

## 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, at that you will get documents.rar folder.
If you unzip that, you will get total of 18828 documents. document name is defined
as 'classlabel_documentNumber' thatLabel'.
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.
```

```
In [ ]: from google.colab import drive
drive.mount('/content/drive')
```

### sample document

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

In article <65882@mimsy.umd.edu> manglee@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", 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?
```

### Preprocessing:

```
useful links: http://www.pyrege.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, test@dm2.d.com]-->[dm1.d.com, dm3.dm4.com]-->[dm1,d,com,dm2,dm3,com]-->
[dm1,d,com,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,
manglee@cs.umd.edu]

preprocessing:
[jcopelan@nyx.cs.du.edu, 65882@mimsy.umd.edu, manglee@cs.umd.edu] => [nyx cs du edu mimsy umd edu cs umd
edu] =>
[nyx.edu mimsy umd.edu cs umd.edu]

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

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\n\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 "< keyword >".
> 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, 's-> is, i've -> 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, as
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 some phrases, and 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 nltk.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://stackoverflow.com/a/11837224/4084039
https://www.nltk.org/howto/tree.html
https://stackoverflow.com/a/44294377/4084039
```

```
In [ ]: print("i am living in the New York -->", list(chunks))
print(" ")
print("50")
print(" ")
print("My name is Mridul Mazumdar -->", list(chunks))
```

```
We did chunking for above two lines and then we got one list where each word is mapped to a
POS_tags of speech and also if you see "New York" and "Mridul Mazumdar",
they got combined and represented as a tree. "New York" was referred as "GPE" and "Mridul Mazumdar"
was referred as "PERSON".
```

```
so now you have to combine the "New York" with "_" i.e "New_York"
```

```
and remove the "Mridul Mazumdar" from the above sentence because it is a person.
```

```
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 words 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_" (dr_berlin ==> berlin). i.e remove the words
which are length less than or equal to 2 after splitting 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.
```

### 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 Average 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 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
```

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/>

### Model-2 : Using 1D convolutions with character embedding

```
Use 1D-convolutions!
```

```
Input tweet sliced by characters
```

```
Input_channel = 1
alpha_beta_size = 1
yo_channels = 1
preprocessed_email = []
preprocessed_subjects = []
preprocessed_text = []
```

```
for x in tqdm(os.listdir('documents/')):
    preprocessed_email.append(x)
```

```
for x in tqdm(os.listdir('documents/')):
    preprocessed_subjects.append(x)
```

```
for x in tqdm(os.listdir('documents/')):
    preprocessed_text.append(x)
```

```
print("preprocessed_email:", len(preprocessed_email))
print("preprocessed_subjects:", len(preprocessed_subjects))
print("preprocessed_text:", len(preprocessed_text))
```

```
df=pd.DataFrame(data)
```

```
df.to_csv('doc_CNN.csv',index=False)
```

```
df=pd.read_csv('doc_CNN.csv').dropna()
df.set_index('id', inplace=True)
```

```
df['text']=df['text'].apply(lambda x: x[1:-1])
```