

Python Project Document

150101027- I.N.Dilip Kumar

150101051- Hareesh Reddi

150101032- L.Sai Shobith

April 10, 2017

Contents

1	Introduction	3
2	Algorithm	4
3	Algorithm Analysis	5
4	Code Description	6
4.1	Section1	6
4.2	Section2	6
4.3	Section3	6
4.4	Section4	7

1 Introduction

In the given attack data file, each category has 10 directories. We select 7 directories for the training and use the other 3 directories for testing.

The problem statement is:

Split the Attack data of each category (Hydra-FTP, Hydra-SSH, Adduser, Java-Meterpreter, Meterpreter and Webshell) into 70% training data and 30 % test data. For instance there are 10 folders in "Adduser" attack. Therefore, 7 of these folders are to be used for training and 3 folders are to be used for testing.

For the Normal data, files in "Training_Data_Master" folder are to be used as training data and files in "Validation_Data_Master" folder are to be used as test data.

Write a python script to find the frequency of occurrences of all unique 3-grams, 5-grams and 7-grams system call sequences in the training data for both Attack data (across all categories of attack) and Normal data.

Perform the same task on files in the Training_Data_Master to obtain all the unique 3-grams, 5-grams and 7-grams.

Once you have obtained the frequencies of all the unique n-grams terms in the training data, use the top 30% n-grams terms with the highest frequency to create a data set.

Apply the same procedure to generate the test dataset from the test files of the attack data (for all attack types) and the normal files in the "Validation_Data_Master" using the top 30% 3-grams terms with highest frequencies obtained during the training phase. The classifier model developed during the training phase will finally be validated on the Test dataset.

2 Algorithm

1. Select the first seven folders for each category and find all the text files in all the directories.
2. Find all the seven grams for each file.
3. Find the frequency of each seven gram and keep them in a dictionary.
4. Using the seven grams, find the 5 grams and 3 grams by taking the first three numbers for 3 gram and first five for 5 gram.
5. Find the frequency for each of the 5 grams and for 3 grams and keep them in a dictionary.
6. Sort the dictionary based on the values as frequencies.
7. Select the top 30 percent of the n-grams and do this for each category and for the training data set.
8. Concat all of the top 30 percent n-grams and keep them in a single dictionary and let us call it feature file.
9. Now for each text file in the test data and validation data master create another text file which contains the 7 gram in the feature file and if it is in the text file write it's frequency and if it is not present write 0 as the frequency.

3 Algorithm Analysis

1. $O(\text{No.Of.TextFiles})$ is the time complexity of first step
2. $O(\text{No.Of.SevenGrams})$ is the time complexity of second step
3. $O(\text{No.Of.SevenGrams})$ is the time complexity for finding the frequencies and $O(N \log N)$ where N is No.Of.SevenGrams for sorting the dictionaries
4. $O(\text{No.Of.ThreeGrams} + \text{No.Of.FiveGrams}) \leq O(\text{No.Of.SevenGrams})$ for the fourth step
5. $O(\text{No.Of.ThreeGrams} + \text{No.Of.FiveGrams})$ is the time complexity of fifth step
6. $O(N \log N)$ where N is No.Of.SevenGrams is time complexity for sorting the dictionaries
7. $O(N)$ where N is No.Of.SevenGrams is time complexity for selecting the top 30% N -Grams from the dictionaries
8. $O(N)$ where N is No.Of.SevenGrams is time complexity for concatenating the dictionaries
9. $O(N * \text{No.Of.TextFiles})$ is the time complexity for the ninth step

4 Code Description

This is the code description

4.1 Section1

The unify(Adduser) will return three dictionaries of 7 tuples, 5tuples, 3tuples as keys and their frequencies as values

```
40 def unify(specific_string):
41     dictionary_7 = collections.OrderedDict()
42     dictionary_5 = collections.OrderedDict()
43     dictionary_3 = collections.OrderedDict()
44     list_of_text_files = get_Attack_Data_Master_files(specific_string)
45     for text_file in list_of_text_files:
46         grams_7_list = List(ngrams(text_file.split(), 7))
47         for item_gram in grams_7_list:
48             if item_gram in dictionary_7:
49                 dictionary_7[item_gram] = dictionary_7[item_gram] + 1
50             else:
51                 dictionary_7[item_gram] = 1
52         grams_5_list = List(ngrams(text_file.split(), 5))
53         for item_gram in grams_5_list:
54             if item_gram in dictionary_5:
55                 dictionary_5[item_gram] = dictionary_5[item_gram] + 1
56             else:
57                 dictionary_5[item_gram] = 1
58         grams_3_list = List(ngrams(text_file.split(), 3))
59         for item_gram in grams_3_list:
60             if item_gram in dictionary_3:
61                 dictionary_3[item_gram] = dictionary_3[item_gram] + 1
62             else:
63                 dictionary_3[item_gram] = 1
64
65     length(dictionary_7, dictionary_5, dictionary_3, specific_string)
66     return dictionary_7, dictionary_5, dictionary_3
67
```

Figure 1: Fig1

4.2 Section2

The Sorted() function will sort the dictionary based on their values in decreasing order and returns the list of keys in that order

```
156 Sorted_List_Training = sorted(dictionary_Training, key=dictionary_Training.__getitem__, reverse = True)
157 Sorted_List_Adduser = sorted(dictionary_Adduser, key=dictionary_Adduser.__getitem__, reverse = True)
158 Sorted_List_HydraFTP = sorted(dictionary_HydraFTP, key=dictionary_HydraFTP.__getitem__, reverse = True)
159 Sorted_List_HydraSSH = sorted(dictionary_HydraSSH, key=dictionary_HydraSSH.__getitem__, reverse = True)
160 Sorted_List_JavaMeterpreter = sorted(dictionary_JavaMeterpreter, key=dictionary_JavaMeterpreter.__getitem__, reverse = True)
161 Sorted_List_Meterpreter = sorted(dictionary_Meterpreter, key=dictionary_Meterpreter.__getitem__, reverse = True)
162 Sorted_List_WebShell = sorted(dictionary_WebShell, key=dictionary_WebShell.__getitem__, reverse = True)
163
164
165
```

Figure 2: Fig2

4.3 Section3

The Sorted_List_Adduser is a list which contains the tuples in sorted order of their frequencies The same process above is done for each type of attack

```
166 Sorted_List_All_Top30=collections.OrderedDict()
167 for s in Sorted_List_Adduser[:int(0.3*len(Sorted_List_Adduser))]:
168     if s in Sorted_List_All_Top30:
169         continue
170     else:
171         Sorted_List_All_Top30[s] = 1
172 for s in Sorted_List_HydraFTP[:int(0.3*len(Sorted_List_HydraFTP))]:
173     if s in Sorted_List_All_Top30:
174         continue
175     else:
176         Sorted_List_All_Top30[s] = 1
177 for s in Sorted_List_HydraSSH[:int(0.3*len(Sorted_List_HydraSSH))]:
178     if s in Sorted_List_All_Top30:
179         continue
180     else:
181         Sorted_List_All_Top30[s] = 1
182 for s in Sorted_List_JavaMeterpreter[:int(0.3*len(Sorted_List_JavaMeterpreter))]:
183     if s in Sorted_List_All_Top30:
184         continue
185     else:
186         Sorted_List_All_Top30[s] = 1
187 for s in Sorted_List_Meterpreter[:int(0.3*len(Sorted_List_Meterpreter))]:
188     if s in Sorted_List_All_Top30:
189         continue
190     else:
191         Sorted_List_All_Top30[s] = 1
192 for s in Sorted_List_WebShell[:int(0.3*len(Sorted_List_WebShell))]:
193     if s in Sorted_List_All_Top30:
194         continue
195     else:
196         Sorted_List_All_Top30[s] = 1
197 for s in Sorted_List_Training[:int(0.3*len(Sorted_List_Training))]:
198     if s in Sorted_List_All_Top30:
199         continue
200     else:
201         Sorted_List_All_Top30[s] = 1
202
203
```

Figure 3: Fig3

4.4 Section4

A dictionary is formed with tuples whose frequencies are top 30 percent in the sorted dictionary for each kind of attack and concatenated into one and is called the feature set. Now the feature set is tested against each file in the Adduser folder and a file is created in which if the tuple in the feature set is present in the n-grams of the selected file then the n-gram and its frequency in the selected file is printed, if it is not present then it will print the n-gram and 0 as the frequency. The same is done for 5 gram and 3 gram

```

309 for file_object in sec_get_Attack_Data_Master_files("Adduser"):
310     new_file_path = get_new_file_ob(file_object,"7")
311     list_new_file = {}
312     file_data = file_object.read()
313     for x in List(ngrams(file_data.split(), 7)):
314         if x in list_new_file:
315             list_new_file[x] += 1
316         else:
317             list_new_file[x] = 1
318     with open(new_file_path, "w") as f:
319         for key,value in Sorted_List_All_Top30.items():
320             if key in list_new_file:
321                 f.write('%s : %d\n' % (str(key), list_new_file[key]))
322             else:
323                 f.write('%s : 0\n' % str(key))
324     new_file_path = get_new_file_ob(file_object,"5")
325     list_new_file = {}
326     file_data = file_object.read()
327     for x in List(ngrams(file_data.split(), 5)):
328         if x in list_new_file:
329             list_new_file[x] += 1
330         else:
331             list_new_file[x] = 1
332     with open(new_file_path, "w") as f:
333         for key,value in Sorted_List_All_Top30_5tuple.items():
334             if key in list_new_file:
335                 f.write('%s : %d\n' % (str(key), list_new_file[key]))
336             else:
337                 f.write('%s : 0\n' % str(key))
338     new_file_path = get_new_file_ob(file_object,"3")
339     list_new_file = {}
340     file_data = file_object.read()
341     for x in List(ngrams(file_data.split(), 3)):
342         if x in list_new_file:
343             list_new_file[x] += 1
344         else:
345             list_new_file[x] = 1
346     with open(new_file_path, "w") as f:
347         for key,value in Sorted_List_All_Top30_3tuple.items():
348             if key in list_new_file:
349                 f.write('%s : %d\n' % (str(key), list_new_file[key]))
350             else:
351                 f.write('%s : 0\n' % str(key))

```

Figure 4: Fig4