

Contents

1	Introduction of Tools	1
1.1	ELK Stack	1
1.1.1	Elastic Search	1
1.1.2	Logstash	1
1.1.3	Kibana	1
1.2	Honeypots	1
1.2.1	Low-interaction honeypots	2
1.2.2	Medium-interaction honeypots	2
1.2.3	High-interaction honeypots	2
1.2.4	Cowrie Honeypot	2
1.3	VirusTotal	3
1.3.1	List of Anti-Virus Engines used by VirusTotal	3
1.4	Github	3
2	Installation of Tools	4
2.1	Elastic Search	4
2.2	Logstash	5
2.3	Kibana	5
2.4	Cowrie Honeypot	6
2.4.1	Install dependencies	6
2.4.2	Create a user account	6
2.4.3	Checkout the code	6
2.4.4	Setup Virtual Environment	6
2.4.5	Port redirection	6
2.4.6	Starting Cowrie	6
3	Working Model	7
4	Monitoring System	8
4.1	Stochastic Process	13
5	Intrusion Detection System	14
5.1	Malicious-URLs Table	16
5.2	Stats Table	20

List of Figures

1	Elastic Search Output	4
2	Kibana Dashboard	5
3	Working Model	7
4	Monitoring System	8
5	Daily Attacks Count Area Graph	11
6	Top 20 Usernames Histogram	11
7	Unique Count of IP's	12

8	Top 20 Passwords Pie	12
9	Global Attacks Map	13
10	Stochastic Process	14
11	Intrusion Detection System	15

List of Tables

1	List of Antivirus Engines	3
2	List of Tables in PostgreSQL Database	15
3	Schema of Malicious-URLs Table	16
4	Schema of Stats Table	20
5	DOS Attacks Clustering Table	27

Listings

1	Cowrie Config File	8
2	Get Malicious URLs From Virus Total	16
3	Add URL's to PostgreSQL Database	17
4	Submit URL's for Scan to VirusTotal	17
5	Get URL Scan results from VirusTotal	18
6	Virus Total API Request Sample	19
7	Get Unique IPs from Elastic Search	20
8	Elastic Search API Request Sample	21
9	Get Commands Executed by IPs from Elastic Search	21
10	Get Login Attempts by IPs from Elastic Search	23
11	Generate Login Attempts in a CSV File	24
12	Applying Spark K-Means Algorithm	24
13	Counting of DOS attacks	25

1 Introduction of Tools

In this section, i will be introducing the tools ,languages and API's which will be used throughout the report.

1.1 ELK Stack

"ELK" is the acronym for three open source projects: Elasticsearch, Logstash, and Kibana.

1.1.1 Elastic Search

Elasticsearch is a search engine based on Lucene. It provides a distributed, multitenant-capable full-text search engine with an HTTP web interface and schema-free JSON documents.

1.1.2 Logstash

Logstash is an open source tool for collecting, parsing, and storing logs for future use.

1.1.3 Kibana

Kibana is an open source data visualization plugin for Elasticsearch. It provides visualization capabilities on top of the content indexed on an Elasticsearch cluster. Users can create bar, line and scatter plots, or pie charts and maps on top of large volumes of data.

1.2 Honeypots

In computer terminology, a honeypot is a trap set to detect, deflect or in some manner counteract attempts at unauthorized use of information systems.

In other words, A server that is configured to detect an intruder by mirroring a real production system. It appears as an ordinary server doing work, but all the data and transactions are phony. Located either in or outside the firewall, these are used to learn about an intruder's techniques as well as determine vulnerabilities in the real system."

Based on design criteria, honeypots can be classified as

1. Low-interaction honeypots
2. Medium-interaction honeypots
3. High-interaction honeypots

1.2.1 Low-interaction honeypots

Low-interaction honeypots simulate only the services frequently requested by attackers. Since they consume relatively few resources, multiple virtual machines can easily be hosted on one physical system, the virtual systems have a short response time, and less code is required, reducing the complexity of the security of the virtual systems.

Low-interaction honeypots present the hacker emulated services with a limited subset of the functionality they would expect from a server, with the intent of detecting sources of unauthorized activity. For example, the HTTP service on low-interaction honeypots would only support the commands needed to identify that a known exploit is being attempted.

1.2.2 Medium-interaction honeypots

Medium-interaction honeypots might more fully implement the HTTP protocol to emulate a well-known vendor's implementation, such as Apache. However, there are no implementations of a medium-interaction honeypots and for the purposes of this paper, the definition of low-interaction honeypots captures the functionality of medium-interaction honeypots in that they only provide partial implementation of services and do not allow typical, full interaction with the system as high-interaction honeypots.

1.2.3 High-interaction honeypots

High-interaction honeypots imitate the activities of the real systems that host a variety of services. It let the hacker interact with the system as they would any regular operating system, with the goal of capturing the maximum amount of information on the attacker's techniques. Any command or application an end-user would expect to be installed is available and generally, there is little to no restriction placed on what the hacker can do once he/she comprises the system.

According to recent researches in high interaction honeypot technology, by employing virtual machines, multiple honeypots can be hosted on a single physical machine. Therefore, even if the honeypot is compromised, it can be restored more quickly. Although high interaction honeypots provide more security by being difficult to detect, but it has the main drawback that it is costly to maintain.

1.2.4 Cowrie Honeypot

Cowrie is used for our research purposes. Cowrie is a medium interaction SSH and Telnet honeypot designed to log brute force attacks and the shell interaction performed by the attacker.

Interesting Features of Cowrie Honeypot

1. Fake filesystem with the ability to add/remove files. A full fake filesystem resembling a Debian 5.0 installation is included.
2. Possibility of adding fake file contents so the attacker can cat files such as /etc/passwd. Only minimal file contents are included.
3. Logging in JSON format for easy processing in log management solutions.

1.3 VirusTotal

VirusTotal is a website created by the Spanish security company Hispasec Sistemas. Launched in June 2004, it was acquired by Google Inc. in September 2012. The company's ownership switched in January 2018 to Chronicle, a subsidiary of Alphabet Inc. (Google's parent company).

VirusTotal aggregates many antivirus products and online scan engines to check for viruses that the user's own antivirus may have missed, or to verify against any false positives. Files up to 256 MB can be uploaded to the website or sent via email. Anti-virus software vendors can receive copies of files that were flagged by other scans but passed by their own engine, to help improve their software and, by extension, VirusTotal's own capability.

1.3.1 List of Anti-Virus Engines used by VirusTotal

AegisLab	Agnitum	AhnLab	Anity
Aladdin	Avast	AVG	Avira
BluePex	Baidu	BitDefender	Bkav
ByteHero	Quick Heal	CMC Antivirus	CYREN
ClamAV	Comodo	CrowdStrike	Doctor Web Ltd.
Emsisoft	Endgame	Eset Software	Fortinet
F-Prot	F-Secure	G Data	Hacksoft
Hauri	IKARUS	nProtect	Invincea
Jiangmin	K7AntiVirus	Kaspersky	Kingsoft
Malwarebytes	McAfee	Microsoft	eScan
Nano Security	Norman	Panda	Rising
Symantec	VIPRE	TotalDefense	TrendMicro

Table 1: List of Antivirus Engines

1.4 Github

GitHub (originally known as Logical Awesome LLC) is a web-based hosting service for version control using git. It is mostly used for computer code. It

offers all of the distributed version control and source code management (SCM) functionality of Git as well as adding its own features. It provides access control and several collaboration features such as bug tracking, feature requests, task management, and wikis for every project.

The code which is written for this project is entirely pushed to Github where it is bug free and version safe.

2 Installation of Tools

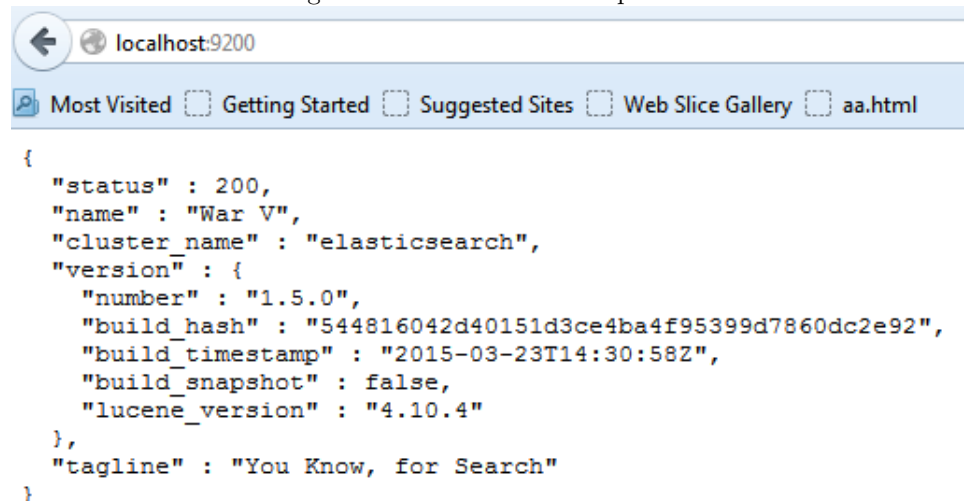
In this section, we will be mentioning about the installation of various tools which will be used in the project.

2.1 Elastic Search

1. Download the Debian Installation file from the below link
2. Install it using the below command
iotsys3@iotsys3-Precision-Tower-3420: dpkg -i elastic-search-6.2.0.deb
3. Now start the service using below command
iotsys3@iotsys3-Precision-Tower-3420: sudo service elasticsearch start

Open the link <http://localhost:9200> in the browser and check for the output which should look like below.

Figure 1: Elastic Search Output



We can check for the working of elastic search in the browser by typing the url `http://localhost:9200` and verifying the output.

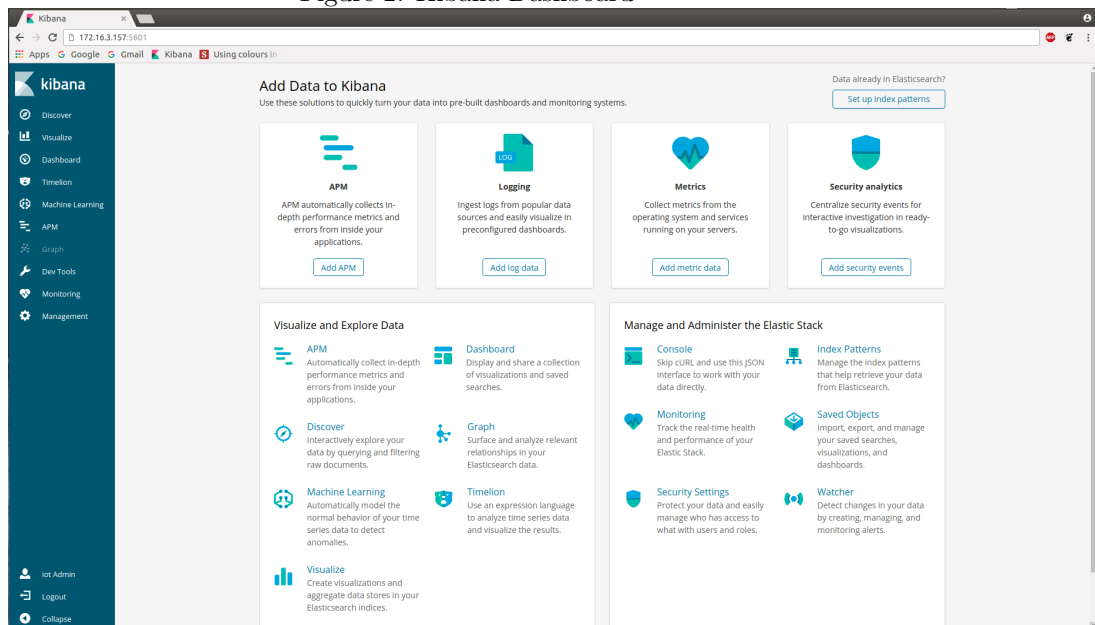
2.2 Logstash

1. Download the logstash debian installation file from below link
`https://www.elastic.co/downloads/logstash`
2. Install it using the below command
`iotsys3@iotsys3-Precision-Tower-3420: dpkg -i logstash-6.2.2.deb`
3. Now start the service using below command
`iotsys3@iotsys3-Precision-Tower-3420: sudo service logstash start`

2.3 Kibana

1. Download the Kibana debian installation file from below link
`https://www.elastic.co/downloads/kibana`
2. Install it using the below command
`iotsys3@iotsys3-Precision-Tower-3420: dpkg -i kibana-6.2.2.deb`
3. Now start the service using below command
`iotsys3@iotsys3-Precision-Tower-3420: sudo service kibana start`

Figure 2: Kibana Dashboard



Open the link (<http://localhost:5601>) in the browser and check for the Kibana Dashboard which looks like above.

2.4 Cowrie Honeypot

Cowrie honeypot is installed in a separate machine with public IP which is connected to the internet. It can be installed by following the steps given below.

2.4.1 Install dependencies

```
iotsys3@iotsys3-Precision-Tower-3420: sudo apt-get install git python-virtualenv  
libssl-dev libffi-dev build-essential libpython-dev python2.7-minimal authbind
```

2.4.2 Create a user account

```
iotsys3@iotsys3-Precision-Tower-3420: sudo adduser --disabled-password cowrie  
iotsys3@iotsys3-Precision-Tower-3420: sudo su - cowrie
```

2.4.3 Checkout the code

```
iotsys3@iotsys3-Precision-Tower-3420: git clone http://github.com/micheloosterhof/cowrie  
iotsys3@iotsys3-Precision-Tower-3420: cd cowrie
```

2.4.4 Setup Virtual Environment

```
iotsys3@iotsys3-Precision-Tower-3420: pwd  
iotsys3@iotsys3-Precision-Tower-3420: virtualenv cowrie-env
```

2.4.5 Port redirection

```
iotsys3@iotsys3-Precision-Tower-3420: sudo iptables -t nat -A PREROUTING  
-p tcp -dport 22 -j REDIRECT --to-port 2222
```

In the same folder, create a copy of *cowrie.cfg.dist* and modify the file according to your needs like changing the name, SSH port number etc.

2.4.6 Starting Cowrie

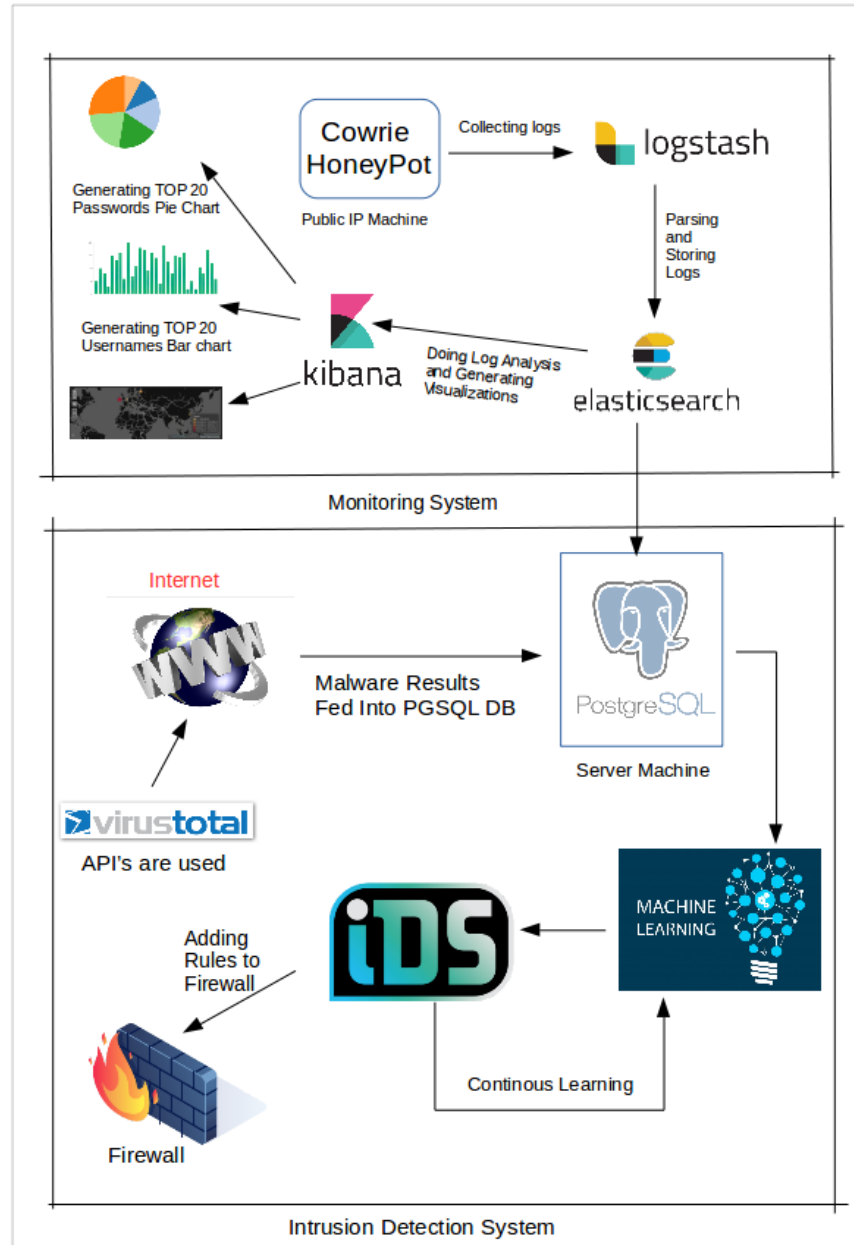
```
iotsys3@iotsys3-Precision-Tower-3420: bin/cowrie start
```

Cowrie will be up and running on port 22 which will monitor all SSH logins and commands executed by the attackers or intruders.

3 Working Model

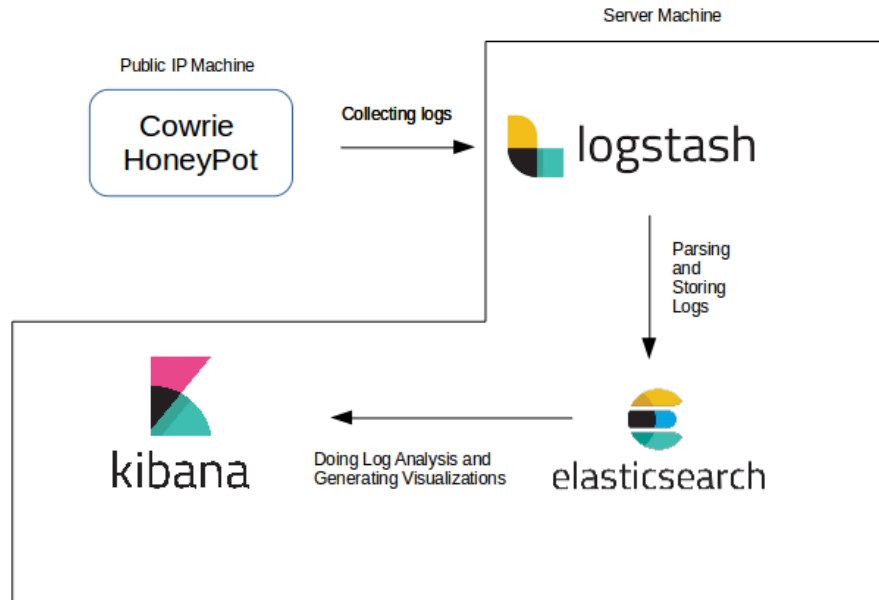
All the tools listed above are integrated to form a complete Working Model which looks like below.

Figure 3: Working Model



4 Monitoring System

Figure 4: Monitoring System



Cowrie will be generating the logs in the form of JSON format which will be available in public IP machine as specified in the figure-3. In order to stream the logs from the public IP machine to Server machine, filebeat can be used which is an open source tool. Unfortunately we did not get to use that tool as there was no physical connection between public IP machine and server machine due to security reasons.

So all the logs are copied from public IP machine to server machine on daily basis. These log files are supplied as input to Logstash which will parse, add metadata and push them to Elastic Search database. Logstash works on config files. Each config file has three parameters namely input, filter and output. Our config file is given below.

Listing 1: Cowrie Config File

```
1 input {
2   file {
3     path => "/home/iotsys3/Documents/Cowrie_Logs/*"
4     start_position => "beginning"
5     ignore_older => 0
6   }
}
```

```

7 }
8 filter {
9
10   json {
11     source => message
12   }
13
14   date {
15     match => [ "timestamp", "ISO8601" ]
16   }
17
18   if [src_ip] {
19
20     dns {
21       reverse => [ "src_host", "src_ip" ]
22       action => "append"
23     }
24
25     geoip {
26       source => "src_ip" # With the src_ip field
27       target => "geoip" # Add the geoip one
28       # Using the database we previously saved
29       database => "/home/iotsys3/Downloads/GeoLite2-
City_20180206/GeoLite2-City.mmdb"
30       add_field => [ "[geoip][coordinates]", "%{[geoip][
longitude]}" ]
31       add_field => [ "[geoip][coordinates]", "%{[geoip][
latitude]}" ]
32     }
33
34     # Get the ASN code as well
35     geoip {
36       source => "src_ip"
37       database => "/home/iotsys3/Downloads/GeoLite2-
ASN_20180206/GeoLite2-ASN.mmdb"
38     }
39
40     mutate {
41       convert => [ "[geoip][coordinates]", "float" ]
42     }
43   }
44 }
45 }
46 output {
47
48   elasticsearch {
49     hosts => ["localhost:9200"]
50     sniffing => true
51     manage_template => false
52     index => "logstash-cowrie-%{+YYYY.MM.dd}"
53     user => "elastic"
54     password => "63TxZjWQNkB4tArR7hri"
55   }
56
57   stdout {
58     codec => rubydebug
59   }

```

1. Input

It contains a parameter file which has options path, start-position and ignore-older. Path is where log files are present and start-position indicates from which point Logstash should start reading the logs from.

2. Filter

Filter is where entire processing happens in logstash. It has many options out of which json option is used as the logs are in JSON format. Some more information is added to the already present logs like country , Continent , latitude and longitude information which are specified in geoip option. Finally latitude and longitude are converted into float variables for plotting on map.

3. Output

It specifies options as to where the data should go after the processing is done. There are many options out of which we are using elastic search where data is stored in json format which will be in key value format.

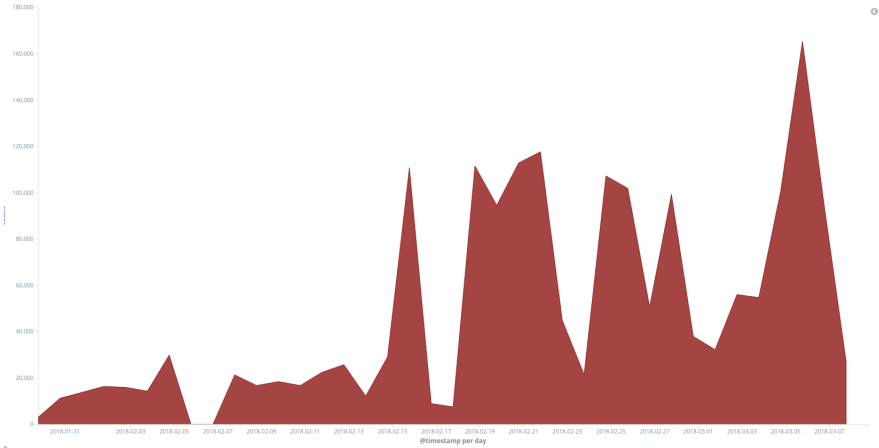
Data collected in the elastic search will be picked up by the Kibana by default as specified in the Kibana config files. Using the data visualizations are created for analyzing and monitoring the regular changes in data. Monitoring system was deployed for 40 days and collected nearly **2 Million Records**.

Using Kibana different visualizations are created which are listed below.

1. Daily Attacks Count Area Graph

An area graph was plotted against the count of attacks on daily basis. X-axis represents the date and Y-axis represents the count of attacks. It can be observed that, attacks kept on increasing at the end.

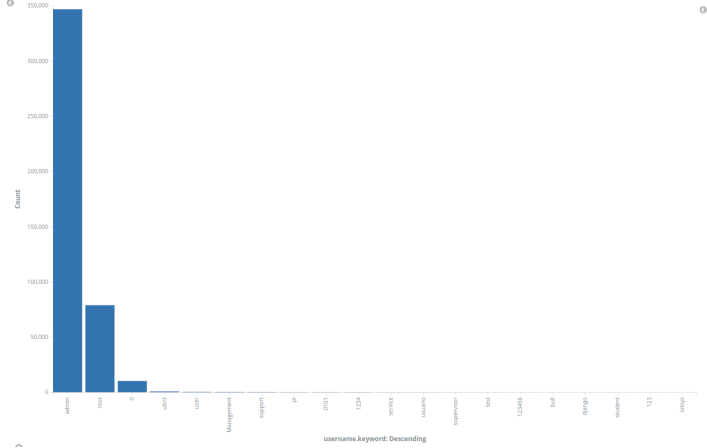
Figure 5: Daily Attacks Count Area Graph



2. Top 20 Usernames Histogram

A Histogram was used to represent the top 20 usernames that have been used by the attackers and intruders.

Figure 6: Top 20 Usernames Histogram

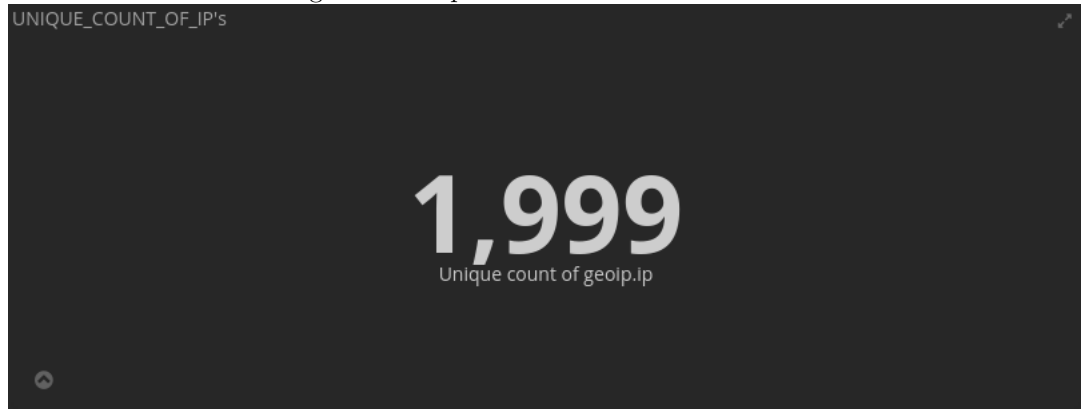


From the graph, it can be observed that usernames "admin", "root" and "ubnt" are mostly used by attackers.

3. Unique Count Of IP's

A Kibana Utility graph named Unique Count has been used for counting the Unique Number of IP's till now monitoring system has captured. It has found to be nearly 2000 IP's.

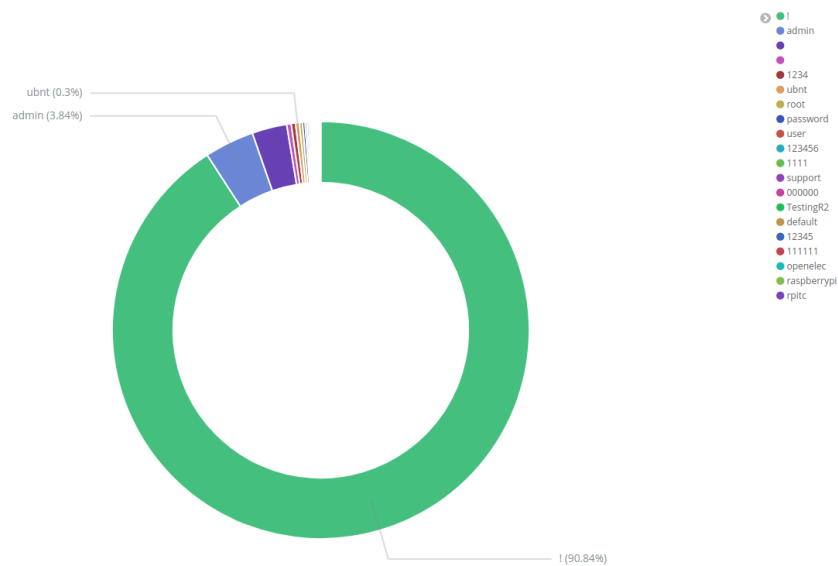
Figure 7: Unique Count of IP's



4. Top 20 Passwords Pie

A pie chart has been drawn against the top 20 passwords used by the attackers and intruders. It can be observed that passwords "ubnt" and "admin" are mostly used.

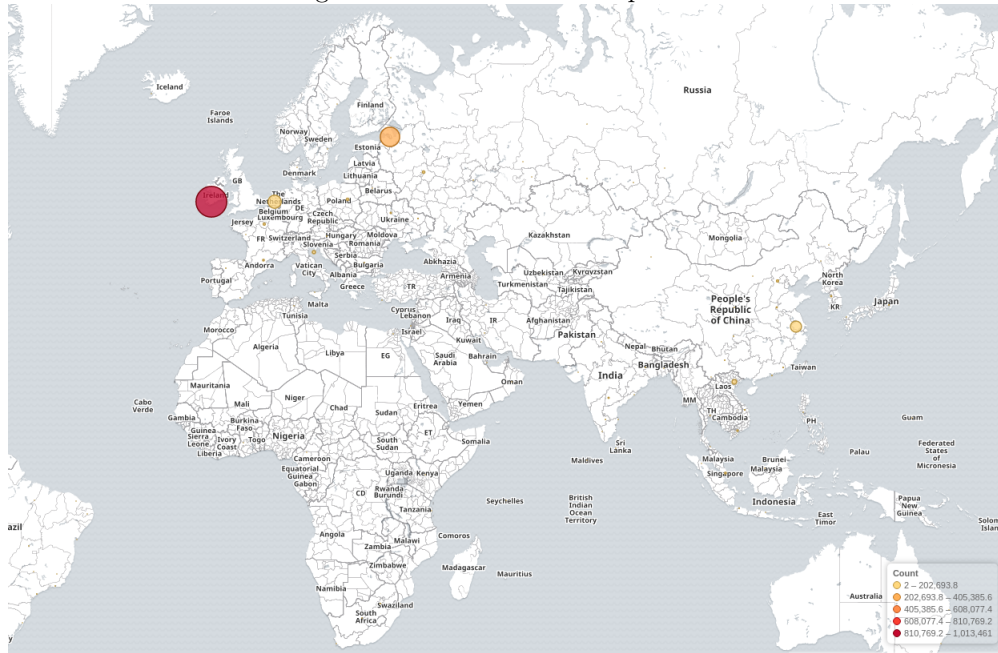
Figure 8: Top 20 Passwords Pie



5. Global Attacks Map

A map is also plotted against the attacks receiving from different countries and continents using the latitude and longitude values. Most of the attacks are received from European continent and China country.

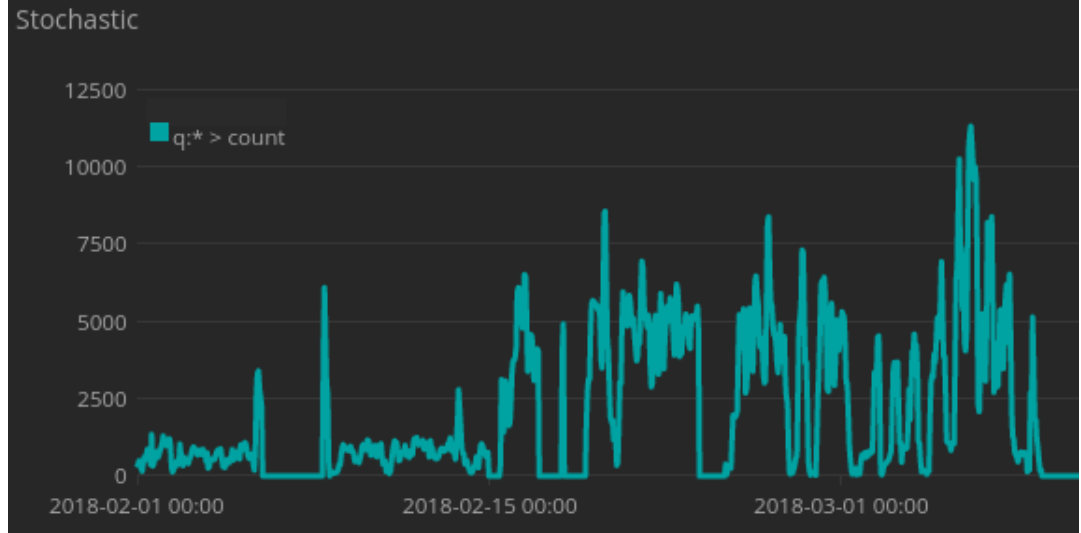
Figure 9: Global Attacks Map



4.1 Stochastic Process

Stochastic process is also called random process. It contains random variables with probabilities associated to it. While exact path cannot be inferred, assumptions can be made based on the probabilities of random variables. Attacks pattern observed is following random process or stochastic process as number of attacks on daily basis is not following any pattern.

Figure 10: Stochastic Process



5 Intrusion Detection System

An Intrusion Detection System is a software application which regularly monitors the network for malicious activities and reports them to central authority for fail-safe measures. After the monitoring system is setup and data is collected, task is now to design an Intelligent Intrusion Detection System which applies machine learning algorithms on the data and comes up with firewall rules to block malicious activities and IP's. For designing, tools like VirusTotal API's, PostgreSQL database, Spark K-Means Algorithm etc. are used which is given below in detail.

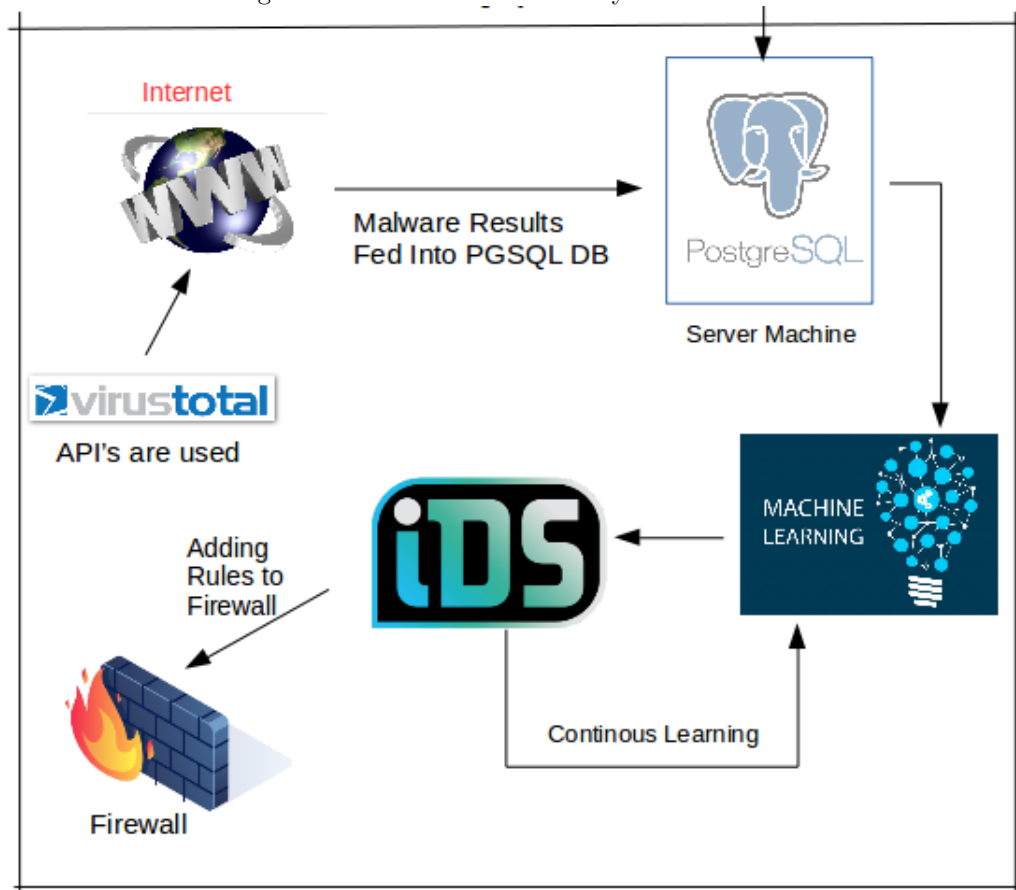
Virus Total is a website where details about different malwares and IP's can be fetched by using their web interface or REST API's. For our convenience, API's are used to integrate with our python programs. Virus Total gives a normal user on an average of 1000 API requests per day combining malware samples and IP's. We have sent a mail to VirusTotal requesting Research User privileges which has access to 10,000 API requests per day and all malware samples present in their database and got their approval.

In PostgreSQL, database named "cowrie" has been created and below tables were created to store the required data from Elastic Search and VirusTotal.

S.No.	Table Name	Schema
1	Dos-Attacks	public
2	Malicious-URLs	public
3	Stats	public

Table 2: List of Tables in PostgreSQL Database

Figure 11: Intrusion Detection System



5.1 Malicious-URLs Table

S.No.	Column Name	Column Datatype
1	url	text
2	scan-url	text
3	positives	integer
3	total	integer
3	identified-engines	text

Table 3: Schema of Malicious-URLs Table

Column "url" is used to store the malicious url and "scan-url" is the ID to identify the job submitted to VirusTotal for malicious url. "positives" contains the number of antivirus engines identified the url as malicious and "total" indicated the number of antivirus engines present in VirusTotal. "identified-engines" contains the names of all anti-virus engines. Malicious-urls table gets the data from VirusTotal API's. Python codes and their explanations are given below.

Listing 2: Get Malicious URLs From Virus Total

```
1 import requests
2 import json
3 import psycopg2
4
5 apikey = "
6     f0aa0772e518487aa5daf1dcf2be2d37c9b9b361b1a0d8339faf972210b36d41
7     "
8
9 scan_url = 'https://www.virustotal.com/vtapi/v2/ip-address/report?
10     apikey='+apikey+'&ip='
11 uniqueURLs = set()
12
13 dbName = "cowrie"
14 offset = 1813
15
16 try:
17     conn = psycopg2.connect("dbname="+str(dbName)+" user='
18     postgres' host='localhost' password='postgres'")
19 except Exception as e:
20     print e
21 cur = conn.cursor()
22
23 #cur.execute("CREATE TABLE IF NOT EXISTS MALICIOUS_URLS(URL TEXT)")
24 #cur.execute("TRUNCATE TABLE MALICIOUS_URLS")
25 #conn.commit()
26
27 cur.execute("SELECT IP FROM STATS OFFSET "+str(offset));
28 rows = cur.fetchall()
29 i = offset
30 for row in rows:
31     print "*****"
32     print i
33     i += 1
34     print row[0]
```

```

29     r = requests.get(str(scan_url)+str(row[0]))
30     resp = json.loads(r.content)
31     #print resp
32     if "detected_urls" in resp:
33         if resp["detected_urls"] != []:
34             #print resp["detected_urls"]
35             for url in resp["detected_urls"]:
36                 print url["url"]
37                 cur.execute("INSERT INTO MALICIOUS.URLS(URL) VALUES
('\" + url[\"url\"] + '\"')")
38                 conn.commit()

```

Listing 3: Add URL's to PostgreSQL Database

```

1 import psycopg2
2
3 dbName = "cowrie"
4
5 try:
6     conn = psycopg2.connect("dbname="+str(dbName)+" user='
postgres' host='localhost' password='postgres'")
7 except Exception as e:
8     print e
9 cur = conn.cursor()
10
11 with open("urls.txt") as fp:
12     url = fp.readline()
13     while url:
14         print url
15         cur.execute("INSERT INTO MALICIOUS.URLS(URL) VALUES('\" + str
(url) + '\"')")
16         conn.commit()
17         url = fp.readline()

```

Listing 4: Submit URL's for Scan to VirusTotal

```

1 import json
2 import psycopg2
3 from subprocess import check_output
4
5 dbName = "cowrie"
6 apikey = "
f0aa0772e518487aa5daf1dcf2be2d37c9b9b361b1a0d8339faf972210b36d41
"
7 scan_url = 'https://www.virustotal.com/vtapi/v2/url/scan '
8
9 offset = 805
10
11 try:
12     conn = psycopg2.connect("dbname="+str(dbName)+" user='
postgres' host='localhost' password='postgres'")
13 except Exception as e:
14     print e
15 cur = conn.cursor()
16
17
18 cur.execute("SELECT URL FROM MALICIOUS.URLS ORDER BY URL OFFSET "+
str(offset))

```

```

19 urls = cur.fetchall()
20 i=offset
21 for url in urls:
22     print "*****"
23     print i
24     i += 1
25     res = json.loads(check_output(["curl", "--request", "POST",
26                                   "--url", "https://www.virustotal.
27                                   com/vtapi/v2/url/scan",
28                                   "--data", "apikey=" + apikey,
29                                   "--data", "url=" + url[0]]))
30     print res['scan_id']
31     cur.execute("UPDATE MALICIOUS_URLS SET SCAN_URL='"+str(res['
scan_id'])+"' WHERE URL='"+str(url[0])+"'")
32     conn.commit()

```

Listing 5: Get URL Scan results from VirusTotal

```

1 import requests
2 import psycpg2
3 from collections import Counter
4 import json
5
6 dbName = "cowrie"
7 apikey = "
f0aa0772e518487aa5daf1dcf2be2d37c9b9b361b1a0d8339faf972210b36d41
"
8 offset = 3
9
10
11 def getScanResults(scan_id):
12     print scan_id
13     softwares = ""
14     report_url = 'https://www.virustotal.com/vtapi/v2/url/report?'
15     \
16         'apikey=' + str(apikey) + '&' \
17         'resource=' + str(scan_id) + '&' \
18         'allinfo=false'
19
20     response = requests.get(report_url)
21     res = response.json()
22     #print json.dumps(res['scans'], indent=4, sort_keys=True)
23     for r in res['scans']:
24         if res['scans'][r]['detected'] == True:
25             softwares += str(r) + ";"
26     print softwares
27     softwares = softwares[0:len(softwares)-1]
28     positives = res['positives']
29     total = res['total']
30     cur.execute("UPDATE MALICIOUS_URLS SET "
31               "POSITIVES="+str(positives)+" , "
32               "TOTAL="+str(total)+" , "
33               "IDENTIFIED_ENGINES="+str(softwares)+" WHERE
SCAN_URL='"+scan_id+"'")
34     conn.commit()
35
36

```

```

37 try:
38     conn = psycopg2.connect("dbname='"+str(dbName)+"' user='
        postgres' host='localhost' password='postgres'")
39 except Exception as e:
40     print e
41 cur = conn.cursor()
42
43 cur.execute("SELECT SCAN.URL FROM MALICIOUS.URLS OFFSET "+str(
        offset))
44 scan_ids = cur.fetchall()
45 i = offset
46 for id in scan_ids:
47     print "
        *****"
48     print i
49     i = i+1
50     if(id[0]!=None):
51         getScanResults(id[0])

```

Above python codes uses Virus Total API and sends requests for all the malicious urls identified for the specific IP. The response is decoded and then stored into the table. response looks like below.

Listing 6: Virus Total API Request Sample

```

1 {
2   'response_code': 1,
3   'verbose_msg': 'Scan finished , scan information embedded in this
        object',
4   'scan_id': '1
        db0ad7dbcec0676710ea0eaacd35d5e471d3e11944d53bcbd31f0cbd11bce31
        -1390467782',
5   'permalink': 'https://www.virustotal.com/url/_urlsha256_/_
        analysis/1390467782/',
6   'url': 'http://www.virustotal.com/',
7   'scan_date': '2014-01-23 09:03:02',
8   'filescan_id': null,
9   'positives': 0,
10  'total': 51,
11  'scans': {
12    'CLEAN MX': {
13      'detected': false,
14      'result': 'clean site'
15    },
16    'MalwarePatrol': {
17      'detected': false,
18      'result': 'clean site'
19    }
20  }
21 }

```

This response is parsed and requires fields are extracted to update the table in the database. These fields include "positives" which is the no of engines identified it as positive, "total" which is the count of the engines in total and "identified-engines" which contains the names of all the engines in the response.

5.2 Stats Table

S.No.	Column Name	Column Datatype
1	ip	character varying(20)
2	commands	text
3	countofcommands	integer
4	loginattempts	integer
5	doscluster	text
6	sentiment	integer

Table 4: Schema of Stats Table

Column "IP" contains the ip of the attacker or intruder and "commands" contains all the commands executed by the attacker or the intruder till now in cowrie honeypot. "Count-of-commands" contains the number of commands attacker executed. "login-attempts" contains the number of times attacker tried to login into the cowrie. "dos-cluster" contains the cluster number to which ip is clustered which will be updated dynamically. "sentiment" contains the value of ip which is updated after running machine learning algorithms.

All the python codes used for filling the table are given below one by one.

Listing 7: Get Unique IPs from Elastic Search

```
1 import json
2 import requests
3 import psycopg2
4 from Globals import *
5
6
7 uniqueIPs = set()
8 dbName = "cowrie"
9
10 try:
11     conn = psycopg2.connect("dbname='"+str(dbName)+"' user='
12     postgres' host='localhost' password='postgres'")
13 except Exception as e:
14     print e
15 cur = conn.cursor()
16
17 def createTables():
18     cur.execute("CREATE TABLE IF NOT EXISTS STATS("
19                 "IP VARCHAR(20) PRIMARY KEY NOT NULL,"
20                 "LOGINATTEMPTS INT,"
21                 "COUNTOFCOMMANDS INT,"
22                 "DOSCLUSTER INT,"
23                 "SENTIMENT INT,"
24                 "COMMANDS TEXT)")
25
26     cur.execute("CREATE TABLE IF NOT EXISTS DOS_ATTACKS("
27                 "IP VARCHAR(20) PRIMARY KEY NOT NULL,"
28                 "LOGINATTEMPTS INT,"
```

```

28         "COUNTOFCOMMANDS INT,"
29         "COMMANDS TEXT")
30
31     conn.commit()
32     cur.execute("TRUNCATE TABLE STATS")
33     conn.commit()
34
35 createTables()
36
37
38 url = ELASTIC_URL + str("?size=1")
39 response = requests.get(url)
40 countOfData = json.loads(response.content)["hits"]["total"]
41 print countOfData
42
43 size = 10000
44
45 for i in range(0, countOfData+1, size):
46     url = ELASTIC_URL + str("?_source=geoip.ip&size=")+str(size)+"&
47     from="+str(i)
48     response = requests.get(url)
49     IPs = json.loads(response.content)["hits"]["hits"]
50     print len(IPs)
51     for IP in IPs:
52         if IP["_source"] != {}:
53             ip = IP["_source"]["geoip"]["ip"]
54             uniqueIPs.add(ip)
55             print "Completed : from = "+str(i)+" and end = "+str(i+size)
56
57 for ip in uniqueIPs:
58     cur.execute("INSERT INTO STATS(IP) VALUES('"+ip+"')")
59     conn.commit()
60
61 print "Total Number of Unique IP's = " + str(len(uniqueIPs))

```

The above python code fetches all the IP's from elastic search, puts them in a set to avoid duplicates and copies all the IP's to the column named "ip" in Stats table. It uses the Elastic Search API's which looks like below.

Listing 8: Elastic Search API Request Sample

```

1 POST localhost:9200/logstash-*/_search
2 {
3     "size" : 0,
4     "aggs" : {
5         "distinct_ip" : {
6             "cardinality" : {
7                 "field" : "geoip.ip"
8             }
9         }
10    }
11 }
12 }

```

Listing 9: Get Commands Executed by IPs from Elastic Search

```

1 import json

```

```

2 import requests
3 import psycopg2
4 from Globals import *
5
6
7 dbName = "cowrie"
8 IPsArray = []
9
10 try:
11     conn = psycopg2.connect("dbname='"+str(dbName)+"' user='
12     postgres' host='localhost' password='postgres'")
13 except Exception as e:
14     print e
15 cur = conn.cursor()
16
17 def getNoOfUniqueIPsFromDB():
18     cur.execute("SELECT COUNT(*) FROM STATS");
19     rows = cur.fetchall()
20     for row in rows:
21         return row[0]
22
23 def getIPsFromDB():
24     cur.execute("SELECT IP FROM STATS");
25     rows = cur.fetchall()
26     for row in rows:
27         IPsArray.append(row[0])
28
29 print getNoOfUniqueIPsFromDB()
30 getIPsFromDB()
31
32 for i in range(0, getNoOfUniqueIPsFromDB()):
33     tillNowCommandsExecuted = ""
34     data = {
35         "query": {
36             "query_string": {
37                 "fields": ["src_ip", "eventid"],
38                 "query": ""+str(IPsArray[i])+" AND cowrie.command.
39             input"
40             },
41             "_source": ["geoip.ip", "eventid", "message"]
42         }
43     }
44     headers = {"Content-Type": "application/json"}
45     url = ELASTIC_URL + str("?size=10000")
46     response = requests.post(url, data=json.dumps(data), headers=
47     headers)
48     countOfCommands = json.loads(response.content)["hits"]["total"]
49     cur.execute("UPDATE STATS SET COUNTOFCOMMANDS=" + str(
50     countOfCommands) + "WHERE IP=" + IPsArray[i] + " ;")
51     conn.commit()
52     commands = json.loads(response.content)["hits"]["hits"]
53     if commands != []:
54         for command in commands:
55             tillNowCommandsExecuted += str(command["_source"] + "
56             message"].replace("CMD:", "")) + " "
57     print tillNowCommandsExecuted

```



```

54 cur.execute("UPDATE STATS SET COMMANDS=" + str(
    tillNowCommandsExecuted.replace(" ", ";").strip()) + " ' WHERE IP
55 =" + IPsArray[i] + "';")
    conn.commit()

```

The above python code fetches all the commands executed per IP , appends them together separated by semicolon and puts them in the column named "commands" in Stats table. It also updates the "countofcommands" column.

Listing 10: Get Login Attempts by IPs from Elastic Search

```

1 import json
2 import requests
3 import psycopg2
4 from Globals import *
5
6
7 dbName = "cowrie"
8 IPsArray = []
9
10 try:
11     conn = psycopg2.connect("dbname="+str(dbName)+" ' user='
        postgres ' host='localhost ' password='postgres '")
12 except Exception as e:
13     print e
14 cur = conn.cursor()
15
16 def getNoOfUniqueIPsFromDB():
17     cur.execute("SELECT COUNT(*) FROM STATS");
18     rows = cur.fetchall()
19     for row in rows:
20         return row[0]
21
22 def getIPsFromDB():
23     cur.execute("SELECT IP FROM STATS");
24     rows = cur.fetchall()
25     for row in rows:
26         IPsArray.append(row[0])
27
28
29 print getNoOfUniqueIPsFromDB()
30 getIPsFromDB()
31
32 for i in range(0, getNoOfUniqueIPsFromDB()):
33     tillNowCommandsExecuted = ""
34     data = {
35         "query": {
36             "query_string": {
37                 "fields": ["src_ip", "eventid"],
38                 "query": str(IPsArray[i])+" AND cowrie.login.*"
39             }
40         },
41         "_source": ["geoip.ip", "eventid", "message"]
42     }
43     headers = {"Content-Type": "application/json"}
44     url = ELASTIC_URL + str("?size=1")

```

```

45     response = requests.post(url, data=json.dumps(data), headers=
46         headers)
47     noOfLoginAttempts = json.loads(response.content)[ "hits" ][ "total
48         " ]
49     cur.execute("UPDATE STATS SET LOGINATTEMPTS=" + str(
50         noOfLoginAttempts) + "WHERE IP=" + IPsArray[i] + "';")
51     conn.commit()

```

The above python code gets the number of login attempts done by each IP from the elastic search. It counts the attempts and stores them in column named "loginattempts" in Stats table.

5.2.1 DOS Attacks

This "loginattempts" login is used to predict whether IP is performing DOS attack or not. For this, only those IP's are considered where after they log in, they don't execute any commands and just leave. Spark K-Means Machine learning algorithm has been used to divide the IP's into 2 categories "High" and "Low". CSV file is generated which is supplied as input to K-Means algorithm. Pseudo code is given below.

Listing 11: Generate Login Attempts in a CSV File

```

1
2 import psycopg2
3
4 dbName = "cowrie"
5
6
7 try:
8     conn = psycopg2.connect("dbname="+str(dbName)+" ' user='
9         postgres' host='localhost' password='postgres'")
10 except Exception as e:
11     print e
12 cur = conn.cursor()
13
14
15
16 cur.execute("SELECT LOGINATTEMPTS FROM STATS WHERE COUNTOFCOMMANDS
17     = 0 ORDER BY LOGINATTEMPTS DESC")
18 rows = cur.fetchall()
19 print "Login-Attempts"
20 for row in rows:
21     print(row[0])

```

Listing 12: Applying Spark K-Means Algorithm

```

1 from pyspark.ml.linalg import Vectors
2 from pyspark.ml.feature import StandardScaler
3 from pyspark.ml.feature import VectorAssembler
4 from pyspark.ml.clustering import KMeans
5 from pyspark.sql import SparkSession
6

```

```

7
8 spark = SparkSession.builder.appName("TEST").getOrCreate()
9
10 Create Dataframe From CSV
11
12 cluster_df = spark.read.csv("/home/iotsys3/PycharmProjects/krishna/
    Cowrie-Scripts/loginattempts.csv",
13     header=True,
14     inferSchema=True,
15 )
16
17
18
19 Removing Columns that are not needed from DataFrame
20
21 #cluster_df = cluster_df.drop("id","date")
22
23
24
25 Convert Data to Vector as Standardised Scaler
26 and Kmeans can only operate on that data
27
28 vectorAssembler = VectorAssembler(
29     inputCols=["Login_Attempts"],
30     outputCol="features"
31 )
32
33
34 vcluster_df = vectorAssembler.transform(cluster_df)
35 #vcluster_df.show()
36 test_df = vectorAssembler.transform(cluster_df)
37
38
39 kmeans = KMeans().setK(3) # set number of clusters
40 kmeans.setSeed(1) # set start point
41 kmodel = kmeans.fit(vcluster_df)
42
43 centers = kmodel.clusterCenters()
44
45 pred_df = kmodel.transform(test_df)
46
47 print("\t**OUTPUT**\n\n\n")
48 print(vcluster_df.show())
49 print("Centers RK = " + str(centers))
50 print(kmodel.summary.clusterSizes)
51 print("\n\n\n")
52 pred_df.show(10)
53 print "Center -1 " + str(centers[0])
54 print "Center -2 " + str(centers[2])
55 print "Center -3 " + str(centers[1])

```

Listing 13: Counting of DOS attacks

```

1 import psycopg2
2
3 dbName = "cowrie"
4
5

```

```

6 cluster1center = 25.21769815
7 cluster2center = 11354.0
8 cluster3center = 20776.75
9
10 try:
11     conn = psycopg2.connect("dbname="+str(dbName)+" user='
12     postgres' host='localhost' password='postgres'")
13 except Exception as e:
14     print e
15 cur = conn.cursor()
16
17 def selectCluster(dis1, dis2, dis3):
18     if (dis1 < dis2) and (dis1 < dis3):
19         return 1
20     elif (dis2 < dis1) and (dis2 < dis3):
21         return 2
22     else:
23         return 3
24
25 cur.execute("SELECT IP, LOGINATTEMPTS FROM STATS WHERE
26     COUNTOFCOMMANDS = 0 ORDER BY LOGINATTEMPTS DESC")
27 rows = cur.fetchall()
28 for row in rows:
29     clusterNo = selectCluster(abs(row[1]-cluster1center), abs(row
30     [1]-cluster2center), abs(row[1]-cluster3center))
31     cur.execute("UPDATE STATS SET DOSCLUSTER="+str(clusterNo)+"
32     WHERE IP="+str(row[0])+" ;")
33     conn.commit()
34
35 countOfSevereDosAttacks = 0
36 countOfMediumDosAttacks = 0
37 countOfLowDosAttacks = 0
38
39 cur.execute("SELECT COUNT(IP) FROM STATS WHERE DOSCLUSTER=1")
40 rows = cur.fetchall()
41 for row in rows:
42     countOfLowDosAttacks = row[0]
43     break
44
45 cur.execute("SELECT COUNT(IP) FROM STATS WHERE DOSCLUSTER=2")
46 rows = cur.fetchall()
47 for row in rows:
48     countOfMediumDosAttacks = row[0]
49     break
50
51 cur.execute("SELECT COUNT(IP) FROM STATS WHERE DOSCLUSTER=3")
52 rows = cur.fetchall()
53 for row in rows:
54     countOfSevereDosAttacks = row[0]
55     break
56
57 print "+++++-----+"
58 print "+++++| Severe DOS Attacks = "+str(countOfSevereDosAttacks)
59 print "+++++| Medium DOS Attacks = "+str(countOfMediumDosAttacks)
60 print "+++++| Low DOS Attacks = "+str(countOfLowDosAttacks)
61 print "+++++-----+"

```

After the above python file is executed, DOS attacks table looks like below.

S.No.	DOS Cluster	No. of IP's
1	Severe	16
2	Medium	5
3	Low	1842
	Total	1863

Table 5: DOS Attacks Clustering Table

5.2.2 Modified Naive Bayes With K-Means for Sentiment Analysis

Based on the commands executed by the

Algorithm 2 Modified Naive Bayes For Sentiment Analysis Algorithm

1: **Probability**(*IP is Malicious / Commands executed by IP*) =

$$\frac{\mathbf{Probability}(IP \text{ is Malicious}) * \mathbf{Probability}(\text{Commands executed by IP} / IP \text{ is Malicious})}{\mathbf{Probability}(\text{Commands executed by IP})}$$

2: **Note** :As we are working on Unclassified Data, We are assuming below values.

3: **Probability**(*IP is Malicious*) = 1

4: **Probability**(*Commands executed by IP*) = 1

5: **Probability**(*Commands executed by IP / IP is Malicious*) =

$$\frac{\text{Frequency of Word} + 1 \text{ (Smoothing Factor)}}{\text{Total Number of Words} + \text{Unique Number of Words}}$$

6: **Note** : *Commands executed by each IP is split into words and their frequencies are stored in Database*

7: *Total No of Words* $\leftarrow N$

8: *Unique No of Words* $\leftarrow \text{dict}$

9: **while** *IP List not Empty* **do**

10: *No of Words* \leftarrow *Splitting Commands into words*

11: **Probability**(*IP is Malicious / Commands executed by IP*)

12: = **Probability** (*Commands executed by IP / IP is Malicious*)

13: = $\prod_{i=1}^{\text{No of Words}}$ **Probability**(*Word_i in Command Executed/ IP is Malicious*)

14: = $\prod_{i=1}^{\text{No of Words}} \frac{n_i + 1}{N + \text{dict}}$

15: **done**

16: **Note** : *Each IP Final Probability will be stored in Database*

17: **Note** : *After Storing the IP's Probabilities in Database, K-Means*

Clustering algorithm is applied to classify the IP's into two classes as high and Low.

K-Means is given below

18: **Input** : *Set of Data Points D and No. of Clusters K.*

19: **Output** : *Cluster Centers that minimizes the squared error distortion.*

20: **Algorithm** :

1. *Pick K Data points randomly from D to form cluster centers.*
 2. *Assign each data point to its nearest cluster center by calculating and taking the minimum Euclidean Squared distance metric.*
 3. *After all data points are assigned to their clusters, move each cluster center to mean of its assigned data points.*
 4. *With new cluster centers, repeat 2 to 3 until there is no convergence.*
-