

FIT5196 Task 1 in Assessment 1

Student Name: Sarthak Sareen

Student ID: 30761182

Date: 13/09/2020

Version: 1.0

- pandas (for dataframe, included in Anaconda Python 2.7)
- re (for regular expression, included in Anaconda Python 2.7)
- os (interacting with os modules included in Anaconda Python 2.7)
- langid (to check string language, included in Anaconda Python 2.7)

1. Indroduction

In the task we have to analyze the textual data that is the extracting the semi structures files. Each file contains the id, text and the date created of the tweet. We have to take out the three things i.e id which is a 19 digit number, text which is the actual tweet and the date created that is the date at which the tweet was posted. After taking out we have to write all the data to a XML file.

1. Import libraries

In []:



```
1
2 #importing the libraies
3 import os
4 import re
5 import langid
```

2. Loading all the files

In []:



```
1 #reading all the file which is present the specified folder path where the file is.txt
2 path = "Ass-1-5196/"
3 all_files = os.listdir(path)
4 new_list = []
5 text2 = ''
6 for root, dirs, files in os.walk(path):
7     for file in files:
8         if file.endswith('.txt'):
9             with open(os.path.join(root, file), 'r',encoding='utf-8') as f:
10                 text = f.read()
11                 #saving the all the text into one string
12                 text2 = text2 + text
13
14
```

3. Function to perform the task

In []:



```

1  #function distribute that will take out the data from the string text2 using regex and
2  def distribute(data):
3      #using regex extracting id into a list
4      list_id = re.findall(r'"id":"[0-9]{19}',data)
5      #using regex extracting created date into a list
6      list_created_date = re.findall(r'"created_at":"[0-9]{4}-[0-9]{2}-[0-9]{2}',data)
7      #using regex extracting text into a list
8      list_text = list_text = re.findall(r'"text":".*?"',data)
9      #defining a new dictionary
10     d={}
11     #adding the first element of the created date list as a key in dictionary (d)
12     #with values first element of id list and text list
13     #first_text = []
14     first_text = list_text[0].replace("//n","/n")
15     first_text = first_text.replace(">","&gt;")
16     first_text = first_text.replace("&","&amp;")
17     first_text = first_text.replace("'", "&quot;")
18     first_text = first_text.replace('"', "&apos;")
19     d[list_created_date[0].split('')[2]] = [[list_id[0].split('')[3],first_text]]
20     #assigning this count for the index which is used to take out the id and text from
21     count = 0
22     #Loop in created date list
23     for key in list_created_date:
24         count= count + 1
25         if count >= len(list_created_date):
26             #breaking the loop if count is greater than length of list created date
27             break
28         else:
29             #assigning the id and text to a variable with the help of the count variable
30             id_element,text_element = list_id[count].split('')[3],list_text[count].split('')[3]
31             text_element = text_element.encode('utf-16', 'surrogatepass').decode('utf-16')
32             text_element = text_element.replace("//n","/n")
33             text_element = text_element.replace("<","&lt;")
34             text_element = text_element.replace(">","&gt;")
35             text_element = text_element.replace("&","&amp;")
36             text_element = text_element.replace("'", "&quot;")
37             text_element = text_element.replace('"', "&apos;")
38             #checking the text is in english or not. Checking id and text are not empty
39             if langid.classify(text_element)[0]=='en' and id_element!="" and text_element!="":
40                 keydic = key.split('')[2]
41                 #checking the key is present in the dictionary or not
42                 if keydic in d:
43                     #appending the id and text to the present key
44                     d[keydic].append([id_element,text_element])
45                 else:
46                     #making the key and adding the values in it
47                     d[keydic] = [[id_element,text_element]]
48             #if the text is not english then continue
49             else:
50                 continue
51     #returning the dictionary
52     return d
53 #Calling function distribute function and saving it into final variable
54 final = distribute(text2)

```

In the above function the steps followed are :

- 1) The function is taking text which is the combined text of all the files as Data.
- 2) Taking out ID from the text using regex - "id":"[0-9]{19}" this means re matches exactly this expression "id:" and then 19 more digits after this. We are taking 19 digits is because we are given that the id is 19 digits. Saving this into list_id.
- 3) Taking out the Created date from the text using regex - created_at":"[0-9]{4}-[0-9]{2}-[0-9]{2}" this means that re matches exactly this expression created_at:" and then 4 digits which are in between 0-9 followed by this expression '-' then 2 digits in between 0-9 followed by this expression '-' then 2 digits in between 0-9. Saving this regex in list_created date
- 4) Taking out the text from the text using regex - "text":".?" *this means it matches this expression "text:" then . means anything ?" means stop when we get 0 or 1 this expression "*. Saving this into list_text.
- 5) defining a empty dictionary
- 6) Adding the first element of created_date_list as key of the dictionary and values are the list of first element of list_id and first element of list_created. Splitting the three of the to remove the unwanted things in the text. Such as the created_date_list has first element as 'created_at' : "2020-03-23 splitting this to remove the created_date and taking out only 2020-02-23 from it. Performing the same task for id and text.
- 7) Initialising count which will be used for the index to take out the id element from the list and text.
- 8) Looping in the created_date_list to take this as key of the dictionary.
 - i) First check is the count should be less than length of the created date list because if the count which is used for index should be in the limit of the created date list length only.
If this happens the loop will break.
 - ii) Then in the else condition id and text is being assigned to variables id_element and text_element respectively using the count as the index and splitting the unwanted text from it.
 - iii) A) check if the text is english and id_element is not empty and text_element is not empty.
 --splitting the created_date_list to take out only the date
 --checking if the key is already present then append the id_element and text_element as values into the key
 --else making the new key and taking id_element and text_element as values into the specific key.

 B) else
 --if the text is english and id_element is not empty and text_element is not empty then continue. Nothing will go into the new dictionary. This will skip the entry.
- 9) returning the dictionary
- 10) calling the function and saving it into a variable final.

In []:



```

1  #making string
2  stri = '<?xml version="1.0" encoding="UTF-8"?>' + '\n' + '<data>\n'
3  #iterating in the final dictionary
4  for key,value in final.items():
5      #adding keys of the dictionary in the string
6      stri = stri + '\t' + '<tweets date="'+str(key) + '">\n'
7      #iterating in the values of the dictionary
8      for i in value:
9          #adding the tweet id and the text in the string
10         stri = stri + '\t\t' + '<tweet id="'+ str(i[0]) + '">' + str(i[1]) + '</tweet>'
11         stri = stri + '</tweets>\n'
12     stri = stri + '</data>'
13     stri= stri.replace("\n","")

```

In the above cell a string is made which contains all the tweets date's tweet id and the text of the specific date and concatenating all of them into one string which is stri. In between the string adding the tags of the data, tweet date and tweets to make the text structured like an XML doc.

In []:



```

1  #writing the string into the xml format
2  with open("30761182.xml", "w", encoding='utf-8') as text_file:
3      text_file.write(stri)

```

Summary

1) First reading all the text files into one string.

the string contains all the text of all the files.

2) Then using regex extracting id, text, created date :

```
list_id = ["id":"1234554412323456789","id":"1234533412341112345,...."]
```

```
list_created_date = ["created date":"2020-03-23","created date":"2020-04-22,...."]
```

```
list_text = ["text":"coronara virus","text":"virus is bad",.....]
```

3) Making a new dictionary with all the handling of the text. The dictionary formed is like :

```
{ "2020-03-22" : [[1233443300987475849.'coronara virus'],[2134443098712345980,'virus is bad'],....] }
```

4) Then making a string according to the XML document

5) writing the string in the output XML file.

The output file is being created as an XML document having tags and is parsable.

