

DDOS DETECTION AND MITIGATION USING MACHINE LEARNING

ARPIT RAMESH GAWANDE

A thesis submitted to the
Graduate School - Camden
Rutgers, The State University of New Jersey

in partial fulfillment of the requirements
for the degree of
Master of Science
Graduate Program in Scientific Computing

Written under the guidance of
Dr. Jean-Camille Birget

May 2018

Contents

1	Existing Systems	3
1.1	What is DDoS attack	3
1.2	Challenges in dealing with DDoS attacks	4
2	DDoS Detection and Mitigation	5
3	Network Functioning	6
4	Our Approach	7
4.1	The router as a point of analysis	7
4.2	Internet packet flow capturing at the router	9
4.3	Analysis techniques to be implemented	10
5	Implementation	11
5.1	Machine Learning	16
5.1.1	Feature Scaling	17
5.1.2	Clustering	17
5.1.3	Anomaly Detection using One Class Support Vector Machine	23
6	Detection	27
7	Conclusion	29

Abstract

Distributed Denial of Service (DDoS) attacks are very common these days [1]. It is evident that the current industry solutions, such as completely relying on Internet Service Provider(ISP) or setting up DDoS defense infrastructure, are not sufficient in detecting and mitigating DDoS attacks, hence consistent research is needed. Most of the current industry solutions involve setting up a centralized expensive hardware system which can analyze the data packets [2] for probable DDoS attacks. Also each router provider has different protocols to communicate between the DDoS attack detection system and the router/networking devices, limiting the reach of DDoS detection systems. In this paper we have discussed a way to detect DDoS attacks using machine learning tools at the routers, instead of setting a centralized analysis system. Also we have proposed a standard communication architecture which can be used across all the networking devices for mitigating DDoS attacks.

1 Existing Systems

1.1 What is DDoS attack

Distributed Denial of Service (DDoS) attack is a way to jam a host network or its resources with a large number of data packets¹ or connections, so that the host becomes disabled. There are different types of DDoS attacks such as : 1. Volume based, e.g. SYN Flood Attacks, in which the victim is flooded with a high volume of packets or connections. 2. Application based, in which an application such as DNS, VOIP or HTTP are attacked. 3. Low rate DDoS attacks, in which the attacker exploits a vulnerability in the application design, e.g. Slowloris. [3]

¹Messages that are sent on the Internet are broken into shorter messages for transmission. These short messages are called packets. Term coined by Donald Watts Davies.

1.2 Challenges in dealing with DDoS attacks

The real challenge in detecting and defending against DDoS attack is because of its dynamic nature. The source² is not a single node or a system on the Internet but there can be many systems participating in DDoS attack, and often these systems are distributed over different regions of the Internet. Also the source of the packet is often spoofed³ [4], which makes harder to know the actual IP address of the system from where attack is originated because original attack source is changed in spoofed data packet. In addition to that, many times the source system itself is not aware that it is compromised and it is being used as a bot [5] by an attacker to launch DDoS attack.

As the source address can't be a reliable way to know the attack source (because of spoofing), detecting and mitigating attack at the destination⁴ is not very useful. Destination may know that the attack is happening but to stop it happening it will have to block all the incoming traffic including the legitimate traffic. To avoid this, many network device producing companies such as Cisco, Netgear have come up with some solutions. Many of the solutions provided by those giants, or the research that is done in this field has been focusing on collecting network traffic flow information [6] at routers(gateways). A flow consist of a number of Internet packets captured during a fixed time interval. Router send that flow information to the central system for the analysis. Central system is a hardware and software infrastructure which is capable of processing and analyzing large flow information.

Some of the major protocols which are widely used for flow collection and analysis are, Internet Protocol Flow Information Export (IPFIX) protocol created by the Internet Engineering Task Force (IETF), Ciscos NetFlow [7] and Sflow(Sampled flow) [8]. These protocols have defined standard way to export the flow information from router and similar devices. All these flow monitoring

²It is a system/device on the Internet that has an IP address and which is involved in DDoS attack

³spoofing is the way to change the source IP address of the message. This is a known issue in the protocol itself not in the implementation

⁴System under DDoS attack

protocols gather information and send the consolidated flow information to the centralized server where user can login and perform functions; such as Security Monitoring, Bandwidth monitoring, Resource Management, Traffic Analysis, Performance Management etc. On such systems, there are some modules which are specifically used for anomaly/DDoS detection.

E.g. Cisco netflow has flow exporter, collector and analysis modules. Flow exporter modules are installed on routers. The routers which are having flow exporter modules, send flow information to the collector module installed on the server. Along with the collector module, server also has analysis module which can be used to detect different patterns in the flow.

These technologies scales well and can be sufficient to indicate trends in network traffic but they have limitations. 1) They are not cross platform, e.g. router with Sflow protocol can't work with Cisco routers. 2) They involve setting up expensive hardware which acts as collector server. 3) Source address is used for flow analysis which is not reliable due to IP spoofing in the case of of DDoS attack.

Now we know that the router based flow analysis can be useful for anomaly detection but it has limitations. We don't want to set up expensive hardware, we want to have a protocol or a system which is compatible with other routers. Also we don't want source IP address for detection analysis. So if we can come up with a way by which we can detect anomaly in the traffic at the intermediate devices on the Internet such as gateway devices(routers) and create a communication protocol between such gateway devices and the destination server or network, then better decisions on regulating the packet flow can be taken.

2 DDoS Detection and Mitigation

DDoS attacks can be detected by checking if there is any anomalous behavior in the network traffic, such as, a sudden increase in the number of packets going to a destination. This can be done at the server by observing all the incoming traffic

or it can be done by observing all the out going/incoming traffic at the ISP's or at every router. Attack can be mitigated if the anomalous packets are blocked from reaching their destination.

3 Network Functioning

A switch creates a network and router connects those networks. A router links computers to the Internet through other routers. Routers are the backbone of the network who helps to forward packet from one point to other point on the Internet. Every packet traveling on the Internet has to go through a router [9]. Router knows where the packet is destined, hence it could serve as first point of knowledge about the change in the packet flow information for a destined network. Each router has interfaces to which hosts or other network are connected. So router is aware to whom it is connected. Router uses protocol to communicate between other networking devices and by that it gathers knowledge about other networks or routers on the Internet. ICMP [10] is one of the most frequently used protocol by routers for communication.

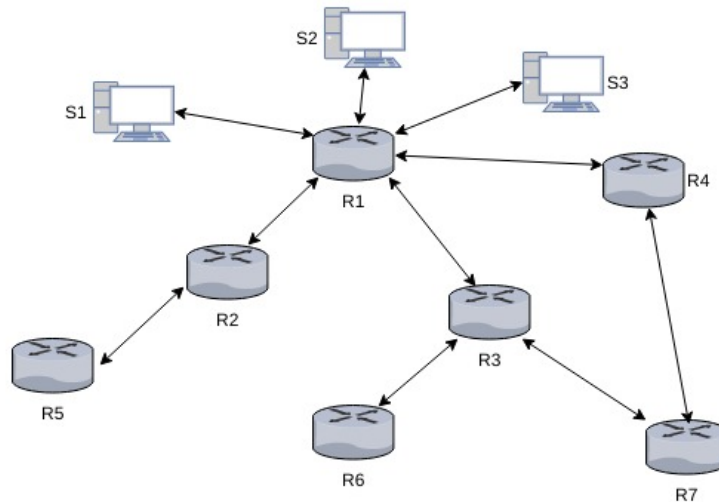


Figure 1: Network Example

Let's illustrate this using an example. In the above figure we can see that

host S1, S2, S3 are connected to a router R1. Router R1 is connected to the Internet through router R2, R3 and R4, thus every packet reaching to the system S2 is coming from either of these three routers. All three routers are located in different geographical region. Most of the websites are regional, either county, state or national (If we leave out few global websites) and hence they are mostly accessed from those region it is meant for. E.g. Rutgers University website is accessed mostly from the eastern region of United States and that too mostly from the New Jersey State or the Philadelphia region.

Using traceroute we can see how many hops⁵ away the destination is. Following is one of the captured traceroute for Rutgers University website.

```
arpit@omega:~$ traceroute camden.rutgers.edu
traceroute to camden.rutgers.edu (128.6.34.90), 30 hops max, 60 byte packets
 1 192.168.0.1 (192.168.0.1) 1.067 ms 1.697 ms 1.684 ms
 2 10.240.177.197 (10.240.177.197) 7.617 ms 9.975 ms 10.302 ms
 3 67.59.225.66 (67.59.225.66) 10.803 ms 12.759 ms 13.074 ms
 4 dstswr1-ge1-2.rh.mhwnj.cv.net (67.83.247.130) 18.962 ms 18.952 ms 18.902 ms
 5 67.59.239.121 (67.59.239.121) 18.844 ms 451be043.cst.lightpath.net (65.19.114.67) 18.314 ms
 6 451be031.cst.lightpath.net (65.19.98.49) 19.762 ms 64.15.3.138 (64.15.3.138) 10.763 ms 17.7
 7 * * *
 8 * * *
 9 RUTGERS-THE.ear3.Newark1.Level3.net (4.14.216.6) 33.338 ms 32.792 ms 33.274 ms
10 * * *
11 * * *
12 * * *
13 * * *
14 * * *
15 web-www.camden.rutgers.edu (128.6.34.90) 22.632 ms 23.859 ms 23.866 ms
```

Figure 2: Trace Route: All the routers in the path to destination

We can see that the packet traverse through the 15 routers to reach to camden.rutgers.edu server. This trace route is taken from a location in the New Jersey State.

4 Our Approach

4.1 The router as a point of analysis

From Figure 1 and 2 we know that routers are located at different geographical locations and also there are specific regions from which a given website or web server is accessed (except few). There are services called as GeoIP services. They

⁵hops are intermediate routers in the communication channel

can detect the geographical location of the system from which the IP packet has originated, but that is just an approximate, based on the source IP and not always correct. In case of DDoS attack this information is unreliable, because the packet source address is often spoofed, it is difficult to know the actual geographical location from which the packet has come, but router through which that packet has traveled can provide its own geographical information. Such information can be useful to understand the path through which a packet has traveled and thus we can know the region from which the packet was originated.

In the normal scenario there is a consistency in which website is accessed from different geographical regions, and this consistency can be found by measuring the number of requests or packets traveled from a router to a destination server. This behavior of accessing different websites from a router can be learned over the time. Thus finding this behavior in a flow⁶ at the router can form the basis of analysis in this paper. If there is a deviation from the learned behavior of accessing a particular destination, then that change in behavior as well as router geographical region information can be communicated to the destination network. The destination network on receiving that information, can decide, whether it wants the reporting router to discard or forward the traffic for it. This is a selective process in which traffic from only specific router is blocked while traffic from other routers remain unaffected

With the advance of electronics and the Internet of Things, processing and storage capacity of the electronic devices has increased. Routers are also not left behind, but storage capacity of routers is always very less compared to server which collect flow data for network traffic analysis. If we use learning techniques that don't need much storage then we don't have to store large chunk of packets on the router. Instead of storing data packets for longer time for analysis, we can learn from a small number of packets and then discard packets once learning is done, leaving behind only the learned information on the router. This is necessary

⁶In this paper we will consider flow only in the sense of destination addresses

because the number of entries on the Internet routing table has steadily grown. Now that the table has passed 500,000 routes [11], so storing each and every flow information for these routes could be difficult.

4.2 Internet packet flow capturing at the router

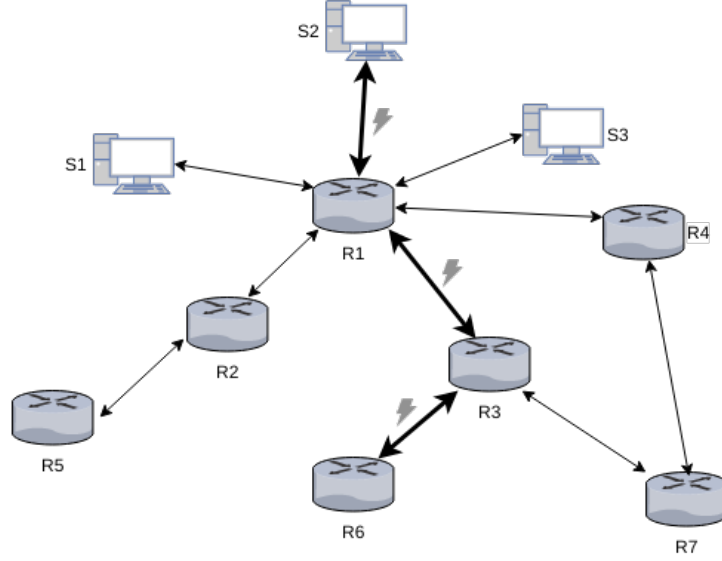


Figure 3: DDoS Attack path

In the above figure an attack is initiated from the region where router R6 is located and from router R6 data packets travel to the victim⁷ system S2. Attack packets traveled through router R3 and R1 to reach system S2. If we can detect an attack at router R6, then router R6 can discard all the packets heading towards system S2. In this process, only traffic from router R6 is affected but traffic from all other routers remain unaffected.

To achieve this, we will gather the flow information during a time window (e.g 300 sec) whose size will be fixed at the beginning. Thus a flow will contain all the packets that traveled from a router to different destinations during a time span of 300 seconds. We can also combine such flows to form a new flow with bigger time windows. E.g. if we combine two flow of 300 seconds then we will have flow of

⁷Victim is a computer system which is under DDoS attack

600 seconds. This flow information can be capture during a particular period or throughout the day. Once we have flow information we can apply learning techniques on each flow iteratively to gain deeper knowledge about normal behavior of the flow, e.g. on average, how many packets of different protocols are destined for a given destination from a given router/region during a particular time of the day.

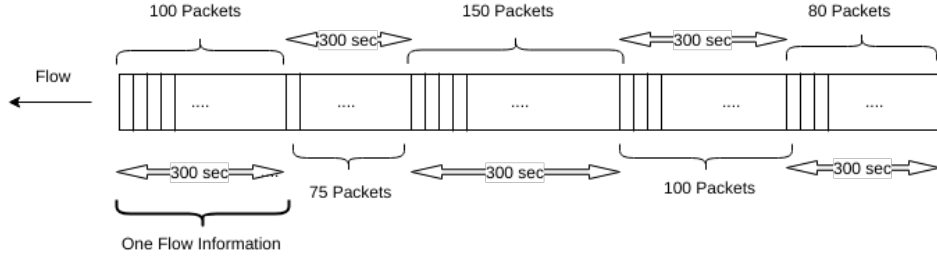


Figure 4: A 300 second time segment corresponds to a flow of information. The number of packets can vary in a flow.

4.3 Analysis techniques to be implemented

In the proposed system, each router will itself act as an analyzer. Each packet will be analyzed and a flow statistic will be created based on the destination IP address. Flow statistic will be used to create feature vectors and those feature vectors will be used for clustering destination IP address. Clustering is a process of examining the collection of points, and grouping the points into clusters according to some distance measure. The goal is to minimize the distance of every point in the cluster to other points in the same clusters. [12].

Once the clustering is learned (i.e. we know which destination IP address is tend to be in which cluster), the learned information will be used as a benchmark for all future flows. The routers will constantly keep clustering destination IP addresses and if there is deviation for the normal traffic at a router for any destination then that will affect the clustering and it will cause the destination IP address to be placed in different cluster. This change in the cluster for a given

destination IP address can be marked as a change in the behavior of the traffic for that destination. Along with the cluster we will also use the Novelty Detection algorithm to achieve more accurate result. This change in traffic behavior will be reported to destination network, which then decides on regulating the traffic coming to itself from the router which has sent the information.

5 Implementation

As discussed previously, there are different types of DDoS attacks, such as volume-based, application-based and also low-rate DDoS attacks. Among these different types of DDoS attacks, the volume-based attacks are most common. In the volume-based attack, a victim is flooded with a high volume of Internet packets (TCP, UDP, HTTP or ICMP), which make it unable to serve the requests.

For the demonstration of the suggested approach we are simulating a volumed-based bot attack. In the real world, the volume-based attacks are orchestrated using bots. Bots use compromised computer systems, controlled by an attacker for launching an attack. They are not bounded by geographical boundaries, so they can be anywhere in the Internet. Botnet (i.e. network of bots) are employed by an attacker to launch a DDoS attack. As we know that the Internet is connection of different computer system that communicate with each other, through different channels such as cables, satellite or radio device and these communication channels run throughout the glob; connecting different computer systems at different locations. For our simulation, instead of any bot program, we are using Low Orbit Ion Canon(LOIC) tool which is a free DDoS attack launching tool. This tool is even used by attackers in the real world to launch DDoS attacks.

We used Wireshark, an open source tool, for capturing Internet packets. Wireshark can capture all digital information received or sent through different devices such as Ethernet or wifi devices, which connect computers to the Internet. It also helps identify different protocol packets (e.g. TCP, UDP) within the wrapper

packet created at Data Link Layer packet.

The image shows the Wireshark 1.12.1 interface. The top menu bar includes File, Edit, View, Go, Capture, Analyze, Statistics, Telephony, Tools, Internals, and Help. Below the menu is a toolbar with various icons. The main window is divided into three panes. The top pane shows a list of captured packets with columns for No., Time, Source, Destination, Protocol, Length, and Info. The middle pane shows the details of the selected packet (No. 1), including Ethernet II, Internet Protocol Version 4, and User Datagram Protocol. The bottom pane shows the raw packet data in hexadecimal and ASCII.

No.	Time	Source	Destination	Protocol	Length	Info
1	0.000000000	192.168.0.7	173.194.175.189	QUIC	187	CID: 18146815167217732570, Seq: 49
2	0.049933000	173.194.175.189	192.168.0.7	QUIC	152	CID: 0, Seq: 56
3	3.105773000	192.168.0.7	192.168.0.18	TCP	223	[TCP segment of a reassembled PDU]
4	3.107726000	192.168.0.18	192.168.0.7	TCP	263	[TCP segment of a reassembled PDU]
5	3.107974000	192.168.0.7	192.168.0.18	TCP	188	52718->8009 [ACK] Seq=116 Ack=118 Win=346 Len=0 TSval=32023603 TSecr=254463
6	4.701349000	192.168.0.7	52.204.61.141	TLSv1.2	396	Application Data
7	4.717465000	52.204.61.141	192.168.0.7	TCP	146	443->40122 [ACK] Seq=1 Ack=289 Win=314 Len=0 TSval=1050903052 TSecr=32024001
8	4.718005000	52.204.61.141	192.168.0.7	TLSv1.2	471	Application Data
9	4.718249000	192.168.0.7	52.204.61.141	TCP	188	40122->443 [ACK] Seq=289 Ack=326 Win=727 Len=0 TSval=32024005 TSecr=1050903052
10	7.255999000	192.168.0.7	54.89.16.99	TCP	188	40110->443 [ACK] Seq=1 Ack=1 Win=237 Len=0 TSval=32024640 TSecr=3880385116
11	8.001888000	65.52.108.76	192.168.0.7	TLSv1.2	1351	Application Data
12	8.013787000	192.168.0.7	65.52.108.76	TCP	1536	[TCP segment of a reassembled PDU]
13	8.014102000	192.168.0.7	65.52.108.76	TLSv1.2	1229	Application Data
14	8.035860000	65.52.108.76	192.168.0.7	TCP	146	443->39116 [ACK] Seq=1206 Ack=2550 Win=514 Len=0 TSval=119147688 TSecr=32024829
15	8.100997000	192.168.0.7	192.168.0.10	TCP	223	[TCP segment of a reassembled PDU]

Frame 1: 187 bytes on wire (856 bits), 187 bytes captured (856 bits) on interface 0
 Radiotap Header v0, Length 14
 IEEE 802.11 QoS Data, Flags: .p.....T
 Logical-Link Control
 Internet Protocol Version 4, Src: 192.168.0.7 (192.168.0.7), Dst: 173.194.175.189 (173.194.175.189)
 User Datagram Protocol, Src Port: 37254 (37254), Dst Port: 443 (443)
 QUIC (Quick UDP Internet Connections)

```

0000 00 00 0e 00 00 00 0a 00 00 00 07 04 07 88 41 .....A
0010 00 00 10 35 0d e3 00 04 ef 18 00 00 73 3c df ...5.....5%
0020 a9 67 c0 39 f0 ef 00 00 34 2f 00 20 00 00 00 00 ...9...4/....
0030 aa aa 03 00 00 00 00 00 45 00 00 33 c1 93 40 00 .....E..3..@.
0040 40 11 5a f7 c0 a0 00 07 ad c2 af bd 91 86 01 bb @.Z.....
0050 00 1f 44 e9 0c da 00 07 60 a2 ba 00 0c 31 d5 00 ..D.....K....1..
0060 1f 60 73 54 c3 d5 ea 93 a9 02 7f .....sT.....

```

Figure 5: Wireshark Tool: snippet of captured packets

To gather data for the demonstration, we have created a small network which has one router and couple of host machine. Each machine can be a victim of a DDoS attack. We have installed Wireshark tool on one of the machine in the this network. For capturing traffic in the network we are using the Promiscuous mode of Wireshark. In Promiscuous mode, a network interface can record not only the traffic that is intended to itself but all of the traffic on the network, so we can see all in/out traffic in our setup network. This setup is similar to any router on the Internet which is connected with different routers and hosts.

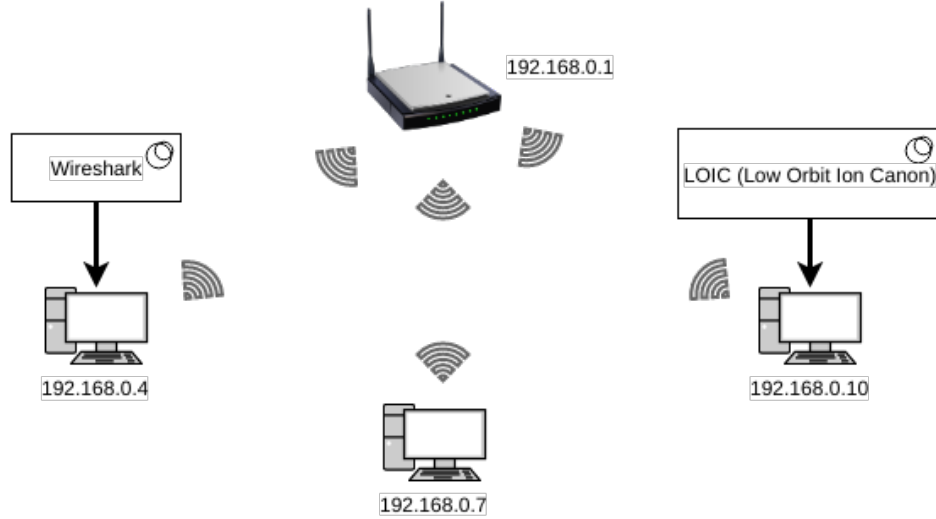


Figure 6: Network for simulating DDoS detection

To collect packets traveling in our network, we start the Wireshark tool. We let it run for a while and then we orchestrate a DDoS attack on one of the host (e.g. 192.168.0.7 in Figure 6) in the network. This attack is engineered using the Low Orbit Ion Canon(LOIC) tool that is installed on one of the host (e.g. 192.168.0.10 in Figure 6) in our network. This tool allow us to launch TCP and UDP flood attack on any destination. In this paper we are analyzing the flow based TCP and UDP attack, which is one of the most common type of DDoS attacks.

In the LOIC tool we need to give the IP address and the port number of the destination where we want to orchestrate attack while rest of the work is done by it. This tool floods destination with packets and if we choose TCP then it will try to create multiple connections and send packets over them. We start flood attack on one of the system(e.g. 192.168.0.7) in demo network and let it continue for few minuets. We launched such attacks few times in between the packet capturing session.

All the traffic, including the normal and the attack traffic, will get captured using Wireshark. Packets are captured for about five hours in the given network. Once the packets are captured they are saved as pcapng file, which is a Wireshark file format for captured packets. Captured packets during the normal operation

and during the attacks are saved separately. The normal packet flow information is used for training and testing the learning algorithms (we will explain it in later sections) while attack packets flow information is used for detecting the attack. Wireshark captures every detail of the packets but we don't need all of the information, we are interested only in the IP layer information of the packet. Most of the routers analyze IP layer of the packet for routing, having said that, there is no reason that other layers of the packet can not be analyzed, but for our demonstration purpose we are analyzing only IP packets.

A data extraction program is written in Python language to extract IP layer information from the captured packets. This program extracts address, port and time information from each of the captured IP packets⁸.

Destination IP Address	Protocol	Time stamp(Sec.)	Sample Number
52.6.129.72	6	1512094785.928596000	1
192.168.0.4	6	1512094785.946987000	1
192.168.0.4	17	1512094786.148488000	1

Table 1: Sample file snippet (row represents an IP packet)

Our data extraction program also divides the captured data into 300 seconds capture window, thus creating a sample data which is a collection of IP packets captured over the time of 300 seconds. It then writes each sample in the separate file for further processing. This sampling of the packets is the continuous process. We then run another program which extracts the flow information from those sample files. One 'flow' contains the number of different packets capture for a given destination. We will store this flow information as sample flow. Then we will train learning algorithms using those sample flows. Once learning is done, those sample files will be discarded and new samples will be created for further training. This training process has to be continuous process in order to correctly reflect the current status of the flow at given router. Whatever the new information learned,

⁸packets containing IP information

is augmented with the previous learning to have the correct understanding of the flow. This learning can be done for the time during a day or during a week of a year. E.g. We can have separate learning information for flow from morning 9 am to 12 pm and also can have information for evening 6 pm to 12 pm.

Flow based model is built, as it is more reliable and fast. Packet analyzing is often difficult due to the size and encryption. Also destination port number is not a reliable information in detecting attacks because of the fact that attacker uses different ports during an attack.

Creating a training set for the learning algorithms is an intermediate step in which IP packet count for each destination for a given protocol(e.g TCP, UDP) is calculated. The training set gives us the flow information for each destination (i.e. how many packets of a particular protocol are recorded for a given destination IP address during a time window e.g. 300 seconds).

Destination IP Address	IP Packet count		
	ICMP	TCP	UDP
172.217.10.134	0	8	12
65.19.96.252	5	0	192
68.67.178.134	0	78	0

Table 2: Training Set with three training examples

We are using Python program to create training sets, as given in ??, from the sample files. Each row in a training set is one training example with destination IP address as label. A training example is \mathbb{R}^3 vector whose elements are the number of packets observed for a particular protocol for a destination during a fixed time (e.g. 300 second). There are around 150 protocols managed and assigned by the Internet Assigned Numbers Authority (IANA) but most commonly used protocols in the DDoS attack are ICMP, TCP and UDP protocols. For the training and analysis purpose we are using only these three protocols as the desired feature. In the larger system such as routers managed by ISP, other protocols can also be

used as features if required.

Each sample file is processed and corresponding train/test set file is created. Train and test sets files are stored on the file systems, so that those can be used by other programs for training and testing the learning algorithms.

5.1 Machine Learning

According to Tom Mitchell (1998), a computer program is said to learn from an experience E with respect to some task T and some performance measure P , if its performance on T , as measured by P , improves with experience.

A learning algorithm builds a hypothesis using a training set as input. Then that hypothesis is used to perform predictions. The most common categorization of machine learning algorithms is Supervised and Unsupervised.

Let f be the function from \mathbb{R}^d to \mathbb{R}^d which we need to guess from an input vectors x^1, x^2, \dots, x^n also called as ‘input variables’ or ‘feature vectors’. Let Ξ be the set of such input vectors. Let be the n number of input vectors in a training set Ξ . Let H be the set of some functions from \mathbb{R}^d to \mathbb{R}^d . Let $h \in H$ be the hypothesis about function $f \in H$. We select h based on a training set $X \in \Xi$, of m input vectors. In Supervised leaning we know the values of f for m samples in the training set X . We assume that if we can find a hypothesis h that closely agrees with f for the members of X , then this hypothesis will be a good guess for f when X is large. In Unsupervised learning, we simply have a training set of vectors without function values for them. The problem in this case, typically, is to partition the training set into subsets, X_1, \dots, X_N , in some appropriate way. [12]

A supervised algorithm such as One Class Support Vector Machine (SVM) [13] is efficient at identifying the anomalies in the data but this algorithm is process and memory intensive, so training the algorithm for each and every IP address is very costly in terms of resources. Because of the resource constraints of the router, our approach is to first cluster the IP addresses based on the features using Unsupervised learning algorithms such as k-means and then apply One Class

SVM on the clusters to decide on the boundaries of those clusters. The k-means algorithm is fast and consumes less resource compared to One Class SVM, that makes it good choice to be used on the devices such as routers which have less processing power and less memory.

5.1.1 Feature Scaling

Before feeding training data, which was acquired in an earlier stage, to a learning algorithm, we have to do feature scaling, also called Standardizing. Feature scaling is done by subtracting the mean and scaling the feature to a unit variance value. It is necessary because different features which are at different scales could cause one feature dominating the others in the algorithm output result. e.g. consider two vectors (1, 2, 3000) (1, 3, 2000). If we calculate the Euclidean distance between these two vectors using the formula $\sqrt{\sum_{i=1}^n (p_i - q_i)^2}$, then the distance will be $\sqrt{(1-1)^2 + (2-3)^2 + (3000-2000)^2}$. From this, it is evident that the larger term is dominating the result.

First we will convert training examples into a vectors. Each vector is one training example and each coordinate in the vector is the feature. To standardize the input vector, the mean and the standard deviation are calculated for the set of input vectors. Then a new vector is created by subtracting the mean from each feature vector and dividing that feature vector by its standard deviation. The new vector created after this step is standardized vector. A set of such standardized vectors is used as input to the learning algorithms.

$$\text{Standardization formula: } x' = \frac{x - \bar{x}}{\sigma}$$

Where x is a feature vector, \bar{x} is the mean and σ is its standard deviation.

5.1.2 Clustering

k-means [14] clustering is one of the most efficient algorithms for creating clusters. This algorithm takes any k randomly chosen points as centroids (also called, cluster centers) $\mu_1, \mu_2, \dots, \mu_k$ as input from the training set $X = \{x^1, x^2, \dots, x^m\}$, $x^i \in \mathbb{R}^d$,

$i = 1, 2, \dots, m$.

The following is Lloyd's algorithm which is most popular heuristic algorithm for k-means clustering. The clustering that we will be doing is of destination IP addresses, such that each cluster will have some number of destination IP addresses.

Algorithm 1 k-means

- 1: **repeat**
 - 2: Calculate the distance for each element (or data point) of the training set to all the centroids
 - 3: For each element, assign the element to the centroid which is nearest
 - 4: Recalculate the new centroid for all the element in one cluster by taking the mean. If $x_1^j, x_2^j, \dots, x_n^j$ are the elements of the cluster j . Then for cluster j , the new centroid will be $\mu_j = \frac{1}{n}[x_1^j + x_2^j + \dots + x_n^j]$ where n is number of point assigned to cluster j . Do this for all the clusters
 - 5: **until** No data point is reassigned to a different centroid
-

For this paper we will be using the k-means++ algorithm [15] which is an improvement of k-means, where an arbitrarily initialization step is replaced by simple, randomized seeding technique. If $D(x)$ is the shortest distance from a data point $x \in X$ to the closest centroid we have already chosen, then the following is the k-means++ algorithm.

Algorithm 2 k-means++

- 1: Take one centroid μ_1 , chosen uniformly at random from X
 - 2: Take a new centroid $\mu_j \neq \mu_1$, by choosing $\mu_j \in X$ with probability $\frac{D(\mu_j)^2}{\sum_{x \in X} D(x)^2}$
 - 3: Repeat Step 2. until we have taken k centers altogether
 - 4: Proceed with the Lloyd's k-means algorithm skipping random the initialization stage
-

For our DDoS attack detection program, we will be using the Scikit-learn libraries. Scikit-learn is the most popular and rich open-source machine learning

software library for the Python programming language and it has implementation for both k-means++ and One Class SVM machine learning algorithms, also it has data preprocessing programs such as feature scaling.

We will take a training set (in the format given in figure 2) and we will cluster training examples from that training set to using the k-means++ clustering algorithm from Scikit-learn library.

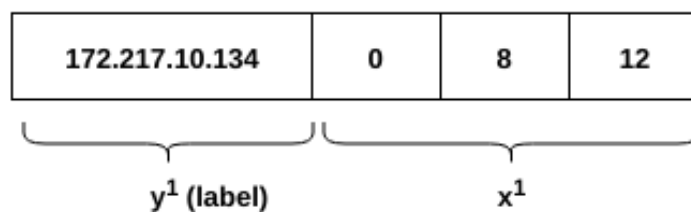


Figure 7: One training example

Before doing the clustering, we have to first determine the number of clusters and centroids.⁹ Deciding on the number of clusters is important, because randomly choosing the number of clusters will not have correct clustering. So, we will use the Elbow method to find the optimal number of clusters in our training set. The Elbow method checks the portion of the variance explained by function of the number of clusters. Following are the Elbow Diagrams for four samples.

⁹Centroid is the vectors which is arithmetic mean of all the point in the cluster

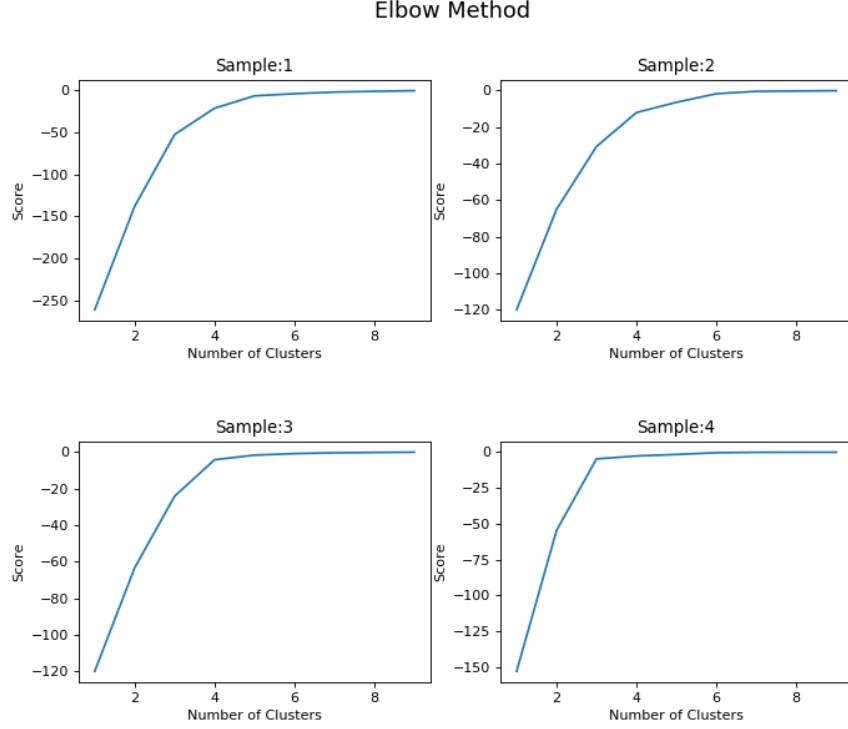


Figure 8: Elbow Method for cluster count detection

Using the elbow method, the variance for each cluster number is calculated and the cluster number which produces less variance for the next cluster number is selected as best choice for the give training set.

To have correct clustering we first run the k-means++ algorithm to determine the central vectors called centroids. Clusters are represented by these centroids. As we have multiple training sets (also called ‘samples’), we have cluster centroids for each sample. But in the future we will require only one set of centroid which would be the more accurate clustering for all our training sets. To achieve this, we will save all the centroids, and then find the median of those centroids after removing the outliers, and that will give us a good estimate of centroids. This estimated centroids are used for clustering the samples in the future. Following is the example of centroid vector, where a rows represents a cluster and a column represents a feature.

Cluster	ICMP	TCP	UDP
0	-0.16815612	-0.14928111	-0.16948046
1	-0.18181818	5.13527652	5.68956244
2	5.08663322	-0.27110845	-0.099885
3	-0.18670401	-0.18804342	-0.018538

Table 3: Example of cluster centroids: A row is a cluster and column is a feature

To draw data points in two dimensional space we had to reduce the three dimensional training example into two dimensions without losing much information. We have achieved this by using Principal-Component Analysis (PCA). PCA is a technique for taking a dataset consisting of a set of tuples representing points in a high-dimensional space and finding the directions along which the tuples line up best. [16]. We have used the Scikit-learn PCA module for this purpose.

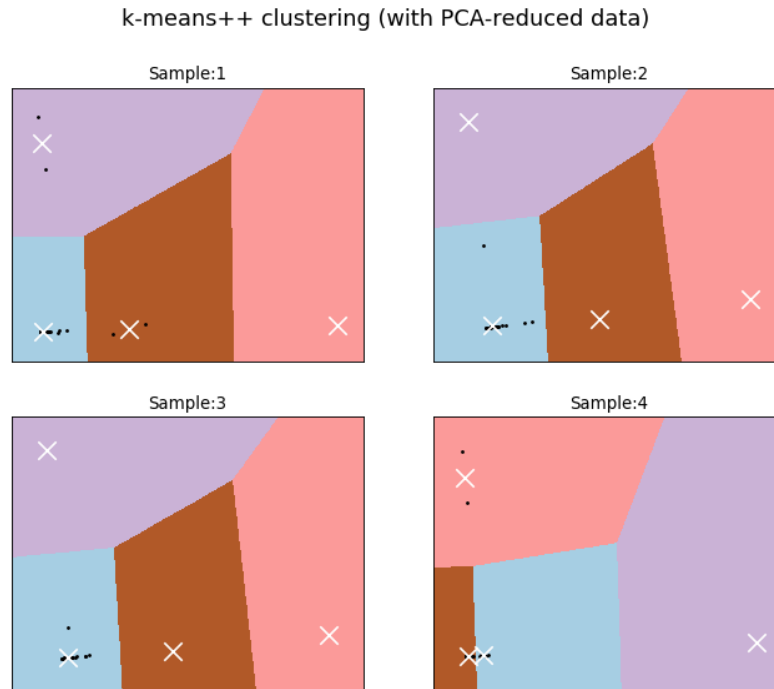


Figure 9: Clustering using k-means++ algorithm

Using the centroids obtained in the previous step, we ran k-means++ algorithm

and then test it for accuracy. k-means++ algorithm gives label to each training example. This label is the cluster number to which that training example has been assigned to. So the new labeled data look like in 4 below.

Destination IP Address	ICMP	TCP	UDP	Cluster
172.217.10.134	0	8	12	1
65.19.96.252	5	0	192	0
68.67.178.134	0	78	0	2

Table 4: Labeled Training Set (with cluster number)

We ran the k-means++ for both training and testing set, and produce labels for both. We then checked how similar the test set clustering is with train set clustering i.e. we are checking, whether an IP address has the same cluster number assigned in both the training and test sets or it differs. For measuring this similarity Rand Index(RI) [17] is used. RI is a measure of how many percent does the test clustering matches with the trained model.

$$RI = \frac{TP+TN}{TP+FP+FN+TN}$$

where TP, TN, FP and FN are the number of true positives, true negatives, false positives and false negatives respectively. We have observed that, with more training sets, the RI index improves.

Because of the training sets contain the flow information for different IP addresses on a router during the window of 300 seconds, there is a very high possibility that the same destination IP address is captured in multiple training set. Our goal is to find the correct cluster for the destination IP address, and to achieve this goal, we check all the cluster assignments for a given destination IP address and we select the cluster which has the highest occurring frequency among all the training sets. We also count the average number of packets going to a given IP address which will help us reduce the error in detecting the anomaly. We will save this labeling information and cluster count information for each IP destination on the file system.

Destination IP Address	Cluster	Packet Count
74.125.141.106	1	113
72.30.2.182	0	16
64.94.191.14	0	22

Table 5: Learned information after clustering

The clustering information tells us the normal behavior of the packets traveling from the router to a given destination. As we have fixed the centroids for the clustering, every time in the future we should expect an IP address to be found in the same cluster mentioned in the clustering information table. But If there is a DDoS attack on a any destination with flooding, then we can expect to see the destination IP address assigned to different cluster.

5.1.3 Anomaly Detection using One Class Support Vector Machine

Support Vector Machine: From the experiments on different data sets, it is found that the destination under attack is labeled with the same cluster number. This happened because there is no other cluster it can be assigned to, so it gets assigned to the cluster whose centroid is nearest.

To avoid such a situation, we will have to create boundaries for the clusters, and if any IP address is out side of the cluster boundries then that can be considered suspicious. This provision will make sure that the attack on the destination IP address will be detected even if the destination is assigned to the same cluster it was assigned in the past.

To create cluster boundaries we have used One Class SVM. One Class SVM is a supervised machine learning algorithm, which means, the training example will have label. I our case, cluster number will be treated as label for a training example (see table:labeled-set).

Because SVM solves a quadratic programming problem (QP), computing and storage requirements of SVM increases rapidly with the number of training vec-

tors. The approach presented in this paper is more efficient because the number of cluster are limited and always be far less than the number of destination IP address. Training one class SVM on clusters will be far less process and memory intensive than training on massive number of IP addresses. E.g. In our simulation training sets we have 4 four clusters, while the unique number of destination IP addresses in all the training sets is 268, which is more than 60 times the count of clusters we have found.

Support Vector Machine is a supervised learning algorithm which tries to classify training examples into two distinct classes. Classification is based on the labels of the training set. Consider the training set $\{(x^1, y^1), \dots, (x^m, y^m)\}$, where $x^i \in \mathbb{R}^d$ is the training example and $y^i \in \{-1, +1\}$ is a label (or a class) for x^i . SVM will separates those classes using a line or a curve.

SVM tries to create a non-linear separation boundary by projecting data points through a non-linear function ϕ (kernel function) to higher dimension space. The data points in space I which can't be separated by a line are projected to the feature space $F \in \mathbb{R}^d$ where there can be a hyperplane that separates data points of one class from another. If that separating hyperplane is projected back on the original space I then we get a non-linear curve. [13]

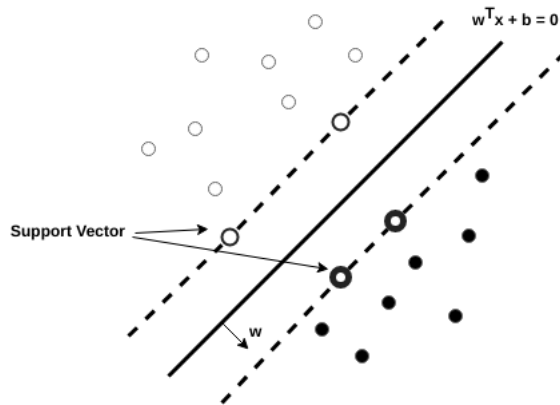


Figure 10: Linear Separation using Support Vector

The hyperplane is represented with the equation $w^T x + b = 0$, where $x \in \mathbb{R}^d$, $w \in F$ and $b \in \mathbb{R}$. ϕ is a kernel function. This hyperplane separates the training

example labeled with -1 and 1 . The position of the hyperplane is such that the distance from the closest point from each class to the hyperplane is same. To avoid the over-fitting, slack variables ξ^i are introduced. Over-fitting happens because the learned hypothesis fits training examples so well that it fails to generalize the new examples. The SVM classification is an optimization problem which is stated as follows. [13] [18]

$$\begin{aligned} & \min_{w,b,\xi^i} \frac{\|w\|^2}{2} + C \sum_{i=1}^m \xi^i \\ & \text{subject to:} \\ & y^i(w^T \phi(x^i) + b) \geq 1 - \xi^i \text{ for all } i = 1, \dots, m \\ & \text{where } y^i \in \{-1, +1\}, x \in \mathbb{R}^d, b \in \mathbb{R} \text{ and } \xi^i \geq 0 \text{ for all } i = 1, \dots, m \end{aligned}$$

The constant $C > 0$ is the regularization parameter. If C is chosen large, miss-classification of training examples can be avoided. If chosen small, then we may miss-classify few examples, but the margin will be large, so most of the points will be far away from the decision boundary. If this minimization problem is solved using Lagrange multipliers then the classification function is:

$$\text{sign}(\sum_{i=1}^m \alpha^i y^i K(x, x^i) + \rho)$$

α^i is the Lagrange multipliers. $\alpha^i y^i$ called the support vectors. The function $K(x, x^i) = \phi(x)^T \phi(x^i)$ is the kernel and $\rho \in \mathbb{R}$ is the intersection.

One Class Support Vector Machine: One Class Support Vector Machine (SVM) is the extension of SVM which detects boundaries of the training set so that every new training example will be classified as belong to training set or not. It separates all the training set data point from feature space F and maximizes the distance of hyperplane from F . This creates a binary function which returns $+1$ for the training example that fits in the trained set region, otherwise it will return -1 .

The minimization function of One Class SVM is slightly different than the SVM. [13]

$$\begin{aligned}
& \min_{w, \rho, \xi^i} \frac{\|w\|^2}{2} + \frac{1}{\nu m} \sum_{i=1}^m \xi^i - \rho \\
& \text{subject to:} \\
& (w \cdot \phi(x^i)) \geq \rho - \xi^i \text{ for all } i = 1, \dots, m \\
& \xi^i \geq 0 \text{ for all } i = 1, \dots, m
\end{aligned}$$

The new parameter $\nu \in (0, 1]$ introduced in place of C in previous SVM equation is used to set upper bound on outliers/anomalies and lower bound on the number of training examples.

We are using Scikit-learn's 'OneClassSVM' library to train the model and create classifier for each cluster. The input vector to the classifier is the set of all the training examples belonging to the same clusters. Thus, if we have four clusters then we will have four classifier. The input vector to the classifier will be of the form shown in table 4 on page 22.

Following is the result of modeling on the training data sets. Each cluster has its own model.

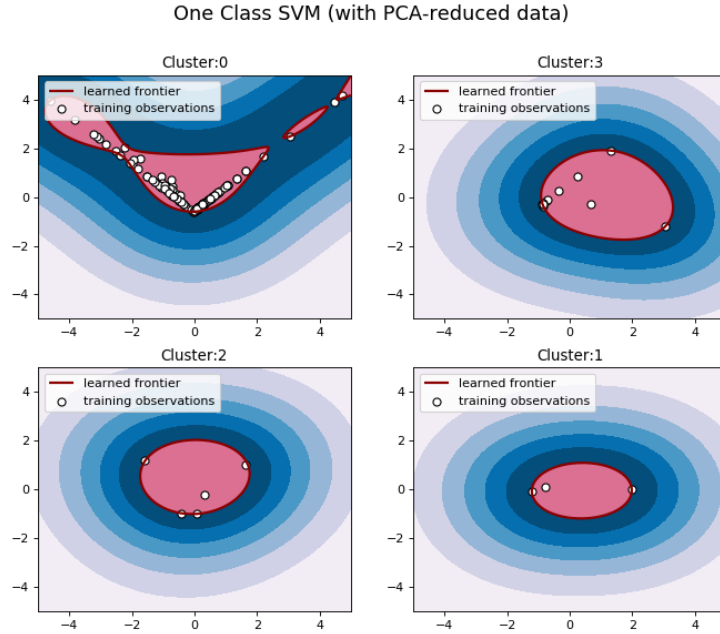


Figure 11: One Class SVM

6 Detection

By this stage, we have all the destination IP address, at the router, labeled with cluster and their count of average number of packets during the flow. Also we have classifier models for each of the cluster.

To detect the attack we capture the flow in the specified interval (e.g. 300 seconds, in which we have trained our models). We transformed the captured flow into a test set. The created test set is then clustered using already centroids which were calculated for the clustering of the training sets. Now we have each IP in the new test set labeled with the cluster number. If there is any destination IP address for which new cluster label does not match with the already known cluster label then that destination address added to the suspect list. If new label matches with the old then the feature vector for that destination address is passed to the cluster classifier to check its subscription to the cluster it was found in. i.e. if an IP address is labeled with cluster number 0 then we use One Class SVM classifier for cluster number 0. If the output of the classifier is -1 then that destination is added to the suspect list. For every destination from this compiled list, the number of packets observed is also recorded. If there is significant difference between the number of capture packet in the past and in the present then that destination is recorded as DDoS attack candidate because failing to find in the same cluster boundary it was in the past is a sign of the change in the behavior of the traffic for destination IP address on the router where this analysis is done.

Using this we have successfully detected attack on the destination IP address 92.168.0.7 in our modeled network.

Next task of the router then will be to communicate the change in the behavior of traffic observed at that router for that destination. The destination address here could be a router or an application server. Destination system i.e. either router or the server will collect the information received. This information can be used by destination to know the nature of the attack.

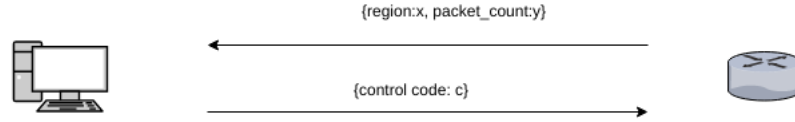


Figure 12: Router Network Communication

To Communicate the destination network about the change in behavior, router can use the existing ICMP protocol. ICMP protocol is used to provide feedback about problems in the communication environment. ICMP messages are sent in several situations: for example 1) when a datagram cannot reach its destination. 2) when the gateway does not have the buffering capacity to forward a datagram. 3) when the gateway can direct the host to send traffic on a shorter route. [10] Similarly we can use ICMP protocol to inform destination system about the change in the traffic. ICMP protocol has many unused type code (there can be 0-255 types but as of now only 0-41 are in use) available. We can create a new ICMP ‘type’ to send DDoS detection information from the router to the destination system and then destination system can send mitigation instructions to the routers.

Depending on the type and the severity of the situation, the destination system can decide whether or not to inform router to block the traffic coming from that router. We will let this decision to be taken at the destination system. There can be different parameter based on which the destination system can decide blocking the flow. Many questions can be asked before making the decision, such as how many routers have reported the change in traffic?, Is the attack information coming from the region which never had traffic in the past?

7 Conclusion

A novel way to identify and mitigate DDoS attack is discussed in this paper. With the advance of NVF(Network Virtual Functional) it is easy to push the learning algorithms to the routers allowing them to detect a DDoS attacks using machine learning algorithms. The DDoS detection and mitigation information can be communicated using existing ICMP protocol making system available to every one who want to use it. This approach has been tested with the small network so further experiments are needed to be performed on the larges networks.

References

- [1] Derek Kortepeter. Destructive ddos attacks increasing at a rapid rate. <http://techgenix.com/ddos-attacks-increasing/>, December 2017. Online; accessed 22-August-2017.
- [2] D. W. Davies. A communication network for real-time computer systems. *Radio and Electronic Engineer*, 37(1):47–51, January 1969.
- [3] Cisco Security portal. A cisco guide to defending against distributed denial of service attacks. <https://www.cisco.com/c/en/us/about/security-center/guide-ddos-defense.html>. Online; accessed 1-August-2017.
- [4] S. M. Bellovin. A look back at "security problems in the tcp/ip protocol suite". In *20th Annual Computer Security Applications Conference*, pages 229–249, Dec 2004.
- [5] Ken Dunham and Jim Melnick. *Malicious Bots: An Inside Look into the Cyber-Criminal Underground of the Internet*, chapter 1. Introduction to Bots. CRC Press, August 2008,.
- [6] Network Working Group. Traffic flow measurement: Architecture. <https://www.ietf.org/rfc/rfc2722.txt>, October 1999.
- [7] Cisco. Introduction to cisco ios netflow. https://www.cisco.com/c/en/us/products/collateral/ios-nx-os-software/ios-netflow/prod_white_paper0900aecd80406232.html, May 2012. Online; accessed 22-August-2017.
- [8] sFlow. Using sflow. http://www.sflow.org/using_sflow/index.php. Online; accessed 16-November-2017.
- [9] Cisco. What is a network switch vs. a router? <https://www.cisco.com/c/en/us/solutions/small-business/resource-center/connect-employees-offices/network-switch-what.html>. Online; accessed 16-November-2017.
- [10] Network Working Group. Icmp. <https://tools.ietf.org/html/rfc792>, September 1981.
- [11] Omar Santos. The size of the internet global routing table and its potential side effects. <https://supportforums.cisco.com/t5/network-infrastructure-documents/the-size-of-the-internet-global-routing-table-and-its-potential/ta-p/3136453>, May 2014. Online; accessed 22-August-2017.
- [12] Jeffrey D. Ullman Jure Leskovec, Anand Rajaraman. *Mining of Massive Datasets*, chapter Clustering. Stanford University, 2014.
- [13] Bernhard Schölkopf, Robert Williamson, Alex Smola, John Shawe-Taylor, and John Platt. Support vector method for novelty detection. In *Proceedings of the 12th International Conference on Neural Information Processing Systems*, NIPS'99, pages 582–588, Cambridge, MA, USA, 1999. MIT Press.

- [14] Nathan S. Netanyahu Christine D. Piatko Ruth Silverman Tapas Kanungo, David M. Mount and Angela Y. Wu. An efficient k-means clustering algorithm: Analysis and implementation. *IEEE*, 24(7):881–892, July 2002.
- [15] David Arthur and Sergei Vassilvitskii. K-means++: the advantages of careful seeding. In *In Proceedings of the 18th Annual ACM-SIAM Symposium on Discrete Algorithms*, 2007.
- [16] Jeffrey D. Ullman Jure Leskovec, Anand Rajaraman. *Mining of Massive Datasets*, chapter 11. Dimensionality Reduction. Stanford University, 2014.
- [17] William M. Rand. Objective criteria for the evaluation of clustering methods. *Journal of the American Statistical Association*, 66(336):846–850, 1971.
- [18] Jeffrey D. Ullman Jure Leskovec, Anand Rajaraman. *Mining of Massive Datasets*, chapter 12. Large-Scale Machine Learning. Stanford University, 2014.