
Date : 01.11.2022

Name : Monisha P S

Roll no : 611219106048

Assignment 4 - SMS SPAM Classification

▼ 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 the dataset and do Preprocessing

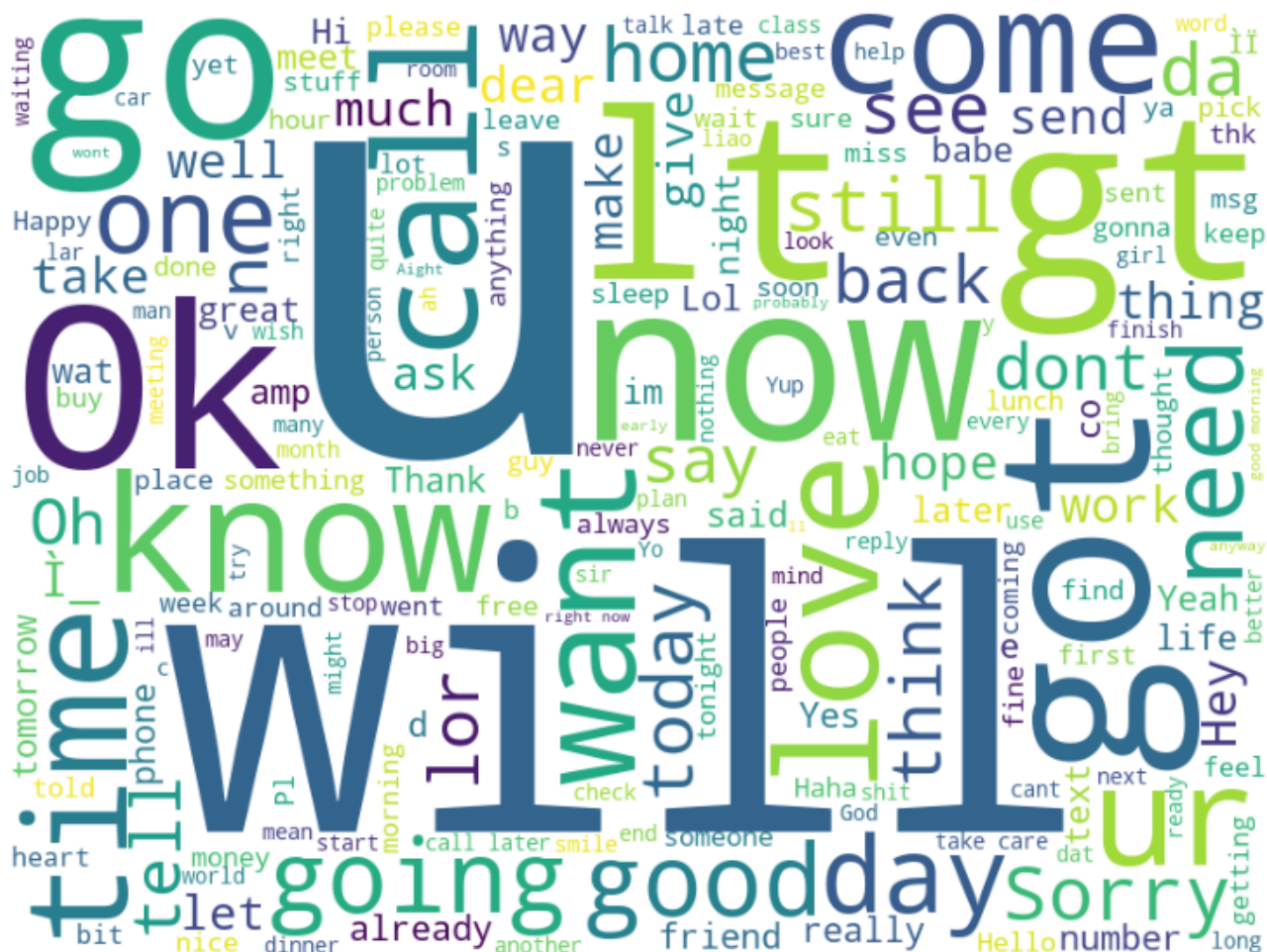
```
df = pd.read_csv("/content/spam.csv", encoding='ISO-8859-1')

df = df.iloc[:, :2]
df.columns = ['label', 'message']
df.head()
```

	label	message
0	ham	Go until jurong point, crazy.. Available only ...
1	ham	Ok lar... Joking wif u oni...

```
df.info()
```

```
ms1 = pd.Series((df.loc[df['label']=='ham', 'message']).tolist()).astype(str)
wordcloud = WordCloud(stopwords=STOPWORDS, width=800, height=600, background_color='white').g
plt.figure(figsize=(20,10))
plt.imshow(wordcloud)
plt.axis('off')
```




```
review = ' '.join(review)
corpus.append(review)
```

▼ 4. Creating the 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
```

```
array([[ 0,  0,  0, ..., 16, 3551,  70],
       [ 0,  0,  0, ..., 359,   1, 1610],
       [ 0,  0,  0, ..., 218,  29,  293],
       ...,
       [ 0,  0,  0, ..., 7042, 1095, 3547],
       [ 0,  0,  0, ...,  842,   1,   10],
       [ 0,  0,  0, ..., 2198,  347,  152]], dtype=int32)
```

```
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=
```

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

```



```
model.summary()
```

Model: "sequential"

Layer (type)	Output Shape	Param #
=====		
embedding (Embedding)	(None, 77, 32)	225408
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


```
model.fit(X_train, y_train,validation_data=(X_test,y_test), epochs=10)
```

```

Epoch 1/10
131/131 [=====] - 45s 307ms/step - loss: 0.3592 - accuracy:
Epoch 2/10
131/131 [=====] - 39s 293ms/step - loss: 13934574.0000 - acc
Epoch 3/10
131/131 [=====] - 36s 278ms/step - loss: 241499.6562 - accur


```

```
Epoch 4/10
131/131 [=====] - 36s 279ms/step - loss: 0.0626 - accuracy:
Epoch 5/10
131/131 [=====] - 39s 298ms/step - loss: 0.0377 - accuracy:
Epoch 6/10
131/131 [=====] - 37s 284ms/step - loss: 0.0327 - accuracy:
Epoch 7/10
131/131 [=====] - 36s 278ms/step - loss: 0.0206 - accuracy:
Epoch 8/10
131/131 [=====] - 37s 279ms/step - loss: 0.0141 - accuracy:
Epoch 9/10
131/131 [=====] - 35s 270ms/step - loss: 0.0096 - accuracy:
Epoch 10/10
131/131 [=====] - 37s 285ms/step - loss: 0.0072 - accuracy:
<keras.callbacks.History at 0x7f7e305dc110>
```



```
model.evaluate(X_test,y_test)
```

```
44/44 [=====] - 1s 24ms/step - loss: 0.1467 - accuracy: 0.98
[0.14665858447551727, 0.9813352227210999]
```



▼ 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

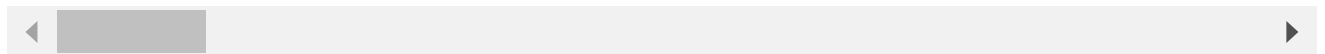
```
def preprocess(raw_mess):
    review = re.sub('[^a-zA-Z]', ' ',raw_mess)
    review = review.lower()
    review = review.split()
    review = [lemmatizer.lemmatize(i) for i in review if not i in set(stopwords.words('eng
    review = ' '.join(review)
    return review
```

```
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_seq = pad_sequences(text_to_seq, maxlen=77, padding="pre")
pred = classifier.predict(padded_seq)
return pred
```

```
msg = input("Enter a input message: ")
predi = predict(msg)
if predi >= 0.6:
    print("It is a spam")
else:
    print("Not a spam")
```

```
Enter a input message: Its a part of checking IQ
WARNING:tensorflow:Layer lstm will not use cuDNN kernels since it doesn't meet the cr
WARNING:tensorflow:Layer lstm_1 will not use cuDNN kernels since it doesn't meet the
WARNING:tensorflow:5 out of the last 5 calls to <function Model.make_predict_function
1/1 [=====] - 0s 317ms/step
Not a spam
```



```
msg = input("Enter a input message: ")
predi = predict(msg)
if predi >= 0.6:
    print("It is a spam")
else:
    print("Not a spam")
```

```
Enter a input message: England v Macedonia - dont miss the goals/team news. Txt ur na
WARNING:tensorflow:Layer lstm will not use cuDNN kernels since it doesn't meet the cr
WARNING:tensorflow:Layer lstm_1 will not use cuDNN kernels since it doesn't meet the
WARNING:tensorflow:6 out of the last 6 calls to <function Model.make_predict_function
1/1 [=====] - 0s 296ms/step
It is a spam
```

