

Performance of Machine Learning and Big Data Analytics Paradigms in Cybersecurity and Cloud Computing Platforms

Professor Gabriel Kabanda

Received: 9 September 2021 Accepted: 4 October 2021 Published: 15 October 2021

Abstract

The purpose of the research is to evaluate Machine Learning and Big Data Analytics paradigms for use in Cybersecurity. Cybersecurity refers to a combination of technologies, processes and operations that are framed to protect information systems, computers, devices, programs, data and networks from internal or external threats, harm, damage, attacks or unauthorized access. The main characteristic of Machine Learning (ML) is the automatic data analysis of large data sets and production of models for the general relationships found among data. ML algorithms, as part of Artificial Intelligence, can be clustered into supervised, unsupervised, semi-supervised, and reinforcement learning algorithms.

Index terms— cybersecurity, artificial intelligence, machine learning, deep learning, big data analytics, cloud computing

1 Introduction a) Background

he era of the Internet of Things (IoT) generates huge volumes of data collected from various heterogeneous sources which may include mobile devices, sensors and social media. This Big Data presents tremendous challenges on the storage, processing and analytical capabilities. Cloud Computing provides a cost-effective and valid solution in support of Big Data storage and execution of data analytic applications. IoT requires both cloud computing environment to handle its data exchange and processing; and the use of artificial intelligence (AI) for data mining and data analytics. However, AI provides value-adding contributions in improving the traditional cybersecurity challenged by both the cloud vulnerability and the networking of IoT devices. Sadly, AI is also being used by hackers to threaten cybersecurity. A hybrid cybersecurity model which uses AI and Machine Learning (ML) techniques may mitigate against IoT cyber threats on cloud computing environments. As the number of IoT devices increases phenomenally, the volumes of cloud-based data and the degree of cybersecurity vulnerability increases astronomically with a high degree of complexity. The situation is exacerbated by the IoT devices that come with inadequate cybersecurity safeguards. Vulnerabilities in the IoT devices opens a window of opportunity for cyber crimes and other forms cybersecurity risks, especially among interconnected devices at now at household level.

The research paper is focused on the Performance of Machine Learning and Big Data Analytics paradigms in Cybersecurity and Cloud Computing platforms. The purpose of the research is to evaluate Machine Learning and Big Data Analytics paradigms for use in Cybersecurity. This is relevant due to the rapid advances in machine learning (ML) and deep learning (DL) as we explore the potency of efficient and cost-effective cloud computing platforms and services. Evaluation of the attacks and defenses using ML and Big Data paradigms is the key subject of this research paper. However, ML and DL techniques are resource intensive and require huge volumes of training data with excellent performance, as is often provided by computational resources such as high-performance graphics processing units (GPUs) and tensor processing units. Security issues related to virtualisation, containerization, network monitoring, data protection and attack detection are interrogated whilst strengthening AI/ML/DL security solutions that involve encryption, access control, firewall, authentication and intrusion detection and prevention systems at the appropriate Fog/Cloud level.

Cybersecurity consolidates the confidentiality, integrity, and availability of computing resources, networks, software programs, and data into a coherent collection of policies, technologies, processes, and techniques

2 FIGURE 1: TYPICAL INTRUSION DETECTION SYSTEM

46 to prevent the occurrence of an attack [1]. Cybersecurity refers to a combination of technologies, processes
47 and operations that are framed to protect information systems, computers, devices, programs, data and
48 networks from internal or external threats, harm, damage, attacks or unauthorized access [2]. The major
49 cybersecurity applications are intrusion detection and malware detection. The rapid advances in mobile
50 computing, communications and mass storage architectures have precipitated the new phenomena of Big Data
51 and Internet of Things (IoT).

52 The transformation and expansion of the cyberspace has resulted in an exponential growth in the amount,
53 quality and diversity of data generated, stored and processed by networks and hosts. These changes have
54 necessitated a radical shift in the technology and operations of cybersecurity to detect and eliminate cyber
55 threats so that cybersecurity remains relevant and effective in mitigating costs arising from computers, networks
56 and data breaches [2].

57 The Network Intrusion Detection Systems (NIDS) is a category of computer software that monitors system
58 behaviour with a view to ascertain anomalous violation of security policies and distinguishes between malicious
59 users and the legitimate network users [3]. The two taxonomies of NIDS are anomaly detectors and misuse
60 network detectors. According to [4], the components in Intrusion Detection and Prevention Systems (IDPSs) can
61 be sensors or agents, servers, and consoles for network management. Data over networks may be secured through
62 the use of antivirus software, firewall, encryption, secure protocols, etc. However, hackers can always devise
63 innovative ways of breaking into the network systems. An intrusion detection and prevention system (IDPS),
64 shown on Figure ?? below, is placed inside the network to detect possible network intrusions and, where possible,
65 prevent the cyber attacks. The key functions of the IDPSs are to monitor, detect, analyze, and respond to cyber
66 threats.

67 The strength of the overall security in Cybersecurity is determined by the weakest link [5]. Access controls
68 and security mechanisms should be encapsulated in the company objectives. Firewall protection has proved to
69 be inadequate because of gross limitations against external threats [6].

70 2 Figure 1: Typical Intrusion detection system

71 Computers are instructed to learn through the process called Machine Learning (ML), a field within artificial
72 intelligence (AI). Artificial intelligence (AI) is the simulating of human intelligence in machines, through
73 programming computers to think and act like human beings [7]. The main characteristic of ML is the automatic
74 data analysis of large data sets and production of models for the general relationships found among data. ML
75 algorithms require empirical data as input and then learn from this input. However, the amount of data provided
76 is often more important than the algorithm itself. Deep Learning (DL), as a special category of ML, brings us
77 closer to AI. ML algorithms as part of Artificial Intelligence (AI) can be clustered into supervised, unsupervised,
78 semi-supervised, and reinforcement learning algorithms. The three classes of ML are as illustrated on Figure 2
79 below [8], and these are: Supervised learning: where the methods are given inputs labeled with corresponding
80 outputs as training examples;

81 Unsupervised learning: where the methods are given unlabeled inputs;

82 Reinforcement learning: where data is in the form of sequences of observations, actions, and rewards.
83 Supervised learning models are grounded on generating functions that maps big datasets (features) into desired
84 outputs [10]. Unsupervised learning is seen as a mathematical process of minimizing redundancy or categorizing
85 huge datasets based on likeness [7]. It is important to note that Machine Learning is a technique of big data
86 analytics that includes programming analytical model construction [10]. The output of a machine learning model
87 often includes perceptions and/or decisions. Big data analytics has emerged as a discipline of ways to analyze,
88 systematically extract/mine information from, or otherwise deal or work with enormous or complex datasets
89 which are large to be dealt with by traditional data-processing methodologies [7].

90 The transformation and expansion of the cyber space has led to the generation, use, storage and processing
91 of big data, that is, large, diverse, complex, multidimensional and usually multivariate datasets [11]. According
92 to [12], Big Data refers to the flood of digital data from many digital sources. The data types include images,
93 geometries, texts, videos, sounds and combinations of each. [13] explained big data as the increase in volume
94 of data that offers difficulty in storage, processing and analysis through the traditional database technologies.
95 Big Data came into existence when the traditional relational database systems were not able to handle the
96 unstructured data generated by organizations, social media, or from any other data generating source [14].
97 The characteristics of big data are volume, velocity, variety, veracity, vocabulary and value [11]. Big data has
98 necessitated the development of big data mining tools and techniques widely referred to as big data analytics.
99 Big data analytics makes use of analytic techniques such as data mining, machine learning, artificial learning,
100 statistics, and natural language processing. In an age of transformation and expansion in the Internet of Things
101 (IoT), cloud computing services and big data, cyber-attacks have become enhanced and complicated [15], and
102 therefore cybersecurity events become difficult or impossible to detect using traditional detection systems [16],
103 [17]. Big Data has also been defined according to the 5Vs as stipulated by [18] where:

104 ? Volume refers to the amount of data gathered and processed by the organisation

105 ? Velocity referring to the time required to do processing of the data ? Variety refers to the type of data
106 contained in Big Data ? Value referring to the key important features of the data. This is defined by the added-

107 value that the collected data can bring to the intended processes. ? Veracity means the degree in which the
108 leaders trust the information to make a decision.

109 Big Data Analytics (BDA) can offer a variety of security dimensions in network traffic management, access
110 patterns in web transactions, configuration of network servers, network data sources, and user credentials. These
111 activities have brought a huge revolution in the domains of security management, identity and access management,
112 fraud prevention and governance, risk and compliance. However, there is also a lack of in-depth technical
113 knowledge regarding basic BDA concepts, Hadoop, Predictive Analytics, and Cluster Analysis, etc. With these
114 limitations in mind, appropriate steps can be taken to build on the skills and competences on security analytics.
115 There is lack of infrastructure to support such innovations, lack of skilled data scientists and lack of policies or
116 legislation that promote such innovations.

117 The supervised machine learning algorithm which can be used for both classification or regression challenges is
118 called the Support Vector Machine (SVM). The original training data can be transformed into a higher dimension
119 where it becomes separable by using the SVM algorithm which searches for the optimal linear separating
120 hyperplane. Estimations of the relationships among variables depends mainly on the statistical process of
121 regression analysis. The independent variables determine the estimation target. The regression function can
122 be linear as in linear regression, or a common sigmoid curve for the logistic function.

123 The easiest and simplest supervised machine learning algorithm which can solve both classification and
124 regression problems is the k-nearest neighbors (KNN) algorithm. Both the KNN and SVM can be applied
125 to finding the optimal handover solutions in heterogeneous networks constituted by diverse cells. Given a set
126 of contextual input cues, machine learning algorithms have the capability to exploit the user context learned.
127 The Hidden Markov Model (HMM) is a tool designed for representing probability distributions of sequences of
128 observations. It can be considered a generalization of a mixture-based model, rather than being independent of
129 each other. The list of supervised learning algorithms includes Regression models, Knearest neighbors, Support
130 Vector Machines, and Bayesian Learning [39].

131 Common examples of generative models that may be learned with the aid of Bayesian techniques include, but
132 are not limited to, the Gaussians mixture model (GM), expectation maximization (EM), and hidden Markov
133 models (HMM) [3, p. 445]. In Table ??, we summarize the basic characteristics and applications of supervised
134 machine learning algorithms.

135 Table ??: Various attack descriptions (Source: [7]) Attack Type Description DoS Denial of service: an attempt
136 to make a network resource unavailable to its intended users; temporarily interrupt services of a host connected
137 to the Internet

138 3 Scan

139 A process that sends client requests to a range of server port addresses on a host to find an active port.

140 4 Local Access

141 The attacker has an account on the system in question and can use that account to attempt unauthorized tasks.

142 5 User to root

143 Attackers access a user account on the system and are able to exploit some vulnerability to gain root access to the
144 system Data Attackers involve someone performing an action that they may be able to do on a given computer
145 system, but that they are not allowed to do according to policy.

146 6 b) Statement of the problem

147 Firewall protection has proved to be inadequate because of gross limitations against external threats The fact
148 is that the most network-centric cyberattacks are carried out by intelligent agents such as computer worms and
149 viruses; hence, combating them with intelligent semi-autonomous agents that can detect, evaluate, and respond
150 to cyberattacks has become a requirement [25]. The rapid development of computing and digital technologies,
151 the need to revamp cyberdefense strategies has become a necessity for most organisations [6]. As a result, there
152 is an imperative for security network administrators to be more flexible, adaptable, and provide robust cyber
153 defense systems in real-time detection of cyber threats.

154 The key problem is to evaluate Machine Learning (ML) and Big Data Analytics (BDA) paradigms for use in
155 Cybersecurity.

156 7 c) Purpose of study

157 The research is purposed to evaluate Machine Learning and Big Data Analytics paradigms for use in
158 Cybersecurity.

159 8 d) Research objectives

160 The research objectives are to:

9 Literature Review a) Overview

Computers, phones, internet and all other information systems developed for the benefit of humanity are susceptible to criminal activity [5]. Cybercrimes consist of offenses such as computer intrusions, misuse of intellectual property rights, economic espionage, online extortion, international money laundering, non-delivery of goods or services, etc. [13]. Intrusion detection and prevention systems (IDPS) include all protective actions or identification of possible incidents, and analysing log information of such incidents [4]. [6] recommends the use of various security control measures in an organisation. Various attack descriptions from the outcome of the research by [7] are shown on Table ???. The monotonic increase in an assortment of cyber threats and malwares amply demonstrates the inadequacy of the current countermeasures to defend computer networks and resources. To alleviate the problems of classical techniques of cyber security, research in artificial intelligence and more specifically machine learning is sought after [1], [2]. To enhance the malware and cyber-attack detection rate, one can apply deep learning architectures to cyber security.

10 b) Classical Machine Learning (CML)

Machine Learning (ML) is a field in artificial intelligence where computers learn like people. We present and briefly discuss the most commonly used classical machine learning algorithms.

11 i. Logistic Regression (LR)

As an idea obtained from statistics and created by [17], logistic regression is like linear regression, yet it averts misclassification that may occur in linear regression. Unlike linear regression, logistic regression results are basically either '0' or '1'. The efficacy of logistic regression is mostly dependent on the size of the training data.

ii. Naive Bayes (NB) Naive Bayes (NB) classifier is premised on the Bayes theorem which assumes independence of features. The independence assumptions in Naive Bayes classifier overcomes the curse of dimensionality.

12 iii. Decision Tree (DT)

A Decision tree has a structure like flow charts, where the root node is the top node and a feature of the information is denoted by each internal node. The algorithm might be biased and may end up unstable since a little change in the information will change the structure of the tree.

13 iv. K-Nearest Neighbor (KNN)

K-Nearest Neighbor (KNN) is a non-parametric approach which uses similarity measure in terms of distance function classifiers other than news cases. KNN stores the entire training data, requires larger memory and so is computationally expensive.

14 v. Ada Boost (AB)

Ada Boost (AB) learning algorithm is a technique used to boost the performance of simple learning algorithms used for classification. Ada Boost constructs a strong classifier using a combination of several weak classifiers. It is a fast classifier and at the same time can also be used as a feature learner. This may be useful in tasks that use imbalanced data analysis.

15 vi. Random Forest (RF)

Random forest (RF), as an ensemble tool, is a decision tree derived from a subset of observations and variables. The Random Forest gives better predictions than an individual decision tree. It uses the concept of bagging to create several minimal correlated decision trees.

16 vii. Support Vector Machine (SVM)

Support Vector Machine (SVM) belongs to the family of supervised machine learning techniques, which can be used to solve classification and regression problems. SVM is a linear classifier and the classifier is a hyper plane. It separates the training set with maximal margin. The points near to the separating hyper plane are called support vectors and they determine the position of hyper plane.

17 c) Modern Machine Learning

Deep learning is a modern machine learning which has the capability to take raw inputs and learns the optimal feature representation implicitly. This has performed well in various long standing artificial intelligence tasks [3]. Most commonly used deep learning architectures are discussed below in detail.

18 i. Deep Neural Network (DNN)

An artificial neural network (ANN) is a computational model influenced by the characteristics of biological neural networks. The family of ANN includes the Feed forward neural network (FFN), Convolutional neural network

211 and Recurrent neural network (RNN). FFN forms a directed graph in which a graph is composed of neurons
212 named as mathematical unit. Each neuron in i th layer has connection to all the neurons in $i + 1$ th layer.

213 Each neuron of the hidden layer denotes a parameter h that is computed by $h_i(x) = f(w_i^T x + b_i)$ (1)
214 $h_i: \mathbb{R}^d \rightarrow \mathbb{R}$ (2) $f: \mathbb{R} \rightarrow \mathbb{R}$ (3)

215 Where $w_i \in \mathbb{R}^{d \times d}$, $b_i \in \mathbb{R}^d$, d_i denotes the size of the input, f is a non-linear activation function,
216 ReLU.

217 The traditional examples of machine learning algorithms include Linear regression, Logistic regression, Linear
218 discriminant analysis, classification and regression trees, Naïve bayes, K-Nearest Neighbour (K-NN), Kmeans
219 clustering Learning Vector Quantization (LVQ), Support Vector Machines (SVM), Random Forest, Monte Carlo,
220 Neural networks and Qlearning. Take note that:

221 i. Supervised Adaptation is carried out in the execution of system at every iteration. Unsupervised
222 Adaptation follows trial and error method. Based on the obtained fitness value, computational model is
223 generalized to achieve better results in an iterative approach.

224 ii. The future of AI in the fight against cybercrimes Exeriments showed that NeuroNet is effective against
225 low-rate TCP-targeted distributed DoS attacks. [19] presented the Intrusion Detection System using Neural
226 Network based Modeling (IDS-NNM) which proved to be capable of detecting all intrusion attempts in the
227 network communication without giving any false alerts [20].

228 The characteristics of NIC algorithms are partitioned into two segments such as swarm intelligence and
229 evolutionary algorithm. The Swarm Intelligence-based Algorithms (SIA) are developed based on the idea of
230 collective behaviours of insects in colonies, e.g. ants, bees, wasps and termites. Intrusion detection and prevention
231 systems (IDPS) include all protective actions or identification of possible incidents and analysing log information
232 of such incidents [4].

233 19 d) Big Data Analytics and Cybersecurity

234 Big Data Analytics requires new data architectures, analytical methods, and tools. Big data environments ought
235 to be magnetic, which accommodates all heterogeneous sources of data. Instead of using mechanical disk drives, it
236 is possible to store the primary data-base in silicon-based main memory, which improves performance. According
237 to [21], there are four critical requirements for big data processing. The first requirement is fast data loading.
238 The second requirement is fast query processing. The «Map» function in Hadoop accordingly partitions large
239 computational tasks into smaller tasks, and assigns them to the appropriate key/value pairs.

240 Behavioral analytics provide information about the behavioral patterns of cybersecurity events or malicious
241 data [2]. Forensics analytics locate, recover and preserve reliable forensic artefacts from specifically identified
242 cybersecurity events or attacks [22]. Forecast analytics attempt to predict cybersecurity events using forecast
243 analytics models and methodologies [23]. Threat intelligence helps to gather threats from big data, analyze
244 and filter information about these threats and create an awareness of cybersecurity threats [2].

245 The situation awareness theory postulated by [24] posits that the success of a cybersecurity domain depends
246 on its ability to obtain real-time, accurate and complete information about cybersecurity events or incidents [20].
247 The situation awareness model consists of situation awareness, decisions and action performance as shown in
248 Figure 3.

249 There is consensus in prior literature that cyber security has evolved to become a problem for big data
250 analytics. This is due to the understanding that the transformation and expansion of the cyberspace [16] has
251 rendered traditional intrusion detection and malware detection systems obsolete. Further, even the data mining
252 models that have been used in the past are no longer sufficient for the challenges in cyber security [16].

253 20 Figure 3: Simplified Theoretical Model Based on Situation 254 Awareness

255 A big data analytics model for cybersecurity can be evaluated on the basis of its agility and robust [16].

256 According to [25], Big Data is defined not only by the amount of the information that it delivers but also
257 by its complexity and by the speed that it is analyzed and delivered. With reference to [26], Big Data can be
258 defined as a multi-faced data which combines the following characteristics: Veracity, Variety, Volume, Velocity
259 and Value.

260 21 e) Advances in Cloud Computing

261 Cloud computing is about using the internet to access someone else's software running on someone else's hardware
262 in someone else's data center [27]. Cloud Computing is essentially virtualized distributed processing, storage,
263 and software resources and a service, where the focus is on delivering computing as a on-demand, pay-as-you-go
264 service. On-demand self-service: A consumer can unilaterally provision computing capabilities such as server
265 time and network storage as needed automatically, without requiring human interaction with a service provider.

22 Broad network access: Heterogeneous client platforms available over the network come with numerous capabilities that enable provision of network access.

Resource pooling: Computing resources are pooled together in a multi-tenant model depending on the consumer demand in a location independent manner.

23 Situational Awareness

Decisions ?Capabilities ?Interface ?Mechanisms ?Stress & Workload ?Complexity ?Workload ?Automation Action Performance ?Capabilities ?Interface ?Mechanisms ?Stress & Workload ?Complexity ?Workload ?Automation

Rapid elasticity: This is when unlimited capabilities are rapidly and elastically provisioned or purchased to quickly scale out; and rapidly released to quickly scale in.

Measured service: A transparent metering capability can be automatically controlled and optimized in cloud systems at some level of abstraction appropriate to the type of service. ii. Platform as a service (PaaS) PaaS provides to the consumer infrastructure for third-party applications. Just like in SaaS the consumer does not manage or control the underlying cloud infrastructure including network, servers, operating systems, or storage, but does have control over the deployed applications and possibly configuration settings for the application-hosting environment [29]; [32]. Examples include Windows Azure, Apache Stratos, Google App Engine, Cloud Foundry, Heroku, AWS (Beanstalk) and Open Shift [33] & [34]. PaaS provides faster and more frequent delivery of functionality for the sake of direct support for business agility. PaaS provides an enabling environment for a consumer to run applications. A PaaS Cloud should be able to support various programming models for different types of Programming. PaaS is a Cloud Computing service that offers a computing platform and solution stack for users, and this may include the following: ? Language ? Operating System (OS) ? Database ? Middleware ? Other applications

iii. Infrastructure as a service (IaaS) This provisions processing, networks, storage, and other essential computing resources on which the consumer is then able to install and run arbitrary software, that can include operating systems (Virtual machines (VM), appliances, etc.) and applications [29]; [32]. Common global examples include Amazon Web Services (AWS), Cisco Metapod, Microsoft Azure, Rackspace and the local ones include TelOne cloud services and Dandemutande [33]. IaaS is a Cloud service that allows existing applications to run on its hardware. It rents out resources dynamically wherever they are needed.

24 Services include:

? Compute Servers ? Data Storage ? Firewall ? Load Balancer f) Cloud Deployment Models

The three commonly-used cloud deployment models are private, public, and hybrid. An additional model is the community cloud. However, this is less commonly used. In a Cloud context the term deployment basically refers to where the software is made available, in other words where it is running.

25 i. Private Cloud

The private cloud is normally either owned or exclusively used by a single organization. The services and infrastructure are permanently kept on a private network, the hardware and software are dedicated solely to the particular organisation. The service provider or the particular organization may manage the physical infrastructure. The major advantage of this model is the improved security as resources are not shared with others thereby allowing for higher levels of control and security [35]. The cloud infrastructure is provisioned for use by the general public. The public cloud is sold to the public, as a mega-scale infrastructure, and is available to the general public. [12] further clarifies that cloud services are provided on a subscription basis to the public. It is typically based on a pay-per-use model. The advantages include lower costs, near-unlimited scalability and high reliability [35].

Examples include Amazon (EC2), IBM's Blue Cloud, Sun Cloud, Google App Engine and Windows Azure [37].

26 iii. Hybrid Cloud

A hybrid cloud model is a mix of two or more cloud deployment models such as private, public or hybrid [36]; [38]. This model requires determining the best split between the public and private cloud components. The advantages include control over sensitive data (private cloud), flexibility, i.e. ability to scale to the public cloud whenever needed and lastly allows for ease transitioning to the cloud through gradual migration [35]. The use of standardized or proprietary technology allows for data and application portability [39].

27 iv. Community Cloud

This model is provisioned for exclusive use by a particular community of consumers bound by shared interests (e.g., policy and compliance considerations, mission and security requirements). A community cloud shares computing resources among several organizations, and can be managed by either organizational IT resources or third-party

321 providers [29]. A typical example is the U.S.-based exclusive IBM Soft Layer cloud which is dedicated for use
322 by federal agencies only. This approach builds confidence in the platform, which cloud consumers will use to
323 process their sensitive workloads [37].

324 v. Cloud computing benefits Cloud computing services are delivered when they are needed in the quantity
325 needed at a certain time. Cloud computing has many benefits for the organizations and these include cost savings,
326 scalability, anytime anywhere access, use of latest software versions, energy saving and quick rollout of business
327 solutions. The cost effectiveness and efficiency of the cloud platforms is tempting most organizations to migrate
328 to the cloud and enjoy a wide range of general benefits [40] which according to [41] include: Traditionally, these
329 networks rely on dedicated hardware-based network equipment and their functions to provide communication
330 services. However, this reliance is becoming increasingly inflexible and inefficient, especially in dealing with traffic
331 bursts for example during large crowd events. NFV strives to overcome current limitations by (1) implementing
332 network functions in software and (2) deploying them in a virtualized environment. The resulting virtualized
333 network functions (VNFs) require a virtual infrastructure that is flexible, scalable and fault tolerant.

334 The growing maturity of container-based virtualization and the introduction of production-grade container
335 platforms promotes containers as a candidate for the implementation of NFV infrastructure (NFVI). Containers
336 offer a simplified method of packaging and deploying applications and services. Virtualization is basically making
337 a virtual image or "version" of something usable on multiple machines at the same time. This is a way of
338 managing the workload by transforming traditional computing to make it more scalable, efficient and economical.
339 Virtualization can be applied to hardware-

340 28 h) Why Virtualization?

341 With virtualization, one can attain better utilization rate of the resources of the service providers, increased
342 ROI for both the service providers and the consumers, and promotes the green IT by reducing energy wastage.
343 Virtualization technology has the drawbacks of the chance of a single point of failure of the software achieving the
344 virtualization and the performance overhead of the entire system due to virtualization. Virtualization in general
345 has tremendous advantages. The advantages of virtual machines are as follows:

346 ? Where the physical hardware is unavailable, run the operating systems, ? Easier to create new machines,
347 backup machines, etc., ? Use of "clean" installs of operating systems and software for software testing ? Emulate
348 more machines than are physically available, ? Timeshare lightly loaded systems on one host, ? Debug problems
349 (suspend and resume the problem machine), ? Easy migration of virtual machines, ? Run legacy systems! Two or
350 more CPUs can work together on the same chip in multicore technology as a single integrated circuit (IC). These
351 single ICs are called a die. Multicore technology can be used to speed up the processing in a multitenant cloud
352 environment. Multicore architecture has become the recent trend of high-performance processors, and various
353 theoretical and case study results illustrate that multicore architecture is scalable with the number of cores.

354 Most of the software vendors raised a complaint that their application is not supported in a virtual state or
355 will not be supported if the end-user decides to virtualize them. To accommodate the needs of the industry
356 and operating environment, to create a more efficient infrastructure -virtualization process has been modified
357 as a powerful platform, such that the process virtualization greatly revolves around one piece of very important
358 software. This is called a hypervisor. Thus, a VM must host an OS kernel.

359 29 i) Compare and Contrast Between Virtualization and Con- 360 tainerization

361 Virtualization allows the running of multiple operating systems on a single physical system and share the
362 underlying hardware resources. Virtualization entails abstraction and encapsulation. However, Clouds rely heavily
363 on virtualization, whereas Grids do not rely on virtualization as much as clouds. In Virtualization, a hypervisor
364 is a piece of computer software that creates and runs virtual machines.

365 Instead of installing the operating system as well as all the necessary software in a virtual machine, the
366 docker images can be easily built with a Dockerfile since the hardware resources, such as CPU and memory,
367 will be returned to the operating system immediately. Therefore, many new applications are programmed into
368 containers. Cgroups allow system administrators to allocate resources such as CPU, memory, network, or any
369 combination of them, to the running containers. This is illustrated in Figure 7 below.

370 30 Virtualization

371 Virtualization is the optimum way to enhance resource utilization in efficient manner. It refers to the act of
372 creating a virtual (similar to actual) variations of the system. Physical hardware is managed with the help of
373 software and converted into the logical resource that will be in a shared pool or can be used by the privileged
374 user. This service is known as VMs we can say Infrastructure as a service. Virtualization is the base of any
375 public and private cloud development. Most of the public cloud providers such as Amazon EC2, Google Compute
376 Engine and Microsoft Azure leverage virtualization technologies to power their public cloud infrastructure [1].
377 The core component of virtualization is Hypervisors.

31 Hypervisor

It is a software which provides isolation for virtual machines running on top of physical hosts. The thin layer of software that typically provides capabilities to virtual partitioning that runs directly on hardware, It provides a potential for virtual partitioning and responsible for running multiple kernels on top of the physical host. This feature makes the application and process isolation very expensive. There will be a big impact if computer resources can be used more efficiently. The most popular hypervisors today are VMware, KVM, Xen, and HyperV.

Basically, a container is nothing but more than a virtual file system which are isolated with some Linux kernel features, such as namespaces and process groups, from the main physical system. Through containers framework it offers an environment as close as desirable one as we want from a VM but without the overhead that comes with running on an another kernel and simulating all the hardware. Due to lightweight nature of containers, more containers can run per host than virtual machines per host. Unlike containers, virtual machine require emulation layers (either software or hardware), which consume more resources and add additional overhead.

Containers are different from Virtualization with respect to the following aspects: 1. Simple: Easy sharing of a hardware resources clean command line interface, simple REST API. 2. Fast:-Rapid provisioning, instant guest boot, and no virtualization overhead so as fast as bare metal. 3. Secure: Secure by default, combine all available kernel security feature with App Armor, user namespaces, SECCOMP. 4. Scalable: The quality-of-service may be broadcast from the from a single container on a developer laptop to a container per host in a data centre. This is also includes remote image services with Extensible storage and networking. 5. Control groups (cgroups): This is a kernel-provided mechanism for administration, grouping and tracking through a virtual file system. Docker containers share the operating system and important resources, such as depending libraries, drivers or binaries, with its host and therefore they occupy less physical resources.

32 III.

33 Research Methodology a) Presentation of the methodology

The Pragmatism paradigm was used in this research and this is intricately related to the Mixed Methods Research (MMR).

Philosophers inclined to the pragmatic paradigm subscribe to the worldview that says it is impossible to access the truth of the real world by employing a single scientific method as supported by the Positivist paradigm or construct social reality under Interpretivist paradigm. In this research, an Interpretivist or Constructivist paradigm was used, as is illustrated on Table 2 below. Cybersecurity is a huge area for consideration and in order to address problems within it, there is need for contextualisation. This is a clear indication that there are multiple realities out there in the world of cybersecurity as supported by the Interpretivist paradigm.

The Research methodology is a way of solving a research problem thoroughly and meticulously and includes steps followed in carrying out the research and the reasoning behind [52]. Research methodology can also be viewed as a procedural or step by step outline or framework within which research is done. Research methodology can be quantitative, qualitative or mixed. Table ?? below shows the differences between qualitative and quantitative research methodologies. The Mixed Methods Research methodology was used. In a mixed methods methodology the researcher mixes both qualitative and quantitative data and employs the practices of both qualitative and quantitative research. It is also underpinned by the pragmatic paradigm.

34 Table 3: Differences between qualitative and quantitative methodologies

Difference with respect to:
Quantitative methodology Qualitative methodology

35 Supporting philosophy

Rationalism. Humans acquire knowledge through their capacity to reason.

36 Empiricism. Humans acquire knowledge through sensory experiences

37 Approach to inquiry

Structured or rigid /predetermined methodology Unstructured /flexible methodology

38 Main purpose of investigation To quantify the extend of variation in a situation or phenomenon

To describe variation in a phenomenon or situation

430 **39 Measurement of variables** Emphasis is on some form of either
431 measurement or classification of variables

432 **40 Emphasis is on the description of variables** Sample size
433 Emphasis is put on a greater sample size Fewer cases

434 **41 Focus of inquiry**

435 Narrows focus in terms of extent of inquiry but draws together required information from a bigger number of
436 respondents Covers multiple issues but draws together required information from a smaller number of respondents

437 **42 Data analysis**

438 Variables are put into frequency distributions or other statistical procedures
439 Responses or observational data is used to identify themes and their descriptions

440 **43 Communication of findings**

441 Organisation is more analytic in nature, drawing inferences and conclusions and testing strength between variables
442 and their relationship
443 Organisation is more descriptive and narrative in nature Source: [51] i.

444 **44 Research approach and philosophy** Research approach

445 The researcher adopts a qualitative approach in form of focus group discussion to research. Since the analysis is
446 done to establish differences in data analytics models for cybersecurity without the necessity of quantifying the
447 analysis [42] .

448 **45 Research philosophy**

449 The researcher adopts a postmodern philosophy to guide the research. Firstly the researcher notes that the
450 definition, scope and measurement of cybersecurity differs between countries and across nations [15]. Further,
451 the post-modern view is consistent with descriptive research designs which seek to interpret situations or models
452 in their particular contexts [43].

453 **46 ii. Research design and methods**

454 **47 Research design**

455 The researcher adopts a descriptive research design since the intention is to systematically describe the facts and
456 characteristics of big data analytics models for cybersecurity. The purpose of the study is essentially an in-depth
457 description of the models [42].

458 **48 Research methods**

459 A case study research method was adopted in this study. In this respect each data analytics model for
460 cybersecurity is taken as a separate case to be investigated in its own separate context [43]. Prior research
461 has tended to use case studies in relation to the study of cybersecurity [15]. However, the researcher develops a
462 control case that accounts for an ideal data analytics model for cybersecurity for comparative purposes.

463 **49 b) Population and sampling i. Population**

464 The research population for the purpose of this study consists of all data analytics models for cybersecurity that
465 have been proposed and developed in literature, journals, conference proceedings and working papers. This is
466 consistent with previous research which involves a systematic review of literature [21].

467 **50 ii. Sample**

468 The researcher identified two data analytics models or frameworks from a review of literature and the sample size
469 of 8. Eight participants in total were interviewed. However, while this may be limited data, it will be sufficient
470 for the present needs of this study. Research in future may review more journals to identify more data analytics
471 models which can be applied to cybersecurity.

472 **51 c) Sources and types of data**

473 The researcher uses secondary data in order to investigate the application of data analytics models in
474 cybersecurity.

52 d) Model for analysis

475
476 In analyzing the different data analytics models for cybersecurity the researcher makes reference to the
477 characteristics of an ideal data analytics model for cybersecurity. In constructing an ideal model, the researcher
478 integrates various literature sources. The basic framework for big data analytics model for cybersecurity consists
479 of three major components which are big data, analytics, and insights [16]. However, a fourth component may
480 be identified as prediction (or predictive analytics) [21]. This is depicted in Figure 8 below:

53 Big data

481
482 The first component in the big data analytics framework for cybersecurity is the availability of big data about
483 cybersecurity. Traditional sources of big data are systems logs and vulnerability scans [16]. However, sources
484 of big data about cybersecurity have extended to include computer-based data, mobile-based data, physical
485 data of users, human resources data, credentials, one-time passwords, digital certificates, biometrics, and social
486 media data [11]. Some authors identify sources of big data about cybersecurity as business mail, access control
487 systems, CRM system and human resources system, a number of pullers in linked data networks, intranet/
488 internet and IIoT/IoT, collectors and aggregators in social media networks and external news tapes [23]. Big
489 data about cybersecurity should be imported from multiple sources to ensure effectiveness in detection and
490 prediction of possible threats [17]. Further, some authors specify the characteristics of security data as consisting
491 of heterogeneous format, diverse semantic and correlation across data sources and classify them into categories
492 for example non-semantic data, semantic data and security knowledge data [17].

54 Big data analytics

493
494 The address the concerns of big data about cybersecurity, more robust big data analytics models for cybersecurity
495 have been developed in data mining techniques and machine learning [16]. Big data analytics employ data mining
496 reactors and algorithms, intrusion and malware detection techniques and vector machine learning techniques for
497 cybersecurity [16]. However, it has been observed that adversarial programs have tended to modify their behavior
498 by adapting to the reactors and algorithms designed to detect them [16]. Further, intrusion detection systems are
499 faced with challenges such as unbounded patterns, data nonstationarity, uneven time lags, individuality, high false
500 alarm rates, and collusion attacks [21]. This necessitates a multi-layered and multi-dimensional approach to big
501 data analytics for cybersecurity [17], [2]. In other words an effective big data analytics model for cybersecurity
502 must be able to detect intrusions and malware at every layer in the cybersecurity framework.

55 Insights for action

503
504 Big data analytics for cybersecurity should be able provide prioritized and actionable insights to cybersecurity
505 personnel. For example setting up effective network defenders that are able to detect flaws in the network and
506 be able to trace the source of threats or attacks [16]. Alternatively, cybersecurity personnel may update existing
507 network defenders in light of new prioritized insights about the cybersecurity system [16].The goal of analysts
508 should be to maximize utility derived from the cybersecurity system.

56 Predictive analytics

509
510 Predictive analytics refer to the application of a big data analytics model for cybersecurity to derive, from current
511 cybersecurity data, the likelihood of a cybersecurity event occurring in future [21]. In essence, a data analytics
512 model for cybersecurity should be able to integrate these components if it is to be effective in its major functions
513 of gathering big data about cybersecurity, analyzing big data about cybersecurity threats, providing actionable
514 insights and predicting likely future cybersecurity incidents.

57 e) Validity and Reliability

515
516 The researcher solicited comments from peers on the emerging findings and also feedback to clarify the biases
517 and assumptions of the researcher to ensure internal validity of the study [43]. The researcher also reliability or
518 consistency in research findings by explaining in detail the assumptions and theories underlying the study [43].

58 f) Summary of research methodology

519
520 In section 3, the researcher developed appropriate methodology for investigating the ideal data analytics models
521 for cybersecurity.

59 g) Possible Outcomes

522
523 The expected accuracy rate for the research should be according to Table 4 below, which shows the international
524 benchmark. The IDS is divided into either as a Host IDS (HIDS) or as a Network IDS (NIDS). Analysis of
525 the network traffic can be handled by a NIDS which distinguishes the unlicensed, illegitimate and anomalous
526 behavior on the network. Packets traversing through the network should generally be captured by the IDS using
527 network taps or span port in order to detect and flag any suspicious activity [10]. Anomalous behavior on the

528 specific device or malicious activity can be effectively detected by a device specific IDS. The vulnerability of
529 networks and susceptibility to cyber attacks is exacerbated by the use of wireless technology [12].

530 The gross inadequacies of classical security measures have been overtly exposed. Therefore, effective solutions
531 for a dynamic and adaptive network defence mechanism should be determined. Neural networks can provide
532 better solutions for the representative sets of training data [12]. [12] argues for the use of ML classification
533 problems solvable with supervised or semi-supervised learning models for the majority of the IDS. However, the
534 one major limitation of the work done by [12] is on the informational structure in cybersecurity for the analysis
535 of the strategies and the solutions of the players.

536 Autonomous robotic vehicles attract cyber attacks which prevent them from accomplishing the intrusion
537 prevention mission. Knowledge-based and vehicle-specific methods have limitations in detection which is
538 applicable to only specific known attacks [3]. The attack vectors of the attack scenarios used by [3] is shown on
539 Figure 11 below. In this experiment, the system is allowed to undertake several missions by the robotic vehicle
540 which diverts the robotic vehicle testbed. The practical experimental setup for the attack vectors used is shown
541 on Figure 12 below. Table ?? shows a comparison of the data mining techniques that can be used in intrusion
542 detection.

543 Intrusion attack classification requires optimization and enhancement of the efficiency of data mining
544 techniques. The pros and cons of each algorithm using the NSL-KDD dataset are shown on Table ?? below.

545 Table ??: Performance of Support Vector Machines, Artificial Neural Network, K-Nearest Neighbour, Naive-
546 Bayes and Decision Tree Algorithms An intrusion detection system determines if an intrusion has occurred,
547 and so monitors computer systems and networks, and the IDS raises an alert when necessary [4]. However, [4]
548 addressed the problems of Anomaly Based Signature (ABS) which reduces false positives by allowing a user to
549 interact with the detection engine and raising classified alerts. The advantages and disadvantages of ABSs and
550 SBSs are summarised on table, Table ??, below.

551 60 Table 7:

552 Advantages and disadvantages of ABSs and SBSs models (Source: [4])

553 An IDS must keep up track of all the data, networking components and devices involved. Additional
554 requirements must be met when developing a cloud-based intrusion detection system due to its complexity
555 and integrated services.

556 61 b) Support vector machine

557 Support Vector Machine is a classification artificial intelligence and machine learning algorithm with a set
558 containing of points of two types in X dimensional place. Support vector machine generates a (X-1) dimensional
559 hyperplane for separating these points into two or more groups using either linear kernel or nonlinear kernel
560 functions [7]. Kernel functions provides a method for polynomial, radial and multi-layer perception classifiers such
561 as classification of bank performance into four clusters of strong, satisfactory, moderate and poor performance.
562 The class of bank performance is defined by the function $Performance\ class = (w \cdot x) + b$

563 Where x is the input vector to the support vector classifier and w is the real vector of weights and f is
564 the function that translates the dot product of the input and real vector of weights into desired classes of bank
565 performance. w is learned from the labeled training data set.

566 62 c) KNN algorithm

567 The K-NN algorithm is a non-parametric supervised machine learning technique that endeavors to classify a data
568 point from given categories with the support of the training dataset [7]. Predictions are performed for a new
569 object (y) by searching through the whole training dataset for the K most similar instances or neighbors. The
570 algorithm does this by calculating the Euclidean distance as follows:

571 63 ? =

572 Where X_s is the standardized value, X is the instance measure, mean and s.d are the mean and standard
573 deviation of instances. Lower values of K are sensitive to outliers and higher values are more resilient to outliers
574 and more voters are considered to decide the prediction.

575 64 d) Multi Linear Discriminant Analysis (LDA)

576 The Linear Discriminant Analysis is a dimensionality reduction technique. Dimensionality reduction is the
577 technique of reducing the amount of random variables under consideration through finding a set of principal
578 variables [7] which is also known as course of dimensionality. The LDA calculates the separability between n
579 classes also known as betweenclass variance. Let D_b be the distance between n classes.

580 65 (

581 $\sum_{i=1}^N (x_i - \bar{x})^2 = \sum_{i=1}^N x_i^2 - N\bar{x}^2$

582 Where \bar{x} the overall is mean, i

640 neural network-based deep learning (DL) approaches to predict Zimbabwean Bank's solvency. Future work will
641 focus on identifying more features that could possibly lead to poor bank performance and incorporate these in
642 our models to develop a robust early warning supervisory tool based on big data analytics, machine learning and
643 artificial intelligence.

644 The researcher analyses the two models that have been proposed in literature with reference to an ideal data
645 analytics model for cybersecurity presented in Section 3.

646 **70 Model 1: Experimental/ Prototype Model**

647 In the first case the researcher makes reference to the model presented in [23] which although developed in the
648 context of the public sector can be applied to the private sector organizations. The proposed model, it is to be
649 noted was demonstrated to be effective in integrating big data analytics with cybersecurity in a cost effective
650 way [23].

651 **71 Model 2: Cloud computing/Outsourcing**

652 The second model involves an organization outsourcing its data to a cloud computing service provider. Cloud
653 computing service providers usually have advanced big data analytics models, with advanced detection and
654 prediction algorithms and better state of the art cybersecurity technologies and better protocols because they
655 specialize in data and networks. However, it is to be noted that cloud computing service providers are neither
656 exempt nor immune from cyber-threats and attacks [11].

657 **72 Application of big data analytics models in cybersecurity**

658 There is overwhelming evidence to support this assertion with many infallible proofs that such application is not
659 only necessary in recent times but a means to survival [11], [23]. The researcher demonstrated by identifying
660 the characteristics of an effective data analytics model, the ideal model, that it is possible to evaluate different
661 models. In the third hypotheses the researcher postulated that, there is an appropriate big data analytics model
662 for cybersecurity for every institution. While the review of literature showed that institutions and countries adopt
663 different big data analytics models for cybersecurity, the researcher also demonstrated that beside the unique
664 requirements these models share major common characteristics for example reactors and detection algorithms
665 are usually present in every model but differ in terms of complexity. Further, using the models presented in
666 this Chapter it is worthy of note that many small organizations will usually adopt Model 2 whereas very large
667 organizations and sensitive public sector organizations will adopt Model 1. This may also explain why models
668 used may differ although the framework used in designing a data analytics model for cybersecurity in a cloud
669 computing services provider may share similar characteristics with that developed by an institution on its own.

670 **73 Summary of analysis**

671 In this section the researcher presented two models for adopting data analytics models to cybersecurity. The
672 first experimental or prototype model involves the design, and implementation of a prototype by an institution
673 and the second model involves the use serviced provided by cloud computing companies. The researcher also
674 demonstrated how this study addressed the hypotheses postulated. In the information era we are currently living
675 in, voluminous varieties of high velocity data are being produced daily, and within them lay intrinsic details and
676 patterns of hidden knowledge which should be extracted and utilized.

677 By applying such analytics to big data, valuable information can be extracted and exploited to enhance decision
678 making and support informed decisions. Thus, the support of big data analytics to decision making was depicted.

679 V.

680 **74 Conclusion**

681 Machine learning algorithms as part of Artificial Intelligence can be clustered into supervised, unsupervised,
682 semi-supervised, and reinforcement learning algorithms. The main characteristic of ML is the automatic data
683 analysis of large data sets and production of models for the general relationships found among data.

684 Big data analytics is not only about the size of data but also clinches on volume, variety and velocity of data
685 .Volume denotes big data as massive; velocity denotes the high speed of big data; variety denotes the diversity
686 of big data; veracity denotes the degrees of trustworthiness in big data; vocabulary denotes conformity of big
687 data to different schema, models and ontologies; and value denotes the cost and worth of big data. Big data has
688 necessitated the development of big data mining tools and techniques widely referred to as big data analytics.
689 Big data analytics refer to a combination of well-known tools and techniques for example machine learning, and
690 data mining, that are capable of leveraging useful data usually hidden in big data and creating an interface in
691 the form of linear and visual analytics.

692 The information that is evaluated in Big Data Analytics includes a mixer of unstructured and semi-structured
693 data, for instance, social media content, mobile phone records, web server logs, and internet click stream data.
694 Big data analytics makes use of analytic techniques such as data mining, machine learning, artificial learning,
695 statistics, and natural language processing. Big Data came into existence when the traditional relational database

696 systems were not able to handle the unstructured data generated by organization, social media, or from any other
 697 data generating source.

698 Passive data sources can include: Computer-based data, for example geographical IP location, computer
 699 security health certificates, keyboard typing and clickstream patterns, WAP data. Data over networks may be
 700 secured through the use of antivirus software, firewall, encryption, secure protocols, etc. However, hackers can
 701 always devise innovative ways of breaking into the network systems. An intrusion detection and prevention
 702 system is placed inside the network to detect possible network intrusions and, where possible, prevent the cyber
 attacks.

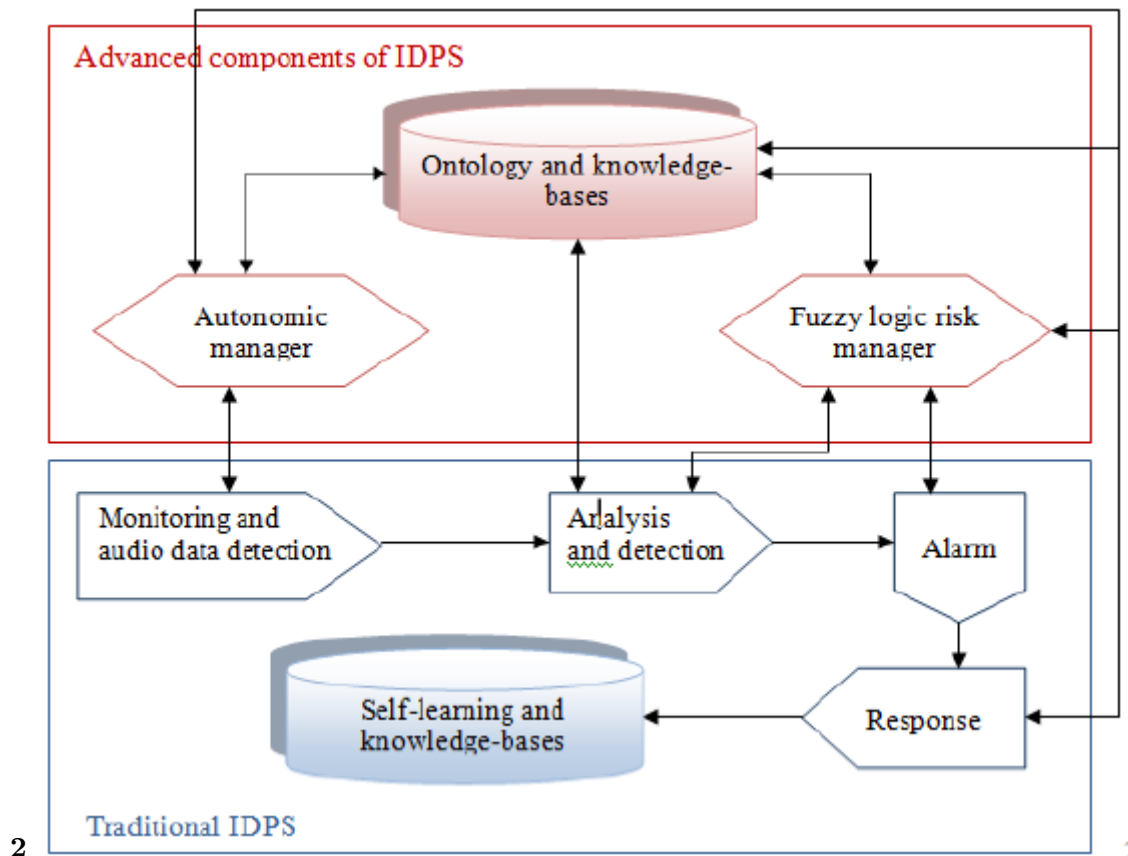


Figure 1: Figure 2 :

703

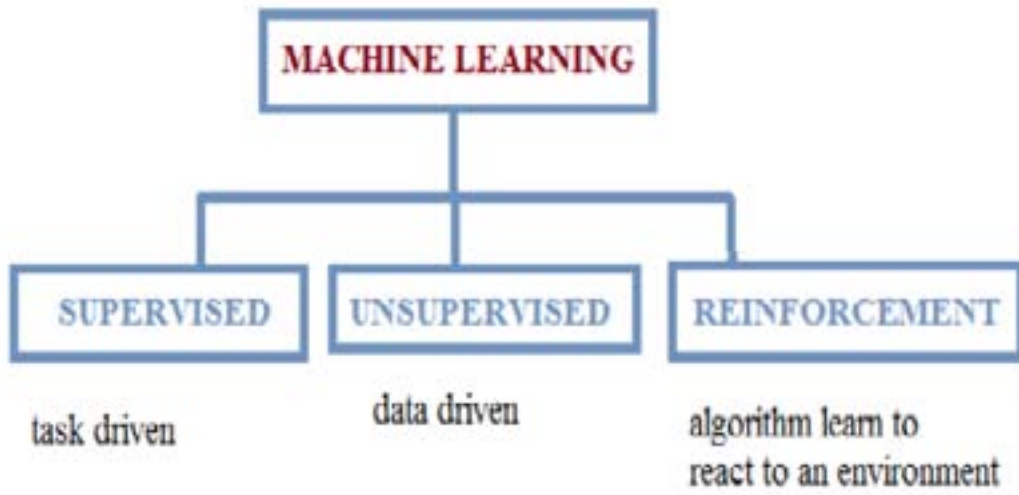


Figure 2:

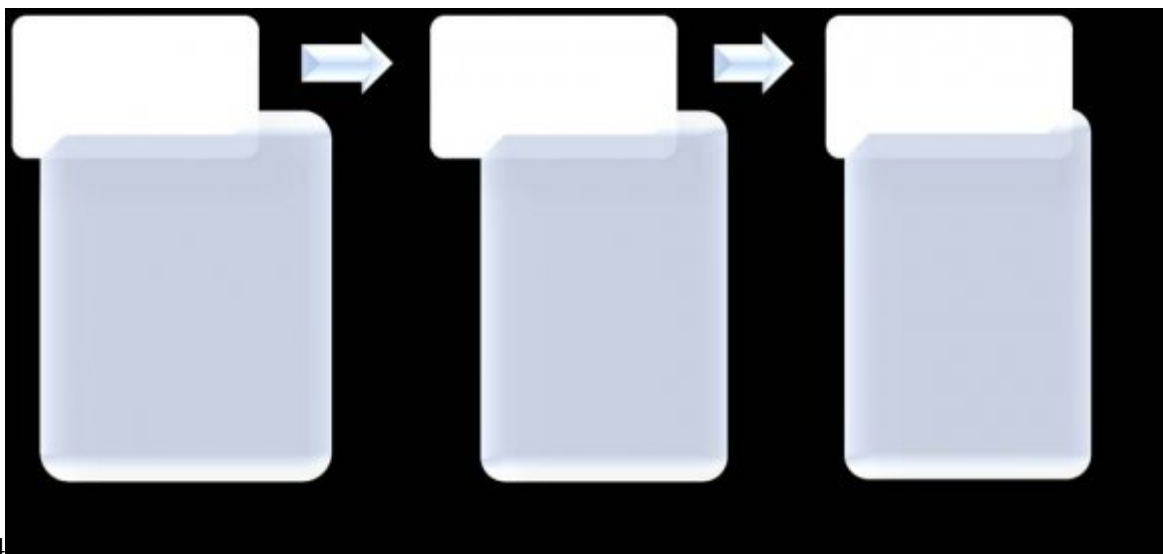
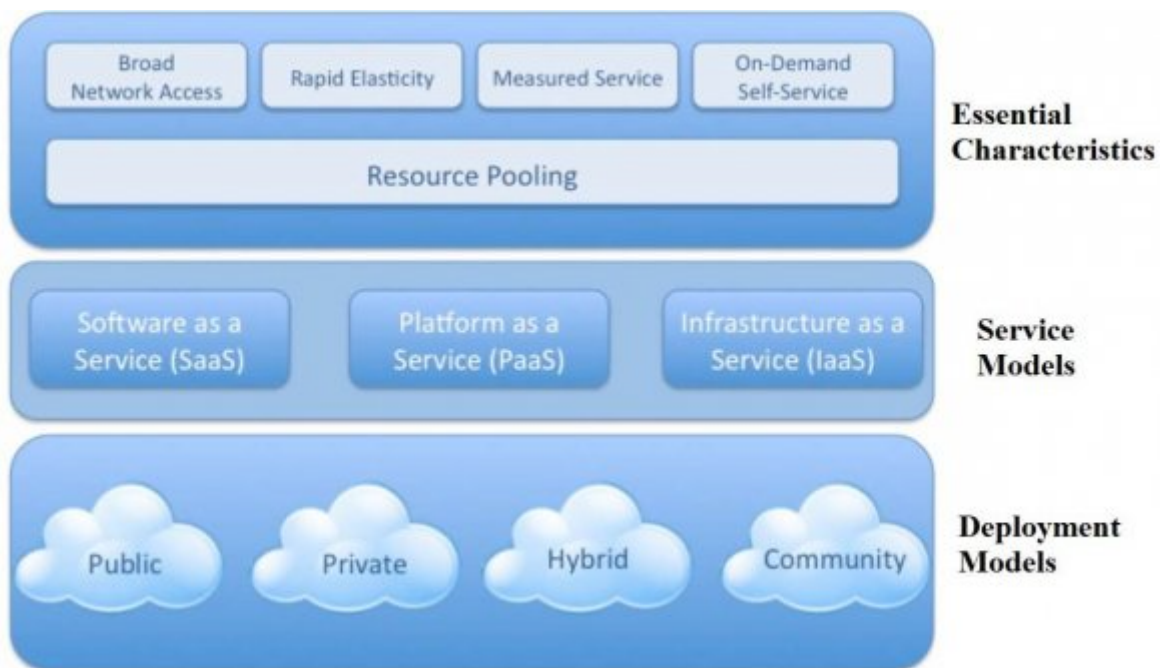


Figure 3: Figure 4 :



5

Figure 4: Figure 5 :



Figure 5: ?

The NIST Cloud Definition Framework

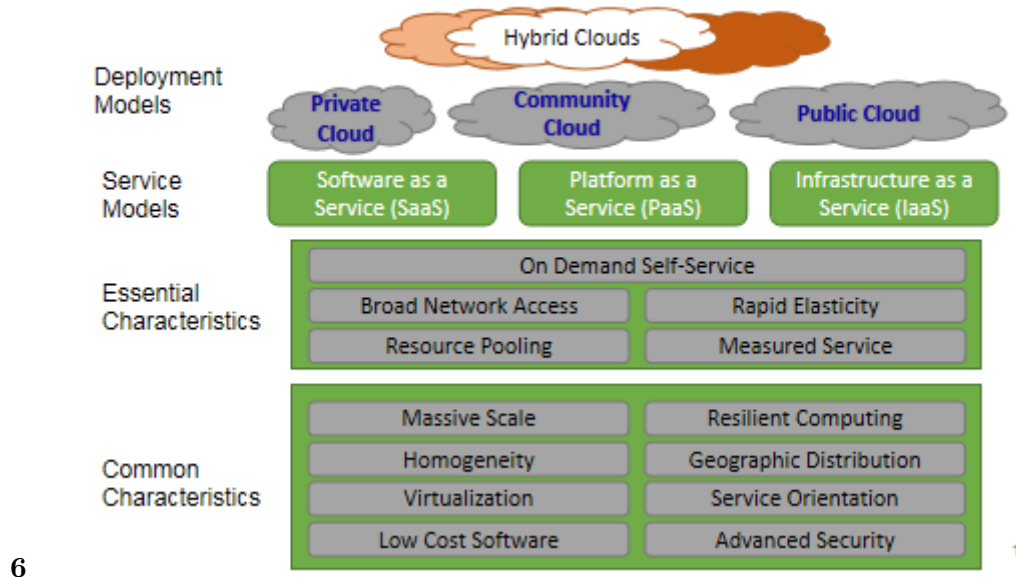
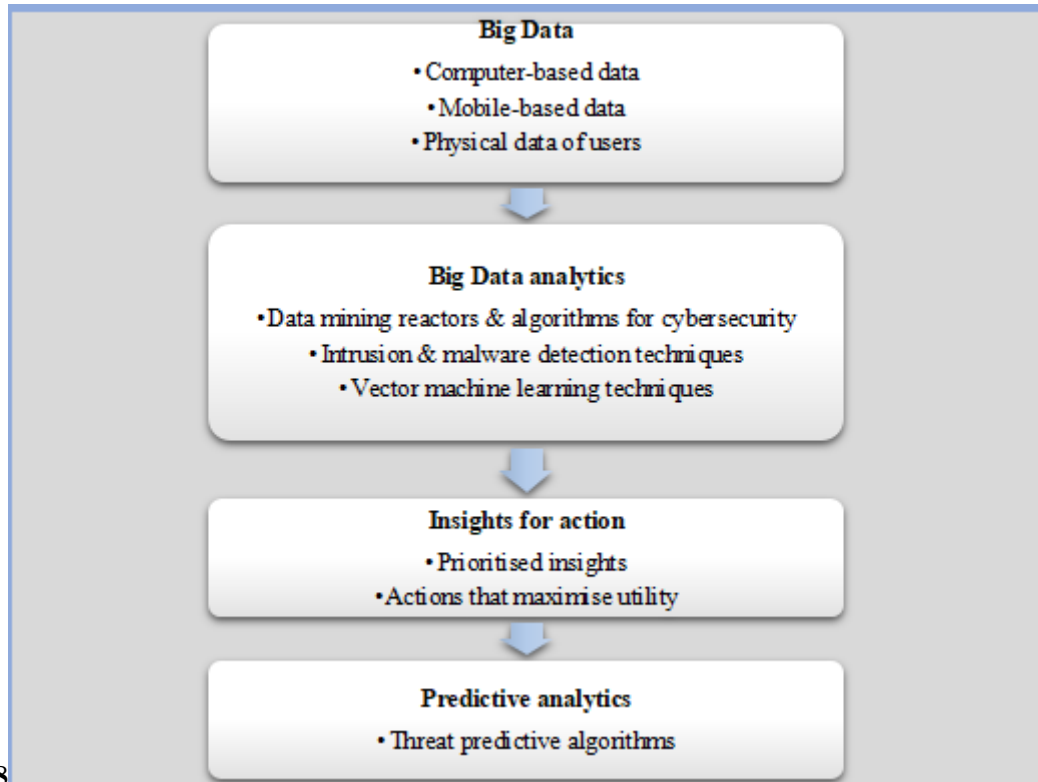


Figure 6: Figure 6 :

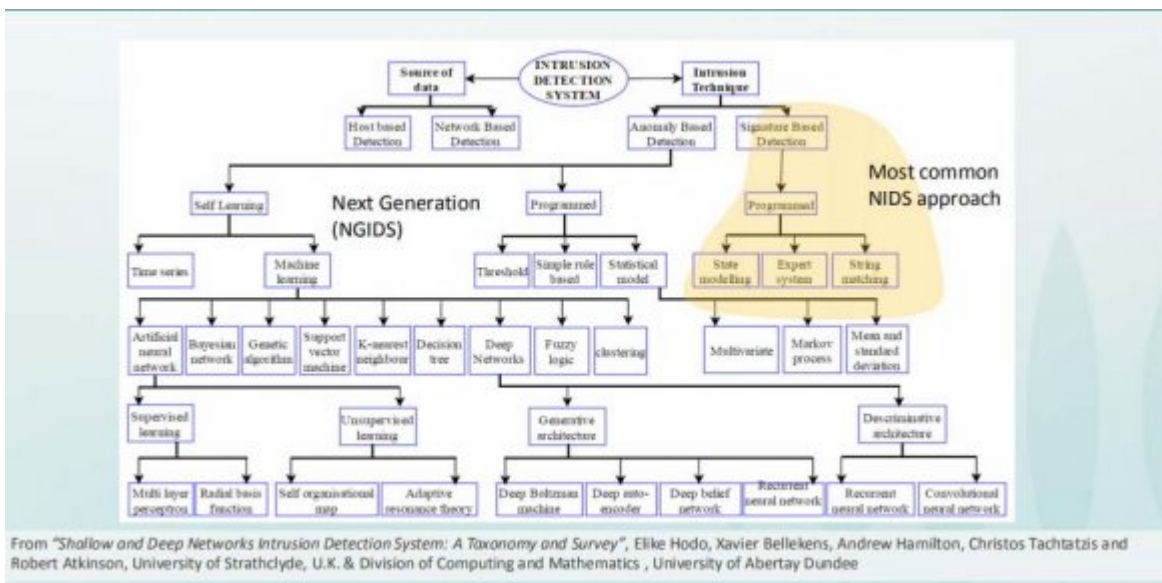


Figure 7: Figure 7 :



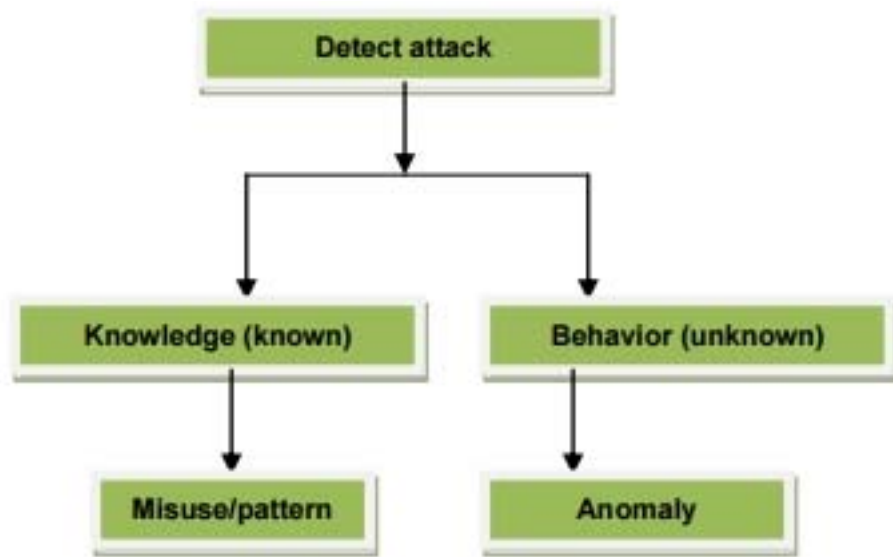
8

Figure 8: Figure 8 :



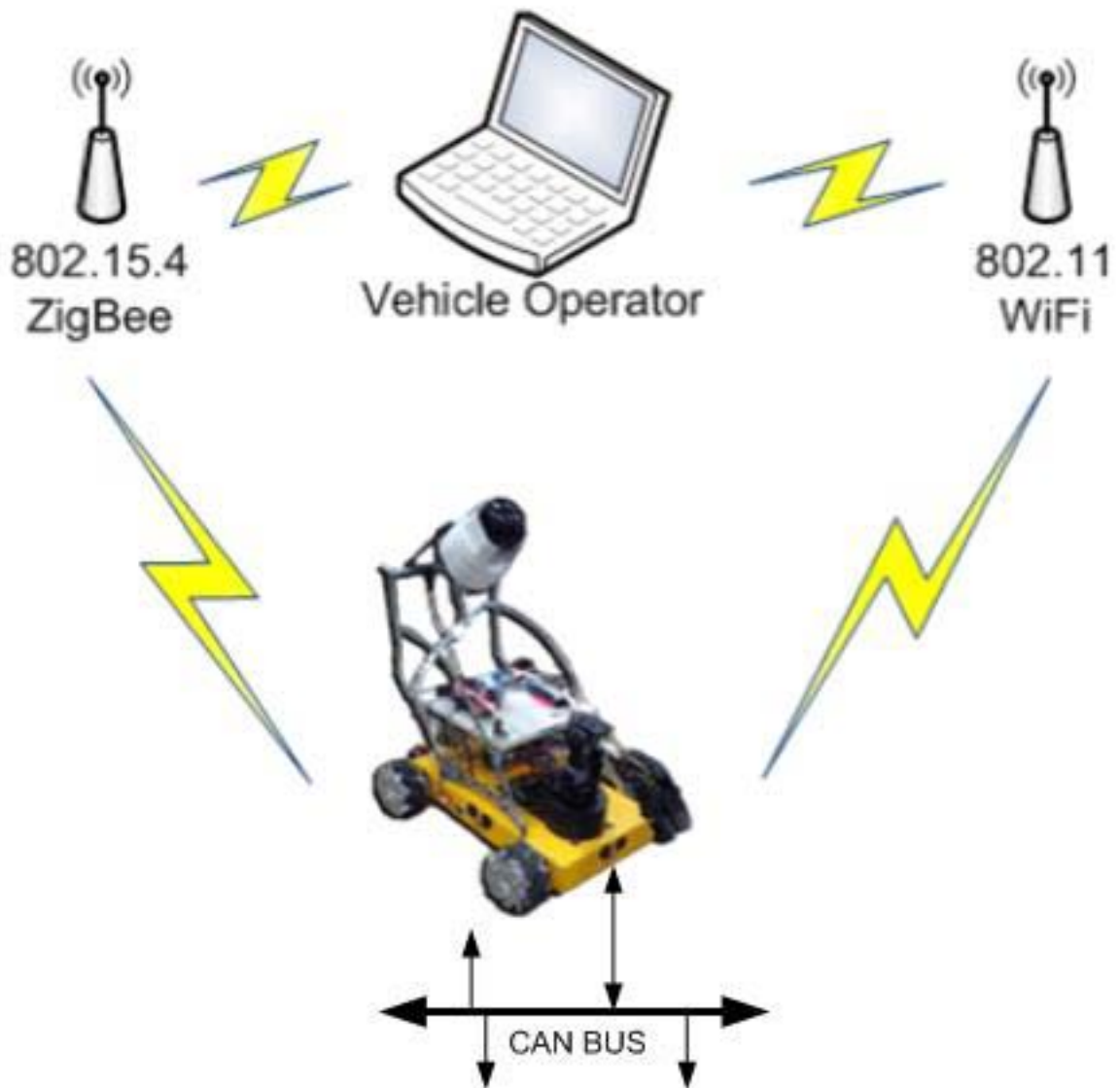
From "Shallow and Deep Networks Intrusion Detection System: A Taxonomy and Survey", Elike Hodo, Xavier Bellekens, Andrew Hamilton, Christos Tachtatzis and Robert Atkinson, University of Strathclyde, U.K. & Division of Computing and Mathematics, University of Abertay Dundee

Figure 9:



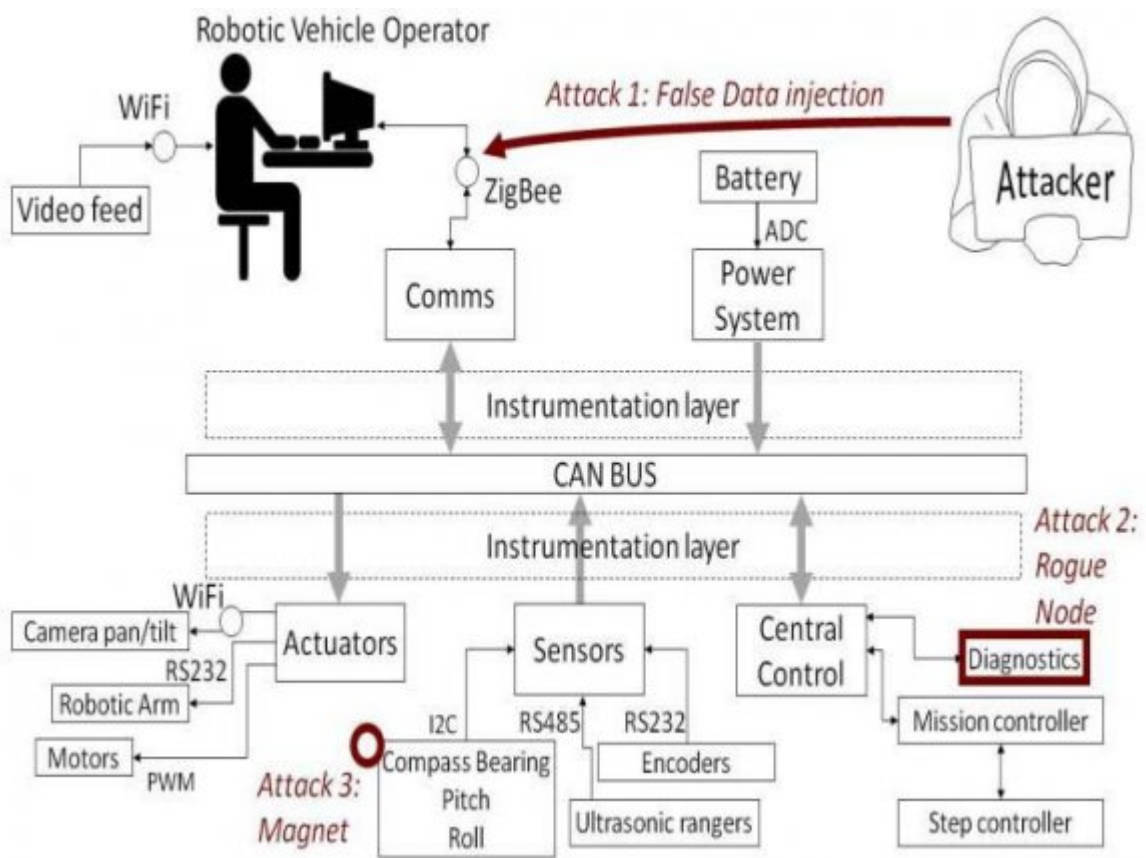
910

Figure 10: Figure 9 :Figure 10 :



11

Figure 11: Figure 11 :



12

Figure 12: Figure 12 :

Technique	Advantages	Disadvantages
Genetic Algorithm	<ul style="list-style-type: none"> - Finding a solution for any optimization problem. - Handling multiple solution search spaces. 	<ul style="list-style-type: none"> - Complexity to propose a problem space. - Complexity to select the optimal parameters - The need to have local searching technique for effective functioning
Artificial Neural Network	<ul style="list-style-type: none"> - Adapts its structure during training without the need to program it. 	<ul style="list-style-type: none"> - Not accurate results with test data as with training data
Naive Bayes Classifier	<ul style="list-style-type: none"> - Very simple structure. - Easy to update. 	<ul style="list-style-type: none"> - Not effective when there are high dependency between features.
Decision tree	<ul style="list-style-type: none"> - Easy to understand - Easy to implement 	<ul style="list-style-type: none"> - Works effectively only with attributes having discrete values. - Very sensitive to training sets, irrelevant features and noise.
K Mean	<ul style="list-style-type: none"> - Very Easy to understand. - Very simple to implement in solving clustering problems. 	<ul style="list-style-type: none"> - Number of clusters is not automatically calculated. - High dependency on initial centroids.

Figure 13:

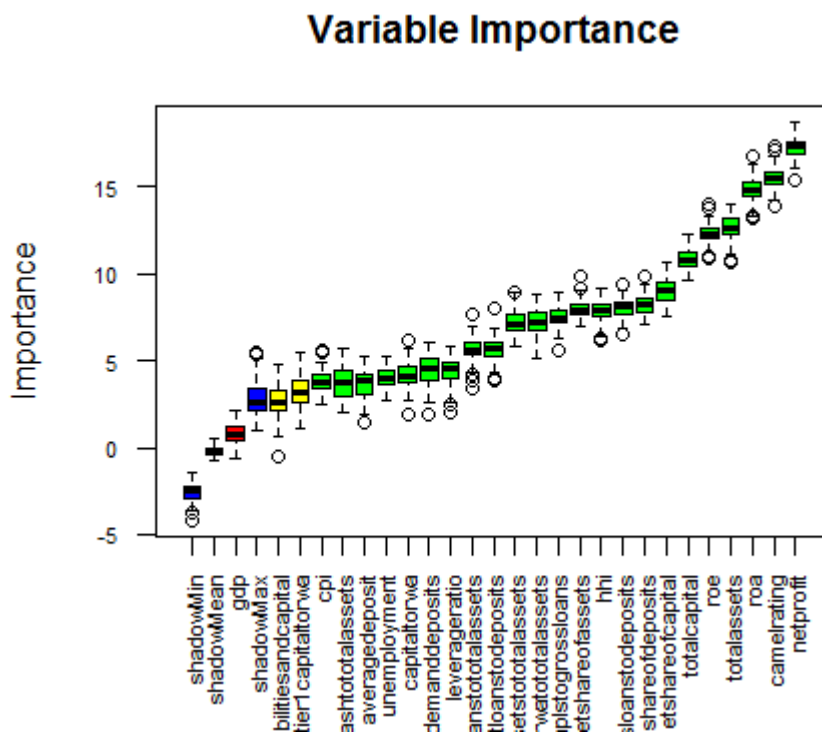
Parameter	SVM	ANN	KNN	NB	DT
Correctly classified instances	24519	24123	25051	22570	25081
Incorrectly classified instances	673	1069	141	2622	111
Kappa Statistic	0.9462	0.9136	0.9888	0.7906	0.9911
Mean Absolute Error	0.0267	0.0545	0.0056	0.1034	0.0064
Root Mean Squared Error	0.1634	0.197	0.0748	0.3152	0.0651
Relative Absolute Error	5.3676%	11.107%	1.1333%	20.7817%	1.2854%

Figure 14: P=

Detection model	Advantages	Disadvantages
Signature-based	Low false positive rate Does not require training Classified alerts	Cannot detect new attacks Requires continuous updates Tuning could be a thorny task
Anomaly-based	Can detect new attacks Self-learning	Prone to raise false positives Black-box approach Unclassified alerts Requires initial training

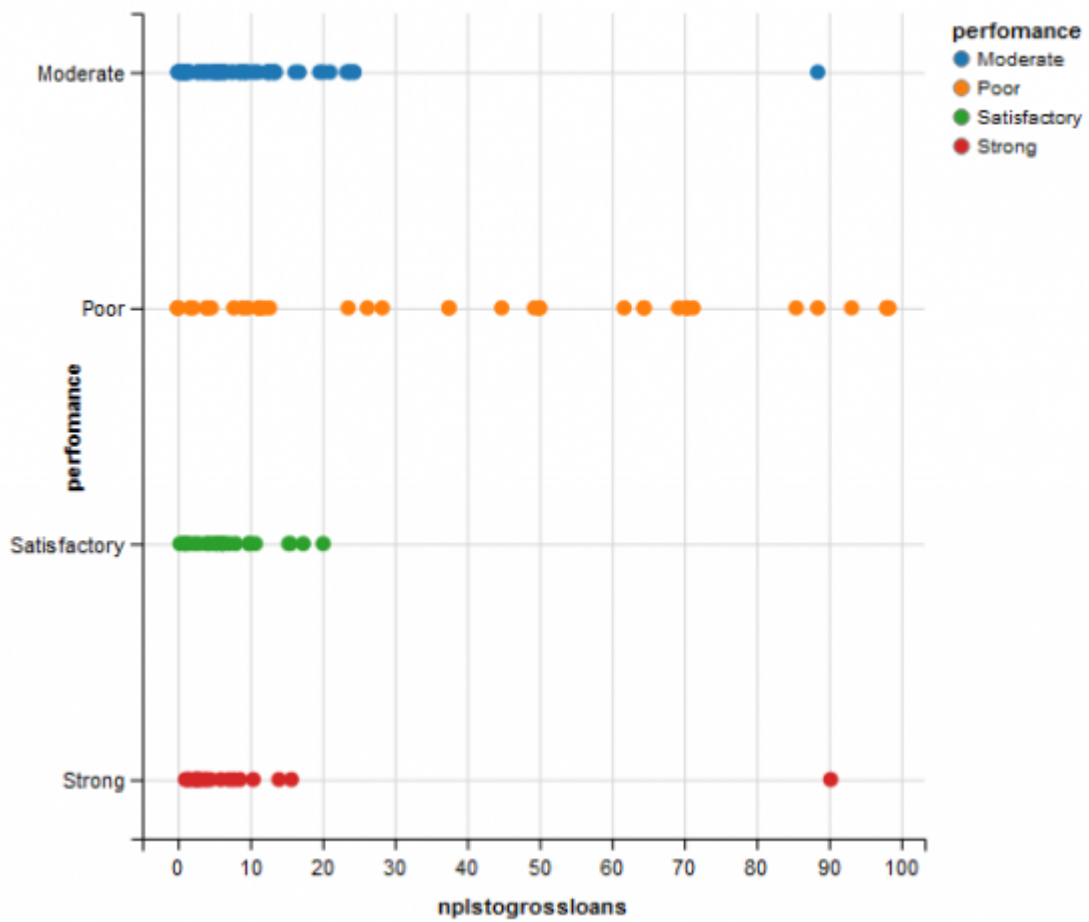
14

Figure 15: Figure 14 :



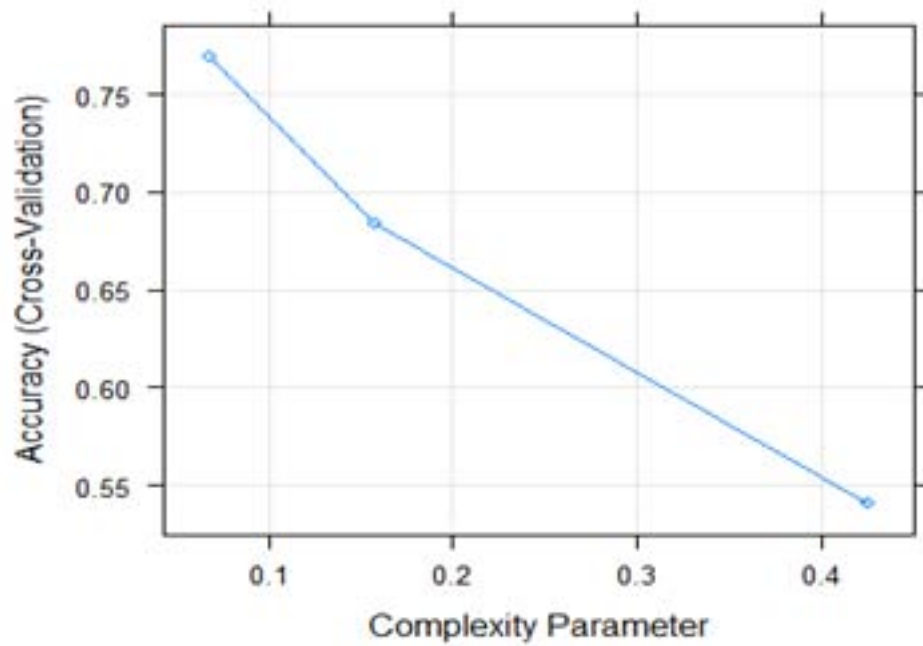
16

Figure 16: Figure 16 :



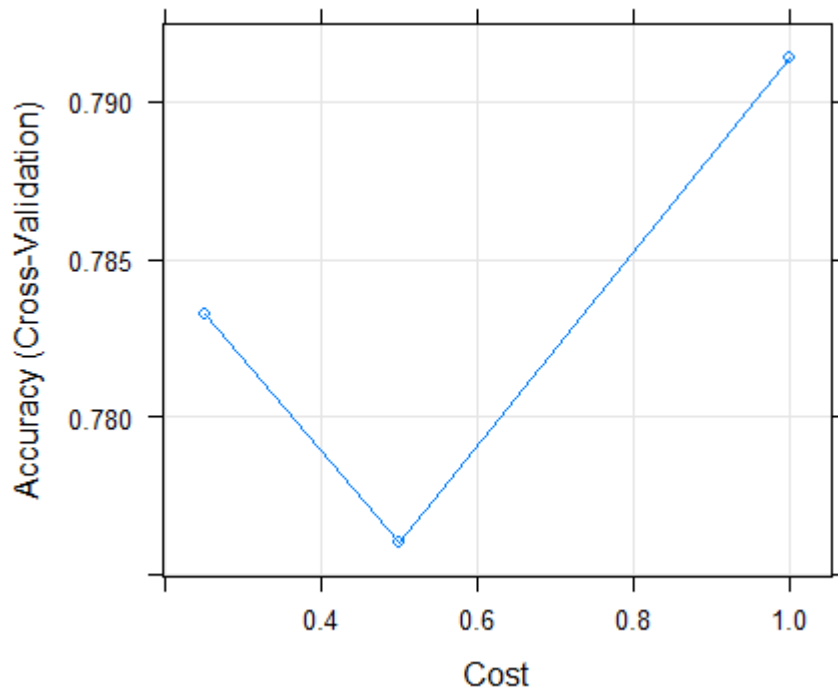
17

Figure 17: Figure 17 :



18

Figure 18: Figure 18 :



19

Figure 19: Figure 19 :

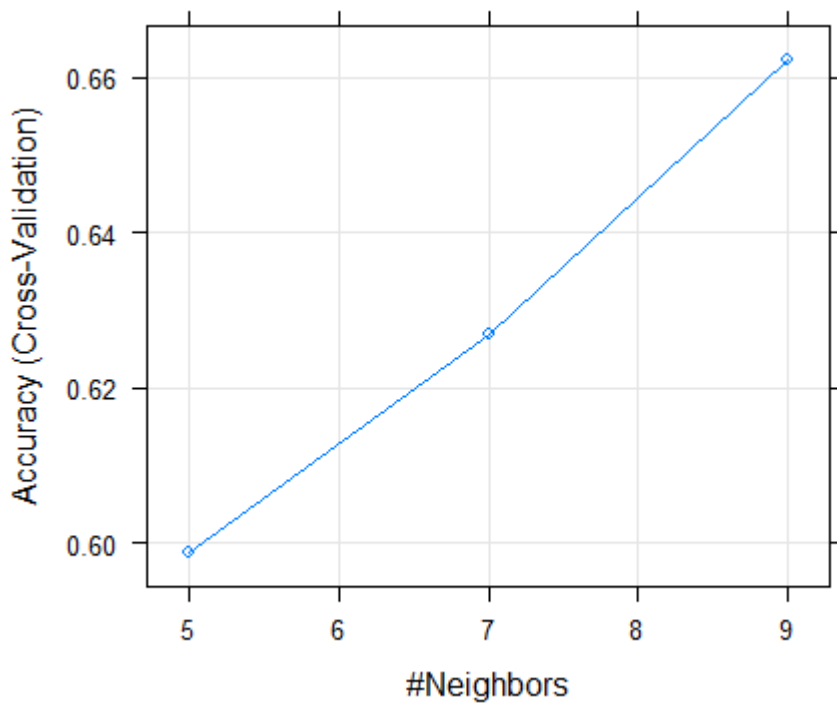


Figure 20:

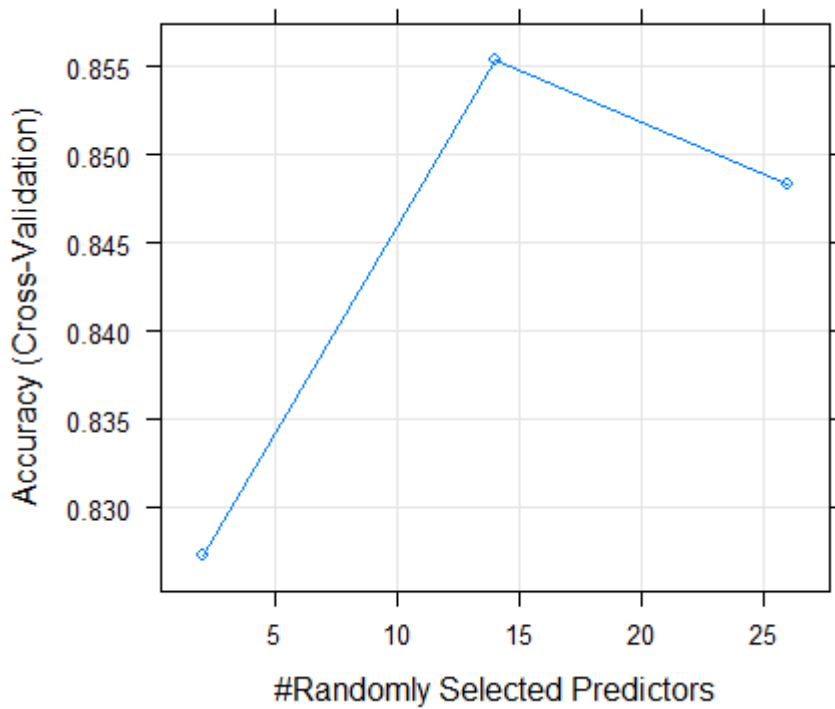


Figure 21:

2

Paradigm component	Explanation
Subjectivist epistemology	Researcher uses his/her own thinking and cognition to derive meaning from the research findings arrived at through interactive processes with the respondents
Relativist ontology	Multiple realities exist in the given setting Meaning is derived from the realities through interactions between the researcher and subjects as well as among participants
Naturalist methodology	Researcher makes use of data collected through text messages, interviews, conversations and reflective sessions as a participant observer
Balanced axiology	Research outcome will reflect the researcher's values, reporting research findings in a balanced manner

Figure 22: Table 2 :

4

Classifier	Detection Accuracy (%)	Time taken to build the Model in seconds	False Alarm rate (%)
Decision Trees (J48)	81.05	**	**
Naive Bayes	76.56	**	**
Random Forest	80.67	**	**
SVM	69.52	**	**
AdaBoost	90.31	**	3.38
Mutlinomal Naive Bayes + N2B	38.89	0.72	27.8
Multinomial Naive Bayes updateable + N2B	38.94	1.2	27.9
Discriminative Multinomial Bayes + PCA	94.84	118.36	4.4
Discriminative Multinomial Bayes + RP	81.47	2.27	12.85
Discriminative Multinomial Bayes + N2B	96.5	1.11	3.0

IV. Analysis and Research Outcomes

Figure 23: Table 4 :

8

Complexity Parameter	Accuracy	Kappa	AccuracySD	KappaSD
0.06849315	0.8275092	0.7519499	0.04976459	0.07072572
0.15753425	0.7783150	0.6683229	0.07720896	0.14039942
0.42465753	0.5222344	0.1148591	0.08183351	0.18732422

Figure 24: Table 8 :

9

sigma	c	Accuracy	Kappa	AccuracySD	KappaSD
0.050398	0.25	0.783223	0.678536	0.095598	0.140312
0.050398	0.50	0.776007	0.661354	0.087866	0.132552
0.050398	1.00	0.791391	0.678694	0.080339	0.126466

j) Linear Discriminant Algorithm

Figure 25: Table 9 :

10

Accuracy	Kappa	AccuracySD	KappaSD
0.8042399	0.7038131	0.1016816	0.159307
On the training dataset, the LDA achieved an accuracy level of 80% as in table 11. The Kappa statistic and the Kappa SD where 70% and 0.16 respectively. On the test dataset, the algorithm achieved an accuracy		level of 90% and a kappa of 84.64%. The algorithm misclassified 4 instance as moderate whose performance is poor in comparison to the CART algorithm.	
k) K-Nearest Neighbor			

Figure 26: Table 10 :

11

K	Accuracy	Kappa	AccuracySD	KappaSD
5	0.5988645	0.3698931	0.1280376	0.2158109
7	0.6268864	0.4072928	0.1564920	0.2703504
9	0.6621978	0.4715556	0.1747903	0.2881390

Figure 27: Table 11 :

12

mtry	Accuracy	Kappa	AccuracySD	KappaSD
2	0.8272527	0.7421420	0.10396454	0.15420079
14	0.8554212	0.7829891	0.06069716	0.09303130
16	0.8482784	0.7718935	0.06455248	0.09881991

Figure 28: Table 12 :

13

Software and Hardware Complex (SHC): Warning-2016	
Model Attributes	Description
	? HBase, a non-relational database, facilitates analytical and predictive operations
HBase working on HDFS (Hadoop Distributed File System)	? Enables users to assess cyber-threats and the dependability of critical infrastructure
	? Processes large amounts of data, interacts with standard configurations servers and is implemented at C language
Analytical data processing module	? Special interactive tools (based on JavaScript/ CSS/ DHTML) and libraries (for example jQuery) developed to work with content of the proper provision of cybersecurity
	? Interactive tools based on JavaScript/ CSS/ DHTML
Special interactive tools and libraries	? Libraries for example jQuery developed to work with content for
	? Designed to ensure the proper provision of cybersecurity
	? Percona Server with the ExtraDB engine
Data store for example (MySQL)	? DB servers are integrated into a multi-master cluster using the Galera Cluster.
Task queues and data caching	? Redis
Database servers balancer	? Haproxy
Web server	? nginx , involved PHP-FPM with APC enabled
HTTP requests balancer	? DNS (Multiple A-records)
Development of special client applications running Apple iOS	? Programming languages are used: Objective C, C++, Apple iOS SDK based on Cocoa Touch, CoreData, and UIKit.
Development of applications running Android OS	? Google SDK
Software development for the web platform	? PHP and JavaScript.
Speed of the service and protection from DoS attacks	? CloudFare (through the use of CDN)

(Source: [23])

Figure 29: Table 13 :

-
- 704 [Truong] , T Truong .
- 705 [Burt et al. ()] , D Burt , P Nicholas , K Sullivan , T Scoles . 2013. Cybersecurity Risk Paradox. Microsoft SIR
- 706 [Umamaheswari and Sujatha ()] , K Umamaheswari , S Sujatha . 2017.
- 707 [Pu and Kitsuregawa ()] , C Pu , M Kitsuregawa . No. GIT-CERCS-13-09. 2019. (Technical Report)
- 708 [Diep and Zelinka ()] , Q B Diep , I Zelinka . 2020.
- 709 [Bloice and Holzinger ()] *A Tutorial on Machine Learning and Data Science Tools with Python*, M Bloice , A
- 710 Holzinger . 2018. Graz, Austria.
- 711 [Kantarcioglu and Xi ()] ‘Adversarial Data Mining: Big data meets cybersecurity’. M Kantarcioglu , B Xi . *CCS*
- 712 2016. p. 16.
- 713 [Moorthy et al. ()] ‘An Analysis for Big Data and its Technologies’. M Moorthy , R Baby , S Senthamaraiselvi .
- 714 <https://ssrn.com/abstract=1118377> *GWU Law School Public Law Research Paper* 2014. (401) .
- 715 [Siti Nurul Mahfuzah et al. ()] ‘An Analysis of Gamification Elements in Online Learning to Enhance Learning
- 716 Engagement’. M Siti Nurul Mahfuzah , S Sazilah , B Norasiken . *6th International Conference on Computing*
- 717 *& Informatics*, 2017.
- 718 [Napanda et al. (2015)] ‘Artificial Intelligence Techniques for Network Intrusion Detection’. K Napanda , H Shah
- 719 , L Kurup . *International Journal of Engineering Research & Technology (IJERT)* 2278-0181. 2015. November-
- 720 2015. 4.
- 721 [Bringas and Santos (ed.) ()] *Bayesian Networks for Network Intrusion Detection*, P B Bringas , I
- 722 Santos . 10.1186/s40537-020-00318-5. [http://www.intechopen.com/books/bayesian-network/](http://www.intechopen.com/books/bayesian-network/bayesian-networks-for-network-intrusion-detection)
- 723 [bayesian-networks-for-network-intrusion-detection](http://www.intechopen.com/books/bayesian-network/bayesian-networks-for-network-intrusion-detection) Bayesian Network, Ahmed Rebai (ed.)
- 724 2010. InTech.
- 725 [Zhang et al. ()] *Bayesian-network-based safety risk analysis in construction projects. Reliability Engineering*
- 726 *and System Safety*, L Zhang , X Wu , M J Skibniewski , J Zhong , Y Lu . 10.1016/j.res.2014.06.006.
- 727 <https://doi.org/10.1016/j.res.2014.06.006> 2014.
- 728 [Berman et al. ()] D S Berman , A L Buczak , J S Chavis , C L Corbett . 10.3390/info10040122. *Survey of Deep*
- 729 *Learning Methods for Cyber Security*, 2019. 2019. 10 p. 122.
- 730 [Mazumdar and Wang ()] ‘Big Data and Cyber security: A visual Analytics perspective in’. S & Mazumdar , J
- 731 Wang . *Guide to Vulnerability Analysis for Computer Networks and Systems*, S Parkinson (ed.) 2018.
- 732 [Fehling et al. ()] ‘Cloud Computing Patterns’. C Fehling , F Leymann , R Retter , W Schupeck , P Arbitter
- 733 . 10.1007/978-3-7091-1568-8. <https://doi.org/10.1007/978-3-7091-1568-8> *Cloud Computing Pat-*
- 734 *terns*, 2014.
- 735 [Pai ()] ‘Corporate Social Responsibility of TSL sector: attitude analysis in the light of research’. Pai . *Logistyka*
- 736 2017. 2015. 2014. (5) p. . (The basis of social responsibility in management, Poltext, Warszawa. 37.
- 737 Marzantowicz)
- 738 [Sarker et al. ()] ‘Cybersecurity data science: an overview from machine learning perspective’. I H Sarker ,
- 739 A S M Kayes , S Badsha , H Alqahtani , P Watters , A Ng . 10.1186/s40537-020-00318-5. <https://doi.org/10.1186/s40537-020-00318-5> *Journal of Big Data* 2020.
- 740
- 741 [Menzes et al. ()] ‘Data Classification with Binary Response through the Boosting Algorithm and Logistic
- 742 Regression’. F S D Menzes , G R Liska , M A Cirillo , M J F Vivanco . 10.1016/j.eswa.2016.08.014.
- 743 <https://doi.org/10.1016/j.eswa.2016.08.014> *Expert Systems with Applications* 2016. 69 p. .
- 744 [Kobielus ()] *Deploying Big Data Analytics Applications to the Cloud: Roadmap for Success*, J Kobielus . 2018.
- 745 Cloud Standards Customer Council.
- 746 [Bezemskij et al. (2017)] ‘Detecting cyber-physical threats in an autonomous robotic vehicle using Bayesian
- 747 Networks’. A Bezemskij , G Loukas , D Gan , Anthony , RJ . [https://ieeexplore.ieee.org/](https://ieeexplore.ieee.org/document/8276737)
- 748 [document/8276737](https://ieeexplore.ieee.org/document/8276737) *IEEE International Conference on Internet of Things (iThings) and IEEE Green*
- 749 *Computing and Communications (GreenCom) and IEEE Cyber, Physical and Social Computing (CPSCom)*
- 750 *and IEEE Smart Data*, 2017. 2017. June 2017. p. . (Smart Data)
- 751 [Gheyas and Abdallah ()] ‘Detection and prediction of insider threats to cyber security: A systematic Literature
- 752 Review and Meta-Analysis’. I A Gheyas , A E Abdallah . *Big Data Analytics*, 2016. 2016. 1 p. 6.
- 753 [Wilson et al. (2015)] ‘Enablers and barriers of cloud adoption among Small and Medium Enterprises in Tamil
- 754 Nadu’. B M R Wilson , B Khazaei , L Hirsch . *2015 IEEE International Conference on Cloud Computing in*
- 755 *Emerging Markets (CCEM)*, 2015. November. IEEE. p. .
- 756 [Iafrate ()] *From Big Data to Smart Data*, F Iafrate . 2015. -848-21755-3 March, 2015. 190. Wiley-ISTE. (Pages)
- 757 [Lee ()] ‘HACKING INTO CHINA’ S CYBERSECURITY LAW’. J Lee . *IEEE International Conference on*
- 758 *Distributed Computing Systems*, 2017. 2017.

- 759 [Impregnable Defence Architecture using Dynamic Correlation-based Graded Intrusion Detection System for Cloud Defence Scien
760 'Impregnable Defence Architecture using Dynamic Correlation-based Graded Intrusion Detection System for
761 Cloud'. DOI: 10.14 429/dsj.67.11118. *Defence Science Journal* November 2017. 67 (6) p. .
- 762 [Almutairi ()] *Improving intrusion detection systems using data mining techniques*, A Almutairi . 2016. 2016.
763 Loughborough University (Ph.D Thesis)
- 764 [Karimpour et al. ()] 'Intrusion detection in network flows based on an optimized clustering criterion'. J
765 Karimpour , S Lotfi , A T Siahmarzkoooh . 17.07. <http://journals.tubitak.gov.tr/elektrik> *Turkish*
766 *Journal of Electrical Engineering & Computer Sciences* 2016. 2016.
- 767 [Kpmsg ()] Kpmsg . *Clarity on Cybersecurity. Driving growth with confidence*, 2018.
- 768 [Stefanova ()] 'Machine Learning Methods for Network Intrusion Detection and Intrusion Prevention Systems'. Z
769 S Stefanova . <https://scholarcommons.usf.edu/etd/7367> *Graduate Theses and Dissertations*, 2018.
770 2018.
- 771 [Suryavanshi ()] 'Magnesium oxide nanoparticle-loaded polycaprolactone composite electrospun fiber scaf-
772 folds for bone-soft tissue engineering applications: in-vitro and in-vivo evaluation'. A Suryavanshi
773 . 10.1088/1748-605X/aa792b/pdf. [https://iopscience.iop.org/article/10.1088/1748-605X/](https://iopscience.iop.org/article/10.1088/1748-605X/aa792b/pdf)
774 [aa792b/pdf](https://iopscience.iop.org/article/10.1088/1748-605X/aa792b/pdf) *Biomed. Mater* 2017. 2017. 12 p. 55011.
- 775 [Nielsen (2015)] R Nielsen . *CS651 Computer Systems Security Foundations 3d Imagination Cyber Security*
776 *Management Plan*, (Los Alamos National Laboratory, USA) 2015. January 2015. (Technical Report)
- 777 [Stallings (2015)] *Operating System Stability*, W Stallings . [https://www.unf.edu/public/cop4610/ree/](https://www.unf.edu/public/cop4610/ree/Notes/PPT/PPT8E/CH15-0S8e.pdf)
778 [Notes/PPT/PPT8E/CH15-0S8e.pdf](https://www.unf.edu/public/cop4610/ree/Notes/PPT/PPT8E/CH15-0S8e.pdf) 2015. Accessed on 27th March, 2019.
- 779 [Hammond ()] *Practical Artificial Intelligence for Dummies®*, *Narrative Science Edition*, K Hammond . 2015.
780 Hoboken, New Jersey: John Wiley & Sons, Inc.
- 781 [Cox and Wang ()] *Predicting the US bank failure: A discriminant analysis*, R Cox , G Wang . 2014. p. .
782 (Economic Analysis and Policy)
- 783 [Proko et al. ()] E Proko , A Hyso , D Gjylapi . <http://www.CEURS-WS.org/Vol-2280/paper-32.pdf>
784 *Machine Learning Algorithms in Cybersecurity*, 2018.
- 785 [Kothari ()] *Research Methodology Methods and Techniques 2 nd Revised Edition*, C R Kothari . 2004. New Age
786 International Publishers.
- 787 [Kumar ()] *Research Methodology: A step by step guide for beginners 3 rd ed*, R Kumar . 2011. London: Sage
788 Publishers.
- 789 [the Cyber Domain: Offense and Defense. Symmetry 2020] *the Cyber Domain: Offense and Defense. Symmetry*
790 *2020*, 12 p. 410.
- 791 [Fernando and Dawson ()] 'The health information system security threat lifecycle: An informatics theory'. J
792 I Fernando , L L Dawson . 10.1016/j.ijmedinf.2009.08.006. [https://doi.org/10.1016/j.ijmedinf.](https://doi.org/10.1016/j.ijmedinf.2009.08.006)
793 [2009.08.006](https://doi.org/10.1016/j.ijmedinf.2009.08.006) *International Journal of Medical Informatics* 2009.
- 794 [Hashem et al. ()] 'The rise of "big data" on cloud computing: Review and open research issues'. I A T
795 Hashem , I Yaqoob , N B Anuar , S Mokhtar , A Gani , S Khan . 10.1016/j.is.2014.07.006. <https://doi.org/10.1016/j.is.2014.07.006>
796 [Information Systems](https://doi.org/10.1016/j.is.2014.07.006), 2015.
- 797 [Yang et al. ()] *Utilizing Cloud Computing to address big geospatial data challenges. Computers, Environment*
798 *and Urban Systems*, C Yang , M Yu , F Hu , Y Jiang , Y Li . 10.1016/j.compenvurbsys.2016.10.010.
799 <https://doi.org/10.1016/j.compenvurbsys.2016.10.010> 2017.