

Text Classification:

Data

1. we have total of 20 types of documents(Text files) and total 18828 documents(text files).
2. You can download data from this [link \(https://drive.google.com/open?id=1rxD15nyeIPIAZ-J2VYPrDRZI66-TBWvM\)](https://drive.google.com/open?id=1rxD15nyeIPIAZ-J2VYPrDRZI66-TBWvM), in that you will get documents.rar folder.
If you unzip that, you will get total of 18828 documents. document name is defined as 'ClassLabel_DocumentNumberInThatLabel'.
so from document name, you can extract the label for that document.
4. Now our problem is to classify all the documents into any one of the class.
5. Below we provided count plot of all the labels in our data.

Assignment:

sample document

Subject: A word of advice
From: jcopelan@nyx.cs.du.edu (The One and Only)

In article < 65882@mimsy.umd.edu > mango@cs.umd.edu (Charley Wingate) writes:
>

>I've said 100 times that there is no "alternative" that should think you
>might have caught on by now. And there is no "alternative", but the point
>is, "rationality" isn't an alternative either. The problems of metaphysical
>and religious knowledge are unsolvable-- or I should say, humans cannot
>solve them.

How does that saying go: Those who say it can't be done shouldn't interrupt
those who are doing it.

Jim

--

Have you washed your brain today?

```
In [2]: 1 import pandas as pd
        2 import numpy as np
        3 import re
        4
        5 from sklearn.preprocessing import label_binarize
        6
        7
        8 import pdb
```

```
In [3]: 1 ## extracting zip file (in current director)
        2
        3 # from zipfile import ZipFile
        4 # with ZipFile('documents2.zip', 'r') as zipObj:
        5 #     # Extract all the contents of zip file in current directory
        6 #     zipObj.extractall()
```

loading files

```
In [4]: 1 # reading file_name from folder
        2
        3 import rarfile
        4 rf = rarfile.RarFile('documents.rar' )
        5
        6 rf.namelist()[ :2]
```

```
Out[4]: ['documents/alt.atheism_49960.txt', 'documents/alt.atheism_51060.txt']
```

```
In [5]: 1 # reading file from folder "documents"
        2
        3 file = []
        4 for f in rf.namelist():
        5     # f is filename we directly opening file from 'documents' folder
        6     try:
        7         e = open(f, 'r')
        8         # print(f)
        9         file.append(e.read())
       10     except:
       11         print("error on loading all file")
       12         print(len(file))
       13
```

```
error on loading all file
18828
```

Preprocessing:

useful links: <http://www.pyregex.com/> (<http://www.pyregex.com/>).

1. Find all emails in the document and then get the text after the "@". and then split those texts by '.' after that remove the words whose length is less than or equal to 2 and also remove 'com' word and then combine those words by space.

In one doc, if we have 2 or more mails, get all.

Eg: [test@dm1.d.com, test2@dm2.dm3.com] --> [dm1.d.com, dm3.dm4.com] --> [dm1,d,com,dm2,dm3,com] --> [dm1,dm2,dm3] --> "dm1 dm2 dm3"
append all those into one list/array. (This will give length of 18828 sentences i.e one list for each of the document).
Some sample output was shown below.

> In the above sample document there are emails [jcopelan@nyx.cs.du.edu, 65882@mimsy.umd.edu, mangoe@cs.umd.edu]

preprocessing:

[jcopelan@nyx.cs.du.edu, 65882@mimsy.umd.edu, mangoe@cs.umd.edu] ==> [nyx cs du edu mimsy umd edu cs umd edu] ==> [nyx edu mimsy umd edu umd edu]

2. Replace all the emails by space in the original text.

```
In [6]: 1 text1 = """Subject: A word of advice
2 From: jcopelan@nyx.cs.du.edu (The One and Only)
3
4 In article < 65882@mimsy.umd.edu > mangoe@cs.umd.edu (Charley Wingate) writes:
5 >
6 >I've said 100 times that there is no "alternative" that should think you
7 >might have caught on by now. And there is no "alternative", but the point
8 >is, "rationality" isn't an alternative either. The problems of metaphysical
9 >and religious knowledge are unsolvable-- or I should say, humans cannot
10 >solve them. new york
11
12 How does that saying go: Those who say it can't be done shouldn't interrupt
13 those who are doing it.
14
15 Jim
16 --
17 Have you washed your brain today?"""
18 print(len(text1))
```

621

instruction

3. Get subject of the text i.e. get the total lines where "Subject:" occur and remove the word which are before the ":" remove the newlines, tabs, punctuations, any special chars.
Eg: if we have sentence like "Subject: Re: Gospel Dating @ \r\r\n" --> You have to get "Gospel Dating"
 Save all this data into another list/array.

4. After you store it in the list, Replace those sentences in original text by space.

5. Delete all the sentences where sentence starts with "Write to:" or "From:".
 > In the above sample document check the 2nd line, we should remove that

6. Delete all the tags like "< anyword >"
 > In the above sample document check the 4nd line, we should remove that "< 65882@mimsy.umd.edu >"

7. Delete all the data which are present in the brackets.
 In many text data, we observed that, they maintained the explanation of sentence or translation of sentence to another language in brackets so remove all those.
Eg: "AAIC-The course that gets you HIRED(AAIC - Der Kurs, der Sie anstellt)" --> "AAIC-The course that gets you HIRED"
 > In the above sample document check the 4nd line, we should remove that "(Charley Wingate)"

8. Remove all the newlines('\n'), tabs('\t'), "-", "\".

9. Remove all the words which ends with ":".
Eg: "Anyword:"
 > In the above sample document check the 4nd line, we should remove that "writes:"

10. Decontractions, replace words like below to full words.
 please check the donors choose preprocessing for this
Eg: can't -> can not, 's -> is, i've -> i have, i'm -> i am, you're -> you are, i'll --> i will

There is no order to do point 6 to 10. but you have to get final output correctly

11. Do chunking on the text you have after above preprocessing.

Text chunking, also referred to as shallow parsing, is a task that follows Part-Of-Speech Tagging and that adds more structure to the sentence. So it combines the some phrases, named entities into single word. So after that combine all those phrases/named entities by separating "_". And remove the phrases/named entities if that is a "Person". You can use `nlk.ne_chunk` to get these. Below we have given one example. please go through it.

useful links:

<https://www.nltk.org/book/ch07.html> (<https://www.nltk.org/book/ch07.html>).

<https://stackoverflow.com/a/31837224/4084039> (<https://stackoverflow.com/a/31837224/4084039>).

<http://www.nltk.org/howto/tree.html> (<http://www.nltk.org/howto/tree.html>).

<https://stackoverflow.com/a/44294377/4084039> (<https://stackoverflow.com/a/44294377/4084039>).

preprocessing

```
In [7]: 1 # 3., .4 - extracting email info
2 def email_extractor(text):
3     l1,l2,l3 =[],[],[];
4
5     email = re.findall(r"[A-Za-z0-9_%+-.]+"
6                        r"@[A-Za-z0-9.-]+"
7                        r"\.[A-Za-z]{2,5}",text)
8
9     for x in email:
10         lst = x.split(sep="@")
11         l1.append(lst[1])
12
13     for x in l1:
14         lst = re.sub('\.',',', x)
15         l2.append(lst)
16
17     for x in l2:
18         lst = x.split(sep=',')
19         lst = [x for x in lst if (len(x)>2 and x not in ["edu","com"] )]
20         str = " ".join(lst)
21         l3.append(str)
22
23     email_info = " ".join(l3)
24
25     return email_info
26
27
28 r = email_extractor(file[1])
29 r
```

Out[7]: 'mantis mantis mantis'

In [8]: 1 file[1]

Out[8]: 'From: mathew <mathew@mantis.co.uk>\nSubject: Alt.Atheism FAQ: Introduction to Atheism\n\nArchive-name: atheism/introduction\nAlt-atheism-archive-name: introduction\nLast-modified: 5 April 1993\nVersion: 1.2\n\n-----BEGIN PGP SIGNED MESSAGE-----\n\nAn Introduction to Atheism\nby mathew <mathew@mantis.co.uk>\n\nThis article attempts to provide a general introduction to atheism. Whilst I have tried to be as neutral as possible regarding contentious issues, you should always remember that this document represents only one viewpoint. I would encourage you to read widely and draw your own conclusions; some relevant books are listed in a companion article.\n\nTo provide a sense of cohesion and progression, I have presented this article as an imaginary conversation between an atheist and a theist. All the questions asked by the imaginary theist are questions which have been cropped repeatedly on alt.atheism since the newsgroup was created. Some other frequently asked questions are answered in a companion article.\n\nPlease note that this article is arguably slanted towards answering questions posed from a Christian viewpoint. This is because the FAQ files reflect questions which have actually been asked, and it is predominantly Christians who proselytize on alt.atheism.\n\nSo when I talk of religion, I am talking primarily about religions such as Christianity, Judaism and Islam, which involve some sort of superhuman divine being. Much of the discussion will apply to other religions, but some of it may not.\n\nWhat is atheism?\n\nAtheism is characterized by an absence of belief in the existence of God. Some atheists go further, and believe that God does not exist. The former is often referred to as the "weak atheist" position, and the latter as "strong atheism".\n\nIt is important to note the difference between these two positions. "Weak atheism" is simple scepticism; disbelief in the existence of God. "Strong atheism" is a positive belief that God does not exist. Please do not fall in to the trap of assuming that all atheists are "strong atheists".\n\nSome atheists believe in the non-existence of all Gods; others limit their atheism to specific Gods, such as the Christian God, rather than making flat-out denials.\n\nBut isn't disbelief in God the same thing as believing he doesn't exist?\n\nDefinitely not. Disbelief in a proposition means that one does not be

In [9]: 1 txt = "From to: FFRF, P.O. Box 750, Madison, WI 53701.\n"
2 txt = txt.lower()
3 n_txt = re.sub(r"from to:" # start with 'From:'
4 r"[-A-Za-z0-9.+, @<>_>()]*", " ", txt, re.I)
5 n_txt

Out[9]: ' \n'


```

In [10]: 1 def preprocessing(txt):
2 #     sub = re.findall("[^'subject:']*w+r'\n'$", text, re.IGNORECASE)
3 #     sub = re.findall("[/n]{2}$", text, re.IGNORECASE)
4
5     # extracting subject
6     sub = re.findall(r"subject:+"
7         r"[-A-Za-z0-9.:, @]+", txt, re.I)
8     sub2 = " ".join(sub)
9     sub2 = re.sub(r"Subject:[ ]*", "", sub2, re.I)
10    subject = re.sub(r"[-A-Za-z0-9_-]+:", "", sub2, re.I)
11
12    # removing Email
13    txt = re.sub(r"[-A-Za-z0-9_%+-.]+"
14        r"@[-A-Za-z0-9.-]+"
15        r"\.[-A-Za-z]{2,5}", " ", txt)
16
17    # 16 convert all to lower
18    txt = txt.lower()
19
20    # 5. removing line start with 'write to:' and 'From : ' "Write to:" or "From:"
21    n_txt = re.sub(r"write\ to:" # start with 'write to:'
22        r"[-A-Za-z0-9.:, @<>_]*", " ", txt, re.I)
23
24    n_txt = re.sub(r"from:" # start with 'From:'
25        r"[-A-Za-z0-9.:, @<>_](*)", " ", n_txt, re.I)
26
27    #6. delete tags like - "< anyword >"
28    n_txt = re.sub(r"<+" # start with
29        r"[-A-Za-z0-9.:, @<>_]+" # selecting every charector
30        ">", " ", n_txt)
31
32    # 7.Delete all the data which are present in the brackets
33    n_txt = re.sub(r"\("
34        r"[-A-Za-z0-9.:, @<>_\"\\n\\t]+"
35        "\)", " ", n_txt)
36    # 8. Remove all the newlines('\n'), tabs('\t'), "-", "\".
37    n_txt = re.sub(r"[\t\n]", ' ', n_txt)
38
39    # 9. Remove all the words which ends with ":".
40    n_txt = re.sub(r"[-A-Za-z0-9_-]+:", '', n_txt)
41
42    # 10. replace can't -> can not, 's -> is, i've -> i have, i'm -> i am, you're -> you are, i'll --> i will
43    n_txt = re.sub("can't ", "can not ", n_txt, re.I)
44    n_txt = re.sub("n't ", " not ", n_txt, re.I)
45    n_txt = re.sub("\'s ", "is", n_txt, re.I)

```

```

46     n_txt = re.sub("i've ", "i have ", n_txt, re.I)
47     n_txt = re.sub("i'm ", " i am ", n_txt, re.I)
48     n_txt = re.sub("you're ", "you are ", n_txt, re.I)
49     n_txt = re.sub("i'll ", "i will ", n_txt, re.I)
50     n_txt = re.sub("it's ", "it is ", n_txt, re.I)
51
52     # self
53     n_txt = re.sub(r"[ -]{2,20}", ' ', n_txt )
54     n_txt = re.sub(r"[\"><?]+", '', n_txt )
55
56
57     #     print(txt)
58     #     print(""*100)
59     #     print(n_txt)
60
61     return n_txt, subject
62
63
64 txt, subject = preprocessing(text1)
65 txt

```

Out[10]: ' a word of advice in article i have said 100 times that there is no alternative that should think you might have caught on by now. and there is no alternative, but the point is, rationality is not an alternative either. the problems of metaphysical and religious knowledge are unsolvable or i should say, humans cannot solve them. new york how does that saying those who say it can not be done sould not interrupt those who are doing it. jim have you washed your brain today'

Chunking

```

In [11]: 1 from nltk import pos_tag
          2 from nltk import RegexpParser
          3 from nltk import ne_chunk, pos_tag, word_tokenize
          4 from nltk.tree import Tree
          5 import time

```

```
In [12]: 1 # geting chunked data
2 # need to remove person name
3 # modified of ref: https://stackoverflow.com/questions/31836058/nltk-named-entity-recognition-to-a-python-list/31837224#31837224
4
5 def get_continuous_chunks(text):
6     chunked = ne_chunk(pos_tag(word_tokenize(text)))
7     continuous_chunk = []
8     continuous_chunk_person = []
9     current_chunk = []
10    for i in chunked:
11
12        if type(i) == Tree and len(i)>1 and i.label() not in ["PERSON", "ORGANIZATION"] : # Label for person use - "PERSON"
13            name = np.array(i.leaves(),)[: ,0]
14            current_chunk.append(list(name))
15        #         pdb.set_trace()
16
17
18        if type(i) == Tree and len(i)>1 and i.label() == "PERSON" : # Label for person use - "PERSON"
19            name = np.array(i.leaves(),)[: ,0]
20            j_name = " ".join(list(name))
21            if j_name not in continuous_chunk:
22                continuous_chunk_person.append(j_name)
23
24
25
26        #         pdb.set_trace()
27        if current_chunk:
28            named_entity = " ".join(current_chunk[0])
29            if named_entity not in continuous_chunk:
30                continuous_chunk.append(named_entity)
31                current_chunk = []
32        else:
33            continue
34    return continuous_chunk, continuous_chunk_person
35
36 get_continuous_chunks(file[0])
```

```
Out[12]: (['North Hollywood', 'New York', 'New England'],  
          ['Laurel Canyon',  
           'Darwin Fish',  
           'Lynn Gold',  
           'Cameron Road',  
           'Holy Horrors',  
           'Glenn Drive',  
           'Norm R. Allen',  
           'United Kingdom Rationalist',  
           'Islington High',  
           'Germany IBKA',  
           'Atheisten Postfach',  
           'Konfessionslosen Postfach',  
           'Fiction THOMAS',  
           'Philip K. Dick Dick',  
           'Joe Fernwright',  
           'Bantam Press',  
           'Temple University Press',  
           'Temple University',  
           'Without Creed',  
           'Without Creed',  
           'Ballantine Books',  
           'Great Men Think',  
           'RICHARD SWINBURNE',  
           'Clarendon Paperbacks',  
           'Holy Horrors'])
```

```
In [13]: 1 # replace with original text
2 def replace_remove_chunk_word(text):
3
4     entity_name, person_name= get_continuous_chunks(text)
5
6     # replacing entity_name with _
7     if len(entity_name)>0:
8         for i in range(len(entity_name)):
9             dual = re.sub("\ ", "_", entity_name[i])
10            text = re.sub(r"{}".format(entity_name[i]), dual, text, re.I)
11
12    # removing person_name
13    if len(person_name)>0:
14        for p in range(len(person_name)):
15            text = re.sub(r"{}".format(person_name[p]), " ", text, re.I)
16
17    return text
18
```

testing w.r.t single file

```
In [14]: 1 dual_name = get_continuous_chunks(file[2])
2 new_text = replace_remove_chunk_word(txt[0])
3 dual_name, new_text[:100]
```

```
Out[14]: (([], ['Gospel Dating', 'Charley Wingate', 'James Felder']), ' ')
```

```
In [15]: 1 #chunking and replace with original text
2
3 txt = preprocessing(file[2])
4
5 # dual_name = get_continuous_chunks(txt[0])
6 new_text = replace_remove_chunk_word(txt[0])
7
8 # print(dual_name)
9
10 print(new_text)
```

gospel dating in article well, john has a quite different, not necessarily more elaborated theology. there is some evidence that he must have known luke, and that the content of q was known to him, but not in a 'canonized' form. this is a new argument to me. could you elaborate a little the argument goes as q-oid quotes appear in john, but not in the almost codified way they were in matthew or luke. however, they are considered to be similar enough to point to knowledge of q as such, and not an entirely different source. assuming that he knew luke would obviously put him after luke, and would give evidence for the latter assumption. i do not think this follows. if you take the most traditional attributions, then luke might have known john, but john is an elder figure in either case. we're talking spans of time here which are well within the range of lifetimes. we are talking date of texts here, not the age of the authors. the usual explanation for the time order of mark, matthew and luke does not consider their respective ages. it says matthew has read the text of mark, and luke that of matthew. as it is assumed that john knew the content of luke's text. the evidence for that is not overwhelming, admittedly. earlier manuscripts of john have been discovered. interesting, where and which how are they dated how old are they unfortunately, i have not got the info at hand. it was in the late '70s or early '80s, and it was possibly as old as c. 200. when they are from about 200, why do they shed doubt on the order of putting john after the rest of the three i don't see your point, it is exactly what james felder said. they had no first hand knowledge of the events, and it is obvious that at least two of them used older texts as the base of their account. and even the association of luke to paul or mark to peter are not generally accepted. well, a genuine letter of peter would be close enough, wouldn't it sure, an original together with id card of sender and receiver would be fine. so what is that supposed to say am i missing something and i don't think a one step removed source is that bad. if luke and mark and matthew learned their stories directly from disciples, then i really cannot believe in the sort of big transformation from jesus to gospel that some people posit. in news reports, one generally gets no better information than this. and if john is a disciple, then there's nothing more to be said. that john was a disciple is not generally accepted. the style and language together with the theology are usually used as counterargument. the argument that john was a disciple relies on the claim in the gospel of john itself. is there any other evidence for it one step and one generation removed is bad even in our times. compare that to reports of similar events in our century in almost illiterate societies. not even to speak of that believers are not necessarily the best sources. it is also obvious that mark has been edited. how old are the oldest manuscripts to my knowledge the oldest is quite after any of these estimates, and it is not even complete. the only clear editing is problem of the ending, and it is basically a hopeless mess. the oldest versions give a strong sense of incompleteness, to the point where the shortest versions seem to break off in midsentence. the most obvious solution is that at some point part of the text was lost. the material from verse 9 on is pretty clearly later and seems to represent a synopsis of the end of luke. in other words, one does not know what the original of mark did look like and arguments based on mark are pretty weak. but how is that connected to a redating of john benedikt

```
In [18]: 1 # print(file[2])
```

example for chunking


```
In [16]: 1 #note - get_continuous_chunks() return 2 list - entity names and person names
2
3 #i am living in the New York
4 print("i am living in the New York -->", get_continuous_chunks("i am living in the New York -->"))
5 print(" ")
6 print("-"*50)
7 print(" ")
8 #My name is Srikanth Varma
9 print("My name is Srikanth Varma -->", get_continuous_chunks("My name is Srikanth Varma -->"))
10
11 print("\nafter removing:-", replace_remove_chunk_word("My name is Srikanth Varma -->"))
```

i am living in the New York --> (['New York'], [])

My name is Srikanth Varma --> ([], ['Srikanth Varma'])

after removing:- My name is -->

We did chunking for above two lines and then We got one list where each word is mapped to a POS(parts of speech) and also if you see "New York" and "Srikanth Varma", they got combined and represented as a tree and "New York" was referred as "GPE" and "Srikanth Varma" was referred as "PERSON".

so now you have to Combine the "New York" with "_" i.e "New_York" and remove the "Srikanth Varma" from the above sentence because it is a person.

```
In [17]: 1 preprocessing(text1)
```

```
Out[17]: (' a word of advice in article i have said 100 times that there is no alternative that should think you might have caught on by no
w. and there is no alternative, but the point is, rationality is not an alternative either. the problems of metaphysical and religio
us knowledge are unsolvable or i should say, humans cannot solve them. new york how does that saying those who say it can not be don
e should not interrupt those who are doing it. jim have you washed your brain today',
'A word of advice')
```

```
In [18]: 1 text1 = text1 + " _word_ "  
2 print(text1)
```

Subject: A word of advice

From: jcopelan@nyx.cs.du.edu (The One and Only)

In article < 65882@mimsy.umd.edu > mangoe@cs.umd.edu (Charley Wingate) writes:

>

>I've said 100 times that there is no "alternative" that should think you
>might have caught on by now. And there is no "alternative", but the point
>is, "rationality" isn't an alternative either. The problems of metaphysical
>and religious knowledge are unsolvable-- or I should say, humans cannot
>solve them. new york

How does that saying go: Those who say it can't be done shouldn't interrupt
those who are doing it.

Jim

--

Have you washed your brain today? _word_

part-2 preprossesing

13. Replace all the digits with space i.e delete all the digits.

> In the above sample document, the 6th line have digit 100, so we have to remove that.

14. After doing above points, we observed there might be few word's like

"_word_" (i.e starting and ending with the _), "_word" (i.e starting with the _),

"word_" (i.e ending with the _) remove the _ from these type of words.

15. We also observed some words like "OneLetter_word"- eg: d_berlin,

"TwoLetters_word" - eg: dr_berlin , in these words we remove the "OneLetter_" (d_berlin ==> berlin) and

"TwoLetters_" (de_berlin ==> berlin). i.e remove the words

which are length less than or equal to 2 after spliiting those words by "_".

16. Convert all the words into lower case and lowe case

and remove the words which are greater than or equal to 15 or less than or equal to 2.

17. replace all the words except "A-Za-z_" with space.

18. Now You got Preprocessed Text, email, subject. create a dataframe with those.

Below are the columns of the df.

```

In [19]: 1 # 13- 17, preprocesses_2 with regular expression
2
3 def preprocessing2(text):
4
5     # 5. removing line start with 'write to:' and 'From : ' "Write to:" or "From:"
6     # n_txt = re.sub(r"write\ to:" # start with 'write to:'
7     #               r"[-A-Za-z0-9.:, @<>_]*", " ", txt)
8
9     # n_txt = re.sub(r"From:" # start with 'From:'
10    #               r"[-A-Za-z0-9.:, @<>_]*", " ", n_txt, re.I)
11
12    # 6. delete tags like - "< anyword >"
13    # n_txt = re.sub(r"<+" # start with
14    #               r"[-A-Za-z0-9.:, @<>_]+" # selecting every charector
15    #               ">", " ", n_txt)
16
17    # 13 replacing digits with space
18    n_txt = re.sub(r"[0-9]+", " ", text)
19
20    # 14 removing _ in " _word_, _word, word_ " these type of words
21    n_txt = re.sub(r"\ _[-A-Za-z0-9.:, @<>_]+", " ", n_txt)
22    n_txt = re.sub(r"[-A-Za-z0-9.:, @<>_]+\ _", " ", n_txt)
23
24    # 15 remove start short word(length less than or equal to 2) "TwoLetters_word" - eg: (d_berlin ==> berlin)
25    n_txt = re.sub(r"\ [A-Za-z0-9]{1,2}\ _[-A-Za-z0-9]+", " ", n_txt)
26
27    # 16 convert all to lower
28    # n_txt = n_txt.lower()
29
30    # 17 replace all the words except "A-Za-z_" with space.
31    n_txt = re.sub(r"[^A-Za-z_ ]", " ", n_txt)
32
33    # removing extra space
34    n_txt = re.sub(r"[ ]{2,5}", " ", n_txt)
35
36    return n_txt
37
38 txt,d = preprocessing("""After de_bfl doing above dee_bfl poin_tkks po_tkks poin_tk, we observed d_berlin there might be few w
39 word_ (i.e ending with the _) remove the from these type of words.""")
40 txt

```

Out[19]: 'after de_bfl doing above dee_bfl poin_tkks po_tkks poin_tk, we observed d_berlin there might be few word /islike _word_ , _word , word_ remove the from these type of words.'

test preprocessing2 w.r.t. single file

```
In [20]: 1 # txt,subject = preprocessing(file[4])
          2 # txt1 = replace_remove_chunk_word(txt)
          3 # dual = get_continuous_chunks(txt1)
          4
          5 # txt = preprocessing2(txt)
          6
          7 # dual, txt
          8
```

```
In [21]: 1 # txt1
```

```
In [22]: 1 # file[4]
```

```
In [23]: 1 # dual = get_continuous_chunks(file[4])
          2 # dual
```

To get above mentioned data frame --> Try to Write Total Preprocessing steps in One Function Named Preprocess as below.

```
In [24]: 1 def preprocess_all(Input_Text):
2         """Do all the Preprocessing as shown above and
3         return a tuple contain preprocess_email,preprocess_subject,preprocess_text for that Text_data"""
4         list_text = []
5         subjects = []
6         list_of_preprocessed_emails = []
7         for i in range(len(Input_Text)):
8
9             # prepossing + chucking
10            email = email_extractor(Input_Text[i])
11            txt, sub1 = preprocessing(Input_Text[i])
12            subject = preprocessing2(sub1)
13            new_text = replace_remove_chunk_word(txt)
14            new_text = preprocessing2(new_text)
15
16
17            # appending
18            list_text.append(new_text)
19            subjects.append(subject)
20            list_of_preprocessed_emails.append(email)
21        #         pdb.set_trace()
22
23        return (list_of_preprocessed_emails,subjects,list_text)
```

Training The models to Classify:

Creating Dataframe

```
In [ ]: 1
```

```
In [299]: 1 %%time
          2 # took more than 1 h 40min
          3
          4 # processing all file data
          5 list_of_preprocessed_emails, subjects, list_texts = preprocess_all(file)
```

Wall time: 52min 21s

```
In [148]: 1 # extracting class/label from file name
          2
          3
          4 file_name = rf.namelist()
          5 classes = []
          6 for i in file_name[:-1]:
          7
          8     l = re.findall(r"\/[A-Za-z.]*", i)
          9     l = re.sub(r"\/", "", " ".join(l))
         10     classes.append(l)
         11
         12 print(len(classes))
```

18828

```
In [288]: 1 classes
```

...

```
In [149]: 1 # reshaping data
          2 # list_of_preprocessed_emails = np.reshape(list_of_preprocessed_emails, (-1,1))
          3 # subjects = np.reshape(subjects, (-1,1))
          4 # list_texts = np.reshape(list_texts, (-1,1))
          5 # label1 = np.reshape(classes, (-1,1))
          6
          7 list_of_preprocessed_emails.shape, subjects.shape, list_texts.shape, label1.shape
```

```
Out[149]: ((18828, 1), (18828, 1), (18828, 1), (18828, 1))
```

```
In [150]: 1 (np.unique(label1))
```

```
Out[150]: array(['alt.atheism', 'comp.graphics', 'comp.os.ms',
                'comp.sys.ibm.pc.hardware', 'comp.sys.mac.hardware',
                'comp.windows.x', 'misc.forsale', 'rec.autos', 'rec.motorcycles',
                'rec.sport.baseball', 'rec.sport.hockey', 'sci.crypt',
                'sci.electronics', 'sci.med', 'sci.space',
                'soc.religion.christian', 'talk.politics.guns',
                'talk.politics.mideast', 'talk.politics.misc',
                'talk.religion.misc'], dtype='<U24')
```

```
In [283]: 1 # label[label=='alt.atheism'].shape
          2 # label[label=='soc.religion.christian'].shape
          3 # sns.distplot(label, kde=True, label= np.unique(label))
          4 # label.shape
          5
```

```
In [303]: 1 # list_texts
```

```
Out[303]: [' alt atheism atheist resources atheism resources resources december  atheist resources addresses of atheist organizations usa fr
eedom from religion foundation darwin fish bumper stickers and assorted other atheist paraphernalia are available from the freedom
from religion foundation in the us  evolution designs evolution designs sell the darwin fish itisa fish symbol like the ones chris
tians stick on their cars but with feet and the word darwin written inside the deluxe moulded d plastic fish is  postpaid in the u
s  north hollywood ca people in the san francisco bay area can get darwin fish from lynn gold try mailing for net people who go to
lynn directly the price is  per fish american atheist press aap publish various atheist books critiques of the bible lists of bibl
ical contradictions and so on one such book the bible handbook by w p ball and g w foote american atheist press pp isbn  nd editi
on bible contradictions absurdities atrocities immoralities contains ball the bible contradicts itself aap based on the king james
version of the bible write american atheist press p o box austin tx  cameron road austin tx  prometheus books sell books includin
g haughtisholy horrors write east amherst street buffalo new york  an alternate address prometheus books glenn drive buffalo ny  a
frican americans for humanism an organization promoting black secular humanism and uncovering the history of black freethought the
y publish a quarterly newsletter aah examiner write norm r allen jr african americans for humanism p o box buffalo ny united kingd
om rationalist press association national secular society islington high street holloway road london n ew london n nl  british hu
manist association south place ethical society lamb s conduit passage conway hall london wc r rh red lion square  london wc r rl f
ax  the national secular society publish the freethinker a monthly magazine founded in germany ibka e v internationaler bund der
konfessionslosen und atheisten postfach d berlin germany ibka publish a miz miz vertrieb postfach d berlin germany for atheist boo
ks write ibdk internationaler bucherdienst der konfessionslosen postfach d hannover germany  books fiction thomas m disch the sant
a claus compromise short story the ultimate proof that santa exists all characters and events are fictitious any similarity to liv
ing or dead gods uh well walter m miller jr a canticle for leibowitz one gem in this post atomic doomsday novel is the monks who s
eat their lives saving blunders from saint leibowitz filling the sheets of paper with ink and leaving white lines and letters
```

```
In [311]: 1 # classes
```



```
In [312]: 1 e1 = ["Geeks", "For", "Geeks"]
          2 e2= ["Geeks", "For", "Geeks"]
          3 e3 = ["Geeks", "For", "Geeks"]
          4
          5 df = pd.DataFrame({
          6     'a':classes[:3],
          7     'b':e2,
          8     'c':e3
          9 })
         10 df
```

Out[312]:

	a	b	c
0	alt.atheism	Geeks	Geeks
1	alt.atheism	For	For
2	alt.atheism	Geeks	Geeks

```
In [313]: 1 # creating dataframe
          2
          3 df = pd.DataFrame({
          4     'label':classes,
          5     'processed_text': list_texts,
          6     'procesossed_subject': subjects,
          7     'procesossed_emails' : list_of_preprocessed_emails })
          8
          9 #stacking data
         10 # processed_data = np.hstack((label1, list_texts, subjects,list_of_preprocessed_emails))
         11
         12 #converting to dataframe
         13 # df = pd.DataFrame(processed_data, columns= ["label","procesossed_text", "procesossed_subject", "procesossed_emails" ])
```

```
In [284]: 1 # np.array2string(list_texts[1])
          2 list_texts[4]
```

Out[284]: array([' soc motss et al princeton axes matching funds for boy scouts in article however i hate economic terrorism and political co
rrectness worse than i hate this policy a more effective approach is to stop donating to any organizing that directly or indirect
ly supports gay rights issues until they end the boycott on funding of scouts can somebody reconcile the apparent contradiction betw
een and rob strom ibm research saw mill river road p o box yorktown heights ny '],
dtype='<U99350')]

```
In [315]: 1 df['procesossed_subject']
```

```
Out[315]: 0 Alt Atheism Atheist Resources
1 Alt Atheism Introduction to Atheism
2 Gospel Dating
3 university violating separation of church
4
...
18823 Religion and marriage
18824 A Message for you Mr How do you know what hap...
18825 Why did they behave as they did
18826 Info about New Age
18827 I
Name: procesossed_subject, Length: 18828, dtype: object
```

```
In [300]: 1 # df2[0]
```

```
In [317]: 1 # 1. combine columns
2 df2 = df['procesossed_subject'] + df['procesossed_emails'] + df['processed_text']
3 # df2 = df.loc[:, ['procesossed_subject', 'procesossed_emails', 'processed_text']]
4
5 # labelizing
6 label = label_binarize(classes, np.unique(df[['label']]))
7
8 df2.shape, label.shape
```

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\validation.py:67: FutureWarning: Pass classes=['alt.atheism' 'comp.graphics' 'comp.os.ms' 'comp.sys.ibm.pc.hardware' 'comp.sys.mac.hardware' 'comp.windows.x' 'misc.forsale' 'rec.autos' 'rec.motorcycles' 'rec.sport.baseball' 'rec.sport.hockey' 'sci.crypt' 'sci.electronics' 'sci.med' 'sci.space' 'soc.religion.christian' 'talk.politics.guns' 'talk.politics.mideast' 'talk.politics.misc' 'talk.religion.misc'] as keyword args. From version 0.25 passing these as positional arguments will result in an error
warnings.warn("Pass {} as keyword args. From version 0.25 ")

```
Out[317]: ((18828,), (18828, 20))
```

Code checking:

After Writing preprocess function. call that functoin with the input text of 'alt.atheism_49960' doc and print the output of the preprocess function

This will help us to evaluate faster, based on the output we can suggest you if there are any changes.

```
In [36]: 1 # # prepossing + chucking
2 # email = email_extractor(file[2])
3 # txt, subject = preprocessing(file[2])
4 # new_text = replace_remove_chunk_word(txt)
5 # new_text = preprocessing2(new_text)
6 # subject = preprocessing2(subject)
```

```
In [37]: 1 # file[2] # file name - alt.atheism_49960
```

```
In [38]: 1 # print(email, '\n', subject, '\n', new_text[:100])
2 # new_text
```

```
In [228]: 1 import tensorflow as tf
2 from tensorflow.keras.models import Model
3
4 from tensorflow.keras.layers import Dense, Input, Activation
5
6 from tensorflow.keras.preprocessing.text import Tokenizer
7 import seaborn as sns
8
```

Splitting Data

```
In [319]: 1 # df2 = np.array(df2)
2
```

Out[319]: (18828,)

```
In [341]: 1 # 2. split the data into Train and test
2 from sklearn.model_selection import train_test_split
3 df2 = np.array(df2)
4 X_train, X_test, y_train, y_test = train_test_split(df2, label, test_size=0.2, stratify= label )
5
```

3. Analyze your text data and pad the sequence

```

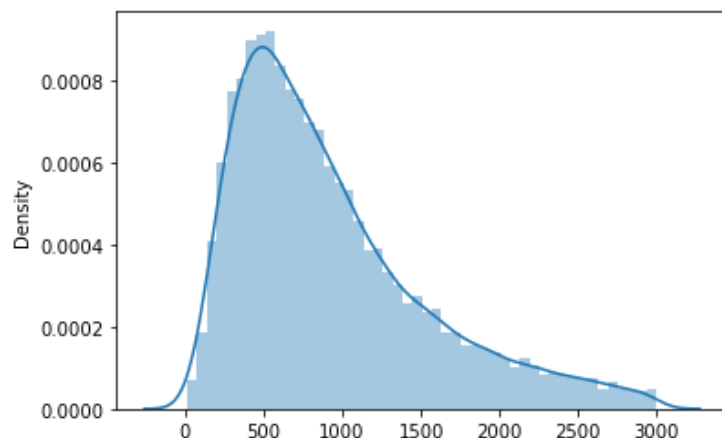
In [323]: 1 # 3. Analyze your text data and pad the sequence
          2
          3 lens_txt = np.array([len(x) for x in X_train])
          4
          5 print(np.median([len(x) for x in X_train]) ,",", np.mean([len(x) for x in X_train]))
          6 sns.distplot(lens_txt[lens_txt<=3000], kde = True)
          7
          8 # pad the text size

```

825.0 , 1435.4373257203558

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)

Out[323]: <AxesSubplot:ylabel='Density'>



observation: max_len for text = 1000, 1000 is good middle value

```

In [324]: 1
          2 lens_txt = np.array([len(x) for x in X_train])
          3 # lens_txt[lens_txt >= 1000]
          4 lens_txt

```

Out[324]: array([1850, 570, 432, ..., 1607, 1195, 2134])

Do Tokenizer And creating embedding matrix through glove file

i.e convert text(word and char) into numbers.

After writing Preprocess function, call the function for each of the document(18828 docs) and then create a dataframe as mentioned above.

```
In [325]: 1 # 4.1 Do Tokenizer i.e convert text into numbers.
2 # https://stackoverflow.com/questions/51956000/what-does-keras-tokenizer-method-exactly-do
3
4
5 t = Tokenizer(num_words = 1000, filters='!"#$%&()*+,-./:;<=>?@[\\]^`{|}~\t\n')
6 t.fit_on_texts(X_train)
7
8 # convert each sentence to matrix(1000 x 1)
9 encoded_docs = t.texts_to_sequences(X_train)
10
11 print("sequences : ",len(encoded_docs),'\n')
12
13 print("word_index : ",len(t.word_index))
14
15
```

sequences : 15062

word_index : 112379

```

In [326]: 1 # 4.2 Do Tokenizer i.e convert char into numbers.
          2
          3 t2 = Tokenizer(num_words = 1000, filters='!"#$%&()*+,-./:;<=>?@[\\]^`{|}~\t\n', char_level = True)
          4 t2.fit_on_texts(X_train)
          5
          6 # convert each sentence to matrix(1000 x 1)
          7 encoded_docs2 = t2.texts_to_sequences(X_train)
          8 encoded_docs_test = t2.texts_to_sequences(X_test)
          9 print("sequences : ",len(encoded_docs2),'\n')
         10
         11 print("word_index : ",(t2.word_index))
         12

```

sequences : 15062

word_index : {' ': 1, 'e': 2, 't': 3, 'a': 4, 'o': 5, 'i': 6, 'n': 7, 's': 8, 'r': 9, 'h': 10, 'l': 11, 'd': 12, 'c': 13, 'u': 14, 'm': 15, 'p': 16, 'f': 17, 'g': 18, 'y': 19, 'w': 20, 'b': 21, 'v': 22, 'k': 23, 'x': 24, 'j': 25, 'q': 26, 'z': 27, '_': 28, '1': 29, '-': 30, '2': 31, '0': 32, '3': 33, '4': 34, '6': 35, '5': 36, '8': 37, '7': 38, '9': 39}

```

In [327]: 1 # char_list = list(t2.word_index.keys())
          2 # char_ohe = Label_binarize(char_list,char_list )
          3 # char_ohe.shape

```

```

In [328]: 1 # padding every encoded txt file to (n x 1000)
          2 from tensorflow.keras.utils import pad_sequences
          3
          4 print(len(encoded_docs[5]))
          5
          6 padded_docs = pad_sequences(encoded_docs, maxlen=1000, padding='post')
          7 padded_docs_test = pad_sequences(encoded_docs_test, maxlen=1000, padding='post')
          8 padded_docs_char = pad_sequences(encoded_docs2, maxlen=1000, padding='post')
          9
         10 print(padded_docs.shape, padded_docs_char.shape, padded_docs_test.shape)

```

77

(15062, 1000) (15062, 1000) (3766, 1000)

```
In [329]: 1 %%time
2 # Load the whole glove file
3 embeddings_index = dict()
4 f = open('glove.6B.100d.txt', encoding="utf8")
5 for i,line in enumerate(f):
6     try:
7         # pdb.set_trace()
8         values = line.split(' ')
9         word = values[0]
10        coefs = np.asarray(values[1:])
11        #pdb.set_trace()
12        # storing vector of word into dict
13        embeddings_index[word] = coefs
14    #     print(i)
15    except:
16        pdb.set_trace()
17
18 f.close()
19 print('Loaded %s word vectors.' % len(embeddings_index))
```

Loaded 400000 word vectors.

Wall time: 26.4 s

```
In [330]: 1 # create a weight matrix for words in training docs
2 vocab_size = 1000
3 embedding_matrix_word = np.zeros((vocab_size, 100))
4 for word, i in t.word_index.items():
5     embedding_vector = embeddings_index.get(word)
6     if i >= 1000:
7         # pdb.set_trace()
8         break
9     if embedding_vector is not None:
10        embedding_matrix_word[i] = embedding_vector
11
12
```

```
In [ ]: 1
```

```
In [331]: 1 %%time
2 # Load the whole glove file for char vectorization
3 char_embed_val = dict()
4 f = open('glove.840B.300d-char.txt', encoding="utf8")
5 for i,line in enumerate(f):
6     try:
7         # pdb.set_trace()
8         values = line.split(' ')
9         char = values[0]
10        coefs = np.asarray(values[1:])
11        #pdb.set_trace()
12        # storing vector of word into dict
13        char_embed_val[char] = coefs
14    #     print(i)
15    except:
16        None
17
18 f.close()
19 print('Loaded %s word vectors.' % len(char_embed_val))
```

Loaded 94 word vectors.

Wall time: 40 ms

```
In [332]: 1 # create a weight matrix for char in training docs
2 char_vocab_size = len(char_embed_val.keys())
3 char_embed_matrix = np.zeros((char_vocab_size, 300))
4
5 for word, i in t2.word_index.items():
6     embedding_vector = char_embed_val.get(word)
7     if i >= 1000:
8         #
9         break
10    if embedding_vector is not None:
11        #     pdb.set_trace()
12        char_embed_matrix[i] = embedding_vector
13
14 print(char_embed_matrix.shape)
```

(94, 300)

Training The models to Classify:

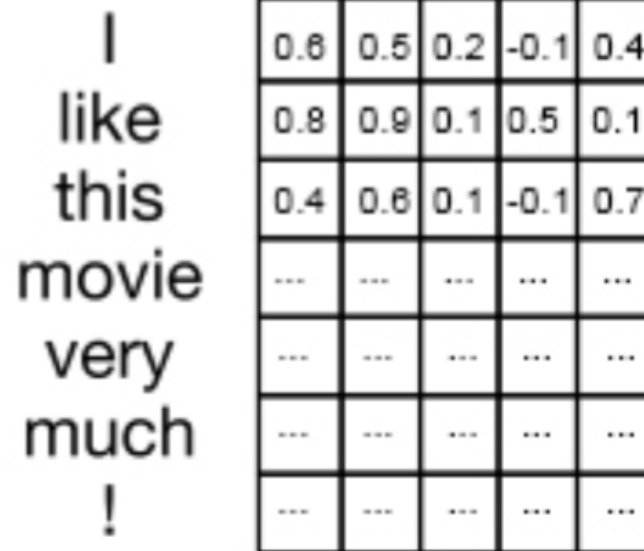
1. Combine "preprocessed_text", "preprocessed_subject", "preprocessed_emails" into one column. use that column to model.
2. Now Split the data into Train and test. use 25% for test also do a stratify split.
3. Analyze your text data and pad the sequence if required.
Sequence length is not restricted, you can use anything of your choice.
you need to give the reasoning
4. Do Tokenizer i.e convert text into numbers. please be careful while doing it.
if you are using tf.keras "Tokenizer" API, it removes the "_", but we need that.
5. code the model's (Model-1, Model-2) as discussed below
and try to optimize that models.
6. For every model use predefined Glove vectors.
Don't train any word vectors while Training the model.
7. Use "categorical_crossentropy" as Loss.
8. Use **Accuracy and Micro Averaged F1 score** as your as Key metrics to evaluate your model.
9. Use Tensorboard to plot the loss and Metrics based on the epoches.
10. Please save your best model weights in to '**best_model_L.h5**' (L = 1 or 2).
11. You are free to choose any Activation function, learning rate, optimizer.
But have to use the same architecture which we are giving below.
12. You can add some layer to our architecture but you **deletion** of layer is not acceptable.
13. Try to use **Early Stopping** technique or any of the callback techniques that you did in the previous assignments.
14. For Every model save your model to image (Plot the model) with shapes
and include those images in the notebook markdown cell,
upload those images to Classroom. You can use "plot_model"
please refer [this \(https://www.tensorflow.org/api_docs/python/tf/keras/utils/plot_model\)](https://www.tensorflow.org/api_docs/python/tf/keras/utils/plot_model). if you don't know how to plot the model with shapes.

Model-1: Using 1D convolutions with word embeddings

Encoding of the Text --> For a given text data create a Matrix with Embedding layer as shown Below.

In the example we have considered $d = 5$, but in this assignment we will get $d =$ dimension of Word vectors we are using.

i.e if we have maximum of 350 words in a sentence and embedding of 300 dim word vector, we result in 350×300 dimensional matrix for each sentence as output after embedding layer



I	0.8	0.5	0.2	-0.1	0.4
like	0.8	0.9	0.1	0.5	0.1
this	0.4	0.6	0.1	-0.1	0.7
movie
very
much
!

Ref: <https://i.imgur.com/kiVQuk1.png>

Reference:

<https://stackoverflow.com/a/43399308/4084039> (<https://stackoverflow.com/a/43399308/4084039>).

<https://missinglink.ai/guides/keras/keras-conv1d-working-1d-convolutional-neural-networks-keras/> (<https://missinglink.ai/guides/keras/keras-conv1d-working-1d-convolutional-neural-networks-keras/>).

[How EMBEDDING LAYER WORKS](https://stats.stackexchange.com/questions/270546/how-does-keras-embedding-layer-work) (<https://stats.stackexchange.com/questions/270546/how-does-keras-embedding-layer-work>).

Go through this blog, if you have any doubt on using predefined Embedding values in Embedding layer -
<https://machinelearningmastery.com/use-word-embedding-layers-deep-learning-keras/>
(<https://machinelearningmastery.com/use-word-embedding-layers-deep-learning-keras/>)

creating and training model

```
In [333]: 1 from tensorflow.keras import layers
          2 from tensorflow.keras.initializers import RandomNormal
          3 import tensorflow as tf
          4 import os
          5 import datetime
          6 import tensorflow_addons as tfa
          7 from sklearn.metrics import f1_score
          8 from tensorflow.keras.layers import Embedding
```

```
In [334]: 1 %load_ext tensorboard
```

The tensorboard extension is already loaded. To reload it, use:
%reload_ext tensorboard

```
In [335]: 1 # callback
2 # Printing costum made accuracy on validation data
3
4 class LossHistory(tf.keras.callbacks.Callback):
5
6     def __init__(self, validation_data):
7         self.x_test = validation_data[0]
8         self.y_test = validation_data[1]
9
10
11     def on_epoch_end(self, epoch, logs={}):
12
13         # # we can get a list of all predicted values at the end of the epoch
14         # # we can use these predicted value and the true values to calculate any custom evaluation score if it is needed for our
15         # # Here we are taking log of all true positives and then taking average of it
16         self.y_pred = self.model.predict(self.x_test)
17         self.y_label_pred = np.argmax(self.y_pred, axis=1)
18
19         #calculating f1_score through sklearn
20         y_pred2 = [1 if x >= 0.5 else 0 for x in self.y_pred[:, 1]]
21         # pdb.set_trace()
22         f1 = f1_score(self.y_test[:, 1], y_pred2, average = "micro")
23
24         print('f1 score', f1)
25
26
```


In [391]:

```
1 def create_model_1():
2     input_layer = Input(shape=(1000))
3     # vocab_size = 1000, embedding_matrix.shape = (1000, 100)
4     embedding = layers.Embedding(vocab_size, 100, weights=[embedding_matrix_word], input_length=1000, trainable=False)(input_layer)
5
6     n1 = layers.BatchNormalization()(embedding)
7
8
9     l2 = layers.Conv1D(100, 4)(n1)
10    l3 = layers.Conv1D(100, 4)(n1)
11    l4 = layers.Conv1D(100, 4)(n1)
12
13    # concatenate 3 output to 1 output
14    l5 = layers.concatenate([l2,l3,l4])
15
16    n2 = layers.BatchNormalization()(l5)
17
18    l6 = layers.MaxPool1D()(n2)
19
20    n3 = layers.Dropout(0.5)(l6)
21
22    l7 = layers.Conv1D(62,5)(n3)
23    l8 = layers.Conv1D(62,5)(n3)
24    l9 = layers.Conv1D(62,5)(n3)
25
26
27    # concatenate 3 output to 1 output
28    L10 = layers.concatenate([l7,l8,l9])
29
30    n4 = layers.BatchNormalization()(L10)
31
32    L11 = layers.MaxPool1D()(n4)
33
34    n5 = layers.Dropout(0.4)(L11)
35
36    L12 = layers.Conv1D(62,3)(n5)
37    L13 = layers.Flatten()(L12)
38    L14 = layers.Dropout(0.4)(L13)
39    L15 = layers.Dense(256, activation= "relu")(L14)
40
41    n6 = layers.BatchNormalization()(L15)
42    n7 = layers.Dropout(0.4)(n6)
43
44    output_layer = layers.Dense(20, activation = 'softmax')(n7)
45
```

```

46     model = Model(input_layer,output_layer, name= "text_df")
47
48     return model
49
50 # tensorboard_callback
51 log_dir = os.path.join("logs", 'fits', datetime.datetime.now().strftime("%Y%m%d-%H%M%S"))
52 tensorboard_callback = tf.keras.callbacks.TensorBoard(log_dir=log_dir,histogram_freq=1,write_graph=True)
53 # custom accuracy callback
54 history_own=LossHistory(validation_data=[padded_docs_test,y_test])
55
56 model = None
57 model = create_model_1()
58
59
60

```

In [392]:

```

1 model.compile(optimizer = tf.keras.optimizers.Adam(learning_rate=1e-3),
2               loss = tf.keras.losses.BinaryCrossentropy(), metrics=[tf.keras.metrics.AUC(multi_label = True, num_labels= 20)
3                               ,tfa.metrics.F1Score(num_classes = 20, average = 'micro')
4                               ])
5
6 model.fit(padded_docs, y_train, batch_size=16, epochs =2, validation_data = [padded_docs_test, y_test], callbacks =[tensorboard_c

```

Epoch 1/2

942/942 [=====] - 733s 776ms/step - loss: 0.2941 - auc_40: 0.4981 - f1_score: 0.0503 - val_loss: 0.3267 - val_auc_40: 0.5042 - val_f1_score: 0.0542

Epoch 2/2

942/942 [=====] - 850s 902ms/step - loss: 0.2138 - auc_40: 0.4988 - f1_score: 0.0504 - val_loss: 0.2970 - val_auc_40: 0.4946 - val_f1_score: 0.0515

Out[392]: <keras.callbacks.History at 0x15228b9b2b0>

In [396]: 1 pip install graphviz-2.38

Note: you may need to restart the kernel to use updated packages.

ERROR: Could not find a version that satisfies the requirement graphviz-2.38 (from versions: none)

ERROR: No matching distribution found for graphviz-2.38

In [393]:

```
1  
2 from tensorflow.keras.utils import plot_model  
3 plot_model(model, "mini_resnet2.png", show_shapes=True)  
4  
5 # y = tf.keras.layers.Conv2D(2, 3, activation='relu', dilation_rate=2, input_shape=input_shape[1:])(x)
```

You must install pydot (`pip install pydot`) and install graphviz (see instructions at <https://graphviz.gitlab.io/download/>) (<https://graphviz.gitlab.io/download/>) for plot_model/model_to_dot to work.

In [369]:

1	<code>model.summary()</code>
---	------------------------------

Model: "text_df"

Layer (type)	Output Shape	Param #	Connected to
input_28 (InputLayer)	[(None, 1000)]	0	[]
embedding_27 (Embedding)	(None, 1000, 100)	100000	['input_28[0][0]']
batch_normalization_4 (Batch Normalization)	(None, 1000, 100)	400	['embedding_27[0][0]']
conv1d_175 (Conv1D)	(None, 998, 164)	49364	['batch_normalization_4[0][0]']
conv1d_176 (Conv1D)	(None, 998, 164)	49364	['batch_normalization_4[0][0]']
conv1d_177 (Conv1D)	(None, 998, 164)	49364	['batch_normalization_4[0][0]']
concatenate_46 (Concatenate)	(None, 998, 492)	0	['conv1d_175[0][0]', 'conv1d_176[0][0]', 'conv1d_177[0][0]']
batch_normalization_5 (Batch Normalization)	(None, 998, 492)	1968	['concatenate_46[0][0]']
max_pooling1d_52 (MaxPooling1D)	(None, 499, 492)	0	['batch_normalization_5[0][0]']
dropout_27 (Dropout)	(None, 499, 492)	0	['max_pooling1d_52[0][0]']
conv1d_178 (Conv1D)	(None, 492, 62)	244094	['dropout_27[0][0]']
conv1d_179 (Conv1D)	(None, 492, 62)	244094	['dropout_27[0][0]']
conv1d_180 (Conv1D)	(None, 492, 62)	244094	['dropout_27[0][0]']
concatenate_47 (Concatenate)	(None, 492, 186)	0	['conv1d_178[0][0]', 'conv1d_179[0][0]', 'conv1d_180[0][0]']
batch_normalization_6 (Batch Normalization)	(None, 492, 186)	744	['concatenate_47[0][0]']
max_pooling1d_53 (MaxPooling1D)	(None, 246, 186)	0	['batch_normalization_6[0][0]']

dropout_28 (Dropout)	(None, 246, 186)	0	['max_pooling1d_53[0][0]']
conv1d_181 (Conv1D)	(None, 244, 62)	34658	['dropout_28[0][0]']
flatten_25 (Flatten)	(None, 15128)	0	['conv1d_181[0][0]']
dropout_29 (Dropout)	(None, 15128)	0	['flatten_25[0][0]']
dense_50 (Dense)	(None, 256)	3873024	['dropout_29[0][0]']
batch_normalization_7 (Batch Normalization)	(None, 256)	1024	['dense_50[0][0]']
dropout_30 (Dropout)	(None, 256)	0	['batch_normalization_7[0][0]']
dense_51 (Dense)	(None, 20)	5140	['dropout_30[0][0]']

```
=====
Total params: 4,897,332
Trainable params: 4,795,264
Non-trainable params: 102,068
```

In [352]: 1 len(padded_docs)

Out[352]: 15062

```
In [355]: 1 model.compile(optimizer = tf.keras.optimizers.Adam(learning_rate=1e-3),
2           loss = tf.keras.losses.BinaryCrossentropy(), metrics=[tf.keras.metrics.AUC(multi_label = True, num_labels= 20)
3           ,tfa.metrics.F1Score(num_classes = 20, average = 'micro')
4           , tf.keras.metrics.Accuracy()])
5
6 model.fit(padded_docs, y_train, batch_size=16, epochs =4, validation_split=0.2, callbacks =[tensorboard_callback] )
```

Epoch 1/4

754/754 [=====] - 209s 275ms/step - loss: 0.2132 - auc_33: 0.4987 - f1_score: 0.0510 - accuracy: 1.3694e-04
- val_loss: 0.2069 - val_auc_33: 0.4980 - val_f1_score: 0.0498 - val_accuracy: 0.0000e+00

Epoch 2/4

754/754 [=====] - 232s 308ms/step - loss: 0.2061 - auc_33: 0.4975 - f1_score: 0.0500 - accuracy: 5.3946e-05
- val_loss: 0.2025 - val_auc_33: 0.4987 - val_f1_score: 0.0561 - val_accuracy: 0.0000e+00

Epoch 3/4

754/754 [=====] - 212s 281ms/step - loss: 0.2034 - auc_33: 0.4946 - f1_score: 0.0493 - accuracy: 0.0000e+00
- val_loss: 0.2003 - val_auc_33: 0.5003 - val_f1_score: 0.0461 - val_accuracy: 0.0000e+00

Epoch 4/4

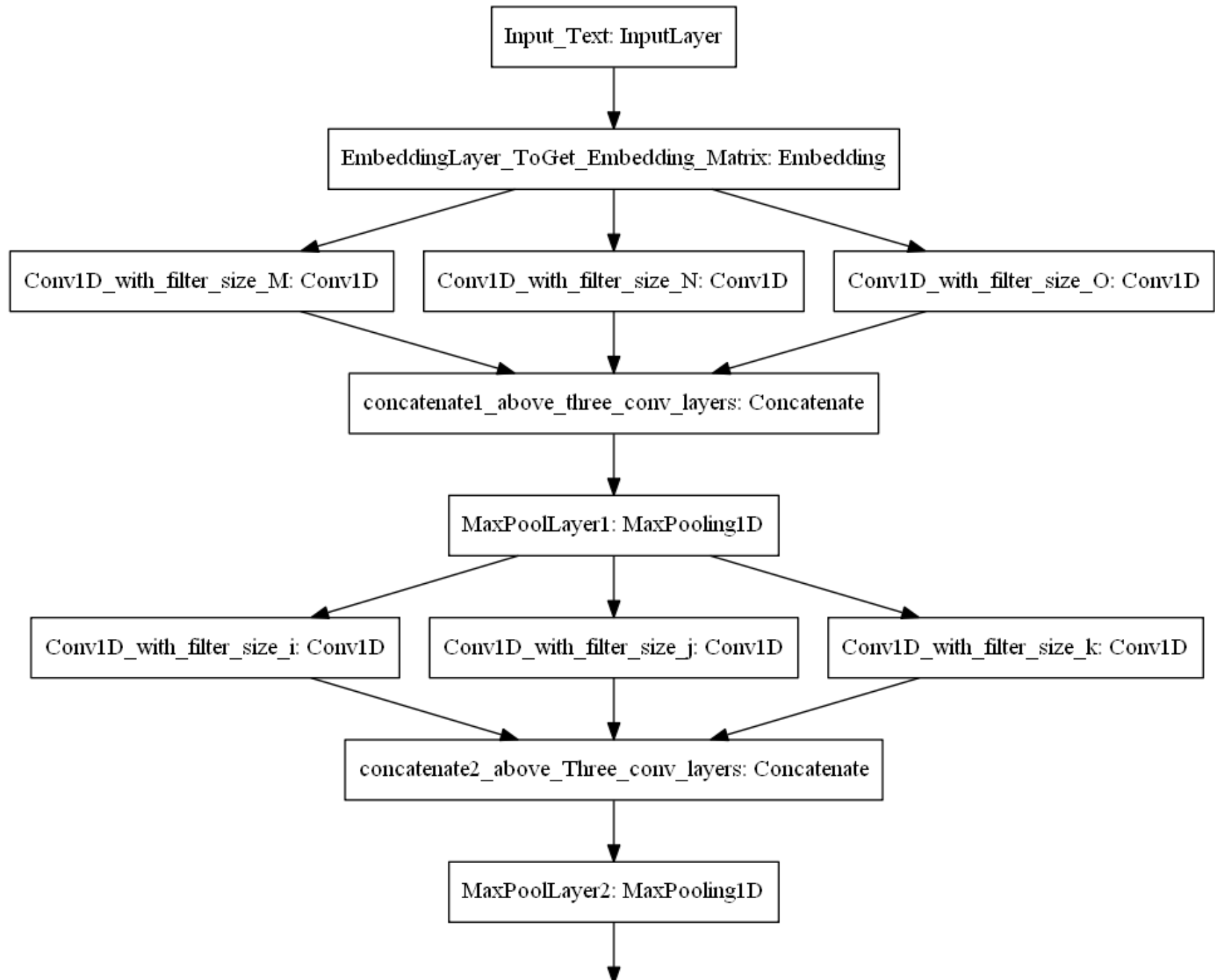
754/754 [=====] - 242s 320ms/step - loss: 0.2014 - auc_33: 0.4960 - f1_score: 0.0499 - accuracy: 0.0000e+00
- val_loss: 0.1991 - val_auc_33: 0.4987 - val_f1_score: 0.0498 - val_accuracy: 0.0000e+00

Out[355]: <keras.callbacks.History at 0x151edce8310>

In [140]: 1 %tensorboard --logdir logs

Reusing TensorBoard on port 6006 (pid 8704), started 0:46:11 ago. (Use '!kill 8704' to kill it.)



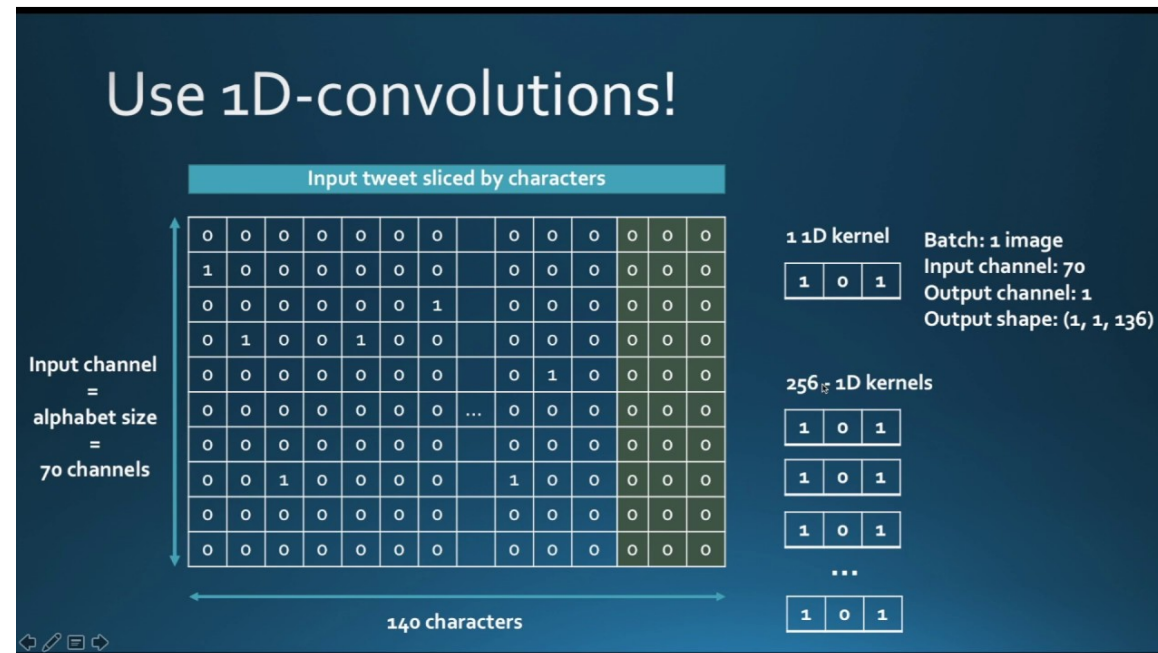


ref: 'https://i.imgur.com/fv1GvFJ.png'

Conv1D with filter size P: Conv1D

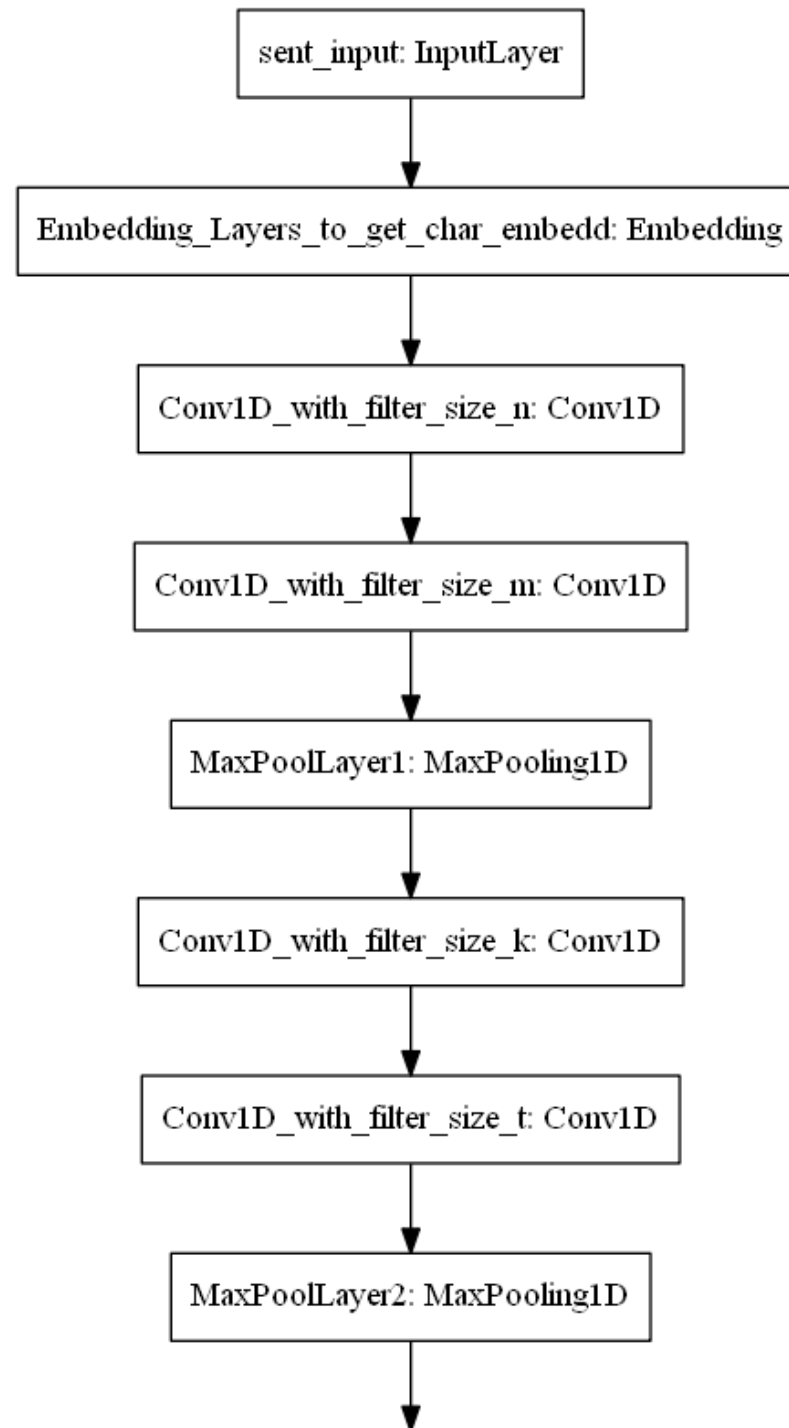
1. all are Conv1D layers with any number of filter and filter sizes, there is no restriction on this.
2. use concatenate layer is to concatenate all the filters/channels.
3. You can use any pool size and stride for maxpooling layer.
4. Don't use more than 16 filters in one Conv layer because it will increase the no of params.
(Only recommendation if you have less computing power)
5. You can use any number of layers after the Flatten Layer.

Model-2 : Using 1D convolutions with character embedding



Here are the some papers based on Char-CNN

1. Xiang Zhang, Junbo Zhao, Yann LeCun. [Character-level Convolutional Networks for Text Classification](http://arxiv.org/abs/1509.01626) (<http://arxiv.org/abs/1509.01626>). NIPS 2015
2. Yoon Kim, Yacine Jernite, David Sontag, Alexander M. Rush. [Character-Aware Neural Language Models](https://arxiv.org/abs/1508.06615) (<https://arxiv.org/abs/1508.06615>). AAAI 2016
3. Shaojie Bai, J. Zico Kolter, Vladlen Koltun. [An Empirical Evaluation of Generic Convolutional and Recurrent Networks for Sequence Modeling](https://arxiv.org/pdf/1803.01271.pdf) (<https://arxiv.org/pdf/1803.01271.pdf>).
4. Use the pretrained char embeddings <https://github.com/minimaxir/char-embeddings/blob/master/glove.840B.300d-char.txt> (<https://github.com/minimaxir/char-embeddings/blob/master/glove.840B.300d-char.txt>).



In [356]:

```
1 import shutil
2 shutil.rmtree(r'C:\Users\ashutosh tiwari\Documents\ML_programs\assignment_aai\21_document with cnn with text data\logs')
```

In [357]:

```
1 vocab_size, char_embed_matrix.shape
```

Out[357]: (1000, (94, 300))

In [360]:

```
1
2 def create_model_2():
3
4     input_layer = tf.keras.Input(shape=(1000,))
5     char_emb = layers.Embedding(char_vocab_size,300,embeddings_initializer=tf.keras.initializers.Constant(char_embed_matrix),trainable=False)(input_layer)
6     convlay1 = layers.Conv1D(32,5,activation='relu',kernel_initializer=tf.keras.initializers.HeNormal(seed=5))(char_emb)
7     convlay2 = layers.Conv1D(32,3,activation='relu',kernel_initializer=tf.keras.initializers.HeNormal(seed=10))(convlay1)
8     maxpool1 = layers.MaxPool1D(2,2)(convlay2)
9     convlay3 = layers.Conv1D(8,3,activation='relu',kernel_initializer=tf.keras.initializers.HeNormal(seed=15))(maxpool1)
10    convlay4 = layers.Conv1D(16,5,activation='relu',kernel_initializer=tf.keras.initializers.HeNormal(seed=20))(convlay3)
11    maxpool2 = layers.MaxPool1D(2,2)(convlay4)
12    flat = layers.Flatten()(maxpool2)
13    drop = layers.Dropout(0.2)(flat)
14    dlay = layers.Dense(40,activation='relu',kernel_initializer=tf.keras.initializers.HeNormal(seed=25))(drop)
15    output_layer = layers.Dense(20,activation='softmax',kernel_initializer=tf.keras.initializers.HeNormal(seed=30))(dlay)
16
17    model = Model(inputs=input_layer,outputs=output_layer)
18    return model
19
20 model = None
21 model = create_model_2()
22
23
24 model.compile(optimizer = tf.keras.optimizers.Adam(learning_rate=1e-3),
25               loss = tf.keras.losses.BinaryCrossentropy(), metrics=[tf.keras.metrics.AUC(),tfa.metrics.F1Score(num_classes = 3,
26
```

In [361]: 1 model.summary()

Model: "model_3"

Layer (type)	Output Shape	Param #
=====		
input_24 (InputLayer)	[(None, 1000)]	0
embedding_23 (Embedding)	(None, 1000, 300)	28200
conv1d_152 (Conv1D)	(None, 996, 32)	48032
conv1d_153 (Conv1D)	(None, 994, 32)	3104
max_pooling1d_46 (MaxPooling1D)	(None, 497, 32)	0
conv1d_154 (Conv1D)	(None, 495, 8)	776
conv1d_155 (Conv1D)	(None, 491, 16)	656
max_pooling1d_47 (MaxPooling1D)	(None, 245, 16)	0
flatten_23 (Flatten)	(None, 3920)	0
dropout_23 (Dropout)	(None, 3920)	0
dense_46 (Dense)	(None, 40)	156840
dense_47 (Dense)	(None, 20)	820
=====		
Total params: 238,428		
Trainable params: 210,228		
Non-trainable params: 28,200		

In [362]:

```
1  
2 model.fit(padded_docs_char, y_train, batch_size=16, epochs =5, validation_split=0.2, callbacks =[tensorboard_callback] )  
3
```

Epoch 1/5

754/754 [=====] - 87s 113ms/step - loss: 0.2044 - auc_35: 0.5099 - f1_score: 0.0532 - val_loss: 0.1995 - val_auc_35: 0.5121 - val_f1_score: 0.0461

Epoch 2/5

754/754 [=====] - 90s 119ms/step - loss: 0.1996 - auc_35: 0.5083 - f1_score: 0.0520 - val_loss: 0.1987 - val_auc_35: 0.5110 - val_f1_score: 0.0514

Epoch 3/5

754/754 [=====] - 85s 112ms/step - loss: 0.1992 - auc_35: 0.5116 - f1_score: 0.0511 - val_loss: 0.1993 - val_auc_35: 0.5151 - val_f1_score: 0.0558

Epoch 4/5

754/754 [=====] - 87s 116ms/step - loss: 0.1990 - auc_35: 0.5183 - f1_score: 0.0524 - val_loss: 0.1991 - val_auc_35: 0.5085 - val_f1_score: 0.0561

Epoch 5/5

754/754 [=====] - 83s 110ms/step - loss: 0.1987 - auc_35: 0.5290 - f1_score: 0.0602 - val_loss: 0.1996 - val_auc_35: 0.5043 - val_f1_score: 0.0478

Out[362]: <keras.callbacks.History at 0x151e99ebb20>

In [363]: 1 %tensorboard --logdir logs

Reusing TensorBoard on port 6006 (pid 8704), started 18:13:22 ago. (Use '!kill 8704' to kill it.)



tensorboard analysis for both model

Link :- https://docs.google.com/document/d/1z74uiJtv1MOyFlt3SJRuVwK8LwjD0teXdg-_4l92T7E/edit?usp=sharing
(https://docs.google.com/document/d/1z74uiJtv1MOyFlt3SJRuVwK8LwjD0teXdg-_4l92T7E/edit?usp=sharing).