## Assignment-4

SMS Spam Classification

Assignment Date	30 September 2022
Student Name	Guru Hari Venkat S
Student Roll Number	513419106012
Maximum Marks	2 Marks

## **SMS SPAM Classification**

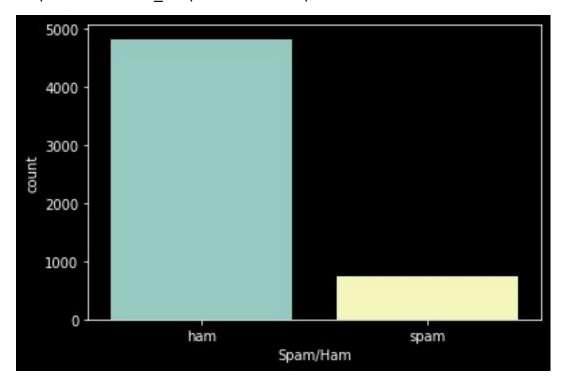
```
Import required libraries
```

```
import pandas as pd
import numpy as np
import re
import collections
import contractions
import seaborn as sns
import matplotlib.pyplot as plt
plt.style.use('dark background')
import nltk
from nltk.stem import WordNetLemmatizer
from nltk.corpus import stopwords
import warnings
warnings.simplefilter(action='ignore', category=Warning)
import keras
from keras.layers import Dense, Embedding, LSTM, Dropout
from keras.models import Sequential
from keras.preprocessing.text import Tokenizer
from keras preprocessing.sequence import pad sequences
Download the dataset
df = pd.read csv("spam.csv", encoding='ISO-8859-1')
df.shape
(5572, 5)
Read dataset and do pre-processing
df.head(10)
     v1 ... Unnamed: 4
0
                    NaN
```

```
v1 ... Unnamed: 4
0 ham ... NaN
1 ham ... NaN
2 spam ... NaN
```

```
3
                     NaN
    ham
4
    ham
                     NaN
5
   spam
                     NaN
6
    ham
                     NaN
         . . .
7
    ham
                     NaN
         . . .
8
   spam
                     NaN
9
                     NaN
   spam
         . . .
[10 rows x 5 columns]
df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5572 entries, 0 to 5571
Data columns (total 5 columns):
     Column
                 Non-Null Count
                                  Dtype
- - -
     _ _ _ _ _ _
                  _____
                                  ----
 0
                 5572 non-null
     v1
                                  object
 1
                 5572 non-null
                                  object
     v2
 2
     Unnamed: 2 50 non-null
                                  object
 3
     Unnamed: 3
                 12 non-null
                                  object
     Unnamed: 4 6 non-null
 4
                                  object
dtypes: object(5)
memory usage: 217.8+ KB
df.isnull().any()
v1
              False
v2
              False
Unnamed: 2
               True
Unnamed: 3
               True
Unnamed: 4
               True
dtype: bool
df.isnull().sum()
                 0
ν1
v2
                 0
Unnamed: 2
              5522
Unnamed: 3
              5560
Unnamed: 4
              5566
dtype: int64
df.drop(["Unnamed: 2", "Unnamed: 3", "Unnamed: 4"], axis=1,
inplace=True)
df.columns = ["Spam/Ham", "Message"]
df.head()
  Spam/Ham
                                                        Message
0
            Go until jurong point, crazy.. Available only ...
                                 Ok lar... Joking wif u oni...
1
       ham
```

```
spam Free entry in 2 a wkly comp to win FA Cup fina...
ham U dun say so early hor... U c already then say...
ham Nah I don't think he goes to usf, he lives aro...
sns.countplot(df["Spam/Ham"])
<matplotlib.axes. subplots.AxesSubplot at 0x7f4de8b9e610>
```



```
def preprocessing(data):
      sms = contractions.fix(data)
      sms = sms.lower()
      sms = re.sub(r'https?://S+|www.S+', "", sms).strip()
      sms = re.sub("[^a-z]", "", sms)
      return sms
X = df["Message"].apply(preprocessing)
from sklearn.preprocessing import LabelEncoder
lb enc = LabelEncoder()
y = lb enc.fit transform(df["Spam/Ham"])
tokenizer = Tokenizer()
tokenizer.fit on texts(X)
text to sequence = tokenizer.texts to sequences(X)
max length sequence = max([len(i) for i in text to sequence])
padded sms_sequence = pad_sequences(text_to_sequence,
maxlen=max_length_sequence,
                                    padding = "pre")
padded sms sequence
```

```
0, ..., 50, 3915, 134],
0, ..., 419, 1, 1715],
                  0,
array([[
           0,
                  0,
           0,
                        0, ..., 2653, 348, 2654],
           0,
                  0,
           0,
                        0, ..., 8472, 222, 8473],
                  0,
           0,
                  0,
                        0, ..., 144, 17, 45],
                        0. .... 3.
                                       61, 233]], dtype=int32)
           0.
                  0.
```

## **Create Model**

```
Add Layers (LSTM, Dense-(Hidden Layers), Output)
TOT_SIZE = len(tokenizer.word_index)+1
def create_model():
    lstm_model = Sequential()
    lstm_model.add(Embedding(TOT_SIZE, 32,
input_length=max_length_sequence))
    lstm_model.add(LSTM(100))
    lstm_model.add(Dropout(0.4))
    lstm_model.add(Dense(20, activation="relu"))
    lstm_model.add(Dropout(0.3))
    lstm_model.add(Dense(1, activation = "sigmoid"))
    return lstm_model
lstm_model = create_model()
```

## **Compile the Model**

```
lstm_model.compile(loss = "binary_crossentropy", optimizer = "adam",
metrics = ["accuracy"])
```

lstm\_model.summary()

Model: "sequential"

Layer (type)	Output Shape	Param #
embedding (Embedding)	(None, 172, 32)	271200
lstm (LSTM)	(None, 100)	53200
dropout (Dropout)	(None, 100)	0
dense (Dense)	(None, 20)	2020
dropout_1 (Dropout)	(None, 20)	0
dense_1 (Dense)	(None, 1)	21

\_\_\_\_\_\_

Total params: 326,441 Trainable params: 326,441

```
Fit the Model
lstm model.fit(padded sms sequence, y, epochs = 5,
validation split=0.2, batch size=16)
Epoch 1/5
0.2117 - accuracy: 0.9361 - val loss: 0.0662 - val accuracy: 0.9821
Epoch 2/5
0.1441 - accuracy: 0.9746 - val loss: 0.0819 - val accuracy: 0.9830
Epoch 3/5
0.0272 - accuracy: 0.9953 - val loss: 0.0638 - val accuracy: 0.9848
Epoch 4/5
0.0135 - accuracy: 0.9980 - val loss: 0.0622 - val accuracy: 0.9865
Epoch 5/5
0.0111 - accuracy: 0.9978 - val loss: 0.0684 - val accuracy: 0.9839
<keras.callbacks.History at 0x7f4de753af10>
Save The Model
lstm model.save('sms.h5')
Test The Model
def predict spam(predict msg):
   new seq = tokenizer.texts to sequences(predict msg)
   padded = pad sequences(new seq, maxlen =max length sequence,
                padding = 'pre')
   return (lstm model.predict(padded))
message = str(input('Enter your message: '))
pred = predict spam(message)
print(f"Prediction: {pred[0][0]}\n")
print(f"The sms is a {'spam' if pred[0][0]<0.5 else 'ham'}")</pre>
Enter your message: you've won something, the IRS is trying to contact
you, you have a refund coming or asking you to verify your bank
account
Prediction: 0.005271350499242544
The sms is a spam
```