

DDOS ATTACK DETECTION BY CONTROL CHARTS

by

Aysun Kalemci

B.S., Electronics and Communications Engineering, Doğuş University, 2012

Submitted to the Institute for Graduate Studies in
Science and Engineering in partial fulfillment of
the requirements for the degree of
Master of Science

Graduate Program in Electrical Electronics Engineering
Boğaziçi University

2019

ACKNOWLEDGEMENTS

This work is supported by the Scientific and Technological Research Council of Turkey (TUBITAK), under Cloud-Based Privileged Access Management System Project, Project No. 117R030.

I would like to thank my thesis advisor, Prof. Emin Anarım for his patience and contributions during this thesis. He has always been sophisticated, kind and ready to help.

Secondly, I want to thank Prof. Mutlu Koca and Assoc. Prof. Şerif Bahtiyar for their precious time to approve this thesis.

Derya Erhan and Ramin Fouladi have great contributions to this thesis. I appreciate for their friendship and assistance.

Thank you for being helpful and supportive.

Finally, I would like to thank my father, Aydın Kalemci and my mother Mualla Kalemci for encouraging me with their best wishes.

Aysun Kalemci, Istanbul, 2019

ABSTRACT

DDOS ATTACK DETECTION BY CONTROL CHARTS

Distributed Denial of Service (DDoS) attacks are considered as the major threats in today's cyberworld. The fact that the source of these threats is often uncertain increases the concerns of many network operators. These types of attacks exhaust the resources to make them unavailable for the legitimate users and they take control over remote hosts. Infrastructure dependent business processes are adversely affected so that companies suffer financial losses. They are violating the security components of information security; confidentiality, integrity and availability. However, many techniques to overcome DDoS attacks have been developing by researchers who have the awareness of these threats.

In this thesis, in order to detect DDoS attacks, we first compared cumulative summation patterns of datasets which have normal and weibull distributions. We applied Tabular CUSUM and V-mask CUSUM methods to two datasets which we maintained at Boğaziçi University by using hping DDoS tool. It was found that these techniques can be applied to detect the anomalies of DDoS attack traffic by analyzing numerical changes of SYN packets during the process. The comparison of the accuracy rates of Tabular CUSUM and V-mask CUSUM techniques was made by Receiver Operating Characteristic (ROC) curves. We made a performance analysis of EWMA and CUSUM control charts evaluating Average Run Length (ARL) approximation. Finally, the Autoregressive Integrated Moving Average (ARIMA) forecasting model was applied in order to obtain the forecasting residuals which are also utilized in the performance evaluation of these two control charts.

ÖZET

KONTROL ÇİZELGELERİ İLE DDoS ATAKLARININ TESPİTİ

DDoS saldırıları, günümüzün siber dünyasında ana tehditler olarak kabul edilmektedir. Bu tehditlerin kaynağının çoğu zaman belirsiz olması, birçok şebeke operatörünün endişelerini artırmaktadır. Altyapıya bağımlı iş süreçleri olumsuz etkilendiğinden, şirketler finansal kayıplara maruz kalmaktadır. Bu tür saldırılar, güvenliğin bileşenlerini; gizlilik, bütünlük ve kullanılabilirliği ihlal ederek kaynakları yasal kullanıcılar tarafından kullanılamaz hale getirmekte ve kaynakların kontrollerini ele geçirmektedir. Ancak, DDoS saldırılarının üstesinden gelmek için, birçok teknik, bu tehditlerin farkındalığına sahip olan araştırmacılar tarafından geliştirilmektedir.

Bu çalışmada, DDoS saldırılarının tespiti için, normal ve weibull dağılımlarına sahip veri kümelerinin kümülatif toplama modellerini karşılaştırdık. Hping DDoS aracını kullanarak elde ettiğimiz iki veri setine Tabular CUSUM ve V-mask CUSUM yöntemlerini uyguladık. Bu tekniklerin, atak trafiğinde TCP syn paketlerinin sayısal değişikliklerinin algılanmasında etkili olduğunu gördük. Tabular CUSUM ve V-mask CUSUM tekniklerinin doğruluk oranlarının karşılaştırılması için Alıcı İşletim Karakteristik (ROC) eğrilerini kullandık. Ortalama Çalışma Uzunluğu (ARL) ile EWMA ve CUSUM kontrol çizelgelerinin performans analizini yaptık. Son olarak, yine bu kontrol çizelgelerinin performansını değerlendirmede kullanılan tahmin kalıntılarını elde etmek için Otoresif Entegre Hareketli Ortalama (ARIMA) tahmin modelini uyguladık.

TABLE OF CONTENTS

ACKNOWLEDGEMENTS	iii
ABSTRACT	iv
ÖZET	v
LIST OF FIGURES	viii
LIST OF SYMBOLS/ABBREVIATIONS	xi
1. INTRODUCTION	1
2. DISTRIBUTED DENIAL OF SERVICE ATTACKS	3
2.1. A Brief History of DDoS Attacks	3
2.2. Geography of DDoS Attacks	4
2.3. Dynamics of the number of DDoS attacks	5
2.4. Types and duration of DDoS attacks	6
3. STATISTICAL PROCESS CONTROL CHARTS	8
3.1. SHEWHART CONTROL CHART	9
3.2. EWMA CONTROL CHART	10
3.3. CUSUM CONTROL CHART	11
3.3.1. Derivations of CUSUM	13
3.3.1.1. Intuitive Derivation	13
3.3.1.2. Repeated Sequential Probability Ratio Test Derivation	15
3.3.2. Tabular Form of CUSUM	16
3.3.3. V-mask Form of CUSUM	18
3.3.3.1. Definition of Chebyshevs Inequality	21
4. CONTROL CHART PERFORMANCE	23
4.1. Average Run Length Approximation	23
4.2. Forecasting Models Approximation	27
5. EXPERIMENTAL RESULTS	29
5.1. Cumulative Sum with Normal Distribution	29
5.2. Cumulative Sum with Weibull Distribution	30
5.3. Applying Tabular CUSUM Method	30
5.4. Applying V-mask CUSUM method	33

5.5. Comparison Analysis of Control Charts	35
5.6. Results and Discussion	38
6. CONCLUSION	45
REFERENCES	47

LIST OF FIGURES

Figure 1.1.	Distribution of DDoS attack types in 2019 Q1	2
Figure 2.1.	DDoS attacks distributions by country in 2018 Q4 and 2019 Q1	5
Figure 2.2.	The number of DDoS attacks in 2019 Q1	6
Figure 2.3.	DDoS attacks distributions by days in 2018 Q4 and 2019 Q1	6
Figure 2.4.	DDoS attacks distribution by time in hours in 2018 Q4 and 2019 Q1	7
Figure 3.1.	Typical behaviors of the processes, (a) Stable, (b) Unstable	8
Figure 3.2.	A typical Shewhart Chart	9
Figure 3.3.	Raw Data and A typical Cusum Chart	12
Figure 3.4.	Typical behavior of the log-likelihood ratio S_k	13
Figure 3.5.	Typical behavior of the CUSUM decision function g_k	13
Figure 3.6.	Repeated use of SPRT	14
Figure 3.7.	Tabular cusum chart	17
Figure 3.8.	The parameters of V-mask	19
Figure 4.1.	ARL Performance of the Tabular Cusum with $k=1/2$ and $h=4$ and $h=5$	23

Figure 4.2.	Values of k and h that give $ARL_0=370$ for Tabular Cusum	23
Figure 4.3.	ARL Performance of the EWMA	23
Figure 4.4.	Shewhart Chart of a sample of data.	25
Figure 4.5.	EWMA Chart of a sample of data.	25
Figure 4.6.	EWMA Chart with different λ values.	26
Figure 5.1.	Data which has normal distribution and its cumulative sum values.	29
Figure 5.2.	Periodic data which has normal distribution and its cumulative sum values.	30
Figure 5.3.	Data which has weibull distribution and its cumulative sum values.	31
Figure 5.4.	Periodic data which has weibull distribution and its cumulative sum values.	32
Figure 5.5.	Bogazici University DDoS attack dataset topology.	33
Figure 5.6.	Distribution of SYN packets in dataset 1.	34
Figure 5.7.	Distribution of SYN packets in dataset 2.	35
Figure 5.8.	Cumsum of SYN packets in dataset 1.	36
Figure 5.9.	Cumsum of SYN packets in dataset 2.	37
Figure 5.10.	Tabular CUSUM Control Chart of number of syn packets in dataset 1.	38

Figure 5.11. Tabular CUSUM Control Chart of number of syn packets in dataset 2.	39
Figure 5.12. Exemplary representation of the V-Mask as an indicator and decision making support	40
Figure 5.13. V-mask cusum applied on cusum of number of syn packets in dataset 2.	40
Figure 5.14. Moving average for smoothing on dataset 2.	41
Figure 5.15. Removed seasonal and trend components of dataset 2.	41
Figure 5.16. Differenced dataset 2.	42
Figure 5.17. Residuals for dataset 2.	42
Figure 5.18. Cusum chart of residuals for dataset 2.	43
Figure 5.19. Ewma chart of residuals for dataset 2.	43
Figure 5.20. ROC analysis of Tabular CUSUM versus V-mask CUSUM.	44

LIST OF SYMBOLS/ABBREVIATIONS

AR	Auto Regressive
ARIMA	Autoregressive Integrated Moving Average
ARL	Average Run Length
ARMA	Auto Regressive Moving Average
CUSUM	Cumulative Sum
DDoS	Distributed Denial of Service attack
DoS	Denial of Service attack
EWMA	Exponentially Weighted Moving Average
FPR	False Positive Rate
HTTP	Hyper-Text Transfer Protocol
ICMP	Internet Control Message Protocol
IP	Internet Protocol
MA	Moving Average
ROC	Receiver Operating Characteristic
SPC	Statistical Process Control
SPRT	Sequential Probability Ratio Test
TCP	Transmission Control Protocol
TPR	True Positive Rate
UDP	User Datagram Protocol

1. INTRODUCTION

Denial of Service (DoS) is a sort of attack that exhausts the sources to make them unavailable for the legitimate users and takes control over remote hosts by sending exceeding number of requests to the sources from invalid return addresses [1].

In a DoS attack, the number of targets and attackers is unrelated. However, Distributed Denial of Service (DDoS) attack pursues having numerous machines each performing a DoS attack towards at least one target in a coordinated manner. Therefore we understand from this definition that DDoS attacks are a subset of DoS attacks. It is worth emphasizing that in a DDoS attack, there must be more than one attacking source which are coordinated.

There are different kinds of DDoS attacks, such as TCP (Transmission Control Protocol) Syn flood attacks, UDP (User Datagram Protocol) flood attacks and ICMP flood attacks. According to general statistics, TCP Syn flood attacks cover the biggest ratio (84.1%) of DDoS attacks as seen from the Figure 1.1 in which distribution of DDoS attack types in the first quarter of 2019 is shown. It is easy to initiate TCP Syn flood attacks and the attacker do not need to send big size of packets. Second most common DDoS attacks are UDP flood attacks. UDP is very common protocol and UDP packets can pass easily from routers and rewalls. From the Figure 1.1 we see that UDP flood attack types are ranked second. In ICMP flood attacks, the attackers send a huge number of ICMP packets of any type. Because of having small size, the ping packets are usually chosen to launch ICMP flood attacks.

In this thesis, we obtained the number of TCP syn packets from two datasets which includes TCP Syn flood attack of DDoS traffic. We obtained CUSUM values of the number of syn packets. In order to evaluate Tabular CUSUM and V-mask CUSUM, we applied these CUSUM methods on two datasets. Our aim was to find which method, Tabular CUSUM or V-mask CUSUM, is the best technique to detect TCP Syn flood attack and which method has the most accurate results. We analyzed

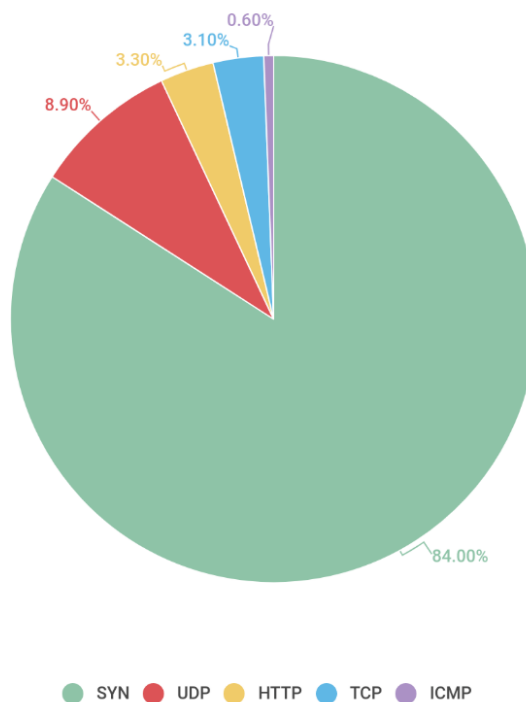


Figure 1.1. Distribution of DDoS attack types in 2019 Q1 [2].

the ROC curve statistics of Tabular CUSUM and V-mask CUSUM and plot the ROC curves in order to evaluate accuracy. We wanted to find the simple detection scheme for TCP Syn flood attack traffics. Therefore, we made a performance analysis of EWMA and CUSUM control charts evaluating Average Run Length (ARL) approximation. Finally, the Autoregressive Integrated Moving Average (ARIMA) forecasting model was applied in order to obtain the forecasting residuals which are also utilized in the performance evaluation of these two control charts. The organization of this thesis is as follow: in the next chapter, DDoS attacks are introduced. A brief history of DDoS attacks is given with some motivations behind DDoS attacks and up-to-date information about these types of attacks. In the third chapter, the statistical process control charts to detect anomalies in unstable processes are discussed. In the fourth chapter, the performance analysis of EWMA and CUSUM control charts is made. In the fifth chapter, the results are discussed. Finally, the last chapter is conclusion which ends this thesis.

2. DISTRIBUTED DENIAL OF SERVICE ATTACKS

2.1. A Brief History of DDoS Attacks

DDoS attacks have expanded considerations on attacks from both system administrators and business in the last thirty years. On the other hand, the system vulnerabilities that cause many types of attacks have been perceived since the beginning of the commercial web.

In the late 1990s, companies were started to worry about denial of service. Groups who were in conflict with social contradictions started utilizing it as a political purpose. In 1998, they used Flood Net tool to make the Frankfurt Stock Exchange and the Pentagon sites unreachable by citizens.

After DoS attacks targeted the websites of eBay, Yahoo, Amazon.com and CNN making them unreachable, denial of service gained new significance in 2000. In order to send spam emails or launch different kinds of attacks to the systems, attackers started to utilize computers controlled by botnets.

The websites of Amazon and Yahoo had a financial loss with an amount of \$1.7bn because of attacks in 2000. Estonia experienced a DDoS attack from Conficker botnet in 2007. The attack was very successful, it totally cut off from the Web for over ten days and the financial institutions of Estonia were compromised. On the other hand, there were DDoS attacks that could not be attributed to any valid cause. Sony and Microsoft gaming services were crashed by a hacking group on Christmas in 2014. One of the member explained the reason behind this attack as a laugh [3].

To launch a DDoS attack is not a complicated process. You can find detailed information on the internet even if you have limited technical ability. Even in order to retain financial power, companies are paying some people to launch DDoS attacks on their competitors. Thus, they cause damage to the infrastructure of competing

companies and indirectly their financial resources [4].

When this is the case, the world of information technology began to look for ways to cope with these types of attacks. Moreover, the damages of DDoS attack mostly gain visibility when attack impacts splash on the economic dimension. However, insufficient and weak technologies and limited experience of system administrators against DDoS attacks has increased attackers' tendency to launch different types of attacks. Especially, smaller sites have much more difficulties on fending off DDoS attacks.

There are many techniques put forwarded to deal with denial of service attacks. Finding traces of the attackers by IP address is an example. Having botnets that exploit vulnerabilities on computers with randomly distributed IP addresses results in difficulties to distinguish the attackers' IP addresses [4].

On the other hand, organizations' security awareness is increasing against DDoS attacks. They have been developing response techniques by establishing incident response, forensic and malware analysis teams responsible for monitoring and analyzing the cyber security situation of an organization to identify security issues.

2.2. Geography of DDoS Attacks

Everyday, thousands of DDoS attacks are happening around the globe. The distribution of DDoS attacks by country in the last quarter of 2018 and the first quarter of 2019 is shown in Figure 2.1.

Kaspersky anti-virus company stated that China is the first target of all DDoS attacks worldwide with a share of 67.89%, and then USA comes second. In third spot, there is Hong Kong with a share of 4.81%. The entire results can be seen in Figure [2].

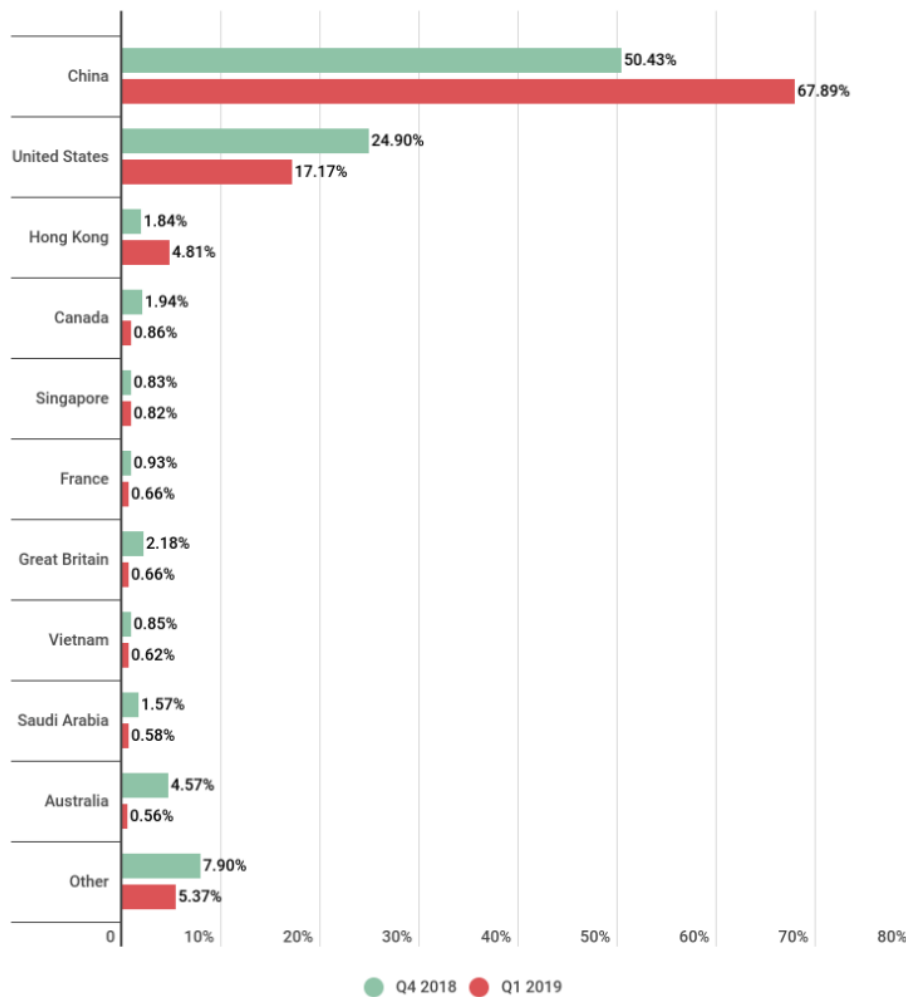


Figure 2.1. DDoS attacks distributions by country in 2018 Q4 and 2019 Q1 [2].

2.3. Dynamics of the number of DDoS attacks

When we scan the Figure 2.2, the number of DDoS attacks reaches the highest level on March 16 with 699 attacks and an important flood is occurred on January 17 with 532 attacks. The entire results can be seen from the Figure 2.2.

We see from the Figure 2.3 that, in the first quarter of 2019 DDoS activity shifted to the weekend compared to the last quarter of 2018. Saturday is the busiest day and Friday is in second spot.

For both Figure 2.2 and 2.3, these time dependent analysis are important for making predictions and taking preventive actions against DDoS attacks.

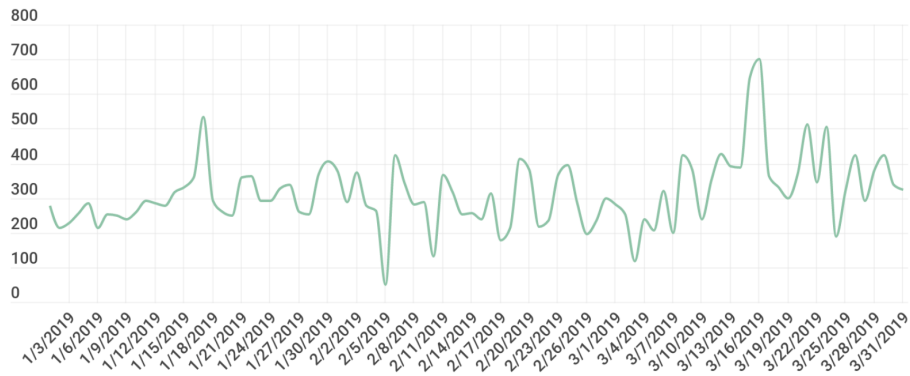


Figure 2.2. The number of DDoS attacks in 2019 Q1 [2].

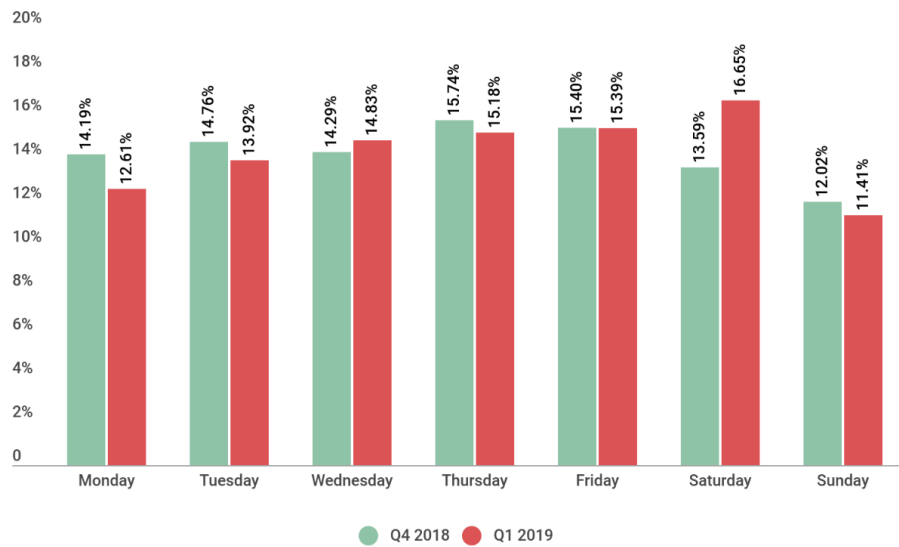


Figure 2.3. DDoS attacks distributions by days in 2018 Q4 and 2019 Q1 [2].

2.4. Types and duration of DDoS attacks

In the first quarter of 2019, duration of the longest attack is approximately 12 days and 289 hours. This was almost 14 days in the last quarter of 2018. As seen from the Figure 2.4, most of the long-duration attacks show an upward trend in the first quarter of 2019 compared to the last quarter of 2018.

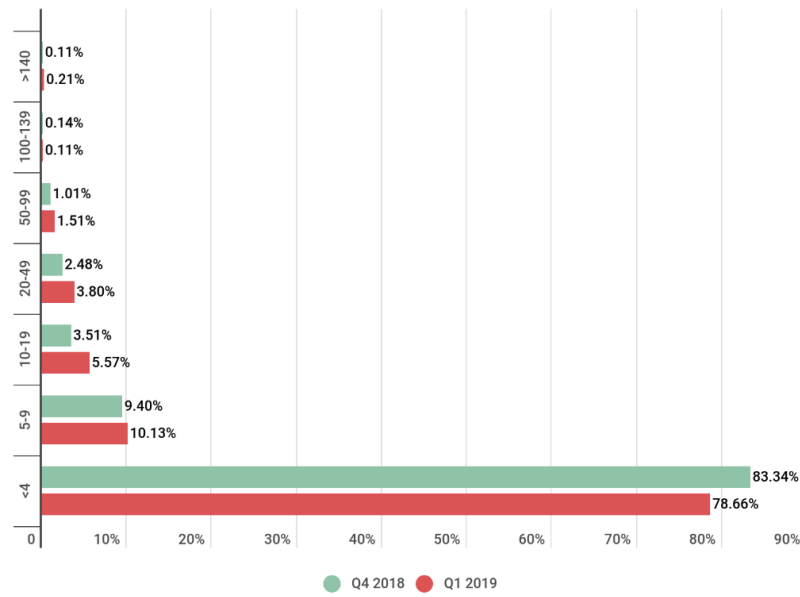


Figure 2.4. DDoS attacks distribution by time in hours in 2018 Q4 and 2019 Q1 [2].

3. STATISTICAL PROCESS CONTROL CHARTS

There are several prevention techniques which researchers put forward to deal with DDoS attacks. These are reactive techniques, proactive techniques and survival techniques [5].

In reactive techniques, when a DDoS attack is occurred, then forensic analysis is made to identify the sources of attacks. In proactive techniques, the purpose is to detect an attack before it's occurrence. Then incident response actions are taken to limit attack traffic. Lastly, in survival techniques, the system which is under DDoS attack has enough resources such as CPU power, memory etc. to maintain its services for its customers [5].

Statistical process control (SPC) is a method for observing and taking control the quality of processes in order to guarantee that these processes work at predetermined levels. SPC has been usually applied to control the manufacturing processes, however it can be applied to any process whose specific features can be extracted [6]. It can therefore be used to detect anomalies in the processes under DDoS attacks.

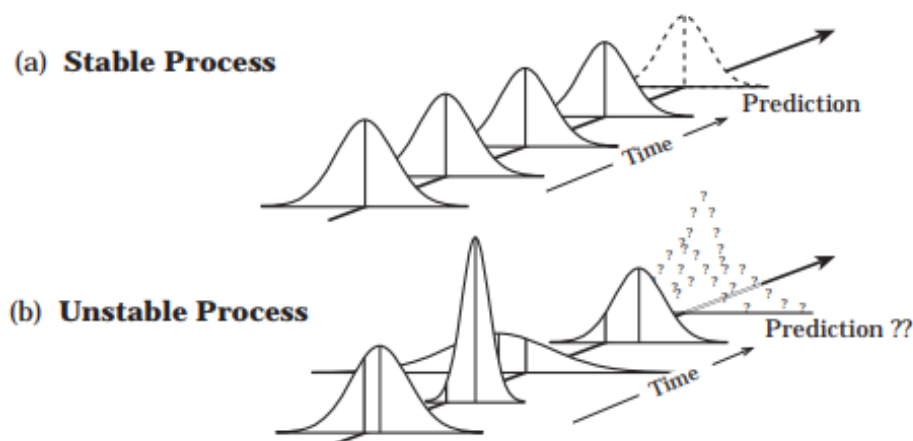


Figure 3.1. Typical behaviors of the processes, (a) Stable, (b) Unstable [7].

Stability is an indicator of prediction on our product in the future, therefore it is an important topic for statistical process control. If we start with such a stable

process shown in (a) in Figure 3.1, control charts signal a change. If there is such an unstable process shown in (b) in Figure 3.1, there is a quality problem and process produces wrong products and we cannot make prediction on our products. Then we seek to eliminate the variability of the process. Such a process is considered to be out of control. Therefore we need some criteria to decide whether the process is starting to go out of control or not. With this approach, control charts which are Shewhart, EWMA and CUSUM control charts are discussed in following sections can be used to detect DDoS attack traffic.

3.1. SHEWHART CONTROL CHART

In normal distributions, if the observations fall within 3 standard deviations from the mean, i.e. between $\mu \pm 3\sigma$, then the process is in control. If observations are started to show up outside of the control limits then the process is said to be out of control.

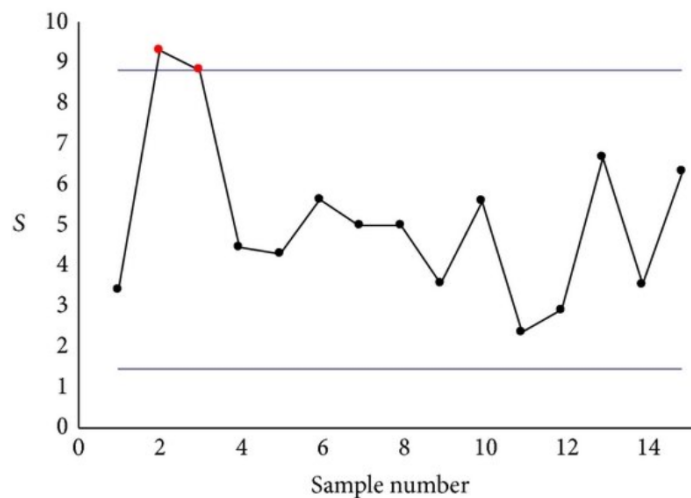


Figure 3.2. A typical Shewhart Chart [8].

For the process observations which have normal distribution and which are uncorrelated, we use 3σ for the control limits' distance. The control limits have an average run length (ARL) of 370. The ARL is calculated as $ARL = \frac{1}{\alpha}$, where α is the probability of any point which appears outside control limits when the process is in control. For the Shewhart control chart, the probability is expressed as:

$$p = Pr(|X - \mu| \geq 3\sigma) = 0.0027,$$

where p is the probability of a point which falls outside control limits when the process is in control. X has a normal distribution with mean μ and standard deviation σ . Then we can calculate the average run length of the Shewhart control chart as:

$$ARL = \frac{1}{\alpha} = \frac{1}{0.0027} = 370.$$

3.2. EWMA CONTROL CHART

The Exponentially Weighted Moving Average (EWMA) control chart consists of plotting a weighted average of measurements, giving heaviest weights to the most recent observations.

EWMA control chart has the statistic z_i at time i defined as:

$$z_i = (1 - \lambda)z_{i-1} + \lambda x_i \quad (3.1)$$

for $i = 1, 2, 3, \dots$, where x_i indicates the observations and the parameter λ , $0 \leq \lambda \leq 1$, is constant. The initial value of statistic z_i is set to the process mean, so that

$$z_0 = \mu_0.$$

To demonstrate that the EWMA control chart statistic z_i , we may substitute z_{i-1} on the right side of equation 3.1 to obtain

$$z_i = \lambda x_i + (1 - \lambda)[\lambda x_{i-1} + (1 - \lambda)z_{i-2}] = \lambda x_i + (1 - \lambda)\lambda x_{i-1} + (1 - \lambda)^2 z_{i-2}.$$

Continuing to substitute for z_{i-j} , $j=2,3,\dots,t$, we obtain

$$z_i = \lambda \sum_{j=0}^{t-1} (1 - \lambda)^j (x_{t-j}) + (1 - \lambda)^t z_0. \quad (3.2)$$

When we have smaller value of λ , there is heavier reliance on past data and we detect small shifts in the process more quickly. If $\lambda = 1$, the EWMA control chart is equivalent to the Shewhart control chart. The comparison of EWMA control charts with different λ values shown in Figure 4.6.

The EWMA chart signals when z_i is outside the control limits. The upper control limit (UCL) and lower control limit (LCL) can be expressed as

$$UCL = \mu_0 + L\sigma_x \sqrt{\frac{\lambda}{2 - \lambda}}$$

and

$$LCL = \mu_0 - L\sigma_x \sqrt{\frac{\lambda}{2 - \lambda}}$$

where L is a positive coefficient. λ with corresponding value of L determines the performance of the EWMA control chart.

The standard deviation of z_i , denoted by σ_{z_i} , is equal to

$$\sigma_{z_i} = \sigma_x \sqrt{\frac{\lambda}{2 - \lambda} [1 - (1 - \lambda)^{2i}]}$$

where σ_x denotes the standard deviation of x_i . The standard deviation z_i is time-varying and it converges, as i increases to

$$\sigma_{z_i} = \sigma_x \sqrt{\frac{\lambda}{2 - \lambda}}.$$

3.3. CUSUM CONTROL CHART

Cumulative sum control chart improves the ability to detect small shifts in the mean of a process at given times by plotting a statistic that uses current and previous

sample values from the process [9].

The simple CUSUM is defined as follows [10]:

$$C_n = \sum_{j=1}^n (\bar{x}_j - \mu_0) \quad (3.3)$$

where \bar{x}_j is the average of the j th sample, μ_0 is the in-control mean or target, n is the sample number. The difference between the \bar{x}_j and μ_0 is summed from $j = 1$ to $j = n$. C_n is called the cumulative sum of the n th observation. The value of C_n is then plotted against n and this forms the CUSUM chart [10] shown Figure 3.3.

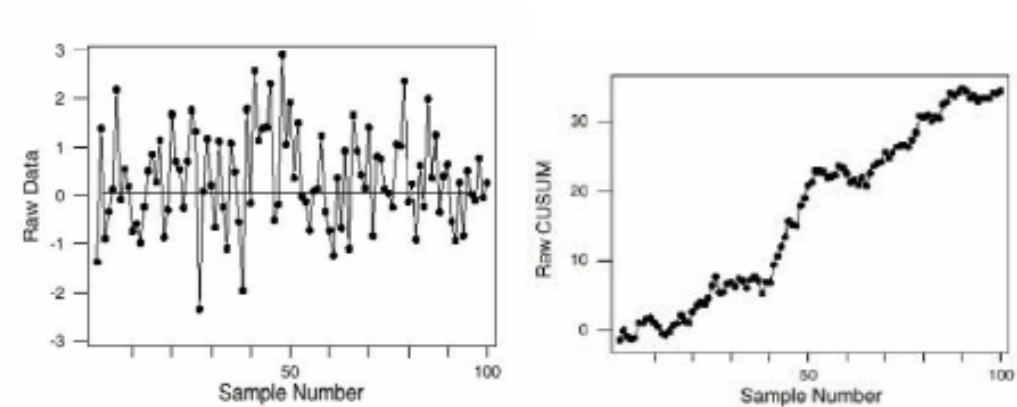


Figure 3.3. Raw Data and A typical Cusum Chart [11].

If the process remains in control about a target value μ_0 then the mean of cumulative is zero. But, if a shift from μ_0 to some value

$$\mu_1 = \mu_0 + k$$

occurs, then an upward drift is said to be present in the cumulative sum. Likewise if a shift from μ_0 to some value

$$\mu_1 = \mu_0 - k$$

occurs, then we say that a downward drift will be present in the cumulative sum. Any trend upward or downward is a possible sign of an out-of-control system.

3.3.1. Derivations of CUSUM

3.3.1.1. Intuitive Derivation. The typical behavior of the log-likelihood ratio S_k shows a negative drift before change, and a positive drift after change. The relevant infor-

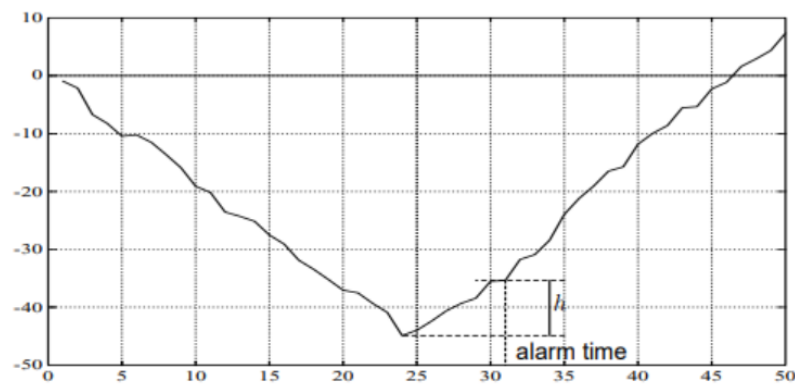


Figure 3.4. Typical behavior of the log-likelihood ratio S_k [12].

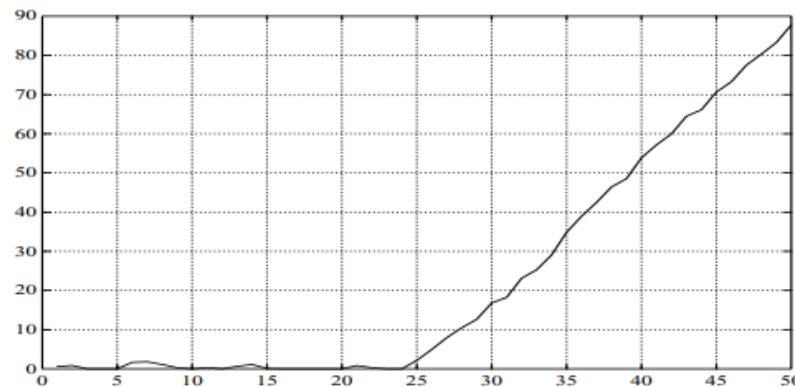


Figure 3.5. Typical behavior of the CUSUM decision function g_k [12].

mation lies in the difference between the value of S_k and its current minimum value m_k . The corresponding decision rule is as follows :

$$g_k = S_k - m_k \geq h \quad (3.4)$$

where

$$m_k = \min_{1 \leq j \leq k} S_j$$

and the log-likelihood ratio is

$$S_k = \sum_{i=1}^k (s_i)$$

where

$$s_i = \ln \frac{p_{\theta_1}(y_i)}{p_{\theta_0}(y_i)}.$$

The typical behavior of g_k is shown in Figure 3.5. The stopping time is

$$t_a = \min\{k : g_k \geq h\},$$

which can be obviously rewritten as

$$t_a = \min\{k : S_k \geq m_k + h\}.$$

The threshold $m_k + h$ keeps the total memory of the whole data contained in the past observations.

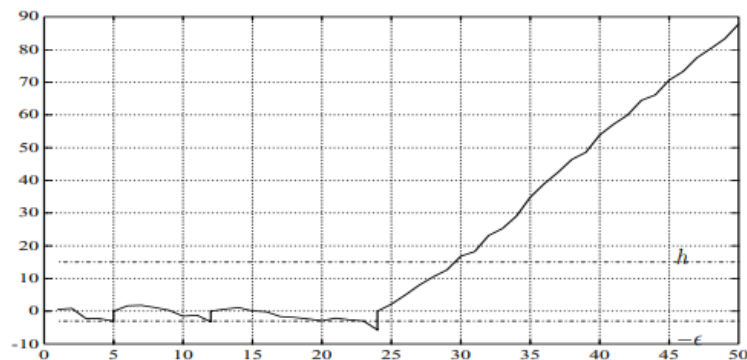


Figure 3.6. Repeated use of SPRT [12].

3.3.1.2. Repeated Sequential Probability Ratio Test Derivation. It is the use of repeated testing of the two simple hypotheses:

$$H_0 : \theta = \theta_0, \quad (3.5)$$

$$H_1 : \theta = \theta_1. \quad (3.6)$$

The Sequential Probability Ratio Test (SPRT) is defined based on the pair d, T where d is the decision rule and T is a stopping time which is the time at which the final decision is taken. The definition of the SPRT is thus

$$d = \begin{cases} 0 & S_1^T \leq -\epsilon, \\ 1 & S_1^T \geq h \end{cases} \quad (3.7)$$

where T is the exit time:

$$T = T_{-\epsilon, h} = \min\{k : S_1^k \geq h \quad \cup \quad S_1^k \leq -\epsilon\} \quad (3.8)$$

where $\epsilon \geq 0$ and $h > 0$ are conveniently chosen thresholds. The typical behavior of this SPRT is shown in Figure 3.6. The first time at which $d = 1$ is the alarm time at which a change is detected.

Starting from the repeated SPRT with the value of lower threshold ϵ , the resulting decision rule can be rewritten as

$$g_k = \begin{cases} g_{k-1} + \ln \frac{p_{\theta_1}(y_k)}{p_{\theta_0}(y_k)}, & g_{k-1} + \ln \frac{p_{\theta_1}(y_k)}{p_{\theta_0}(y_k)} > 0 \\ 0, & \ln \frac{p_{\theta_1}(y_k)}{p_{\theta_0}(y_k)} \leq 0 \end{cases} \quad (3.9)$$

where $g_0 = 0$. Remembering the definition of s_k ,

$$g_k = (g_{k-1} + s_k)^+ \quad (3.10)$$

where $(x)^+ = \sup(0, x)$. Finally, the stopping rule and alarm time are defined by

$$t_a = \min\{k : g_k \geq h\} \quad (3.11)$$

where g_k is given in 4.15. On the other hand, it can also be written as:

$$g_k = (S_{k-N_k+1}^k)^+ \quad (3.12)$$

where

$$N_k = N_{k-1} \cdot 1_{\{g_{k-1} > 0\}} + 1. \quad (3.13)$$

$1_{\{x\}}$ is the indicator of event x , and t_a is defined in 3.11. N_k is the number of observations after restart of the SPRT.

3.3.2. Tabular Form of CUSUM

The tabular CUSUM sum up deviations from μ_0 that are above and below target statistics \bar{C}_n^+ and \bar{C}_n^- respectively. These statistics are [13];

$$\bar{C}_n^+ = \max[0, \bar{x}_n - (\mu_0 + K) + \bar{C}_{n-1}^+], \quad (3.14)$$

$$\bar{C}_n^- = \max[0, \bar{x}_n - (\mu_0 - K) - \bar{C}_{n-1}^-] \quad (3.15)$$

with \bar{C}_n^+ and \bar{C}_n^- initialized to 0. The value K is referred to as the reference value. The statistic \bar{C}_n^+ and \bar{C}_n^- are called one sided lower and one sided upper CUSUM respectively. The K value which is often chosen halfway between the target value μ_0 and an out-of-control value μ_1 , and it is expressed as

$$K = \frac{|\mu_1 - \mu_0|}{2}, \quad (3.16)$$

if the shift is expressed in standard deviations, such as $\mu_1 = \mu_0 + \delta\sigma$ then K can be expressed as

$$K = \frac{\delta\sigma}{2}. \quad (3.17)$$

The system is considered out of control anytime either of the statistics \bar{C}_n^+ or \bar{C}_n^- exceeds the control limits of H , where H is chosen to be 5σ where σ is the standard deviation of sample mean.

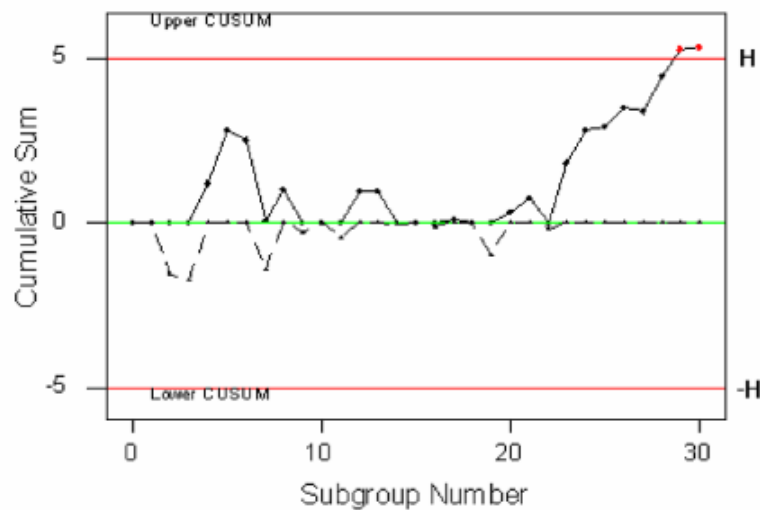


Figure 3.7. Tabular cusum chart [14].

\bar{N}^+ specifies the number of consecutive times $\bar{C}_n^+ > 0$ and \bar{N}^- specifies the number of consecutive times $\bar{C}_n^- < 0$. Then the system is said to have gone out of control at

time N , where

$$N = \min\{N^-, N^+\}. \quad (3.18)$$

If the system exceeds the control limit at some period n then the process probably went out of control at period $n - N$.

In situations where process is out-of-control, in order to make process in control, it will be useful to have an estimate of the new process mean as:

$$\hat{\mu}_0 = \begin{cases} \mu_0 + K + \frac{\bar{C}_n^+}{N^+}, & \bar{C}_n^+ > H, \\ \mu_0 - K + \frac{\bar{C}_n^-}{N^-}, & \bar{C}_n^- < -H. \end{cases} \quad (3.19)$$

3.3.3. V-mask Form of CUSUM

There is another method known as V-mask is at time used to determine if a process is out of control [15]. The performance of V-mask rely on the definition of some design parameters estimated. We must choose these parameters carefully so that there are very few false alarms occur when the process is in control. The V-mask design parameters are as follows [16]: μ_0 is the target value of the process, σ is the standard deviation of the process, n is subgroup size, α is the probability of indicating a process out of control although the process is in control in reality, β is the probability of indicating a process in control although the process is out-of-control in reality, δ is the amount of shift we want to detect, k is the slope of the V-mask arm corresponding to one sampling unit, d is the distance from the origin to the vertex of the mask, h is the vertical distance in the arm corresponding. These parameters are indicated in Figure 3.8.

The V-mask is applied to values of the CUSUM statistic

$$C_n = \sum_{j=1}^n y_j = y_j + C_{n-1} \quad (3.20)$$

where y_j is the standardized observation $y_j = \frac{x_j - \mu_0}{\sigma}$.

In the V-mask procedure, we place V-mask on the cusum control chart with the point O on the last value of C_n . If all the previous cumulative sums, C_1, C_2, \dots, C_n are in between the upper and lower arms of the V-mask, the process is in control. Otherwise, if they are not in between the arms of V-mask, the process is out of control which is shown in Figure 3.8. Practically, the V-mask would be placed each new point on the cusum chart as soon as it was plotted [17].

The performance of the V-mask is determined by the distance d and the angle θ which are shown in 3.8. The tabular cusum and the V-mask scheme are equal if

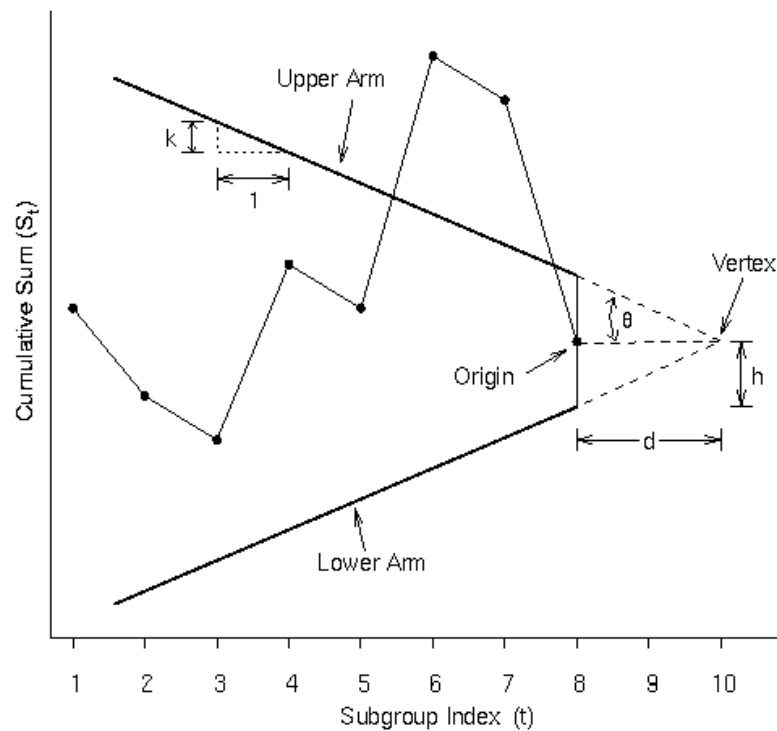


Figure 3.8. The parameters of V-mask [16].

$$k = \tan \theta = \frac{\delta\sigma}{2}$$

and

$$h = dk.$$

The recommended parameters of the V-mask parameters [18]:

$$\theta = \arctan \frac{\delta\sigma}{2}$$

and

$$d = \frac{2}{\delta^2} \ln \frac{1 - \beta}{\alpha}$$

where α is the probability of indicating a process out of control although the process is in control in fact, β is the probability of indicating a process in control although the process is out-of-control in fact. If β is small, then

$$d \approx -2 \frac{\ln \alpha}{\delta}.$$

The details of probabilities α , β and the shift δ will be examined in the next section.

If the distribution of x is not normal, the Chebyshevs inequality shows us that without relevance of the distribution, at least 89% of observations fall within the 3-sigma limits. We can use Chebyshevs inequality theorem to find out what percentage of the data is gathered around the mean.

3.3.3.1. Definition of Chebyshevs Inequality. Chebyshevs inequality states that the probability that the outcome of an experiment with the random variable x will fall more than standard deviations beyond the mean of x is less than $1/k^2$, which is expressed as

$$P(|x - \mu| > k\sigma) < \frac{1}{k^2}. \quad (3.21)$$

Proof. Let S be the sample space for a random variable, x , and let $f_x(x)$ stands for the pdf of x . Let R_1, R_2, R_3 be partition S , such that for every sample point $x \in S$

$$x \in \begin{cases} R_1, x < \mu - k\sigma, \\ R_2, |x - \mu| \leq k\sigma, \\ R_3, x > \mu + k\sigma, \end{cases} \quad (3.22)$$

$$\sigma^2 = \sum_{R_1} (x - \mu)^2 f_x(x) + \sum_{R_2} (x - \mu)^2 f_x(x) + \sum_{R_3} (x - \mu)^2 f_x(x). \quad (3.23)$$

Clearly,

$$\sigma^2 \geq \sum_{R_1} (x - \mu)^2 f_x(x) + \sum_{R_3} (x - \mu)^2 f_x(x) \quad (3.24)$$

since the term that evaluates to the variance in R_2 has been subtracted on the right-hand side. For any sample point $x \in R_1$

$$x < \mu - k\sigma$$

which implies

$$= (x - \mu)^2 < k^2 \sigma^2$$

therefore, for any sample point $x \in R_1$ or $x \in R_3$, it can be said that $(x - \mu)^2 < k^2\sigma^2$, thus

$$\sigma^2 \geq \sum_{R_1} k^2\sigma^2 f_x(x) + \sum_{R_3} k^2\sigma^2 f_x(x). \quad (3.25)$$

Dividing each side of the inequality by $k^2\sigma^2$ results in

$$\frac{1}{k^2} \geq \sum_{R_1} f_x(x) + \sum_{R_3} f_x(x),$$

or, in other terms

$$\frac{1}{k^2} > P(|x - \mu| > k\sigma). \quad (3.26)$$

□

We usually want to choose k as $k = 0.5\delta$, where δ is the size of the shift. This approach is very close to minimizing the ARL_1 while maximizing ARL_0 values which will be examined in details in the next section.

4. CONTROL CHART PERFORMANCE

4.1. Average Run Length Approximation

The performance of a control chart is measured by the α error and β error, the same as evaluating any hypothesis test. Similarly, a control chart is actually a series of sample-by sample hypothesis tests, in order to measure the performance of a chart more directly, average run length ARL is used as a performance measure of control chart.

Shift in Mean (multiple of σ)	$h = 4$	$h = 5$
0	168	465
0.25	74.2	139
0.50	26.6	38.0
0.75	13.3	17.0
1.00	8.38	10.4
1.50	4.75	5.75
2.00	3.34	4.01
2.50	2.62	3.11
3.00	2.19	2.57
4.00	1.71	2.01

Figure 4.1. ARL Performance of the Tabular Cusum with $k=1/2$ and $h=4$ and $h=5$ [11].

k	0.25	0.5	0.75	1.0	1.25	1.5
h	8.01	4.77	3.34	2.52	1.99	1.61

Figure 4.2. Values of k and h that give $ARL_0=370$ for Tabular Cusum [11].

Shift in mean (multiple of σ)	Average Run Length (ARL)				
	$L=3.054$ $\lambda=0.40$	$L=2.998$ $\lambda=0.25$	$L=2.962$ $\lambda=0.20$	$L=2.814$ $\lambda=0.10$	$L=2.615$ $\lambda=0.05$
0	500	500	500	500	500
0.25	224	170	150	106	84.1
0.50	71.2	48.2	41.8	31.3	28.8
0.75	28.4	20.1	20.1	15.9	16.4
1.00	14.3	11.1	11.1	10.3	11.4
1.50	5.9	5.5	5.5	6.1	7.1
2.00	3.5	3.6	3.6	4.4	5.2
2.50	2.5	2.7	2.7	3.4	4.2
3.00	2.0	2.3	2.3	2.9	3.5
4.00	1.4	1.7	1.7	2.2	2.7

Figure 4.3. ARL Performance of the EWMA [11].

When a value of data is outside of the control limits, control chart signals on the sample number and that is Run length (RL). Its expected value is called average run length.

If the process is in fact in control and control chart signals on the sample number, that is $ARL(ARL_0)$, in other words, in-control ARL . If the process is in fact out-of-control and control chart signals on the sample number, that is $ARL(ARL_1)$, in other words, out-of-control ARL . We would like ARL_0 to be as large as possible and ARL_1 to be as small as possible.

In control $ARL(ARL_0)$ is expressed as

$$ARL_0 = E(RL) = \sum_{k=1}^{\infty} kPR(RL = k) = \sum_{k=1}^{\infty} k(1 - \alpha)^{k-1}\alpha \quad (4.1)$$

$$ARL_0 = \frac{1}{\alpha}, \quad (4.2)$$

where α is the probability of error. Out of control $ARL(ARL_1)$ is expressed as

$$ARL_1 = E(RL) = \sum_{k=1}^{\infty} kPR(RL = k) = \sum_{k=1}^{\infty} k(\beta)^{k-1}(1 - \beta), \quad (4.3)$$

where $(1 - \beta)$ is the probability of error. It is worth saying that when sample size n increases, both α error and β error decrease, so ARL_0 is getting larger while ARL_1 is getting smaller.

A Shewhart chart does not react to little changes rapidly it is because it has not memory, specifically that a Shewhart chart consistently utilize the most recent update of process measurements as opposed to using the data in the whole sequence of measurements.

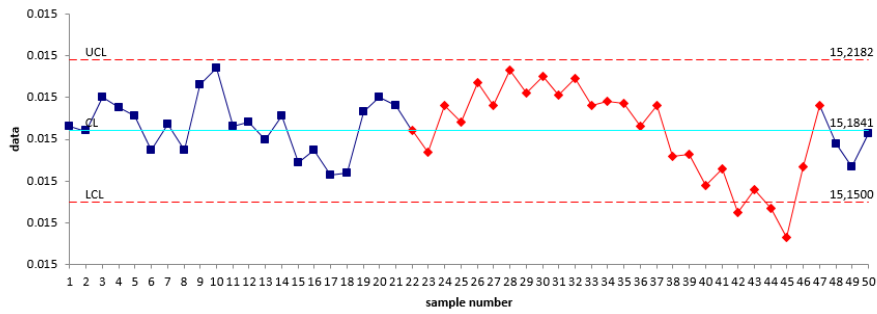


Figure 4.4. Shewhart Chart of a sample of data.

By considering previous samples, a chart can accomplish quicker detection of small mean shifts. Fundamentally, giving a control chart memory is like expanding its sample size.

Design of an EWMA chart based on choosing λ and L parameters. λ controls how much memory the chart has. When λ gets smaller, the chart has longer memory and is good for detecting small changes. When λ gets larger the chart has shorter memory and is good for detecting large changes. ARL performance of EWMA control chart is shown in Figure 4.3.

Figure 4.6 indicates how Ewma chart signals when we change the values of λ . As seen from the Figure 4.6, When λ gets smaller value as 0,02, the chart becomes good for detecting small changes. When λ gets larger value as 1, the chart becomes good for detecting large changes.

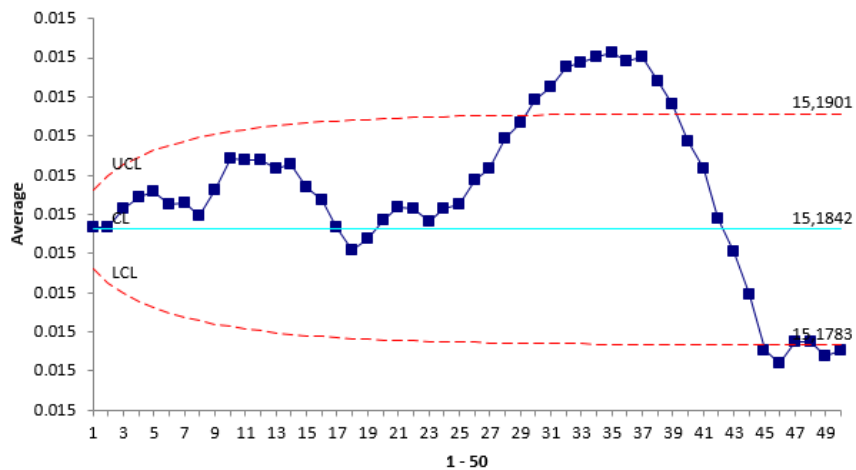


Figure 4.5. EWMA Chart of a sample of data.

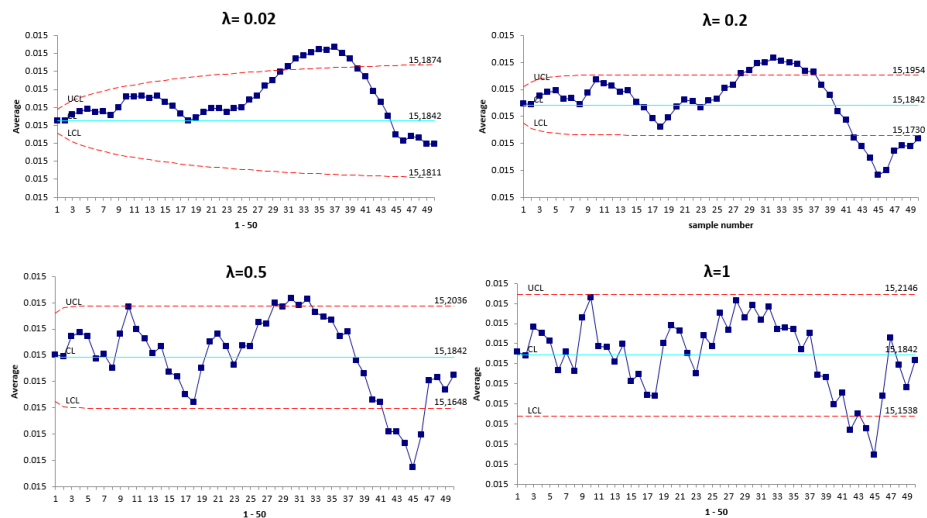


Figure 4.6. EWMA Chart with different λ values.

Design of a CUSUM chart is based on choosing K and H parameters. Selecting appropriate K value and corresponding H value actually means making a proper choice of average run length. Remembering that we would like ARL_0 to be as large as possible and ARL_1 to be as small as possible. The parameters K and H are expressed in terms of σ ,

$$K = k.\sigma, \quad (4.4)$$

$$H = h.\sigma. \quad (4.5)$$

As seen from the equations, choosing K and H means choosing k and h . The combination of $k=0.5$ and $h=5$ produces a good performance and ARL properties.

An approximation formula is developed for calculating the average run length for the CUSUM method. For one sided CUSUM, ARL [19] is given by:

$$ARL^\pm = \begin{cases} \frac{\exp(-2\Delta b) + 2\Delta b}{2\Delta^2}, & \text{if } \Delta \neq 0 \\ b^2, & \text{if } \Delta = 0 \end{cases} \quad (4.6)$$

where

$$\Delta = \begin{cases} \delta^* - k, \text{ for } C_i^+ \\ -\delta^* - k, \text{ for } C_i^- \end{cases} \quad (4.7)$$

$$b = h + 1.166 \quad (4.8)$$

$$\delta^* = \frac{\mu - \mu_0}{\sigma}. \quad (4.9)$$

For two sided CUSUM, ARL is given by:

$$\frac{1}{ARL} = \frac{1}{ARL^+} + \frac{1}{ARL^-}. \quad (4.10)$$

When $\delta^* = 0$, we get ARL_0 . When $\delta^* \neq 0$, we get ARL_1 . In Figure 4.1, $h = 4$ results in an in-control $ARL_0 = 168$ samples and $h = 5$ results in $ARL_0 = 465$ samples. If we choose $h = 4.77$, we have the value $ARL_0 = 370$. We see that it is the same ARL_0 value for a Shewhart control chart with 3σ limits [20].

A shift, δ , would be detected in either 8.38 samples (with $k = 1/2$ and $h = 4$) or 10.4 samples (with $k = 1/2$ and $h = 5$) in a Tabular CUSUM control chart. By comparison, a Shewhart control chart would require 43.96 samples to detect this shift [11].

4.2. Forecasting Models Approximation

There is another research which analyzes the performance of EWMA and CUSUM control charts applied to the residuals of forecasting models processing two datasets.

In [6], firstly, the forecasting residuals between the forecasts and the observed traffic are obtained by using classical forecasting models. Monitoring forecasting errors, the comparison of the performance of the selected control charts is made. Forecasting models accuracy is calculated based on the mean absolute percentage error (MAPE).

We applied the Autoregressive Integrated Moving Average (ARIMA) forecasting model to obtain the forecasting residuals which are important components to make comparison of the performance of CUSUM and EWMA control charts.

5. EXPERIMENTAL RESULTS

5.1. Cumulative Sum with Normal Distribution

We created datasets which have normal distribution and weibull distribution in order to implement cumulative summation algorithms which we explained in fourth section.

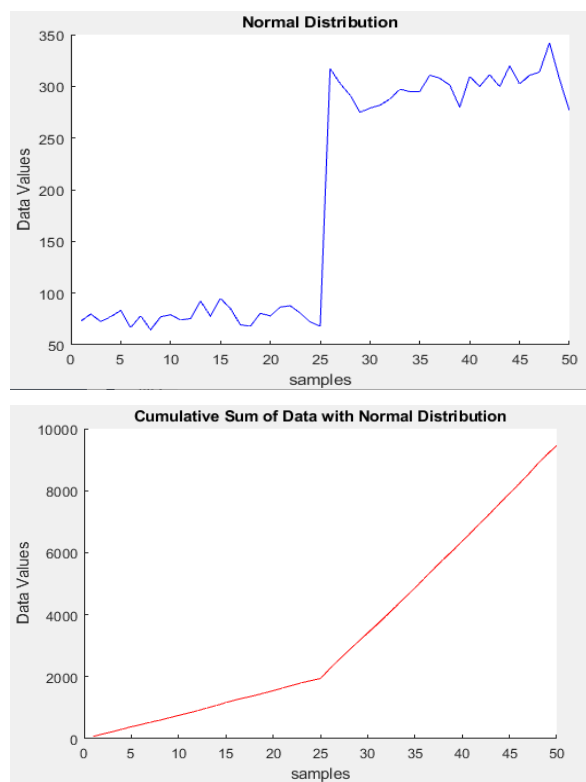


Figure 5.1. Data which has normal distribution and its cumulative sum values.

Figure 5.1 shows the distribution of data and cumulative sum values of this data. The cumulative sum graph indicates that there is a positive drift at the 25th sample when we increased the mean value of the data at the 25th sample which means that the graph shows typical behavior of cusum decision function g_k in Figure 3.5.

Then, by making this dataset periodic with same mean and standard deviation, we saw that cumulative summation graph pattern did not change shown in Figure 5.2.

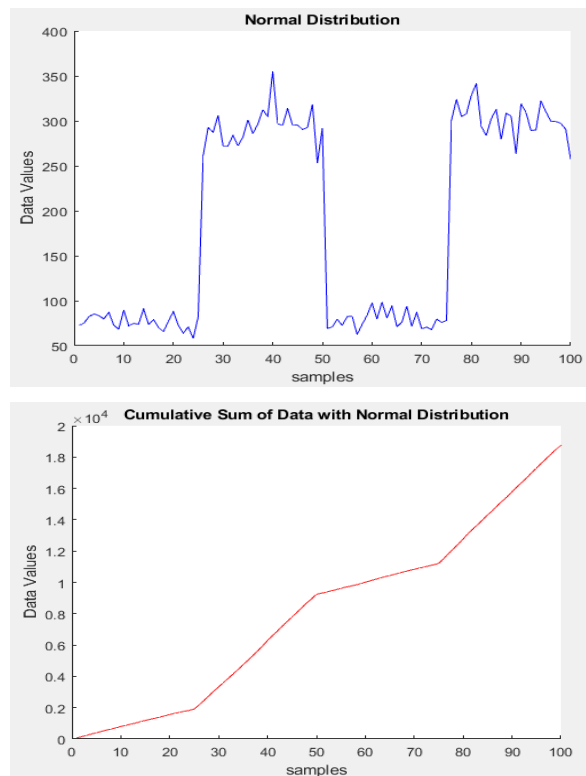


Figure 5.2. Periodic data which has normal distribution and its cumulative sum values.

5.2. Cumulative Sum with Weibull Distribution

In Figure 5.3, cusum graph indicated a positive drift at the 25th sample when we increased the mean value of the data at the 25th sample which also shows a typical behavior of cusum decision function g_k in Figure 3.5.

Then, by making this dataset periodic with same mean and standard deviation, we saw that cumulative summation graph pattern did not change shown in Figure 5.4.

5.3. Applying Tabular CUSUM Method

We applied tabular CUSUM and V-mask CUSUM methods on two datasets. We had 2 TCP syn flood attack datasets which are dataset 1 and dataset 2. We maintained dataset 2 at Bogazici University by using hping DDoS tool. In dataset 2, the attack launched from electrical-electronic engineering department in North Campus. The

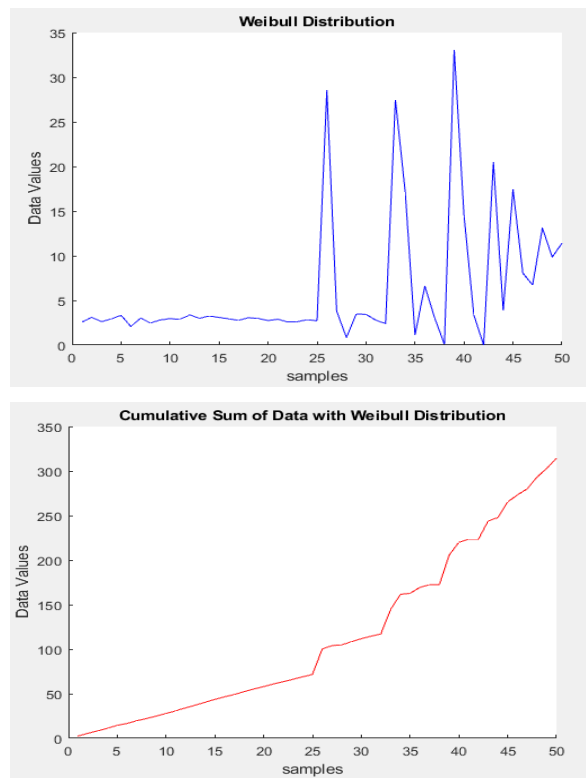


Figure 5.3. Data which has weibull distribution and its cumulative sum values.

victim server was located in Information Technologies office in South Campus and the data trace was recorded between LAN switch and backbone switch, which has over 4000 host connections. Average packet per second was around 14000. The topology of TCP Syn flood attack was given in Figure 5.5.

In dataset 1, approximately 350 seconds attack period applied with 3797 seconds waiting time. The dataset 1 is created 9330 seconds in total. In dataset 2, 20 seconds attack period is applied with 80 seconds waiting period. The dataset 2 is created to be 480 seconds in total. There were different packet rates in each attack period in both datasets. The size of dataset 1 is around 10 mb and the size of dataset 2 is around 8gb. Figure 5.6 shows the distribution of SYN packets in dataset 1. The distribution of the SYN packets of dataset 2 shown in Figure 5.7.

Cumulative sum values of number of syn packets in dataset 1 and dataset 2 are shown in Figure 5.8 and Figure 5.9, respectively.

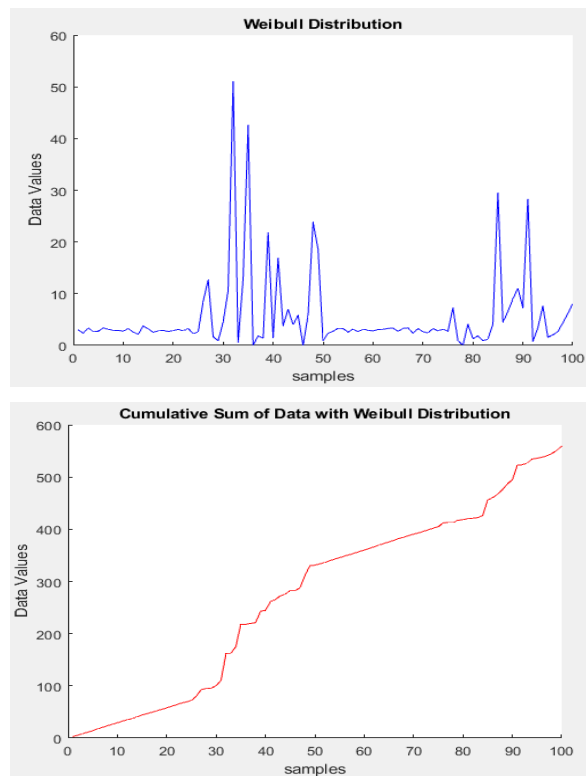


Figure 5.4. Periodic data which has weibull distribution and its cumulative sum values.

The datasets are maintained in 9330 seconds and 480 seconds with time interval $dt=0.000001$ and the number of rows are 76199 and 6028975 in whole syn flood data packets in dataset 1 and dataset 2 respectively. Due to the performance problem of the code that we read data from, we had to set $dt=0.1$ on both datasets. We obtained tabular CUSUM graph for the data in dataset 1 show in Figure 5.10.

We first computed the mean and standard deviation of the number of syn packets and applied tabular CUSUM using these numbers as the target mean and the target standard deviation. Then we highlighted the points (shown in red color) where the cumulative sum drifts more than five standard deviations beyond the target mean and set the minimum detectable mean shift to one standard deviation.

It is possible to see on the chart of syn number packets that some points are out of the H interval, five standard deviations beyond target mean. This fact is observed only in the tabular CUSUM control chart, not being noticeable in individual charts, which

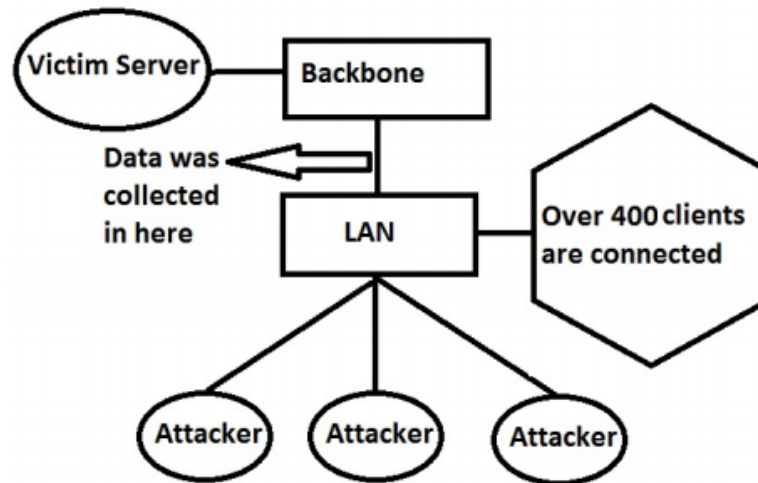


Figure 5.5. Bogazici University DDoS attack dataset topology.

indicates that the process may be considered out of statistical control [21]. The process is in statistical control when the values are close to zero. We commented on that the system becomes unstable in the time periods 3213-3637 seconds, 7358-7805 seconds, as seen from the Figure 5.6. The tabular CUSUM values (highlighted in red) related to this time intervals are discrepant in relation to others considering the cumulative sum which we evaluate this situation as there is syn flood attack in these intervals. It is also seen that cumsum values shows a positive and sudden drift in these time intervals in Figure 5.8. If the process is in a state of statistical control, the cumsum should vary around 0, shown by a horizontal line.

We had also obtained tabular CUSUM graph for the data in dataset 2. We commented on that the system becomes unstable in the time periods 791-1075 seconds, 1788-2098 seconds, 2780-3133 seconds and 3800-4206 seconds as seen from the Figure 5.11.

5.4. Applying V-mask CUSUM method

We know that if the process is in a state of statistical control, the cumulative sum values should vary around 0, shown by a horizontal line. If the entire cumulative sum values are in between the upper and lower arms of the mask, the process is said

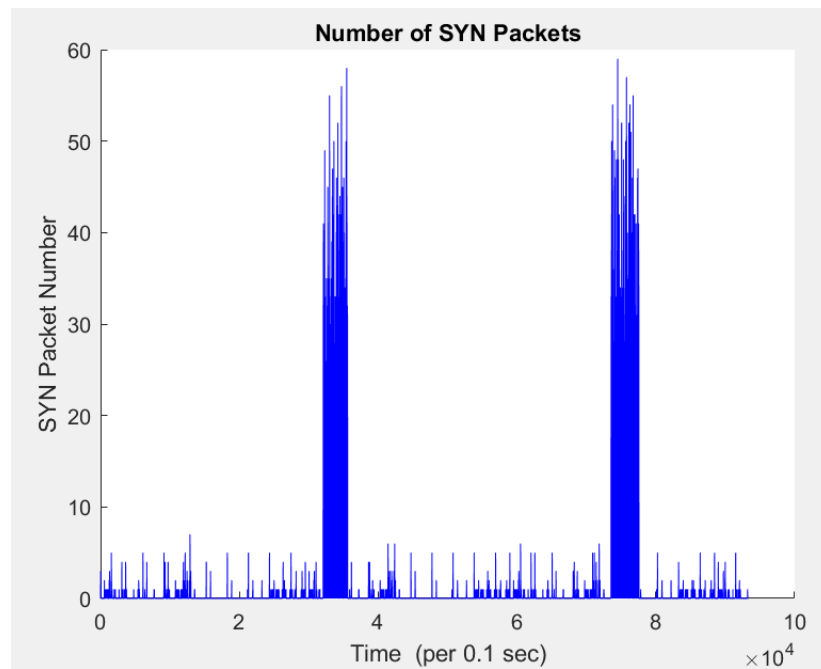


Figure 5.6. Distribution of SYN packets in dataset 1.

to be in under statistical control. If any point is not in between the upper and lower arms of the mask, the process is said to be out-of-control.

As we mentioned in the previous sections, the parameters d and θ of the V-mask CUSUM chart determine the performance of the chart. These parameters are related to the α which is the probability of incorrectly generating an out-of-control signal when the process is in control and the β which is the probability of not generating an out-of-control signal when the process shifts by an amount $\delta\sigma$.

We applied V-mask CUSUM method in order to detect anomalies on dataset 2.

The opening arms of the mask is stated at an angle of $\pm\theta$ from the horizontal at each sample. We followed the procedure which is shown in Figure 5.12.

To determine whether or not the number of syn packets is in a state of statistical control at time i , we placed V-mask at a distance d in front of C_i at each sample. Every single sample value is directly monitored with V-mask.

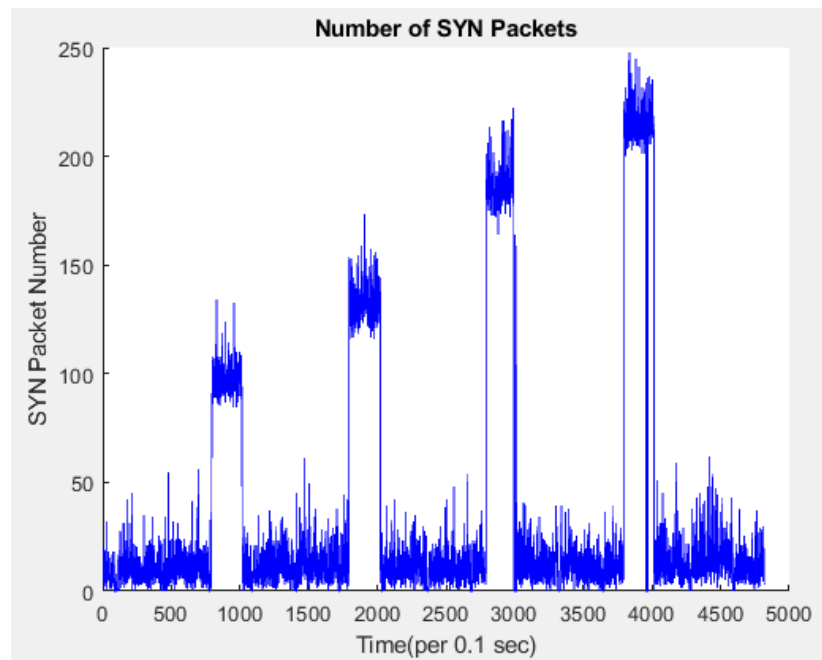


Figure 5.7. Distribution of SYN packets in dataset 2.

We assumed that the process has a mean level μ . Then we calculated the mean number of syn packets as 37.69. The minimum and maximum values of number of syn packets are 0 and 248 respectively. We set α and β parameters as 0.0027 and 0.01 as Chebyshevs inequality theorem says. Then we plotted the cumulative sum of deviations from the mean value with the V-mask procedure.

Figure 5.13 shows when V-mask positioned between the sample values i 787-1150, 1786-2180, 2779-3247 and 3799-4288, the cumsum values are outside the upper and lower arms of the mask. These are the time intervals at which the chart indicated an out-of-control condition. Therefore the system becomes unstable in these periods which means there is an unexpected increase on the number of syn packets in the data.

5.5. Comparison Analysis of Control Charts

In Average Run Length approximation, as mentioned on previous sections, we know from searches on the literature Cusum control chart gives $ARL_0 = 500$ and $ARL_1(1\sigma) = 9.5$ at its best performance for the values $k = 0.5, h = 5, ARL_0 = 500$ and $ARL_1(1\sigma) = 9.5$ whereas EWMA control chart has the best performance for the

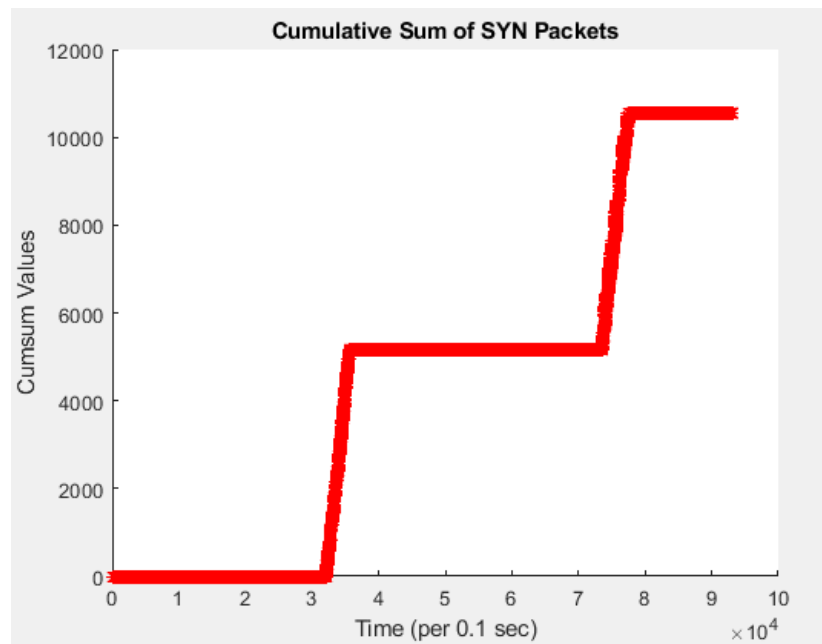


Figure 5.8. Cumsum of SYN packets in dataset 1.

values $\lambda = 0.1L = 2.814$, and it gives the result $ARL_0 = 465$ and $ARL_1(1\sigma) = 10.3$ remembering that we would like ARL_0 to be as large as possible and ARL_1 to be as small as possible. Therefore a CUSUM control chart gives better results based on average run length approximation.

In forecast model approximation, based on the samples of forecasting residuals, we applied the ARIMA(2,1,2) model for dataset 2. We evaluated the Ewma and CUSUM control charts described in previous sections. We firstly examined the dataset 2 whether it has any outliers or irregularities. As seen from the graph of our data, there are many suspected fluctuations with different magnitudes especially in the attack period. We identified and removed outliers using series smoothing and decomposition.

After removing outliers, our data still had unpredictable series. We used moving average for smoothing data because smoothing makes series more stable and predictable. We know that the wider the window of the moving average, the smoother original series becomes. We used moving average order 10 and 50 for smoothing comparison shown in Figure 5.14.

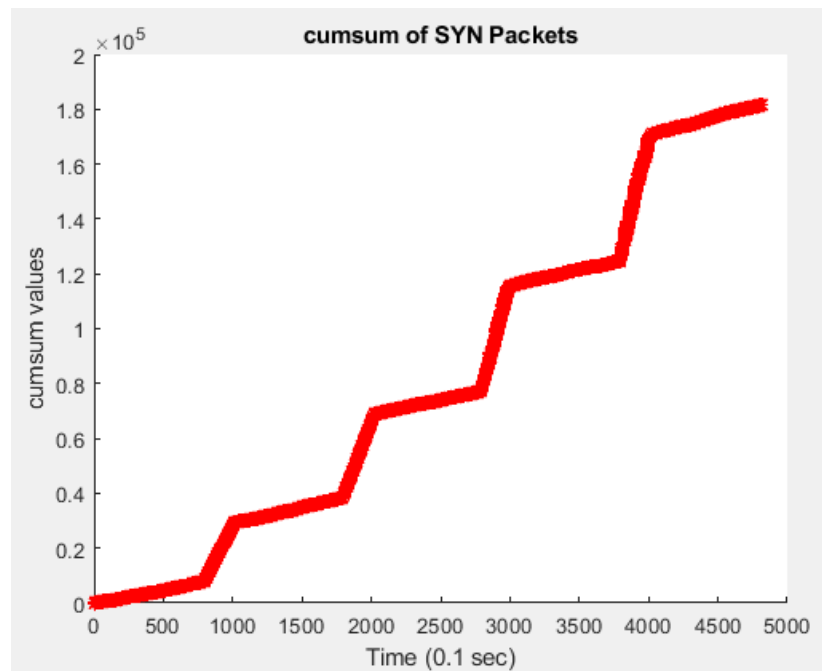


Figure 5.9. Cumsum of SYN packets in dataset 2.

Time series mostly has seasonality and trend components. Seasonal component means fluctuations in the data. Trend component represents variations of low frequency. The decomposition is process of removing these components. We decomposed our data in order to eliminate seasonality and trend components show in Figure 5.15. We removed seasonal and trend components of our data.

To fit an ARIMA model on a time series, the series must be stationary which means that its mean, variance, and autocovariance are time invariant. We used the Augmented Dickey-Fuller (ADF) test which is a formal statistical test for stationarity. Our data is non-stationary, because the average number of syn packets changes through time. In order to make a non-stationary process stationary, differencing method is used which is the parameter represented by the d component of ARIMA. We applied differencing method on our data using differencing order 1. We saw an oscillating pattern around 0 which means there is no remarkable trend that we should focus on to remove. Therefore, differencing of order 1 terms is sufficient and should be used in the model. After applying differencing method to our data, we use the augmented Dickey-Fuller test to see whether our data is stationary or not. Finally, we had a stationary data. The differenced dataset 2 is shown in Figure 5.16.

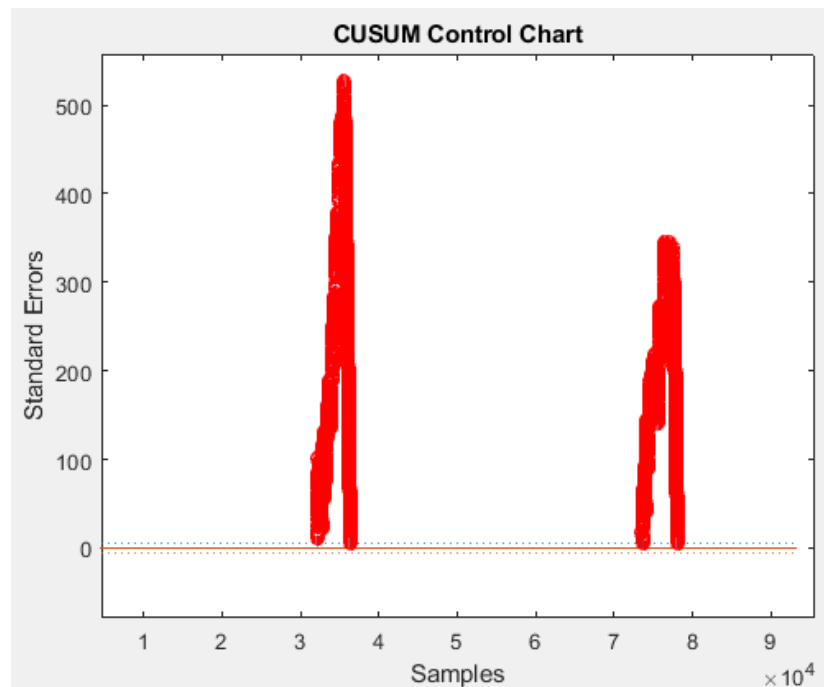


Figure 5.10. Tabular CUSUM Control Chart of number of syn packets in dataset 1.

The parameters in the ARIMA(2,1,2) model includes differencing of degree 1, and uses an autoregressive model of second lag and a moving average model of order 2. We fitted a model that can produce a forecast, we obtained and examined ACF and PACF plots for model residuals. We saw from the Figure 5.17 that there is no significant autocorrelation present.

Figure 5.18 shows the results of CUSUM control chart applied to monitor the ARIMA(2,1,2) forecasting model for dataset 2. The values in y-axis represent the forecasting errors for each prediction. Figure 5.19 show the results of EWMA control chart applied to monitor the ARIMA(2,1,2) forecasting model for dataset 2.

In forecast models based comparison, the graphical analysis show that the CUSUM control chart gives more accurate results compared to the EWMA control chart.

5.6. Results and Discussion

We have realistic datasets that we maintained in our laboratory environment. The test evaluations with ROC curve analysis of tabular CUSUM and V-mask CUSUM

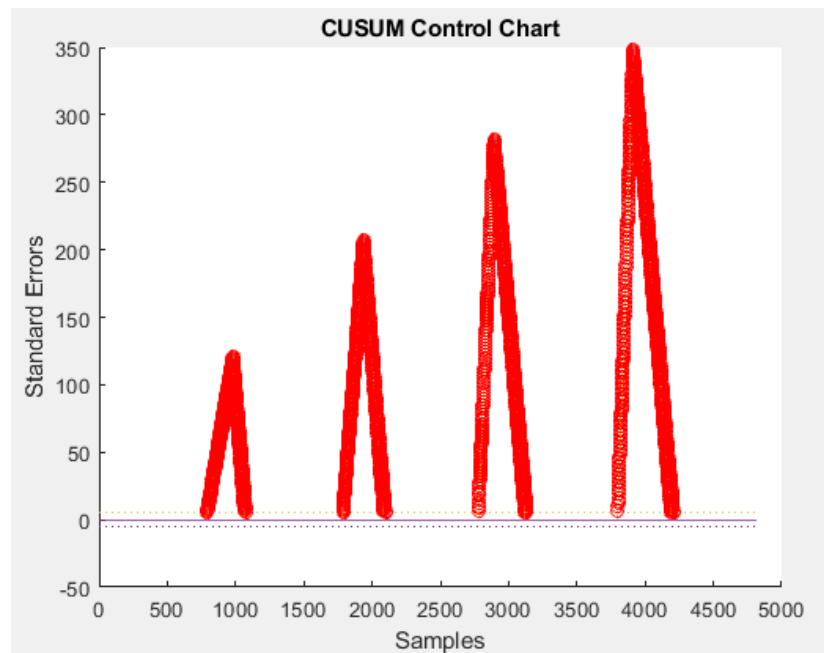


Figure 5.11. Tabular CUSUM Control Chart of number of syn packets in dataset 2.

methods are based on the exact reference values by means of samples which attack occurs. Therefore it was a realistic comparison between two CUSUM methods.

In a ROC curve the true positive rate (TPR) is plotted in function of the false positive rate (FPR) for different cut-off points. In the ROC curve, if the ROC curve is closer to the upper left corner, we obtain higher overall accuracy of the test.

One of our aim was to find which method, tabular CUSUM or V-mask CUSUM, has higher accuracy rate to detect tcp synflood attacks. We analyzed the ROC curve statistics of tabular CUSUM and V-mask CUSUM and plot the ROC curves while evaluating accuracy. The results of the ROC analysis in Figure 5.20 showed that tabular CUSUM has higher TPR and lower FPR than V-mask while detecting an anomaly on our dataset.

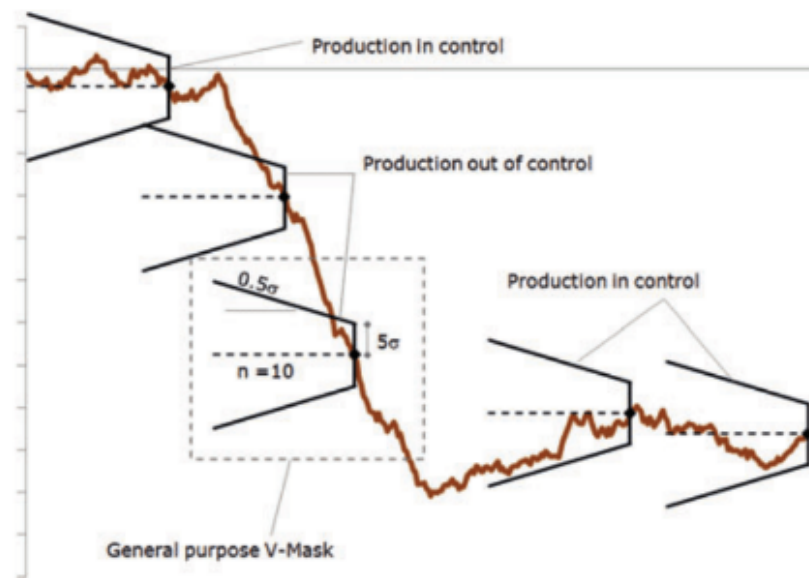


Figure 5.12. Exemplary representation of the V-Mask as an indicator and decision making support [22].

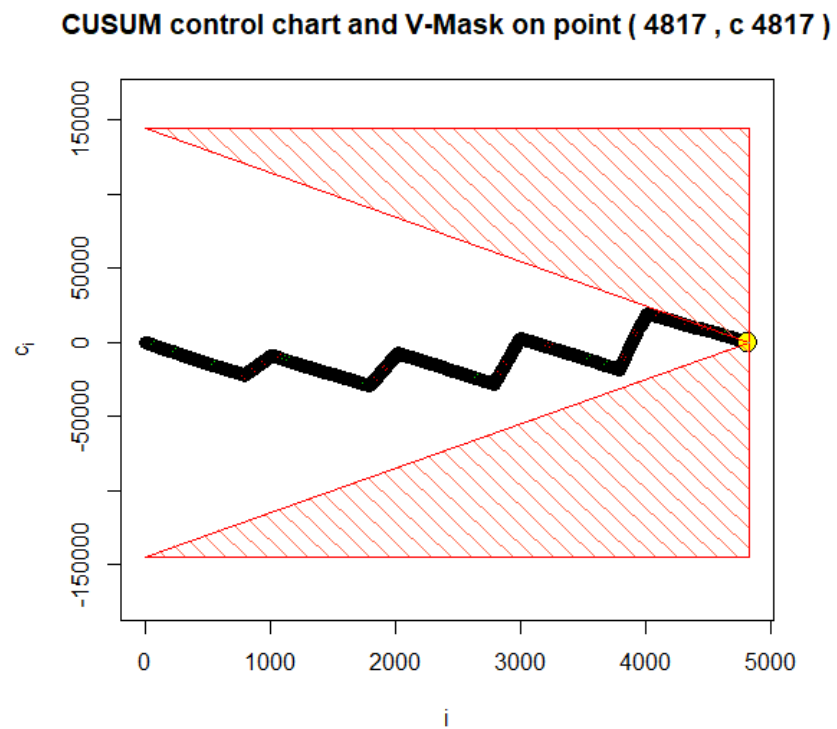


Figure 5.13. V-mask cusum applied on cusum of number of syn packets in dataset 2.

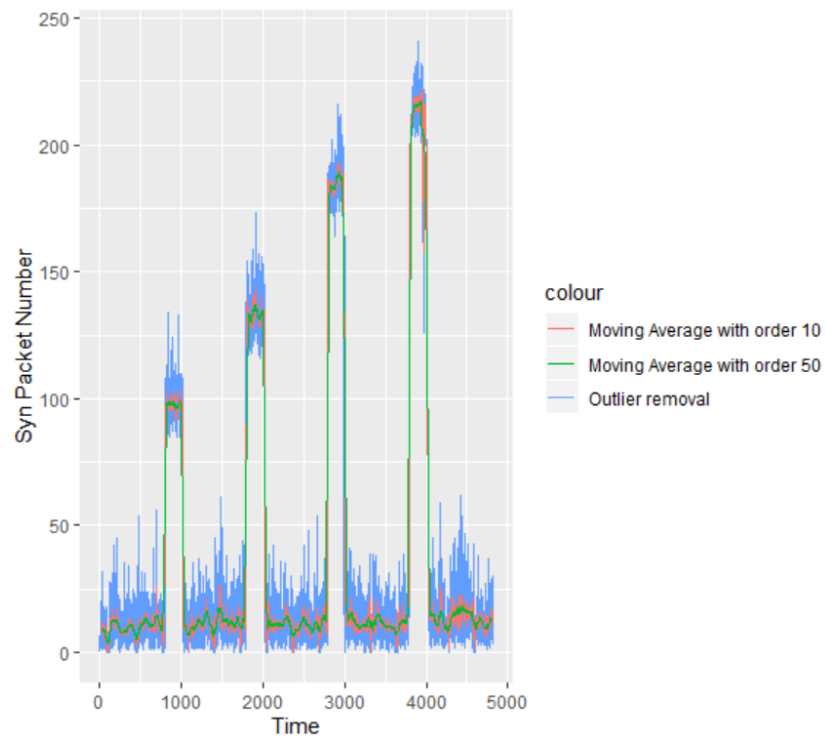


Figure 5.14. Moving average for smoothing on dataset 2.

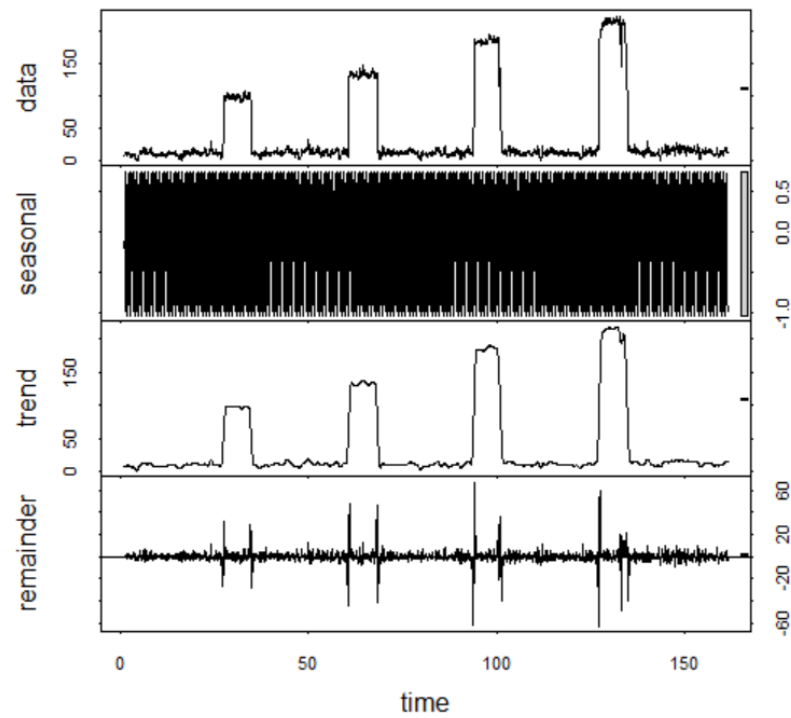


Figure 5.15. Removed seasonal and trend components of dataset 2.

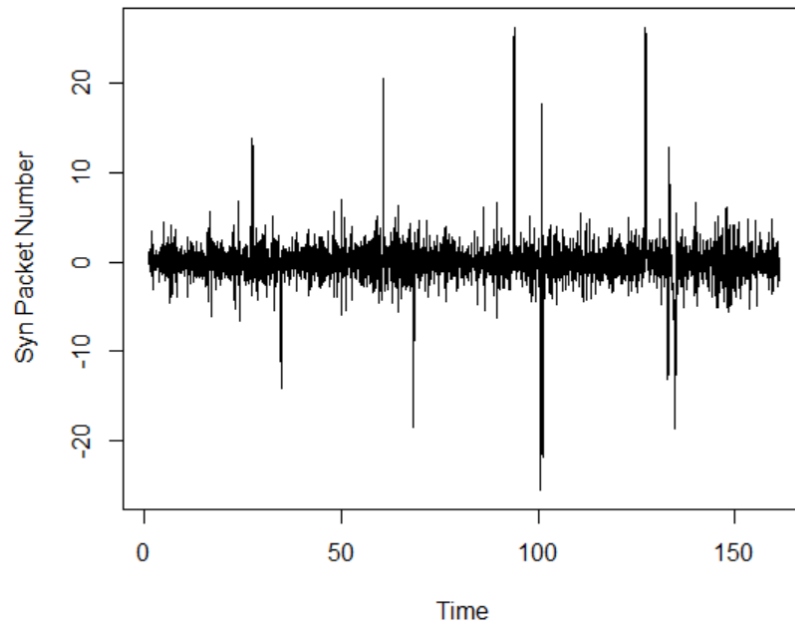


Figure 5.16. Differenced dataset 2.

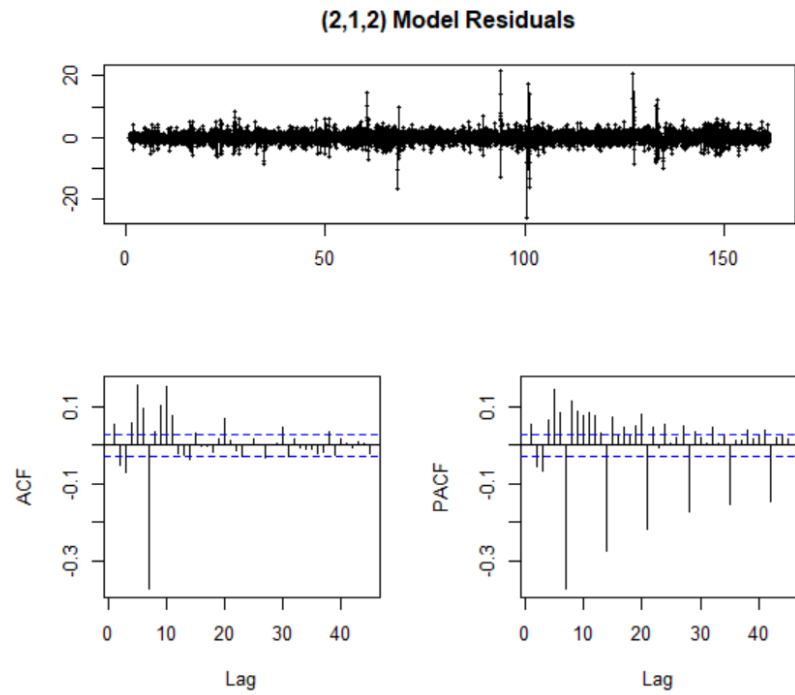


Figure 5.17. Residuals for dataset 2.

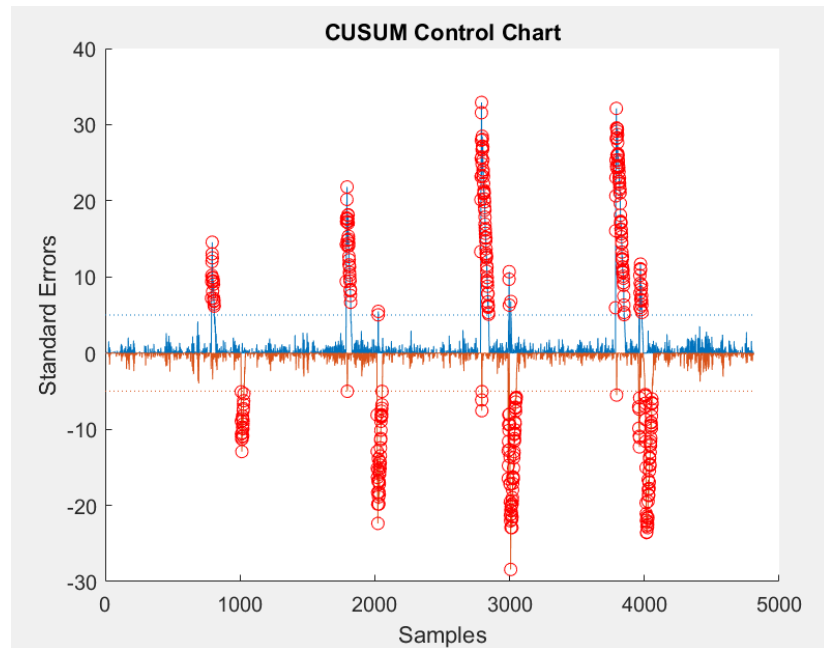


Figure 5.18. Cusum chart of residuals for dataset 2.

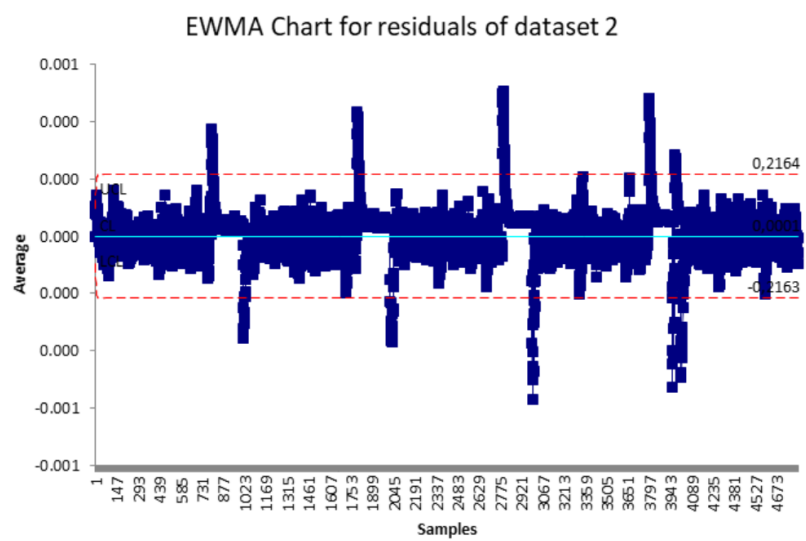


Figure 5.19. Ewma chart of residuals for dataset 2.

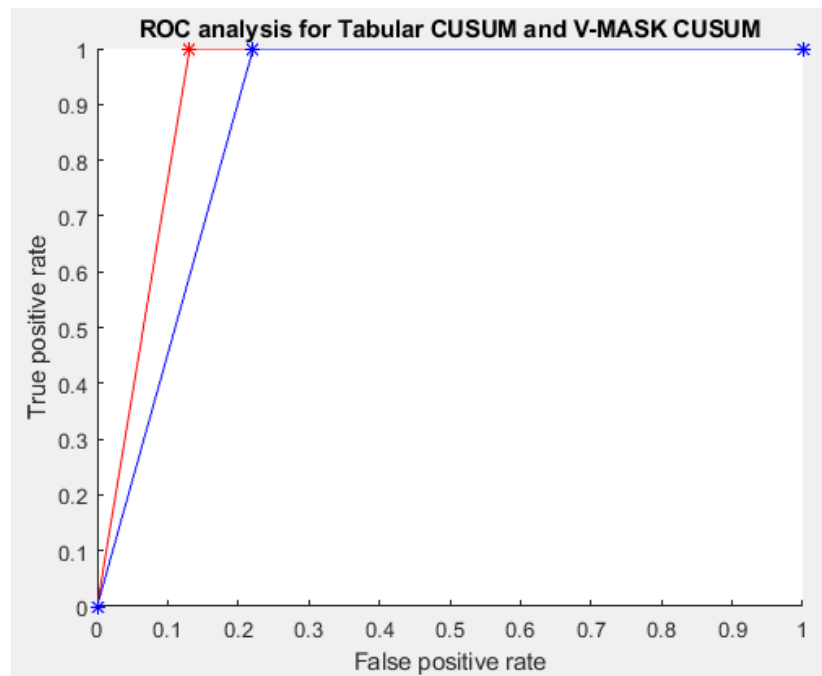


Figure 5.20. ROC analysis of Tabular CUSUM versus V-mask CUSUM.

6. CONCLUSION

In this thesis project, we obtained the number of TCP syn packets from two datasets which includes synflood attack. We obtained CUSUM values of the number of syn packets. We performed Tabular CUSUM and V-mask methods on two datasets to find out which method is the best technique to detect DDoS attacks and to obtain the most accurate results. We analyzed the ROC curve statistics of Tabular CUSUM and V-mask CUSUM and plot the ROC curves in order to evaluate the accuracy of this methods. The ROC curve of the Tabular CUSUM is closer to the upper left corner, therefore it has the higher the overall accuracy than V-mask CUSUM. It is clear that Tabular CUSUM method is the best CUSUM detection scheme for TCP synflood attack traffic when compared to V-mask CUSUM.

Then we made a performance analysis of EWMA and CUSUM charts evaluating average run length approximation. In average run length approximation, CUSUM chart gives $ARL_0 = 500$ and $ARL_1(1\sigma) = 9.5$ at its best performance for the values $k = 0.5, h = 5, ARL_0 = 500$ and $ARL_1(1\sigma) = 9.5$ whereas EWMA control chart has the best performance for the values $\lambda = 0.1L = 2.814$, and it gives the result $ARL_0 = 465$ and $ARL_1(1\sigma) = 10.3$ remembering that we would like ARL_0 to be as large as possible and ARL_1 to be as small as possible. Therefore Cusum control chart gives better results based on average run length approximation.

Finally, we used the Autoregressive Integrated Moving Average (ARIMA) forecasting model in order to obtain the forecasting residuals to be used in the performance analysis of EWMA and CUSUM control charts. We removed trend and seasonal components of our dataset 2 and applied the augmented Dickey-Fuller test to see whether our data is stationary or not. We had a stationary data as a result of transformations. In forecasting ARIMA(2,1,2) model based comparison, the numerical results and graphical analysis for dataset 2 shows us that EWMA control chart gives false alarms and have higher probability of error as seen from the Figure 5.19.

Shewhart control chart considers the last observations on data samples, this results in that at large short-term shifts, it is expected to report false alarms. CUSUM and EWMA charts consider the last data point and past data points. That is the reason makes them sensitive to both short and long term shifts [23]. However, our experimental results indicates that the CUSUM control chart demonstrates the best adequacy when compared to the EWMA control chart.

REFERENCES

1. Carl, G., G. Kesidis, R. R. Brooks and S. Rai, “Denial-of-service attack-detection techniques”, *IEEE Internet computing*, Vol. 10, No. 1, pp. 82–89, 2006.
2. Kaspersky, “Kaspersky DDoS Intelligence Report Q1 2019”, *Kaspersky Lab*, 2019.
3. Jamal, T., Z. Haider, S. A. Butt and A. Chohan, “Denial of Service Attack in Cooperative Networks”, *arXiv preprint arXiv:1810.11070*, 2018.
4. Zuckerman, E., H. Roberts, R. McGrady, J. York and J. Palfrey, “Distributed denial of service attacks against independent media and human rights sites”, *The Berkman Center*, 2010.
5. Aamir, M. and M. A. Zaidi, “A survey on DDoS attack and defense strategies: from traditional schemes to current techniques”, *Interdisciplinary Information Sciences*, Vol. 19, No. 2, pp. 173–200, 2013.
6. Matias, R., A. M. Carvalho, L. B. Araujo and P. R. Maciel, “Comparison analysis of statistical control charts for quality monitoring of network traffic forecasts”, *2011 IEEE International Conference on Systems, Man, and Cybernetics*, pp. 404–409, IEEE, 2011.
7. Wild, C. J. and G. A. Seber, “Chance Encounters: A First Course in Data Analysis and Inference Reviewed by Flavia Jolliffe”, .
8. Yang, S.-F. and B. C Arnold, “A Simple Approach for Monitoring Business Service Time Variation”, *TheScientificWorldJournal*, Vol. 2014, p. 238719, 05 2014.
9. Gan, F. F., “Cumulative Sum (CUSUM) Chart”, *Wiley StatsRef: Statistics Reference Online*, 2014.
10. Hawkins, D. M. and D. H. Olwell, *Cumulative sum charts and charting for quality*

- improvement*, Springer Science & Business Media, 2012.
11. Montgomery, D. C., *Introduction to statistical quality control*, John Wiley & Sons, 2007.
 12. Basseville, M., I. V. Nikiforov *et al.*, *Detection of abrupt changes: theory and application*, Vol. 104, Prentice Hall Englewood Cliffs, 1993.
 13. Baldewijns, G., S. Luca, B. Vanrumste and T. Croonenborghs, “Developing a system that can automatically detect health changes using transfer times of older adults”, *BMC medical research methodology*, Vol. 16, No. 1, p. 23, 2016.
 14. Prajapati, D., “Effectiveness of conventional CUSUM control chart for correlated observations”, *International Journal of Modeling and Optimization*, Vol. 5, No. 2, p. 135, 2015.
 15. Woodall, W. H. and B. M. Adams, “The statistical design of CUSUM charts”, *Quality Engineering*, Vol. 5, No. 4, pp. 559–570, 1993.
 16. Lucas, J. M., “The design and use of V-mask control schemes”, *Journal of Quality Technology*, Vol. 8, No. 1, pp. 1–12, 1976.
 17. Lucas, J. M., “A Modified V Mask Control Schemet”, *Technometrics*, Vol. 15, No. 4, pp. 833–847, 1973.
 18. Taylor, H. M., “The economic design of cumulative sum control charts”, *Technometrics*, Vol. 10, No. 3, pp. 479–488, 1968.
 19. Koshti, V., “Cumulative sum control chart”, *International journal of physics and mathematical sciences*, Vol. ISSN, pp. 2277–2111, 12 2011.
 20. Gan, F., “An optimal design of CUSUM quality control charts”, *Journal of Quality Technology*, Vol. 23, No. 4, pp. 279–286, 1991.

21. Follador, F. A. C., M. A. A. V. Boas, M. Schoenhals, E. Hermes and C. Rech, “Tabular cusum control charts of chemical variables applied to the control of surface water quality”, *Engenharia Agrícola*, Vol. 32, pp. 951 – 960, 10 2012, [http://www.scielo.br/scielo.php?script=sci_arttextpid = S0100 – 69162012000500014nrm = iso](http://www.scielo.br/scielo.php?script=sci_arttextpid = S0100 - 69162012000500014nrm = iso).
22. Schmidt, W., H.-C. Kühne and B. für Materialforschung, “Cusum control charts with V-mask in process control”, *CPI International*, 2015.
23. Cheng, S. W. and K. Thaga, “Single variables control charts: an overview”, *Quality and Reliability engineering international*, Vol. 22, No. 7, pp. 811–820, 2006.