Assignment Date: 03 November 2022

Student Name: R K Tharun Kumar

1. Download the dataset link

- Label Ham or Spam
- Message Message

```
import warnings
warnings.filterwarnings("ignore")
```

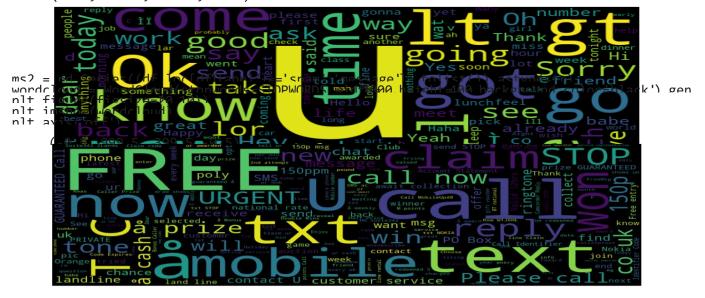
2. Importing Required Library

```
import re import nltk import pandas as
pd import numpy as np import
matplotlib.pyplot as plt from
nltk.stem import WordNetLemmatizer
from nltk.corpus import stopwords
from wordcloud import WordCloud,STOPWORDS,ImageColorGenerator
```

3. Read dataset and do Preprocessing

```
df = pd.read_csv("/content/spam.csv",encoding='ISO-8859-1')
df = df.iloc[:,:2]
df.columns=['label','message']
df.head()
          label
                                                       message
      0
                     Go until jurong point, crazy.. Available only ...
           ham
       1
           ham
                                       Ok lar... Joking wif u oni...
      2 spam
                  Free entry in 2 a wkly comp to win FA Cup fina...
                   U dun say so early hor... U c already then
      3
           ham
```

say...



from nltk.stem.wordnet import
WordNetLemmatizer lemmatizer =
WordNetLemmatizer() corpus = []

import nltk from nltk.corpus
import stopwords
nltk.download('all')

```
for i in range(len(df)):
    review = re.sub('[^a-zA-Z]',' ',df['message'][i])
                                                          review = review.lower()
                                                                                       review
                     review = [lemmatizer.lemmatize(i) for i in review if not i in
= review.split()
                                review = ' '.join(review)
set(stopwords.words('englis
                                                              corpus.append(review)
     [nltk_data]
                        Unzipping corpora/pe08.zip.
                    | Downloading package perluniprops to [nltk_data]
     [nltk data]
     /root/nltk_data...
                        Unzipping misc/perluniprops.zip.
     [nltk data]
                    | Downloading package pil to /root/nltk_data...
     [nltk_data]
     [nltk_data]
                        Unzipping corpora/pil.zip.
     [nltk_data]
                      Downloading package pl196x to /root/nltk_data...
     [nltk_data]
                        Unzipping corpora/pl196x.zip.
                     Downloading package porter_test to /root/nltk_data...
     [nltk_data]
                        Unzipping stemmers/porter_test.zip. [nltk_data]
                                                                           Downloading
     [nltk_data]
     package ppattach to /root/nltk_data...
                        Unzipping corpora/ppattach.zip.
     [nltk_data]
                      Downloading package problem_reports to
     [nltk_data]
                          /root/nltk_data... [nltk_data]
     [nltk data]
                                                                Unzipping
     corpora/problem_reports.zip.
     [nltk_data]
                    | Downloading package product_reviews_1 to [nltk_data]
     /root/nltk_data...
     [nltk data]
                        Unzipping corpora/product reviews 1.zip. [nltk data]
     Downloading package product_reviews_2 to [nltk_data]
                                                                    /root/nltk data...
     [nltk_data]
                        Unzipping corpora/product_reviews_2.zip. [nltk_data]
     Downloading package propbank to /root/nltk_data...
                    Downloading package pros_cons to /root/nltk_data...
     [nltk_data]
                        Unzipping corpora/pros_cons.zip.
     [nltk_data]
     [nltk_data]
                      Downloading package ptb to /root/nltk_data...
                        Unzipping corpora/ptb.zip.
     [nltk_data]
                    Downloading package punkt to /root/nltk_data...
     [nltk_data]
                        Unzipping tokenizers/punkt.zip.
     [nltk_data]
                      Downloading package qc to /root/nltk_data...
     [nltk data]
                        Unzipping corpora/qc.zip.
     [nltk_data]
     [nltk_data]
                    Downloading package reuters to /root/nltk_data...
     [nltk data]
                      Downloading package rslp to /root/nltk data...
     [nltk_data]
                        Unzipping stemmers/rslp.zip.
                    Downloading package rte to /root/nltk_data...
     [nltk_data]
     [nltk_data]
                        Unzipping corpora/rte.zip.
                      Downloading package sample grammars to
     [nltk_data]
     [nltk_data]
                          /root/nltk_data...
                        Unzipping grammars/sample_grammars.zip.
     [nltk_data]
                      Downloading package semcor to /root/nltk_data...
     [nltk_data]
                      Downloading package senseval to /root/nltk data...
     [nltk_data]
     [nltk_data]
                        Unzipping corpora/senseval.zip.
                    Downloading package sentence_polarity to [nltk_data]
     [nltk_data]
     /root/nltk_data...
     [nltk data]
                        Unzipping corpora/sentence polarity.zip. [nltk data]
     Downloading package sentiwordnet to [nltk_data]
                                                              /root/nltk_data...
                        Unzipping corpora/sentiwordnet.zip.
     [nltk_data]
                    Downloading package shakespeare to /root/nltk_data... [nltk_data]
     [nltk data]
         Unzipping corpora/shakespeare.zip.
     [
               ]
                            pp g
                                                        р
```

р

р

```
[nltk_data]
               Downloading package sinica_treebank to
[nltk data]
                     /root/nltk data... [nltk data]
                                                           Unzipping
corpora/sinica_treebank.zip.
[nltk_data]
               Downloading package smultron to /root/nltk data...
[nltk_data]
                   Unzipping corpora/smultron.zip.
                 Downloading package snowball data to
[nltk data]
[nltk_data]
                     /root/nltk_data...
[nltk_data]
                Downloading package spanish grammars to
[nltk_data]
                    /root/nltk_data...
```

▼ 4. Create Model

```
from keras.preprocessing.text import Tokenizer from
keras_preprocessing.sequence import pad_sequences
from keras.layers import Dense, Dropout, LSTM, Embedding
from keras.models import Sequential, load model
token = Tokenizer()
token.fit_on_texts(corpus) text_to_seq =
token.texts_to_sequences(corpus)
max_length_sequence = max([len(i) for i in text_to_seq]) padded_seq =
pad_sequences(text_to_seq, maxlen=max_length_sequence, padding="pre")
padded_seq
                          0, ..., 16, 3551,
    array([[
               0,
                     0, ..., 359,
                                     1, 1610],
                           0, ..., 218,
            0, ..., 7042, 1095, 3547],
               0,
                     0,
                     0,
                          0, ..., 842,
                                            1,
               0,
                          0, ..., 2198, 347, 152]], dtype=int32)
               0,
from sklearn.preprocessing import
LabelEncoder le = LabelEncoder() y =
le.fit_transform(df['label'])
from sklearn.model_selection import train_test_split
X train, X test, y train, y test = train test split(padded seq, y, test size=0.25, random state=42)
X_train.shape
    (4179, 77)
```

5. Add Layers

```
TOT_SIZE = len(token.word_index) + 1
model = Sequential()
#IP Layer
model.add(Embedding(TOT_SIZE,32,input_length=max_length_sequence))
model.add(LSTM(units=50, activation = 'relu',return_sequences=True))
model.add(Dropout(0.2))
#Layer2 model.add(LSTM(units=60, activation =
'relu')) model.add(Dropout(0.3))
#output layer
model.add(Dense(units=1, activation='sigmoid'))

WARNING:tensorflow:Layer lstm will not use cuDNN kernels since it doesn't meet the crite
WARNING:tensorflow:Layer lstm_1 will not use cuDNN kernels since it doesn't meet the crite
```

model.summary()

Model: "sequential"

Layer (type)	Output Shape	Param #
embedding (Embedding)	 (None, 77, 32)	
lstm (LSTM)	(None, 77, 50)	16600
dropout (Dropout)	(None, 77, 50)	0
lstm_1 (LSTM)	(None, 60)	26640
dropout_1 (Dropout)	(None, 60)	0
dense (Dense)	(None, 1)	61

Total params: 268,709
Trainable params: 268,709

Non-trainable params: 0

6 Compile the model

model.compile(optimizer='adam', loss='binary_crossentropy',metrics=['accuracy'])

7 Fit the model

```
Epoch 3/10
 Epoch 4/10
 Epoch 5/10
 Epoch 6/10
 Epoch 7/10
 131/131 [================== ] - 33s 249ms/step - loss: 0.0427 - accuracy: 0.9
 Epoch 8/10
 131/131 [=================== ] - 32s 244ms/step - loss: 0.1001 - accuracy: 0.9
 Epoch 9/10
 Epoch 10/10
 <keras.callbacks.History at 0x7feed0948990>
model.evaluate(X_test,y_test)
```

8. Save the Model

```
from pickle import
dump,load tfid = 'tfid.sav'
lstm = 'lstm.sav'

dump(token,open(tfid,'wb'))
model.save('nlp.h5')
```

9. Test the Model

```
review = re.sub('[^a-zA-Z]',' ',raw_mess)
def preprocess(raw_mess):
                  review = review.split()
                                            review = [lemmatizer.lemmatize(i) for i in
review.lower()
                                                  review = ' '.join(review)
review if not i in set(stopwords.words('englis
review
def predict(mess):
    vect = load(open(tfid, 'rb'))
                                    classifier =
load_model('nlp.h5')
                       clean = preprocess(mess)
                                                     text_to_seq =
token.texts_to_sequences([mess])
                                    padded seg =
```

```
pad_sequences(text_to_seq, maxlen=77, padding="pre")
                                                      pred =
classifier.predict(padded seq)
                                 return pred
msg = input("Enter a message:
") predi = predict(msg) if
predi >= 0.6:
   print("It is a
spam") else:
print("Not a spam")
    Enter a message: "Thanks for your Ringtone Order, Reference T91. You will be charged GBP
    WARNING:tensorflow:Layer lstm will not use cuDNN kernels since it doesn't meet the crite
    WARNING:tensorflow:Layer lstm 1 will not use cuDNN kernels since it doesn't meet the cri
    1/1 [======= ] - 0s 284ms/step It
    is a spam
    msg = input("Enter a message:
") predi = predict(msg) if
predi >= 0.6:
   print("It is a
spam") else:
print("Not a spam")
    Enter a message: Keep my payasam there if rinu brings,,,
    WARNING:tensorflow:Layer lstm will not use cuDNN kernels since it doesn't meet the crite
    WARNING:tensorflow:Layer lstm_1 will not use cuDNN kernels since it doesn't meet the cri
    1/1 [======== ] - 0s 250ms/step Not
    a spam
```