



# Swinburne University of Technology

Faculty of Science, Engineering and Technology

## ASSIGNMENT AND PROJECT COVER SHEET

Unit Code: 30015 Unit Title: IT Security

Assignment number and title: Assignment 1 Research Report

Due date: 8<sup>th</sup> September

Lab group: NA Tutor: Jamie OOI

Lecturer: Lin Li

Family name: Efaz

Identity no: 103494172

Other names: Ekrar Uddin Mohammed

### To be completed if this is an INDIVIDUAL ASSIGNMENT

I declare that this assignment is my individual work. I have not worked collaboratively, nor have I copied from any other student's work or from any other source except where due acknowledgment is made explicitly in the text, nor has any part been written for me by another person.

Signature: Ekrar Efaz

Marker's comments:

Total Mark:

---

# Machine Learning for Intrusion Detection Systems

---

*COS30015 IT Security*

*Assignment 1 Research Report*

---

*Ekrar Efaz (103494172)*

*Submission Date: 8<sup>th</sup> September*

*Submission Due Date-Time: 11:59 pm 8<sup>th</sup> September*

***Abstract: Intrusion Detection is one of the widely studied topics in the space of Cyber-Security. Intrusion Detection Systems are yet to be perfect and even after years of research the community still faces severe difficulties. The uneven data-growth is a lot to handle for manually updated IDSs. The problem of reducing enormous number of false positives is still lingering and there have been little progress into detecting novel attacks. However, current progress in research has shined some light and there may be answers to these dilemmas. This paper postulates an outline of research directions for using Machine Learning to solve Intrusion Detection core problems.***

***Keywords: Intrusion, Machine Learning, Intrusion Detection, Algorithm, Anomaly, Supervised Learning, Unsupervised Learning, Support Vector Machine,***

## ***1. Introduction***

The explosive use of the Internet and Networks has raised a lot of security issues and threats and the study of Intrusion Detection has received a lot of attention as a result.

First, we make a clear distinction among the terms Intrusion, Intrusion Detection, and Intrusion Detection System. Intrusion is the act of compromising a system or its CIA (Confidentiality, Integrity, Authenticity). Intrusion Detection automates the process of monitoring events in a system or network to detect any anomaly or signs of intrusion. Intrusion Detection system is a hardware/software tool to aid in Intrusion Detection. [1]

Early Intrusion detection worked based on two assumptions [2]: one, the detection of anomaly based on rules that would express intrusion in an obvious manner; and two, deviation from what is normal behaviour according to the program or the user or the admin could be a sign of intrusion. The exponential increase of users also increases the types of behaviour that would be normal/abnormal, and it is becoming hard to handle the growing data which is why researchers are looking into Machine Learning for Intrusion detection.

Machine learning is the branch of computer science which makes use of available data and complex algorithms to make machines learn by imitating how humans learn gradually improving its accuracy. The use of Machine Learning trains the Intrusion Detection System to learn the normal behaviour and detect anomalies equipped with the power to handle large amount of system logs or network traffic. Machine Learning techniques have gotten a lot of attention from researchers to address the flaws of knowledge-based detection.

## ***2. Overview***

The paper is organized to give an overview of the problems that come with traditional Intrusion detection viewpoint like; handling growing data and traffic, failure to detect or predict novel attacks, false positives, signature-modified attacks to evade IDS and how the Machine Learning approach today is attempting to deal with the problems.

In the first part of the paper, we talk about more in-depth about general detection methodologies and approaches in IDS. Generally, how IDSs work and their architecture. Elaborate description of detection methodologies and approaches are provided in this part.

In the second part we bring about the how Machine Learning helps us with intrusion detection and then talk about a few Intrusion Detection ML approaches that are widely studied. Discussion on what is Machine Learning, Anomaly detection using Machine Learning, how algorithms in machine learning are trained and how machine learning solves issues with Intrusion detection using various models can be found in this part of the paper.

### **3. Literature Review**

Traditionally, the study of Intrusion Detection is done in two major perspectives. One is anomaly detection which aims to detect all deviated behaviours from the stated normal behaviour [3]. On the other hand, Misuse detection is understood as the detection of precise and specific representable techniques of system exploitation. [4]

We only focus on the anomaly-based detection in this paper for further study and survey.

There are different methodologies to Intrusion Detection such as knowledge-based intrusion detection, behaviour-based detection, and specification-based detection. Intrusion Detection methodologies are classified as three major categories. [5, 6]

#### **3.1.1. Signature-based (knowledge-based)**

It works by matching the traffic to a specific signature dataset of known intrusion patterns. This system is dependent on the signature dataset which makes it time consuming to maintain the knowledge, and hard to keep the signatures updated with increasing threats. It is also ineffective against signature-modified or novel attacks.

#### **3.1.2. Anomaly-based (behaviour-based)**

Anomaly-based detection works by identifying anormal behaviour/traffic from the flow of normal traffic. This system is trained to learn normal behaviour and anything that deviates from its training dataset is flagged as an anomaly. But this method is prone to a lot of false positives and with ever changing “normal traffic behaviour” it’s hard to make out which one is an anomaly.

#### **3.1.3. Stateful protocol analysis (specification-based)**

Stateful protocol analysis studies and weigh up profiles of predetermined accepted benign protocol for each activity and then identifies deviations. Unlike anomaly-detection which identifies threats based on behaviours that deviate from normal behaviour stateful analysis uses vendor developed universal protocol.

### **3.2. Issues with Intrusion Detection Methodologies [7]**

#### **3.2.1. Detecting Novel Attacks**

Novel attacks can be considered as 0-day attacks such that they have never been encountered before and their behaviour although anomalous is completely new to the Intrusion Detection system. It is extremely rare for traditional IDS to be able to detect Novel attacks because they don’t have their signatures in their database even, they are hard to place under normal or abnormal traffic based on how they are programmed.

### ***3.2.2. False Positives & Alarm Management***

IDSs are designed to monitor and report any suspected activity via sounds, email or by paging the *Sysadmin*. These reports are then studied, and action is taken to mitigate the issue. But Most IDSs face the problem of falsely classifying a normal connection or traffic as an anomaly and therefore obstructs or raises false alarms denying legitimate user access and headache for security analysts. False-Positives and alarms reduce the reliability of the IDS. This is because of attacks becoming more and more sophisticated and modified attacks are intentionally made to look like normal traffic.

### ***3.2.3. Updating Signature Database***

Keeping the signature dataset updated is a big issue nowadays because of the rise in attacks and new novel attacks are happening every other week which increases the number of updates that are required to be released and installed making it much more ineffective if an update is failed to release.

## ***3.3. Machine Learning***

Machine learning is the critical study of teaching computers to make smarter performance decisions based on learning and experience. Machine learning has wide range of application domain mainly useful for (a) problem domains which are poorly understood and requires effective algorithms, but little knowledge exists (b) problem domains where large datasets are available from which regularities can be discovered. (c) domains where programs must adapt to evolving situations. To no surprise Intrusion Detection is a fertile field where lots of problems can be formulated as learning problems and Machine Learning approach would be an effective measure to solve them. [6]

### ***3.4. Anomaly Detection using Machine Learning***

Theoretically more intelligence is brought into IDS to output the best performance and it would be possible to avoid false positives and low detection rates if the ML Algorithm is trained with infinite datasets, but it is computationally limited, and IDS would be unable to provide response in real-time.

ML builds a prediction model after it is trained from a specific dataset for anomaly detection. Algorithms are used to learn from data and generate predictions based on it. Machine Learning shines a promising light on the study for anomaly detection and even classification based on many features. Machine Learning Algorithms are of two types: [8]

#### ***3.4.1. Supervised Learning***

Supervised algorithms are predictive models that are specifically trained with labelled dataset where the normal and anomalous samples are labelled. In general, supervised models are supposed to be more effective because they have access to more information, but sometimes it is not possible to overcome the scarce comprehensive training dataset. This results in a more error-prone predictive model.

### 3.4.2. *Unsupervised Learning*

Task of unsupervised learning is ambiguous to some extent. Unsupervised models are provided dataset containing data points that are not labelled. The aim of unsupervised learning is to use AI algorithms to discover pattern or structure in the data. Different models work in very different ways and discover very different structures based on the same dataset.

### 3.5. *Tackling IDS issues with ML*

#### 3.5.1. *False-Positives and Alarm Management* [9]

Anomaly detection usually has a high rate of false-positive. The anomaly detection model depends on its training dataset greatly and it is difficult to curate high quality training data because new attacks are emerging every day, but training datasets are not updated as frequently. However, to deal with the false positives system tuning must be done accordingly which is a strenuous task for the system administrators. Alarm filtering and attack classification can reduce the number of false alarms as well as take care of classifying the type of attacks which will make it easier for system administrators to analyse the reports and filter out false alarms.

A lot of hybrid machine learning algorithms sort and classify the attacks upon detection. Most notable among them is the *Multi-Step Multi-Class SVM-based Anomaly Detection*. MMIDS combines two types of anomaly detection algorithm to detect and classify attack using the *Fuzzy-ART* to put out detailed clustering for each attack.

#### 3.5.2. *Novel Attacks*

Novel attacks are very rare for IDS to pick up because it's an attack that the signature database doesn't have any signature on. It is also hard to train a model to detect novel attacks because it has completely new features that the training dataset had no datapoints on, which makes the trained model oblivious to the attempts made by the attacker. Anomaly detection before an attack i.e., during reconnaissance is the best way to make out a novel attack. Every attack is divided into few phases and the recon phase is similar for most attacks which include probing ports and vulnerabilities. If the IDS is trained to sniff out probing and recon attempts and notifies the system administrator about it actions can be taken for security.

A simple anomaly detection algorithm like a One-Class SVM can be trained with normal traffic introduced to train it. It will pick up on the normal tendencies and any outlier is marked, alerted, and blocked from further actions.

#### 3.5.3. *Updated Signature Database* [10]

Signature Databases a huge problem to keep updated. Taking the signature-based detection approach is also very ineffective because of the new attacks that emerge. The attacks are to be detected and then a signature made and integrated into the database which sounds more like a post-attack counter measure rather than a pre-attack measure. To improve anomaly detection, we can incorporate packet-based anomaly detection.

Improved Live Anomaly Detection System (I-LADS) processes and examines the IP Packet header data to make way for better sorting of detected anomalies. We can also take on the adaptive learning effort to solve the Signature Database problem, which forces user to collect and construct high quality dataset. But we are going focus on I-LADS and how its implemented in this paper.

### ***3.6.Survey of Related Work***

#### ***3.6.1. One-Class SVM anomaly detection*** [11] [12]

*SVM or Support Vector Machines* are supervised learning models that are associated with learning algorithms that are used for classification and analysis. SVMs usually deal with multi-class classification problems where they sort out different datasets based on their learning model. Surprisingly OCSVM or One-Class support vector mechanism solves only single class classification problem which means it only knows a single class of normal data and upon presentation of fresh data it determines if It belongs to the group or deviates from the norm to mark it as anomalous.

#### ***3.6.2. Multi-Step Multi-Class SVM anomaly detection (MMIDS)*** [13]

MMIDS is a hybrid anomaly detection algorithm made up of 3 modules; OCSVM to perceive normal data and attack data; multi-class SVM then catalogues the attack into one of the trained categories and later further clustering of each attack is done by the Fuzzy-ART.

Three-step procedure for MMIDS to work is described below:

*Phase 1:* The OCSVM is trained with normal data to detect deviances. Training only requires normal data which ensures faster training speed. After detection it is pushed into Phase 2.

*Phase 2:* The multi-class SVM provides taxonomy for the type of attack into one of the following: probing, DOS, R2L etc and offers the system more information about the type of attack.

*Phase 3:* For further analysis and study, Fuzzy-ART provides clustering information for the specified attack

#### ***3.6.3. Improved Live Anomaly Detection System*** [14]

Most of the present effort focuses on detecting anomaly on the amount of flow of data or traffic. This although is a very simple solution and somewhat effective poses a serious problem. Predicting anomaly based on variations in traffic flow is now always accurate because huge traffic is not always the worst-case scenario for a lot of businesses, and it might be the by-product of their blooming business or advertising.

Proposed I-LADS examines the IP packet headers and its data to allow for better detection of anomalies. Since the amount of information that I-LADS processes is a lot more than traffic flow deviations it presents a much more detailed and accurate report for anomaly detection.

I-LADS has two versions, and its selection depends on the desired performance and use-case. I-LADSV2 allows working with much larger datasets while I-LADSV1 let us work with filtered IP addresses.

## **4. Discussion**

### **4.1.Pre-Attack Detection**

Attackers prepare for an attack in early stages by Recon which involves gathering information about the system. This pre-attack stage could be sensed by IDS and further action from the specific IPs could be handled and the system administrator could be alerted to act. OCSVM and I-LADS would do a very good job of discovering pre-attack signs.

If the OCSVM is trained with normal daily traffic as its datapoints and deformities in the traffic flow would be flagged. On the other hand, I-LADS examines IP Packet Header data to differentiate between normal and anomalous traffic. Probing packets are a huge giveaway if not masked and can be flagged by the I-LADS instantly. Further measures are required for masked packets designed to look like normal traffic.

### **4.2.Mitigating False Positives and Alarms**

To crack the False-Positive problem in Anomaly based Intrusion detection, we must start by classifying the uncovered attacks and then do further clustering on each type of attack to make it easier for system administrators to go through the attack report. Cataloguing each type of attack would filter out the false alarms because they wouldn't fall into certain attack categories.

MMIDS would be a good alternative to reduce False Positives as well as I-LADS. MMIDS would identify the attack data and then in its 3-step process sort and cluster the type of attack. Whereas I-LADS would avoid false positives because it makes its prediction based on a lot more information about the traffic packet rather than the amount of traffic flow.

### **4.3.Countering Data-Growth [15]**

Which the exponential rise in data and its usage it is nearly impossible to make any sense of the data or drive any insight from it manually and unstructured data makes it much more difficult to work with. We need technologies that aren't afraid of sorting, analysing, and driving insight from huge amounts of data. When it comes to big data sets machine learning produces the best results because they can handle sizable measures of data and in fact huge datasets are great training tools for Machine Learning Algorithms

In my observation and study, I understand that most hybrid machine learning algorithms i.e., mixture of two ML algorithms that could be both supervised or one supervised and the other unsupervised, makes the best use cases for anomaly-based intrusion detection. While a single model is yet to solve all the problems that come with Intrusion Detection Systems, MMIDS and I-LADS attempts to solve the False-Positive problem which makes the IDSs much more reliable, furthermore OCSVM and I-LADS can detect attack vectors before the [6]commencement of the attack from information gathering stage of the attackers.



## 5. *Conclusion*

Machine Learning approach to Intrusion Detection is receiving a lot of attention because of its ability to handle large datasets. Intrusion Detection is still the frontline of defence for many enterprises and systems. The world's growing data demands for a smarter approach to Intrusion Detection to solve its core problems to make it more reliable and performance efficient. Machine Learning algorithms learn and attempt to disentangle the False-Positive problems by classifying and clustering attacks, making smart pre-attack predictions by using specific training datasets and making more informed decisions by examining traffic parameters instead of just the volume of traffic flow. Hybrid ML algorithms that combine two or more algorithms demonstrate the most promising future for Intrusion Detection Systems.

## References

---

- [1] R. B. & P. Mell, Intrusion Detection Systems, NIST Special Publication, 2000.
- [2] G. Bruneau, The History and Evolution of Intrusion Detection, SANS, 2021.
- [3] V. V. R. Prasad, V. Jyothsna and K. M. Prasad, "A Review of Anomaly based Intrusion Detection Systems," *International Journal of Computer Applications* , Vols. 28 - No.7, 2011.
- [4] C. Lawrence, Intrusion Prevention System: The Future of Intrusion Detection?, University of Auckland, 2004.
- [5] H.-J. Liao, C.-H. R. Lin, Y.-C. Lin and K.-Y. Tung, "Intrusion detection system: A comprehensive review," *Journal of Network and Computer Applications*, 2012.
- [6] "Practical Data Science," [Online]. Available: <https://www.datasciencecourse.org/notes/unsupervised/>. [Accessed 08 09 2022].
- [7] M. Aljanabi, M. A. Ismail and A. H. Ali, "Intrusion Detection Systems, Issues, Challenges, and Needs," *International Journal of Computational Intelligence Systems*, vol. 14, no. 1, 2021.
- [8] A. N. H. H. J. Salima Omar, "Machine Learning Tech- niques for Anomaly Detection: An Overview," *International Journal of Computer Applications (0975 – 8887)*, vol. Volume 79 – No.2, 2013.
- [9] Z. Yu and J. Tsai, Intrusion Detection: A Machine Learning Approach, Imperial College Press, 2011.
- [10] V. Kumar and D. O. P. Sangwan, "Signature Based Intrusion Detection System Using SNORT," *International Journal of Computer Applications & Information Technology*, vol. 1, no. 3, 2012.
- [11] R. Khandelwal, One Class Support Vector Machine Anomaly Detection Techniques: Part 2, 2021.
- [12] M. G. a. S. A. Mennatallah Amer, "Enhancing one-class support vector machines for unsupervised anomaly detection," 2013.
- [13] H. S. J. P. D. Lee, "Intrusion Detection System Based on Multi-class SVM.," in *Lecture Notes in Computer Science*, vol 3642. Springer, Berlin, Heidelberg., 2005.
- [14] G. Granadillo, A. Bedoya and R. Diaz, "An Improved Live Anomaly Detection System (I-LADS) based on Deep Learning Algorithms.," Vols. 568-575, 2021.
- [15] E. Efaz, "Anomaly Detection using Machine Learning," 2022.