

Index

1. Introduction.....	03
2. Python Code.....	04
3. Discription of loops.....	10
4. Discription of Functions.....	10
5. Graph of time complexity.....	12

Introduction

Computer Security is one of the most important aspects of Computer Science as in this modern age sensitive and important data is stored in computers. Its knowledge or manipulation could have catastrophic effects. To prevent attacks on computers to hack it, we must first detect attacks. System logs is an important way to detect attacks. A particular combination of numbers in the log file can imply there was an attack.

Example- multiple login/wrong password sequence of numbers can indicate the attempt to guess someone's password.

Thus finding such sequences can be useful in improving security features of computers and websites. One way to do so is find recurring sequences in known attacks. Thus these sequences can then help identify that attack. Thus documenting sequences with highest frequencies, it can help detect that attack, and this project uses this understanding and tries to detect attack.

Note:- Files are merged using command prompt using command merging files in folder :- **type *.txt >merged.txt**

Listing 1: Python code

```
1 import operator
2 import time
3 result=[]
4 def ngrams(input, n):
5     input = input.split(' ')
6     output = {}
7     for i in range(len(input)-n+1):
8         g = ' '.join(input[i:i+n])
9         output.setdefault(g, 0)
10        output[g] += 1
11    return output
12 # https://docs.python.org/3/library/operator.html
13 # Source for itemgetter() function
14 # Return a callable object that fetches item from its operand
15 # using the operands __getitem__() method.
16 # If multiple items are specified, returns a tuple of lookup values. For
    example:
17 # After f = itemgetter(2), the call f(r) returns r[2].
18 # After g = itemgetter(2, 5, 3), the call g(r) returns (r[2], r[5], r[3]).
19 # def itemgetter(*items):
20 #     if len(items) == 1:
21 #         item = items[0]
22 #         def g(obj):
23 #             return obj[item]
24 #     else:
25 #         def g(obj):
26 #             return tuple(obj[item] for item in items)
27 #     return g
28 def sort_ngrams(output):
29     input=sorted(output.items(),key=operator.itemgetter(1),reverse=True)
30     return input
31
32 def top30percent(x):
33     length=len(x)*3/10
34     count=1
35     for i in x:
36         if(count>length):
37             break
38         result.append(i[0])
39         count+=1
40 def main():
41     n=input("Enter n-gram value: ")
42
43     #This assumes that your program takes at least a tenth of second to run.
44     start_time = time.time() +0.1
45
46     Allfiles=[]
47     Adduser=("Adduser.txt","r")
48     Hydra_FTP=("Hydra.FTP.txt","r")
49     Hydra_SSH=("Hydra.SSH.txt","r")
50     Java_Meterpreter=("Java.Meterpreter.txt","r")
51     Meterpreter=("Meterpreter.txt","r")
52     Web_Shell=("Web.Shell.txt","r")
53     NormalData=("NormalData.txt","r")
54
55     Allfiles.append(Adduser)
56     Allfiles.append(Hydra_FTP)
57     Allfiles.append(Hydra_SSH)
58     Allfiles.append(Java_Meterpreter)
```

```

59 Allfiles.append(Meterpreter)
60 Allfiles.append(Web_Shell)
61 Allfiles.append(NormalData)
62
63 train_dataset=open("trained.txt","w")
64
65 name=["Adduser","Hydra_FTP","Hydra_SSH","Java_Meterpreter","Meterpreter","
66       Web_Shell","Normal"]
67
68 names=["Adduser","Hydra_FTP","Hydra_SSH","Java_Meterpreter","Meterpreter",
69        "Web_Shell"]
70
71 Adduser=[]
72 Hydra_FTP=[]
73 Hydra_SSH=[]
74 Java_Meterpreter=[]
75 Meterpreter=[]
76 Web_Shell=[]
77 Normal=[]
78
79 Adduser_vector=[]
80 Hydra_FTP_vector=[]
81 Hydra_SSH_vector=[]
82 Java_Meterpreter_vector=[]
83 Meterpreter_vector=[]
84 Web_Shell_vector=[]
85 Normal_vector=[]
86
87 Adduser_vector_final=[]
88 Hydra_FTP_vector_final=[]
89 Hydra_SSH_vector_final=[]
90 Java_Meterpreter_vector_final=[]
91 Meterpreter_vector_final=[]
92 Web_Shell_vector_final=[]
93 Normal_vector_final=[]
94
95 count=0
96 for i in range(1,65):
97     count+=1
98     s="Adduser"+str(count)+".txt"
99     Adduser.append(s)
100     s=""
101 count=0
102 for i in range(1,115):
103     count+=1
104     s="Hydra_FTP"+str(count)+".txt"
105     Hydra_FTP.append(s)
106     s=""
107 count=0
108 for i in range(1,124):
109     count+=1
110     s="Hydra_SSH"+str(count)+".txt"
111     Hydra_SSH.append(s)
112     s=""
113 count=0
114 for i in range(1,88):
115     count+=1
116     s="Java_Meterpreter"+str(count)+".txt"
117     Java_Meterpreter.append(s)
118     s=""

```

```

117 count=0
118 for i in range(1,54):
119     count+=1
120     s="Meterpreter"+str(count)+".txt"
121     Meterpreter.append(s)
122     s=""
123 count=0
124 for i in range(1,84):
125     count+=1
126     s="Web_Shell"+str(count)+".txt"
127     Web_Shell.append(s)
128     s=""
129 count=0
130 for i in range(1,834):
131     count+=1
132     s="Normaldata"+str(count)+".txt"
133     Normal.append(s)
134     s=""
135
136 for j in Allfiles:                                     #Loop 1
137     output_temp=ngrams(open(j[0],j[1]).read(),n) # output_temp is dict
138     sorted_elements=sort_ngrams(output_temp) #sort is list of tuples
139     top30percent(sorted_elements)
140
141 for i in Adduser:
142     output_temp=ngrams(open(i,"r").read(),n)
143     Adduser_vector.append(output_temp)
144 for i in Hydra_FTP:
145     output_temp=ngrams(open(i,"r").read(),n)
146     Hydra_FTP_vector.append(output_temp)
147 for i in Hydra_SSH:
148     output_temp=ngrams(open(i,"r").read(),n)
149     Hydra_SSH_vector.append(output_temp)
150 for i in Java_Meterpreter:
151     output_temp=ngrams(open(i,"r").read(),n)
152     Java_Meterpreter_vector.append(output_temp)
153 for i in Meterpreter:
154     output_temp=ngrams(open(i,"r").read(),n)
155     Meterpreter_vector.append(output_temp)
156 for i in Web_Shell:
157     output_temp=ngrams(open(i,"r").read(),n)
158     Web_Shell_vector.append(output_temp)
159 for i in Normal:
160     output_temp=ngrams(open(i,"r").read(),n)
161     Normal_vector.append(output_temp)
162
163 temp=[]
164 for k in Adduser_vector: # seven times only                #Loop 2
165     for i in result:
166         flag=k.get(i,0) #O(1)time by get() function
167         t=(i,flag)
168         temp.append(t)
169     Adduser_vector_final.append(temp)
170     temp=[]
171
172 for k in Hydra_FTP_vector: # seven times only
173     for i in result:
174         flag=k.get(i,0) #O(1)time by get() function
175         t=(i,flag)
176         temp.append(t)

```

```

177     Hydra_FTP_vector_final.append(temp)
178     temp=[]
179
180     for k in Hydra_SSH_vector: # seven times only
181         for i in result:
182             flag=k.get(i,0) #O(1)time by get() function
183             t=(i, flag)
184             temp.append(t)
185         Hydra_SSH_vector_final.append(temp)
186         temp=[]
187
188     for k in Java_Meterpreter_vector: # seven times only
189         for i in result:
190             flag=k.get(i,0) #O(1)time by get() function
191             t=(i, flag)
192             temp.append(t)
193         Java_Meterpreter_vector_final.append(temp)
194         temp=[]
195
196     for k in Meterpreter_vector: # seven times only
197         for i in result:
198             flag=k.get(i,0) #O(1)time by get() function
199             t=(i, flag)
200             temp.append(t)
201         Meterpreter_vector_final.append(temp)
202         temp=[]
203
204     for k in Web_Shell_vector:
205         for i in result:
206             flag=k.get(i,0) #O(1)time by get() function
207             t=(i, flag)
208             temp.append(t)
209         Web_Shell_vector_final.append(temp)
210         temp=[]
211
212     for k in Normal_vector:
213         for i in k:
214             flag=k.get(i,0) #O(1)time by get() function
215             t=(i, flag)
216             temp.append(t)
217         Normal_vector_final.append(temp)
218         temp=[]
219
220     feature_vector=[]
221     feature_vector.append(Adduser_vector_final)
222     feature_vector.append(Hydra_FTP_vector_final)
223     feature_vector.append(Hydra_SSH_vector_final)
224     feature_vector.append(Java_Meterpreter_vector_final)
225     feature_vector.append(Meterpreter_vector_final)
226     feature_vector.append(Web_Shell_vector_final)
227     feature_vector.append(Normal_vector_final)
228
229     index=0
230     string_header=str(n)+"-grams feature vectors:\n\n"
231     string_vector_final=""
232     for k in feature_vector: #Loop 3
233         for i in k:
234             string2=""
235             count=0
236             for j in i:

```

```

237         string2+=str(j[1])+" "
238         string2+=name[index]+"\\n\\n"
239         string_vector_final+=string2
240         index+=1
241
242     final_vector=string_header+string_vector_final
243     train_dataset.write(final_vector)
244     print("—— %s seconds ——" % (time.time() - start_time))
245     #print(final_string) # apply one tab here to get optiized format but that
                        will disturb formatting
246
247     main()

```


Loops

Loop	Discription	Complexity
Loop 1	all the log files are cocatenated into 'Allfiles' n-grams of Allfiles is sorted and n-grams of top 30% frequency is stored in join (list of Dictionary)	$O(n \log n)$
Loop 2	top 30% of each type of log file is rechecked in all log file for frequency of their occurence	$O(n)$
Loop 3	All of the frequncies are shown with their feature vector	$O(n)$

Functions

Function	Parameters	Output	Complexity
setdefault	key -This is the key to be searched. default -This is the Value to be returned in case key is not found	This method returns the key value available in the dictionary and Default value if given key is not available	$O(n)$
ngrams	input -values from log file n - n gram value	return n-grams in required format	$O(n)$
sortngram	output -unsorted n-gram	sorted n-grams in decreasing order of frequency	$O(n \log n)$
top30percent	x -sorted n-grams in decreasing order of frequency	30% values from start	$O(n)$

Complexity of algorithm :- $O(n \log n)$

1. ngrams(input,n)
 - Parameter- n -n gram value
 - input- values from log file
 - output- return n-grams in required format
 - complexity- $O(n)$
2. sort_ngrams(output)
 - Parameter- output -unsorted n gram
 - output- sorted n-grams in decreasing order of frequency
 - complexity- $O(n \log n)$
3. top30percent(x)
 - Parameter- x -sorted n-gram in decreasing order of frequency
 - output- 30% values from start
 - complexity- $O(n)$

Time taken for computing n-grams

