```
1   from google.colab import drive
2   drive.mount('/content/gdrive')
```

Mounted at /content/gdrive

```
1   """
2   # This Python 3 environment comes with many helpful analytics libran
3   # For example, here's several helpful packages to load
4
5   import numpy as np # linear algebra
6   import pandas as pd # data processing, CSV file I/O (e.g. pd.read_cs
7
8   # Input data files are available in the read-only "../input/" direct
9   # For example, running this (by clicking run or pressing Shift+Enter
10
11  import os
12  for dirname, _, filenames in os.walk('path goes here..'):
13      for filename in filenames:
14          print(os.path.join(dirname, filename))
15  """
```

'\n# This Python 3 environment comes with many helpful analytics libraries installed\n#
For example, here\'s several helpful packages to load\n\nimport numpy as np # linear al
gebra\nimport pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)\n\n# Inpu
t data files are available in the read-only "../input/" directory\n# For example, runni
ng this (by clicking run or pressing Shift+Enter) will list all files under the input d

```
1   pip install stop_words
```

Collecting stop_words
  Downloading https://files.pythonhosted.org/packages/1c/cb/d58290804b7a4c5daa42abbbe2a9
Building wheels for collected packages: stop-words
  Building wheel for stop-words (setup.py) ... done
  Created wheel for stop-words: filename=stop_words-2018.7.23-cp36-none-any.whl size=329
  Stored in directory: /root/.cache/pip/wheels/75/37/6a/2b295e03bd07290f0da95c3adb9a74ba
Successfully built stop-words
Installing collected packages: stop-words
Successfully installed stop-words-2018.7.23

```
1   import numpy as np
2   import pandas as pd
3   import os
4   import matplotlib.pyplot as plt
5   import re
6   import nltk
```

```
7   #from nltk.corpus import stopwords
8   from sklearn.model_selection import train_test_split
9   from sklearn.metrics import confusion_matrix
10  from mlxtend.plotting import plot_confusion_matrix
11  from sklearn import preprocessing
12  import tensorflow as tf
13  from tensorflow.keras.models import Sequential
14  from tensorflow.keras.layers import LSTM, MaxPool1D, Dropout, Dense,
15  from keras.utils import to_categorical
16  from keras.preprocessing.text import Tokenizer
17  from keras.preprocessing.sequence import pad_sequences
18  from stop_words import get_stop_words
19
```

## Importing dataset , dropping unnamed columns and ceeating a test data

```
1   train_data_path = '/content/gdrive/My Drive/Colab Notebooks/Hate Spe
2   train_data = pd.read_csv(train_data_path + 'train.csv')
3
4   train_data = train_data.loc[:, ~train_data.columns.str.contains('^Un
5
6   print('original train data shape',train_data.shape)
7
8   #test_data = train_data.iloc[137571:,:]
9   #test_data = test_data.reset_index()
10  #train_data = train_data.iloc[:137570,:]
11
12  #print('test data shape',test_data.shape)
13  #print('new train data shape',train_data.shape)
14  #print(test_data.shape[0]+train_data.shape[0])
```

```
original train data shape (159571, 8)
```

```
1   train_data.head(5)
```

| | id | comment_text | toxic | severe_toxic | obscene | threat | insult | ider |
|---|---|---|---|---|---|---|---|---|
| **0** | 0000997932d777bf | Explanation\nWhy the edits made under my usern... | 0 | 0 | 0 | 0 | 0 | |
| | | D'aww! He | | | | | | |

## ▾ Preprocessing the dataset

```
1   def deleteSmallWords(text):
2       return ' '.join([word for word in text.split() if len(word) > 3]
3   def cleanText(text):
4       # clean the text
5       text = re.sub(r"Https?://[A-Za-z0-9./]+","url",text)
6       text = re.sub(r"[^A-Za-z0-9^,!.\/'+-=]"," ",text)
7       text = re.sub(r"what's","what is ",text)
8       text = re.sub(r"\'s"," ",text)
9       text = re.sub(r"\s+[a-zA-Z]\s+", ' ', text) # Single character r
10      text = re.sub(r"\'ve"," have ",text)
11      text = re.sub(r"\n't"," not ",text)
12      #text = re.sub(r"\i'm","i am ",text)
13      text = re.sub(r"\'re"," are ",text)
14      text = re.sub(r"\'d"," would ",text)
15      text = re.sub(r"\'ll"," will ",text)
16      text = re.sub(r"\."," ",text)
17      text = re.sub(r"!"," ",text)
18      text = re.sub(r"\/"," ",text)
19      text = re.sub(r"\^"," ^ ",text)
20      text = re.sub(r"\+"," + ",text)
21      text = re.sub(r"\-"," - ",text)
22      text = re.sub(r"\="," = ",text)
23      text = re.sub(r":"," : ",text)
24      text = re.sub(r"'"," ",text)
25      text = re.sub(r"(\d+)(k)",r"\g<1>000",text)
26      text = re.sub(r" e g "," eg ",text)
27      text = re.sub(r" b g "," bg ",text)
28      text = re.sub(r" u s "," amarican ",text)
29      text = re.sub(r"\0s","0",text)
30      text = re.sub(r" 9 11 ","911",text)
31      text = re.sub(r"e - mail","email",text)
32      text = re.sub(r"j k","jk",text)
33      text = re.sub(r"\s{2,}"," ",text)
```

```
33     text = re.sub(r"\s{2,}","  ",text)
34     text = re.sub(r"@[A-Za_z0-9]+","",text)
35     text = re.sub(r"(\w)\1{2,}",r"\1\1",text)
36     text = re.sub(r"\w(\w)\1{2}","",text)
37     return text
38 def deleteNonAlphaWords(text):
39     return ''.join([word for word in text.split() if word.isalpha()]
40 def deleteStopWords(text):
41     return ' '.join([word for word in text.lower().split() if not wo
```

```
1 train_data['comment_text'] = train_data['comment_text'].apply(lambda
2 train_data['comment_text'] = train_data['comment_text'].apply(lambda
3 train_data['comment_text'] = train_data['comment_text'].apply(lambda
4 train_data['comment_text'] = train_data['comment_text'].apply(lambda
5
6 #test_data['comment_text'] = test_data['comment_text'].apply(lambda
7 #test_data['comment_text'] = test_data['comment_text'].apply(lambda
8 #test_data['comment_text'] = test_data['comment_text'].apply(lambda
9 #test_data['comment_text'] = test_data['comment_text'].apply(lambda
```

## ▾ Tokenize the data

```
1  num_texts = len(train_data.index)
2  print(num_texts)
3  token = Tokenizer(num_words=num_texts)
4  token.fit_on_texts(train_data['comment_text'])
5  text = token.texts_to_sequences(train_data['comment_text'])
6  text = pad_sequences(text)
7
8  tt = 'bitch hate nigga'
9  text_test = token.texts_to_sequences(tt)
10 text_test = pad_sequences(text_test)
11 print(type(text_test))
```

```
159571
<class 'numpy.ndarray'>
```

```
1 columns = train_data.columns
2 columns = list(columns[2:])
3 print(columns)
4 #y = train_data.loc[:,columns].values
```

```
5  y = train_data['toxic']
6  print(text.shape,y.shape)
7  print(text[1],y[0])
8
```

```
['toxic', 'severe_toxic', 'obscene', 'threat', 'insult', 'identity_hate']
(159571, 1) (159571,)
[1015] 0
```

```
1  x_train, x_test, y_train, y_test = train_test_split(text, y, test_si
2  print(x_train.shape,x_test.shape)
3  print(y_train.shape,y_test.shape)
4  print(type(text))
```

```
(127656, 1) (31915, 1)
(127656,) (31915,)
<class 'numpy.ndarray'>
```

```
 1  max_features = num_texts
 2  embedding_dim = 32
 3
 4  model = Sequential()
 5  model.add(Embedding(max_features, embedding_dim))
 6  model.add(Dropout(0.2))
 7  model.add(LSTM(32, return_sequences=True))
 8  model.add(Dropout(0.2))
 9  model.add(Dense(1))
10  model.add(Activation('sigmoid'))
11  model.summary()
```

```
Model: "sequential"
_____
Layer (type)                 Output Shape              Param #
=================================================================
embedding (Embedding)        (None, None, 32)          5106272
_____
dropout (Dropout)            (None, None, 32)          0
_____
lstm (LSTM)                  (None, None, 32)          8320
_____
dropout_1 (Dropout)          (None, None, 32)          0
_____
dense (Dense)                (None, None, 1)           33
_____
activation (Activation)      (None, None, 1)           0
=================================================================
Total params: 5,114,625
Trainable params: 5,114,625
```

```
    Non-trainable params: 0
    _____
```

```
1   # compile and train model
2   print(x_train.shape,y_train.shape)
3   model.compile(loss='binary_crossentropy', optimizer='adam', metrics=
4   history = model.fit(x_train, y_train, validation_data=(x_test, y_tes
```

```
(127656, 1) (127656,)
Epoch 1/10
3990/3990 [==============================] - 227s 57ms/step - loss: 0.3282 - accuracy: (
Epoch 2/10
3990/3990 [==============================] - 238s 60ms/step - loss: 0.0510 - accuracy: (
Epoch 3/10
3990/3990 [==============================] - 227s 57ms/step - loss: 0.0033 - accuracy: (
Epoch 4/10
3990/3990 [==============================] - 229s 57ms/step - loss: 0.0024 - accuracy: (
Epoch 5/10
1377/3990 [========>....................] - ETA: 2:28 - loss: 0.0016 - accuracy: 0.9996
```

```
1   #model.predict(text_test)
```