Assignment-4 (SMS SPAM Classification)

```
import numpy as np import.pandas.as.pd import.seaborn.as.sns import
matplotlib.pyplot as plt
data = pd.read csv('/content/sample data/spam.csv', delimiter=', ',
encoding='latin-1') data.head()
v1
V2 Unnamed: 2 Unnamed: 3 Unnamed: 4
()
ham
NaN
NaN
NaN
Go until jurong point, crazy.. Available only ...
Ok lar... Joking wif u oni...
1
ham
NaN
NaN
NaN
Free entry in 2 a wkly comp to win FA Cup fina...
NaN
NaN
NaN
2 3
```

```
spam ham
U dun say so early hor... U c already then say...
NaN
NaN
4
ham
Nah I don't think he goes to usf, he lives aro...
NaN
NaN
NaN
data.columns
Index(['vi', 'V2', 'Unnamed: 2', 'Unnamed: 3', 'Unnamed: 4'],
dtype='object')
#drop the unamed columns data-data.drop(columns=["Unnamed: 2", "Unnamed:
3", "Unnamed: 4"])
#rename the two relevant columns data-data.rename(
"V1":"Category",
"V2": "Message" } , axis=1)
data.head()
https://colab.research.google.com/drive/1RUT5WjmG_-tjTK299qMNKS96wYxYODMm#scrollTo=47k2F
S4CDrqL&printMode=true
1/14
10/29/22, 11:33 AM
Assignment_4.ipynb - Colaboratory
Category
```

```
Message
```

KB

```
Go until jurong point, crazy.. Available only ...
ham ham
1
Ok lar... Joking wif u oni...
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... #check for null values
data.isnull().sum()
Category 0 message dtype: int64
data.info()
<class 'pandas.core. frame. DataFrame'> Range Index: 5572 entries, 0 to
5571 Data columns (total 2 columns):
# Column Non-Null Count Dtype --- -----
0 Category 5572 non-null object
```

1 Message 5572 non-null object dtypes: object(2) memory usage: 87.2+

```
data["Message Length"]=data["Message"].apply(len)
```

```
fig=plt.figure(figsize=(12,8)) sns.histplote
X=data["Message Length"], hue=data["Category"]
```

plt.title("ham & spam messege length comparision") plt.show()

```
S4CDrqL&printMode=true
2/14
10/29/22, 11:33 AM
Assignment 4.ipynb - Colaboratory
ham & spam messege length comparision
Category O ham
spam
Count
300 +
#Display the description of length of ham and spam messages
seperately on an individual serie
ham_desc=data[data["Category"]=="ham"]["Message
Length"].describe() spam_desc=data[data["Category"]
=="spam"]["Message Length"].describe() print("Ham
Messege Length Description:\n", ham_desc)
```

https://colab.research.google.com/drive/1RUT5WjmG_-tjTK299qMNKS96wYxYODMm#scrollTo=47k2F

```
print("Spam Message Length Description:\n", spam desc)
min
Ham Messege Length Description:
Count 4825.000000 mean
71.023627 std
58.016023
2.000000 25%
33.000000
52.000000 75%
92.000000 max
        910.00000 Name: Message Length, dtype: float64
50%
*********
Spam Message Length Description:
count 747.000000 mean 138.866131 sta
29.183082 min
13.000000 25% 132.500000 50% 149.000000 75% 157.000000 max
224.000000 Name: Message Length, dtype: float64
```

data.describe(include="all")

$https://colab.research.google.com/drive/1RUT5WjmG_-tjTK299qMNK \textbf{S96wYxYODMm} \\ \textit{\#scrollTo=47k2FS4CDrqL\&printMode=true}$

3/14

10/29/22, 11:33 AM

Assignment 4.jpynb - Colaboratory

Category

Message Message Length

count

5572

5**572**

5**572.000000**

unique

5169

NaN

top

ham

Sorry, I'll call later

NaN

freq

4825

30

NaN

mean

NaN

NaN

80.118808

```
NaN
NaN
59.690841
min
NaN
NaN
2.000000 36.000000
25%
NaN
NaN
50%
NaN
NaN
61.000000
data["Category"].value_counts()
ham 4825 spam 747 Name: Category, dtype: int64
sns.countplot(
data=data, x="Category"
plt.title("ham vs spam") plt.show()
ham vs spam
50001
4000
3000
count
2000
```

std

```
ham
spam
```

Category

```
ham_count=data["Category"].value_counts()[0]
spam_count=data["Category"].value_counts()[1]
```

total count=data.shape [0]

https://colab.research.google.com/drive/1RUT5WjmG_-tjTK299qMNKS96wYXYODM m#scrollTo=4762FS4CDqL&printMode=true

10/29/22, 11:33 AM

Assignment_4.ipynb - Colaboratory print("Ham contains:{:.2f}% of total data.".format(ham_count/total_count*100)) print("Spam contains:{:.2f}% of total data.", format(spam_count/total_count*100))

Ham contains:86.59% of total data. Spam contains:13.41% of total data.

#compute the length of majority & minority class minority_len=len (data[data["Category"]=="spam"])
majority_len=len(data[data["Category"]=="ham"])

#store the indices of majority and minority class
minority_indices=data[data["Category"]=="spam"].index
majority_indices=data[data["category"]=="ham"].index

#generate new majority indices from the total majority_indices #with size equal to minority class length so we obtain equivalent number of indices length random_majority_indices=np.random.choice
majority indices, size=minority len, replace=False

```
#concatenate the two indices to obtain indices of new dataframe
undersampled_indices=np.concatenate([minority indices, random_majority indices])
#create df using new indices df=data.loc[undersampled indices]
#shuffle the sample df=df.sample(frac=1)
#reset the index as its all mixed dfsdf.reset index()
#drop the older index df=df.drop
columns=["index"],
df.shape
(1494, 3)
df["Category"].value_counts()
ham 747 spam 747 Name: Category, dtype: int64
https://colab.research.google.com/drive/1RUT5WjmG_-tjTK299qMNKS96wYxYODMm#scrollTo=47k2FS
4CDrqL&printMode=true
5/14
10/29/22, 11:33 AM
Assignment 4.ipynb - Colaboratory
sns.countplot(
data-df, x="Category"
plt.title("ham vs spam") plt.show()
ham vs spam
count
```

ham spam Category

df.head()

Category

Message Message Length

ham

Aah! A cuddle would be lush! I'd need lots of ...

ham

I'm in solihull, do you want anything?

spam

Double Mins & 1000 txts on Orange tariffs. Lat...

ham

spam

No we put party 7 days a week and study light...

URGENT!! Your 4* Costa Del Sol Holiday or å£50...

#created new column Label and encode ham as 0 and spam as 1 df["Label"]=df["Category"].map

"ham":0, "spam":1

```
df.head()
```

 $https://colab.resear ch.google.com/drive/1RUT5WjmG_-tjTK299qMNK\\ \textbf{S96wYxYODMm\#scrollTo=47k2FS4CDrqL\&printMode=true}$

6/14

10/29/22, 11:33 AM

Assignment 4.ipynb - Colaboratory

Category

Message Message Length Label

ham

Aah! A cuddle would be lush! I'd need lots of ...

87

0

1.

ham

I'm in solihull, do you want anything?

40

0

spam

Double Mins & 1000 txts on Orange tariffs. Lat...

151

1

ham

No we put party 7 days a week and study lightl...

1260

import re import nitk

from nlt k.corpus import stopwords from nit k.stem import PorterStemmer

```
stemmer=porterStemmer()
nltk.download('stopwords')
[nltk data] Downloading package stopwords to /root/nltk data... [nltk data]
Unzipping Corpora/StopWords.zip. True
#declare empty list to store tokenized message COrpus=[]
#iterate through the df["Message"] for message in df["Message"]:
#replace every special characters, numbers etc.. with whitespace of
message #It will help retain only letter/alphabets message=re. sub("[^a-zA-Z]", "
", message)
#convert every letters to its lowercase message=message. lower()
#split the word into individual word list message=message.split()
#perform stemming using PorterStemmer for all non-english-stopwords
message=[stemmer.stem(words)
for words in message
if words not in Set(Stopwords.WordS("english"))
#join the word lists with the whitespace message=" ".join(message)
#append the message in corpus list corpus.append(message)
```

https://colab.research.google.com/drive/1RUT5WjmG_-tjTK299qMNKS96wYxYODMm#scrollTo=47k2FS4CDrqL&

Assignment 4.ipynb - Colaboratory from tensorflow.keras.preprocessing.text

printMode=true

10/29/22, 11:33 AM

import one hot vocab size=10000

```
oneHot_doc=[one_hot(words, n=vocab_size)
for words in corpus

df["Message Length"].describe()

count 1494.000000 mean 105.203481
std
61.166448 min
3.000000 25%
48.000000 50% 118.000000 75% 153.000000 max 790.000000 Name: Message
Length, dtype: float64

fig=plt.figure(figsize=(12,8)) sns.kdeplot(
X=df["Message Length"], hue=df["Category"]

plt.title("ham & spam messege length Comparision") plt.show()
```

https://colab.research.google.com/drive/1RUT5WjmG_-tjTK299qMNKS96wYxYODMm#scrollTo=47k2F S4CDrqL&printMode=true 8/14

10/29/22, 11:33 AM

Assignmen<u>t 4</u>.ipynb - Colaboratory

ham & spam messege length comparision

Cate gory

from tensorflow.keras.preprocessing. sequence import pad_sequences sentence_len=200 embedded_doc=pad_sequences (

```
oneHot_doc,
  maxlen=sentence len
  , padding="pre"
extract features=pd.Data
Frame
  data=embedded
  doc
target=df["label"]
     W.Ut i 1 i
df_final=pd.concat([extract_features, target], axis=1)
     0.0
     02
     +
                   77
df_final.head()
       1234 5 6 7 8 9 ... 191 192 193 194 195 196 197 198 1
            000000000000 ... 2090 1632 42897158 478 5808
   200000000000 . . . 1275 702 1694 4114 4162 3935 4162
            8536 72 30000000000 3705 9946 5462 7158
    98834500 8030 8630 29 4 0000000000 . . . 47536414 5018 1953
                                         216 1175 8861 2485 60
   5 rows x 201
   columns
```

X=df_final.drop ("Label", axis=1)

y=df_final["label"]

```
X_{trainval,x_test,y_trainval,y_test=train_te}
   st_split(
       Χ
       у,
       random_state=4
       test_size=0.15
   X train,x val,y train, y val=train tes
   t split(
       X trainval
https://colab.research.google.com/drive/1RUT5WjmG_-tjTK299qMNKS96wYxYODMm#scrollTo=47k2
FS4CDrqL&printMode=true
       Assignment 4.ipynb - Colaboratory
        10/29/22, 11:33 AM
       y_trainval, random state=42, test size=0.15
       from tensorflow.keras.layers import LSTM from tensorflow.keras.layers import Dense
       from tensorflow.keras.layers import Embedding from tensorflow.keras models
       import Sequential model=Sequential()
       feature num=100 model.add(
       Embedding
```

input_dim=vocab size, output_dim=feature_num,

9/1

from sklearn.model_selection import

train_test Split

```
input_length=sentence len
model.add(
LSTMC units=128
model.add( Dense(
units=1, activation="sigmoid"
from tensorflow.keras.optimizers import Adam model.compiled
optimizer=Adam(learning rate=0.001
loss="binary crossentropy", metrics=['accuracy"]
model.fit
X_train, y train, validation data=(
x val, y val
https://colab.research.google.com/drive/1RUT5WjmG -tjTK299qMNKS96wYxYODMm#scrollTo=47k2F
S4CDrqL&printMode=true
10/14
Assignment 4.ipynb - Colaboratory
10/29/22, 11:33 AM
epochs=10
33ms/step - loss: 0.5258 -accuracy: 0.7653 Epoch 2/10 34/34
```

```
accuracy: 0.9453 Epoch 3/10 34/34
accuracy: 0.9842 Epoch 4/10 34/34
[=====] - 1s 15ms/step -
loss: 0.0254 - accuracy: 0.9926 Epoch 5/10 34/34
-accuracy: 0.9954 Epoch 6/10 34/34
-accuracy: 0.9963 Epoch 7/10 34/34
0.0150 - accuracy: 0.9954 Epoch 8/10 34/34
loss: 0.0112 - accuracy: 0.9972 Epoch 9/10 34/34
[======] - 1s 16ms/step -
loss: 0.0062 - accuracy: 0.9981 Epoch 10/10 34/34
[======] - 1s 16ms/step - loss:
0.0050 - accuracy: 0.9991 < keras.callbacks. History at
0 \times 7 = 3263 = 7850 >
accuracy
```

```
from sklearn.metrics import accuracy score, ConfuSion_matrix
score=accuracy_score(y_test, y_pred) print("Test
Score:{:.2f}%", format(score*100))
Test Score:96.00%
cm=confusion_matrix(y test,y pred) fig=plt.
figure(figsize=(12,8)) sns.heatmap
cm , annot=True ,
plt.title("Confusion Matrix")
```

cm

https://colab.research.google.com/drive/1RUT5WjmG_-tjTK299qMNK**S96wYxYODMm**#scrollTo=47k2FS4CDrqL&printMode=true
11/14
10/29/22, 11:33 AM
Assignment 4.ipynb - Colaboratory

```
array([[100, 2], [7, 116]])

Confusion Matrix
```

- 100

le+02

-40

12e+02

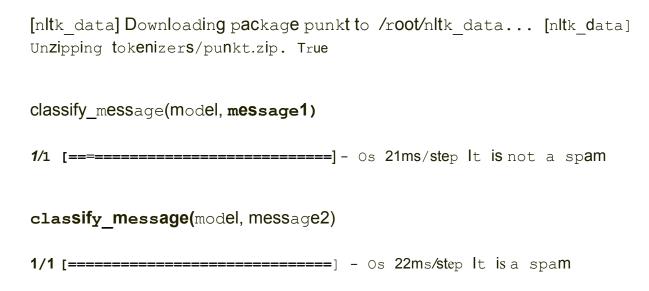
#The function take model and message as parameter def classify message(model, message):

#We will treat message as a paragraphs containing multiple sentences(lines) #we will extract individual lines for sentences in message:

sentences=nltk.sent_tokenize(message)

#Iterate over individual sentences for sentence in sentences:

```
#replace all special characters words=re. sub ("[^a-zA-z]"," ",
sentence)
#perform word tokenization of all non-english-stopwords if words not
in Set(stopwords.words('english')):
word=nltk.word_tokenize(words) word=" ".join(word)
#perform one hot on tokenized word oneHot=[one_hot(word,
n=vocab_size)]
#create an embedded documnet using pad Sequences #this can
be fed to our model text=pad sequences
(oneHot,maxlen=sentence len, padding="pre")
#predict the text using model
https://colab.research.google.com/drive/1RUT5WjmG -tjTK299qMNKS96wYxYODMm
#scrollTo=47k2FS4CDrgL&printMode=true
12/14
Assignment 4.ipynb - Colaboratory
10/29/22, 11:33 AM
predict=model.predict(text)
#if predict value is greater than 0.5 its a spam if predict>0.5:
print("It is a spam") #else the message is not a spam else:
print("It is not a spam")
message1="lam having a bad day and lwould like to have a break
today" message2="This is to inform you had won a lottery and the subscription will
end in a week so
nitk.download('punkt')
```



 $https://colab.resear ch.google.com/drive/1RUT5WjmG_-tjTK299qMNK\\ \textbf{S96wYxYODMm\#scrollTo=47k2FS4CDrqL\&printMode=true}$

13/14

10/29/22, 11:33 AM

Assignment_4.ipynb - Colaboratory

Colab paid products - Cancel contracts here

#scrollTo=47k2FS4CDrqL&printMode=true 14/14