



Sardar Patel Institute of Technology, Mumbai  
Department of Electronics and Telecommunication Engineering  
B.E. Sem-VII (2022-2023) Data Analytics

## Experiment: Exploratory Data Analysis (EDA)

Name: Rahul Alshi

UID: 2019110001

BE ETRX

DA LAB 6

**Aim:** Perform Classification on SPAM on a Dataset in Python

### Dataset Overview

The dataset 'SPAM text\_message 20170820 - Data.csv' contains with 5527 rows and 2 columns contains information on emails wther it is a spam or ham based on the message in the box .:

Code :

```
In [ ]: !pip install scikit-plot

Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Collecting scikit-plot
  Downloading scikit-plot-0.3.7-py3-none-any.whl (33 kB)
Requirement already satisfied: scikit-learn>=0.18 in /usr/local/lib/python3.7/dist-packages (from scikit-plot) (1.0.2)
Requirement already satisfied: matplotlib>=1.4.0 in /usr/local/lib/python3.7/dist-packages (from scikit-plot) (3.2.2)
Requirement already satisfied: joblib>=0.10 in /usr/local/lib/python3.7/dist-packages (from scikit-plot) (1.2.0)
Requirement already satisfied: scipy>=0.9 in /usr/local/lib/python3.7/dist-packages (from scikit-plot) (1.7.3)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.7/dist-packages (from matplotlib>=1.4.0->scikit-plot) (0.11.0)
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.7/dist-packages (from matplotlib>=1.4.0->scikit-plot) (1.4.4)
Requirement already satisfied: numpy>=1.11 in /usr/local/lib/python3.7/dist-packages (from matplotlib>=1.4.0->scikit-plot) (1.21.6)
Requirement already satisfied: pyparsing!=2.0.4,!=2.1.2,!=2.1.6,>=2.0.1 in /usr/local/lib/python3.7/dist-packages (from matplotlib>=1.4.0->scikit-plot) (3.0.9)
Requirement already satisfied: python-dateutil>=2.1 in /usr/local/lib/python3.7/dist-packages (from matplotlib>=1.4.0->scikit-plot) (2.8.2)
Requirement already satisfied: typing-extensions in /usr/local/lib/python3.7/dist-packages (from kiwisolver>=1.0.1->matplotlib>=1.4.0->scikit-plot) (4.1.1)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.7/dist-packages (from python-dateutil>=2.1->matplotlib>=1.4.0->scikit-plot) (1.15.0)
Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.7/dist-packages (from scikit-learn>=0.18->scikit-plot) (3.1.0)
Installing collected packages: scikit-plot
Successfully installed scikit-plot-0.3.7

In [ ]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
sns.set()

import re
import string
from wordcloud import WordCloud
from collections import Counter

import warnings
warnings.filterwarnings('ignore')

from nltk import sent_tokenize, word_tokenize
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer

from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import MultinomialNB
```

```
In [ ]:
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
sns.set()

import re
import string
from wordcloud import WordCloud
from collections import Counter

import warnings
warnings.filterwarnings('ignore')

from nltk import sent_tokenize, word_tokenize
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer

from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import MultinomialNB
from sklearn.ensemble import RandomForestClassifier

from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

from scikitplot.metrics import plot_confusion_matrix, plot_roc
```

```
In [ ]:
data = pd.read_csv('/content/sample_data/SPAM text message 20170820 - Data.csv')
data.head()
```

```
Out[ ]:
  Category      Message
0      ham  Go until jurong point, crazy.. Available only ...
1      ham             Ok lar... Joking wif u oni...
2     spam  Free entry in 2 a wkly comp to win FA Cup fina...
3      ham  U dun say so early hor... U c already then say...
4      ham  Nah I don't think he goes to usf, he lives aro...
```

```
In [ ]:
print(data.shape)

(5572, 2)
```

```
In [ ]:
print(data.shape)

(5572, 2)
```

```
In [ ]:
data.isnull().sum()
```

```
Out[ ]:
Category    0
Message     0
dtype: int64
```

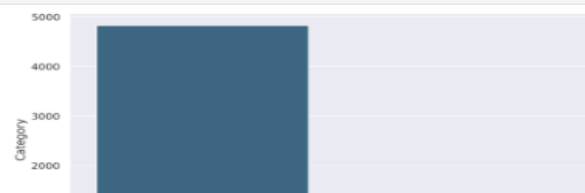
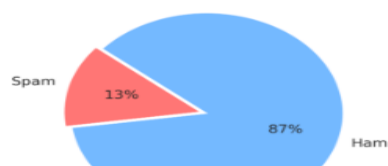
```
In [ ]:
data['Category'].value_counts()
```

```
Out[ ]:
ham      4825
spam      747
Name: Category, dtype: int64
```

```
In [ ]:
labels = ['Spam', 'Ham']
sizes = [747, 4825]
custom_colours = ['#ff7675', '#74b9ff']

plt.figure(figsize=(20, 6), dpi=227)
plt.subplot(1, 2, 1)
plt.pie(sizes, labels = labels, textprops={'fontsize': 15}, startangle=140,
        autopct='%1.0f%%', colors=custom_colours, explode=[0, 0.05])

plt.subplot(1, 2, 2)
sns.barplot(x = data['Category'].unique(), y = data['Category'].value_counts(), palette= 'viridis')
plt.show()
```

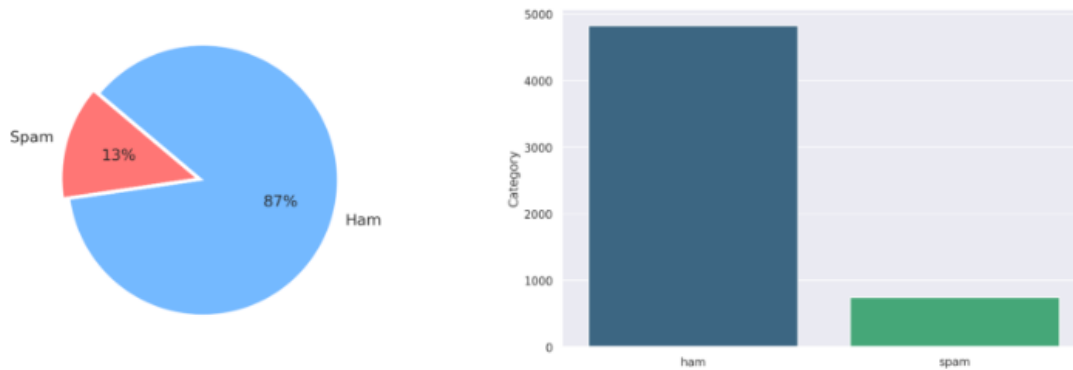


```
In [ ]:
labels = ['Spam', 'Ham']
sizes = [747, 4825]
custom_colours = ['#ff7675', '#74b9ff']

plt.figure(figsize=(20, 6), dpi=227)
plt.subplot(1, 2, 1)
plt.pie(sizes, labels = labels, textprops={'fontsize': 15}, startangle=140,
        autopct='%1.0f%%', colors=custom_colours, explode=[0, 0.05])

plt.subplot(1, 2, 2)
sns.barplot(x = data['Category'].unique(), y = data['Category'].value_counts(), palette= 'viridis')

plt.show()
```



```
In [ ]:
data['Total Words'] = data['Message'].apply(lambda x: len(x.split()))

def count_total_words(text):
    char = 0
    for word in text.split():
        char += len(word)
    return char

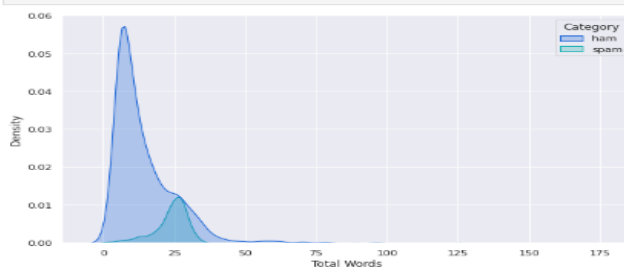
data['Total Chars'] = data['Message'].apply(count_total_words)
```

```
In [ ]:
data.head()
```

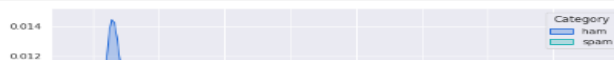
```
In [ ]:
data.head()
```

```
Out[ ]:
   Category  Message  Total Words  Total Chars
0    ham  Go until jurong point, crazy.. Available only ...      20        92
1    ham  Ok lar... Joking wif u oni...                6        24
2   spam  Free entry in 2 a wkly comp to win FA Cup fina...      28       128
3    ham  U dun say so early hor... U c already then say...      11        39
4    ham  Nah I don't think he goes to usf, he lives aro...
```

```
In [ ]:
plt.figure(figsize = (10, 6))
sns.kdeplot(x = data['Total Words'], hue= data['Category'], palette= 'winter', shade = True)
plt.show()
```

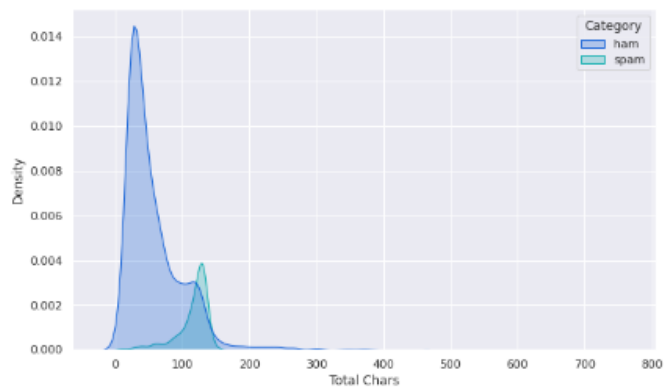


```
In [ ]:
plt.figure(figsize = (10, 6))
sns.kdeplot(x = data['Total Chars'], hue= data['Category'], palette= 'winter', shade = True)
plt.show()
```



In [ ]:

```
plt.figure(figsize = (10, 6))
sns.kdeplot(x = data['Total Chars'], hue= data['Category'], palette= 'winter', shade = True)
plt.show()
```



In [ ]:

```
def convert_lowercase(text):
    text = text.lower()
    return text

data['Message'] = data['Message'].apply(convert_lowercase)
```

In [ ]:

```
def remove_url(text):
    re_url = re.compile('https?://\S+|www\.\S+')
    return re_url.sub('', text)

data['Message'] = data['Message'].apply(remove_url)
```

In [ ]:

```
exclude = string.punctuation

def remove_punc(text):
    return text.translate(str.maketrans('', '', exclude))

data['Message'] = data['Message'].apply(remove_punc)
```

```
In [ ]:
def convert_lowercase(text):
    text = text.lower()
    return text

data['Message'] = data['Message'].apply(convert_lowercase)
```

```
In [ ]:
def remove_url(text):
    re_url = re.compile('https?://\S+|www\.\S+')
    return re_url.sub('', text)

data['Message'] = data['Message'].apply(remove_url)
```

```
In [ ]:
exclude = string.punctuation

def remove_punc(text):
    return text.translate(str.maketrans('', '', exclude))

data['Message'] = data['Message'].apply(remove_punc)
```

```
In [ ]:
!pip install nltk

Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: nltk in /usr/local/lib/python3.7/dist-packages (3.7)
Requirement already satisfied: regex>=2021.8.3 in /usr/local/lib/python3.7/dist-packages (from nltk) (2022.6.2)
Requirement already satisfied: click in /usr/local/lib/python3.7/dist-packages (from nltk) (7.1.2)
Requirement already satisfied: tqdm in /usr/local/lib/python3.7/dist-packages (from nltk) (4.64.1)
Requirement already satisfied: joblib in /usr/local/lib/python3.7/dist-packages (from nltk) (1.2.0)
```

```
In [ ]:
import nltk
nltk.download('punkt')
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Unzipping tokenizers/punkt.zip.
```

```
Out[ ]: True
```

```
In [ ]:
nltk.download('stopwords')
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Unzipping corpora/stopwords.zip.
```

```
Out[ ]: True
```

```
In [ ]:
def remove_stopwords(text):
    new_list = []
    words = word_tokenize(text)
    stopwords = stopwords.words('english')
    for word in words:
        if word not in stopwords:
            new_list.append(word)
    return ' '.join(new_list)

data['Message'] = data['Message'].apply(remove_stopwords)
```

```
In [ ]:
def perform_stemming(text):
    stemmer = PorterStemmer()
    new_list = []
    words = word_tokenize(text)
    for word in words:
        new_list.append(stemmer.stem(word))

    return ' '.join(new_list)

data['Message'] = data['Message'].apply(perform_stemming)
```

```
In [ ]:
data['Total Words After Transformation'] = data['Message'].apply(lambda x: np.log(len(x.split())))
```

```
In [ ]:
data.head()
```

```
Out[ ]:
```

	Category	Message	Total Words	Total Chars	Total Words After Transformation
0	ham	go jurong point crazy avail bugi n great world...	20	92	2.772589
1	ham	ok lar joke wif u oni	6	24	1.791759
2	spam	free entri 2 wkli comp win fa cup final tkt 21...	28	128	3.135494
3	ham	u dun say earli hor u c already say	11	39	2.197225
4	ham	nah dont think goe usf live around though	13	49	2.079442

```
In [ ]:
text = ' '.join(data[data['Category'] == 'spam']['Message'])
plt.figure(figsize = (15, 10))
wordcloud = WordCloud(max_words=500, height= 800, width = 1500, background_color="black", colormap= 'viridis').generate(text)
plt.imshow(wordcloud, interpolation="bilinear")
plt.axis('off')
plt.show()
```

Output :

```
text = " ".join(data[data['Category'] == 'spam']['Message'])
plt.figure(figsize = (15, 10))
wordcloud = WordCloud(max_words=500, height= 800, width = 1500, background_color="black", colormap= 'viridis').generate(text)
plt.imshow(wordcloud, interpolation="bilinear")
plt.axis('off')
plt.show()
```



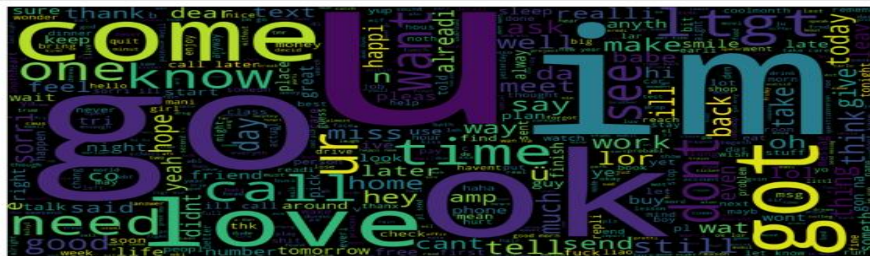
```

text = " ".join(data[data['Category'] == 'ham']['Message'])
plt.figure(figsize = (15, 10))
wordcloud = WordCloud(max_words=500, height= 800, width = 1500, background_color="black", colormap= 'viridis').generate(text)
plt.imshow(wordcloud, interpolation="bilinear")
plt.axis('off')
plt.show()

```



```
In [ ]: text = ".join(data[data['Category'] == 'ham']['Message'])
plt.figure(figsize = (15, 10))
wordcloud = WordCloud(max_words=500, height= 800, width = 1500, background_color="black", colormap= 'viridis').generate(text)
plt.imshow(wordcloud, interpolation="bilinear")
plt.axis('off')
plt.show()
```



```
In [ ]: data['Category'] = data['Category'].replace({'span':0, 'ham':1})
```

```
In [ ]: all_span_words = []
for sentence in data[data['Category'] == 0]['Message'].to_list():
    for word in sentence.split():
        all_span_words.append(word)

df = pd.DataFrame(Counter(all_span_words).most_common(25), columns= ['Word', 'Frequency'])

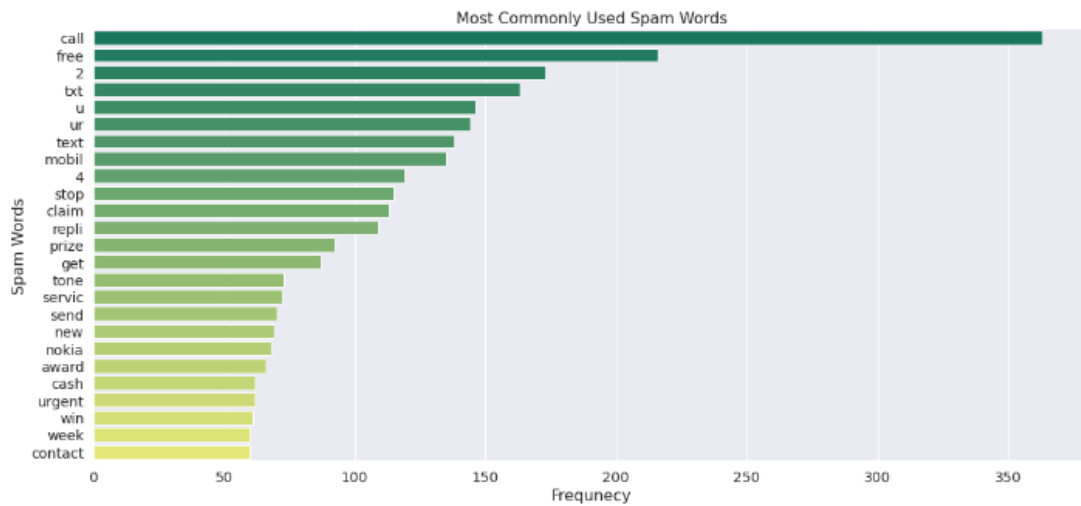
sns.set_context('notebook', font_scale= 1.3)
plt.figure(figsize=(18,8))
sns.barplot(y = df['Word'], x = df['Frequency'], palette= 'summer')
plt.title("Most Commonly Used Span Words")
plt.xlabel("Frequency")
plt.ylabel("Span Words")
plt.show()
```

```
In [ ]: data['Category'] = data['Category'].replace({'spam':0,'ham':1})
```

```
In [ ]: all_span_words = []
for sentence in data[data['Category'] == 0]['Message'].to_list():
    for word in sentence.split():
        all_span_words.append(word)

df = pd.DataFrame(Counter(all_span_words).most_common(25), columns= ['Word', 'Frequency'])

sns.set_context('notebook', font_scale= 1.3)
plt.figure(figsize=(18,8))
sns.barplot(y = df['Word'], x= df['Frequency'], palette= 'summer')
plt.title("Most Commonly Used Spam Words")
plt.xlabel("Frequency")
plt.ylabel("Spam Words")
plt.show()
```

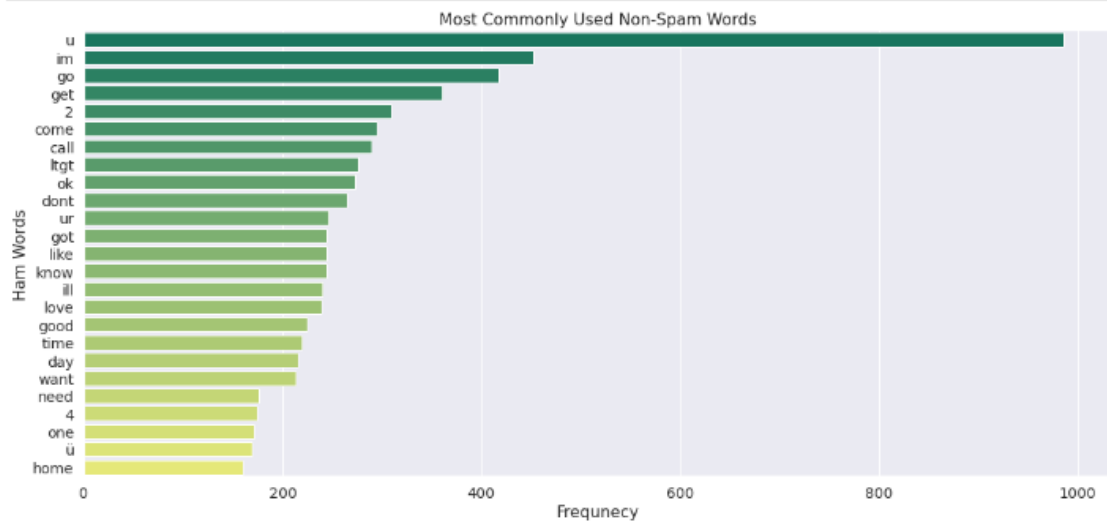


In [ ]:

```
all_ham_words = []
for sentence in data[data['Category'] == 1]['Message'].to_list():
    for word in sentence.split():
        all_ham_words.append(word)

df = pd.DataFrame(Counter(all_ham_words).most_common(25), columns= ['Word', 'Frequency'])

sns.set_context('notebook', font_scale= 1.3)
plt.figure(figsize=(18,8))
sns.barplot(y = df['Word'], x= df['Frequency'], palette= 'summer')
plt.title("Most Commonly Used Non-Spam Words")
plt.xlabel("Frequency")
plt.ylabel("Ham Words")
plt.show()
```



In [ ]:

```
X = data["Message"]
y = data["Category"].values

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size= 0.2, random_state= 42, stratify = y)
```



```
In [ ]: X = data["Message"]
y = data["Category"].values

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size= 0.2, random_state= 42, stratify = y)
```

```
In [ ]: tfidf = TfidfVectorizer(max_features= 2500, min_df= 2)
X_train = tfidf.fit_transform(X_train).toarray()
X_test = tfidf.transform(X_test).toarray()
```

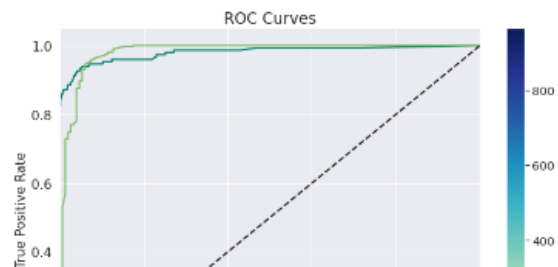
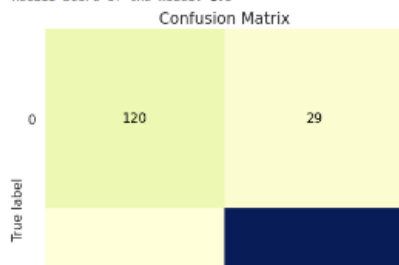
```
In [ ]: def train_model(model):
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    y_prob = model.predict_proba(X_test)
    accuracy = round(accuracy_score(y_test, y_pred), 3)
    precision = round(precision_score(y_test, y_pred), 3)
    recall = round(recall_score(y_test, y_pred), 3)

    print(f'Accuracy of the model: {accuracy}')
    print(f'Precision Score of the model: {precision}')
    print(f'Recall Score of the model: {recall}')

    sns.set_context('notebook', font_scale= 1.3)
    fig, ax = plt.subplots(1, 2, figsize = (25, 8))
    ax1 = plot_confusion_matrix(y_test, y_pred, ax= ax[0], cmap= 'YlGnBu')
    ax2 = plot_roc(y_test, y_prob, ax= ax[1], plot_macro= False, plot_micro= False, cmap= 'summer')
```

```
In [ ]: rf = RandomForestClassifier(n_estimators= 300)
train_model(rf)
```

Accuracy of the model: 0.974  
Precision Score of the model: 0.971  
Recall Score of the model: 1.0

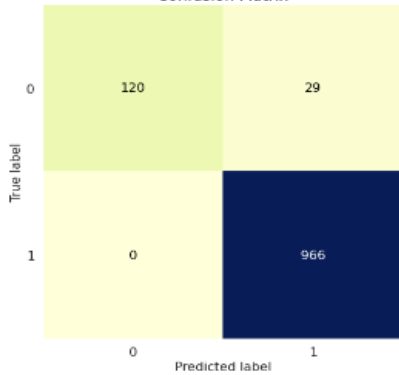


In [ ]:

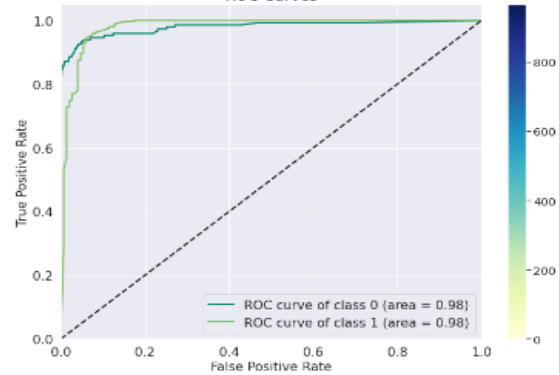
```
rf = RandomForestClassifier(n_estimators= 300)
train_model(rf)
```

Accuracy of the model: 0.974  
Precision Score of the model: 0.971  
Recall Score of the model: 1.0

Confusion Matrix



ROC Curves



## Conclusion:

1. Performed classification analysis in Python for spam classifier dataset and build a model
2. Few insights we found from the dataset:
  - The accuracy calculated of the model was to be 97.4%
  - The ROC curve found to be constant at 1 after increasing at a constant from 0.2 .