: Text processing is easy.

Unique Words: ['and','easy','important','is','necessary','processing','text']

```
[[0 0 0 1 1 1 1]
 [1 0 1 1 1 1 1]
 [0 1 0 1 0 1 1]]
```

} Count Vectorizer array

```
[[0.          0.          0.          0.46333427 0.59662724 0.46333427
  0.46333427]
 [0.52523431 0.          0.52523431 0.31021184 0.39945423 0.31021184
  0.31021184]
 [0.          0.69903033 0.          0.41285857 0.          0.41285857
  0.41285857]]
```

} TF-IDF Array

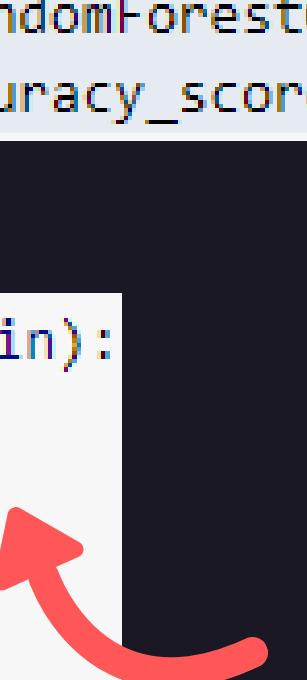
This will be fitted into different classifier models.

Let's try out different classifier ML algorithms

We first split the dataset into features train and test, labels train and test, which will be used for classifier.fit and classifier.predict methods.

```
[ ] from sklearn.linear_model import LogisticRegression  
from sklearn.svm import SVC  
from sklearn.naive_bayes import MultinomialNB  
from sklearn.tree import DecisionTreeClassifier  
from sklearn.neighbors import KNeighborsClassifier  
from sklearn.ensemble import RandomForestClassifier  
from sklearn.metrics import accuracy_score
```

```
def train_classifier(clf, feature_train, labels_train):  
    clf.fit(feature_train, labels_train)  
  
def predict_labels(clf, features):  
    return (clf.predict(features))
```

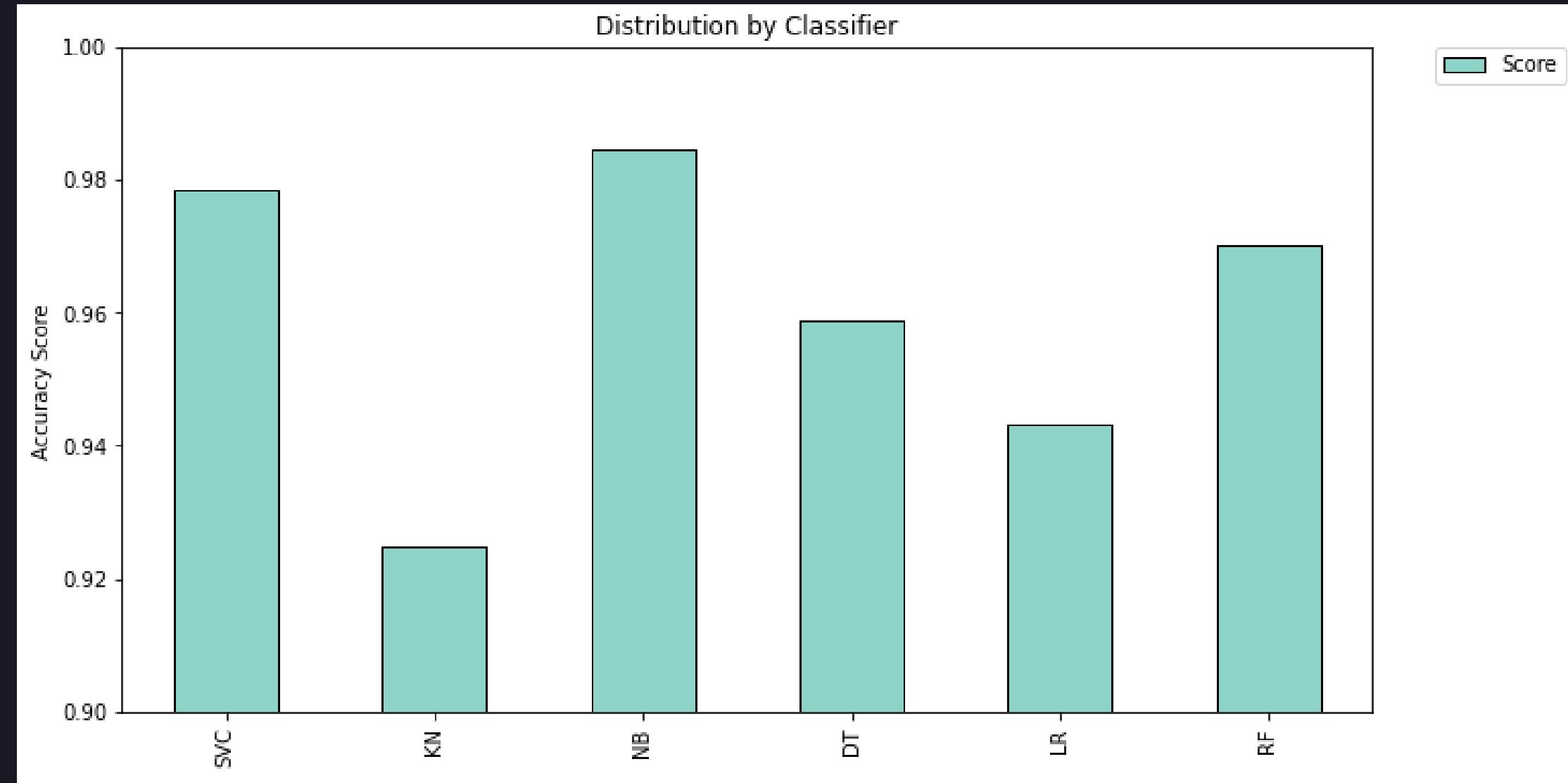


After importing the classifier algorithms and creating their object, I convert the objects into the dictionary for easier accessibility and then use a for loop to call these fit and predict user-defined functions.

Accuracy and Graphs (Part 1)

I create a list, having the accuracy value of different classifiers and then convert this whole thing into a data frame for easier plotting.

Score	
SVC	0.978469
KN	0.924641
NB	0.984450
DT	0.958732
LR	0.943182
RF	0.970096



Ensemble Algorithms like **Decision Trees** and **Random Forest** have performed relatively poor in comparision to **Support Vector Machine** and **Naive Bayes**.

A study of twists and turns

Now let me introduce you to Stemming, it is Text Normalization, these techniques in the field of Natural Language Processing are used to prepare text, words, and documents for further processing.

Playing, Plays, Played are having a common root "play".
Similarly Cars, Car's, car have a common root "car".

In our case Stemming helps us to achieve the root forms.

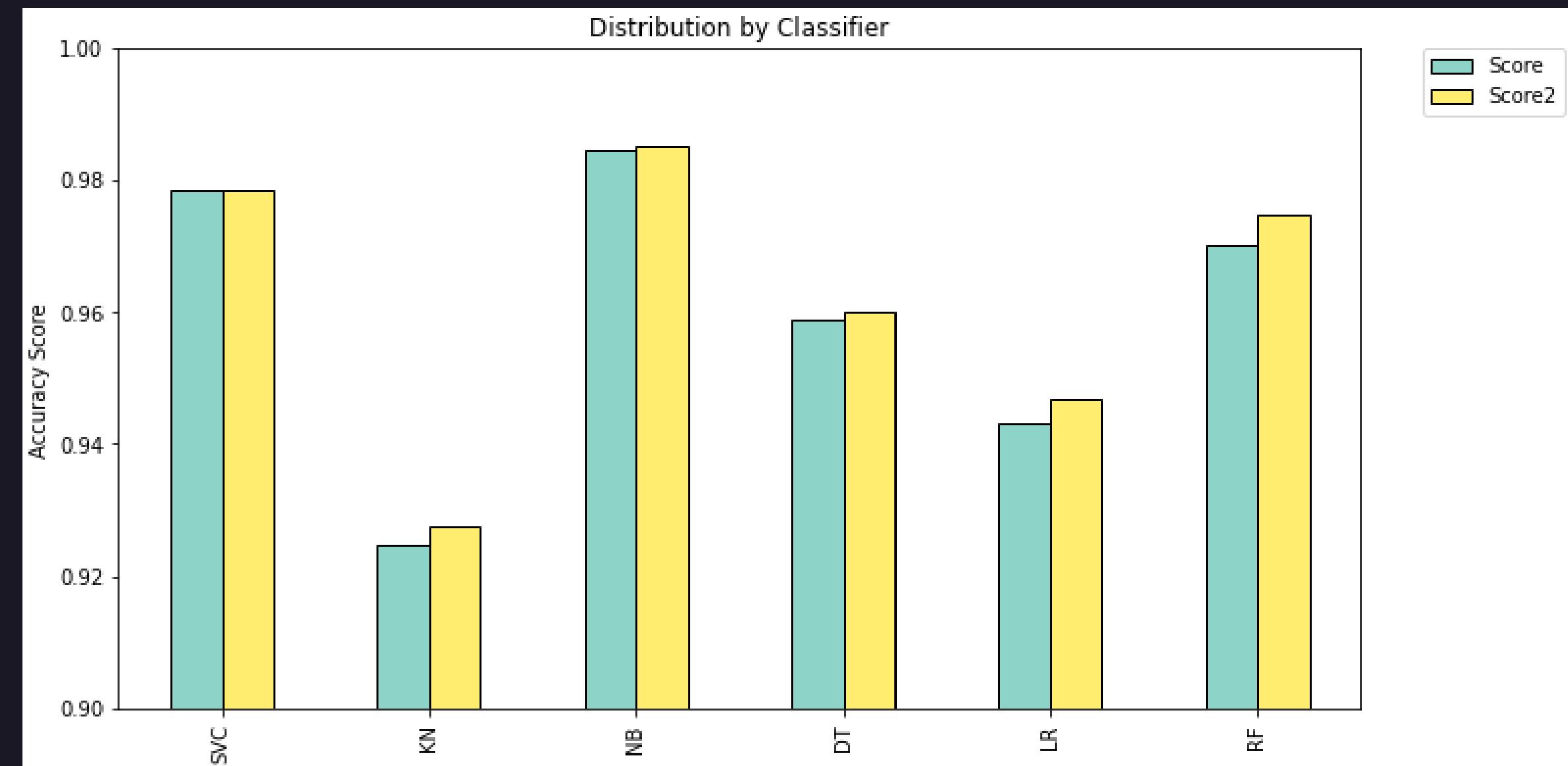
```
0 Go jurong point crazy Available bugis n great ...
1                               Ok lar Joking wif u oni
2 Free entry 2 wkly comp win FA Cup final tkts 2...
3                               U dun say early hor U c already say
4                               Nah dont think goes usf lives around though
5 ...
5567 2nd time tried 2 contact u U å£750 Pound prize...
5568           i b going esplanade fr home
5569           Pity mood Soany suggestions
5570 guy bitching acted like id interested buying s...
5571           Rofl true name
```

```
0 go jurong point crazi avail bugi n great world...
1                               ok lar joke wif u oni
2 free entri 2 wkli comp win fa cup final tkts 2...
3                               u dun say earli hor u c alreadi say
4                               nah dont think goe usf live around though
5 ...
5567 2nd time tri 2 contact u u å£750 pound prize 2...
5568           i b go esplanad fr home
5569           piti mood soani suggest
5570 guy bitch act like id interest buy someth els ...
5571           rofl true name
```

Accuracy and Graphs (Part 2)

I vectorized the stemmed data, split it to train and test, ran a for loop so every classifier algorithm could run fit and predict methods on the modified 'stemmed' data. I then obtain the accuracy scores inside a data frame for visualization.

	Score	Score2
SVC	0.978469	0.978469
KN	0.924641	0.927632
NB	0.984450	0.985048
DT	0.958732	0.959928
LR	0.943182	0.946770
RF	0.970096	0.974880

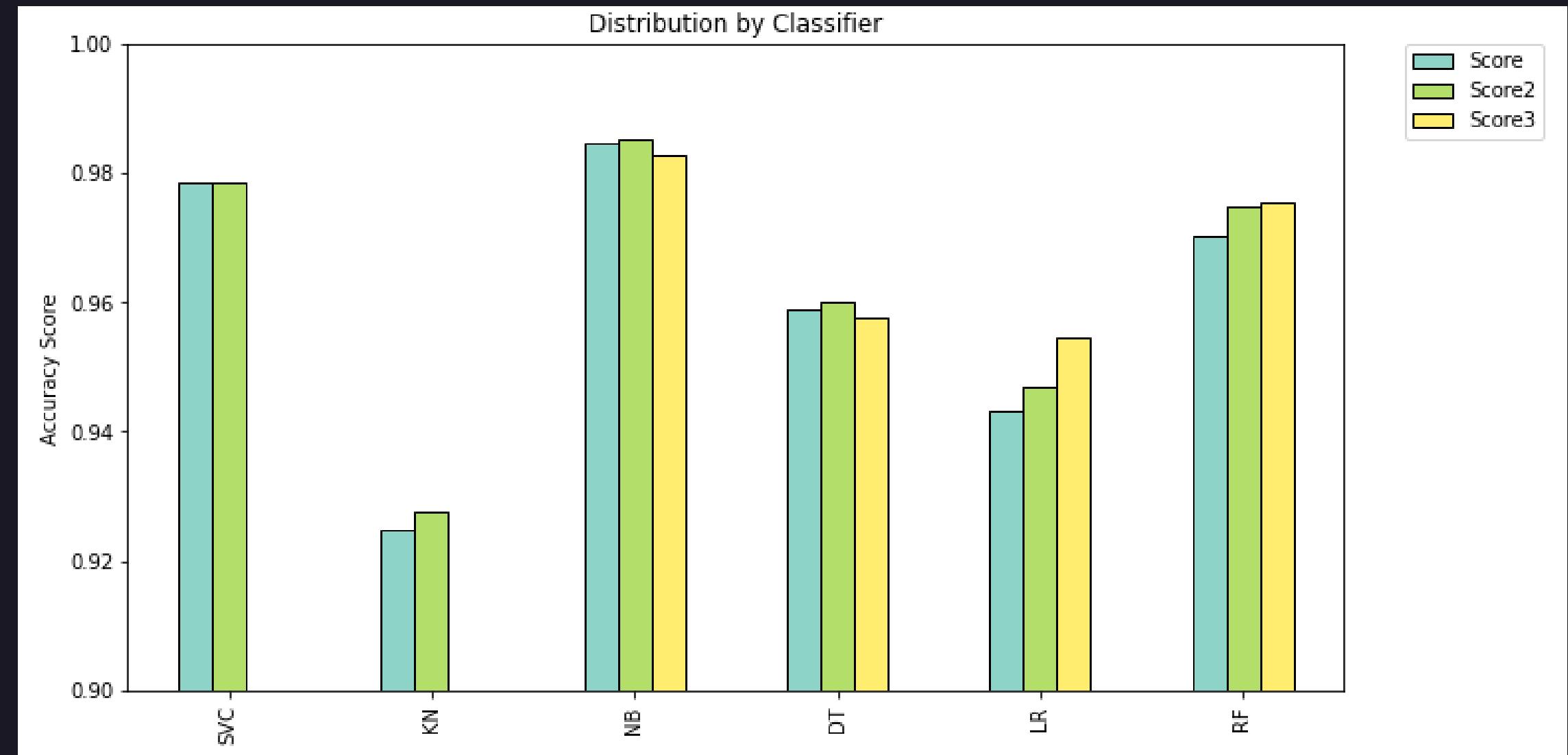


The scores are slightly better than previous check.

Accuracy and Graphs (Part 3)

After considering a normal message with stop words removed and shortening the length of the word to just its root (Stemming), we will now check the accuracy of these classifiers considering a key factor that is, length of the messages.

	Score	Score2	Score3
SVC	0.978469	0.978469	0.861244
KN	0.924641	0.927632	0.881579
NB	0.984450	0.985048	0.982656
DT	0.958732	0.959928	0.957536
LR	0.943182	0.946770	0.954545
RF	0.970096	0.974880	0.975478



In this study, Multinomial Naive Bayes comes out as the undisputed best algorithm for SPAM classification, having the best accuracy regardless of the conditions, but why?

Understanding Naive Bayes Algorithm

Naive Bayes algorithm is a supervised learning algorithm, which is based on Bayes theorem and used for solving classification problems.

We will be using multinomial classification as it suits best for the discrete values like word counts. It is very useful in text classification. Let's understand why.

Let's take the example of normal and spam message:

The probabilities of certain words from a friendly message and a spam. (example)



$P('Dear'|NM) : 0.47$
 $P('Friend'|NM) : 0.29$
 $P('Lunch'|NM) : 0.18$
 $P('Money'|NM) : 0.06$



$P('Dear'|SPAM) : 0.29$
 $P('Friend'|SPAM) : 0.14$
 $P('Lunch'|SPAM) : 0$
 $P('Money'|SPAM) : 0.57$

Let's assume 8 out of 12 messages
are normal messages so
 $P(\text{Normal Messages}) : 0.67$
 $P(\text{Spam Messages}) : 0.33$

Probabilities or Likelihoods?

Understanding Naive Bayes Algorithm

We get a message "Lunch Money Money Money"

We take an initial guess that it's a SPAM message, the probability of our initial guess is called Prior Probability. From the first look we can say that this is SPAM, because of the higher probability of 'money' in SPAM.

Score of the message being Normal = 0.000002

Score for the message being a SPAM = $P(SPAM) \times P(Lunch|SPAM) \times 4 \times P(Money|SPAM)$
= $0.33 \times 0 \times 4 \times 0.57 = 0$

By the logic of high score we will conclude **wrongly** that this message is **NORMAL**. ✓

To avoid this 0 factor coming into play for the probability, we increase the count of words by 1 as a buffer.

This increasing count by x value is termed as **ALPHA**.

Understanding Naive Bayes Algorithm



$P(\text{'Dear'}|\text{SPAM})$: abc
 $P(\text{'Friend'}|\text{SPAM})$: abc
 $P(\text{'Lunch'}|\text{SPAM})$: 0.09
 $P(\text{'Money'}|\text{SPAM})$: abc

Every probability will be updated accordingly. But more importantly, we won't have a 0 to deal with.

Score of the message being Normal = 0.00001
Score for the message being a SPAM = 0.00122

Because of the alpha value we got a greater score for the message being spam than for the score for the message being normal, we can conclude that the message we have received is a SPAM.

The thing about Naive Bayes is it treats all the words equally any order of the same set of words will consequently return the same score.

And so for this reason of ignoring order and just considering importance of each message, Multinomial Naive Bayes does comparatively well against other classifier algorithms for text classification.

Thank You

Link to the repository:

