



Building a K-Medoids Clustering technique to Detect Cloud Security-Related Anomaly

Hauwa Kulu Haruna, Kabiru Ibrahim Musa, Fatima Umar Zambuk, Abuzairu Ahmad
Department of Mathematical Sciences,
Abubakar Tafawa Balewa University, Bauchi

ABSTRACT

The showiest technological advancement of the twenty-first century is perhaps cloud computing. This is because, compared to other technologies in the field, it has had the fastest widespread acceptance. The proliferation of smartphones and other mobile devices with internet connectivity has been the main driver of this acceptance. The typical person can also benefit from cloud computing; it's not simply good for businesses and organizations. The increasing volume of data exchanged between businesses and cloud service providers creates risk factors for purposeful and unintentional leaks of private information to unreliable third parties. The majority of cloud service data breaches are caused by criminal activities, insider threats, malware, weak credentials, and human mistake. The most crucial security mechanisms against complex and expanding network threats are intrusion detection systems (IDSs) and intrusion prevention systems (IPSs). The two known IDS are the Signature base capable of detecting only known attack, while the anomaly detection alerts you to suspicious behavior that is both known and unknown. Issues from the main three anomaly detection techniques were challenges in chosen threshold in statistical approach, difficulties in maintenance from Data mining approach and challenges of expanding information in machine learning. Among the three techniques categorized for anomaly detection in Cloud Machine learning is the most essential strategy the detection, since it allowed machine learning and get better with time. From the Machine learning algorithm Recursive Feature Elimination (RFE) and K-medoids Clustering will be used to build a sophisticated model that will help achieve two separate objectives. SMO will obtain the desired nature for the model, while the K-medoids was used to build the model. The Model will be implemented using the Standard benchmark datasets, KDD CUP 99 and CICIDS2017. The performance evaluation of the proposed system was measured with an accuracy of 54.4% and 84.4% CICIDS2017 and KDDCUP 99 respectively. This shows that KDDCUP 99 is still more acceptable when clustering anomalies in cloud than the proposed CICIDS2017 dataset.

ARTICLE INFO

Article History

Received: August 2023

Received in revised form: October, 2023

Accepted: November, 2023

Published online: December, 2023

KEYWORDS

K-Medoids, Cloud computing, Clustering technique, Cloud security detection, Cloud anomaly detection, Recursive Feature Elimination, RFE, Intrusion detection system.

INTRODUCTION

Anything that includes offering hosted services via the internet is referred to as cloud computing. Infrastructure as a service (IaaS),

platform as a service (PaaS), and software as a service (SaaS) are the three basic forms of cloud computing services. It is possible to have a private or public cloud. Anyone on the internet

Corresponding author: Hauwa Kulu Haruna.

✉ hauwakuluharuna@gmail.com

Department of mathematical Sciences, Abubakar Tafawa Balewa University, Bauchi.

© 2022. Faculty of Technology Education. ATBU Bauchi. All rights reserved



can buy services from a public cloud. A private cloud is a closed network or data center that provides hosted services to a small group of people with specific access and permissions. The purpose of cloud computing, whether private or public is to give quick, scalable access to computer resources and IT services. The hardware and software components required for efficient cloud computing implementation are referred to as cloud infrastructure (Wesley & Stephen, 2021). Cloud computing is a next-generation technology based on the network and internet that delivers IT services like file transfers, software for hire, secure storage, servers for rent, databases, analytics, and intelligence. These flexible resources are enabling faster innovation, and shaping economies.

According to the "Right Scale 2019 State of the Cloud Report," public cloud is the top priority for over 30% of company IT decision-makers in 2019. Nonetheless, organizational adoption of the public cloud, particularly for mission-critical applications, has lagged behind expectations, according to several experts. Today, however, mission-critical workloads are increasingly likely to be moved to public clouds. One reason for this trend is that corporate executives who want their organizations to compete in the new world of digital transformation are seeking the public cloud. Business executives are also looking to the public cloud for its elasticity, to upgrade internal computer systems, and to empower essential business divisions and their DevOps teams (Wesley & Stephen, 2021). Networks play important roles in modern life, and cyber security has become a dynamic research area (Abuzairu et al, 2023). Cloud computing has revolutionized the way businesses and individuals use and manage their data and applications. However, like any technology, it also comes with its share of challenges and issues. Some common issues with cloud computing include: Security and Privacy, Downtime and Reliability, Vendor Lock-in, Data Transfer Bottlenecks, Compliance and Legal Issues, Cost Management, and Dependency on Service Providers (Binod, 2023). Despite these challenges, cloud computing

continues to evolve, and many organizations find the benefits far outweigh the potential issues. With proper planning, security measures, and risk management strategies, businesses can mitigate these challenges and leverage the advantages offered by cloud computing (Sehgal, 2013). Cloud computing security issues are a major concern for organizations utilizing cloud services. Some common security issues associated with cloud computing, are; Data Breaches, Insider Threats, Insecure Interfaces and APIs, Data Loss, Shared Infrastructure Vulnerabilities, Compliance, and Legal Issues, Lack of Transparency and Control, Cloud Service Provider Vulnerabilities and Malicious insiders (David, 2023). To address these security issues, organizations adopt a comprehensive approach to cloud security, including implementing strong access controls, encrypting data, conducting regular security assessments, monitoring for suspicious activities, and establishing clear policies and agreements with cloud providers regarding security responsibilities and incident response protocols (John & Mello, 2022).

Statement of the problem

An anomaly is an abnormal behavior or departure from the usual. The process of removing these abnormal or anomalous behaviors from the data or services is known as anomaly detection. Cloud service providers must deliver services to users that behave normally. Anomaly detection becomes a crucial and engaging topic of research in order to give services to users in a proper and normal manner. Anomaly detection methods come in a variety of forms, but they can be broadly categorized into three groups: statistical, data mining-based, and machine learning-based methods. It's crucial to make sure the cloud is reliable (Nassif, et, al., 2021). However, the reality is that today's cloud solutions are insufficiently trustworthy. The S3 service was interrupted for 4 hours on February 28th, 2017 in the Northern Virginia (US-EAST-1) Region, which quickly spread to other online service providers who rely on the S3 service. Amazon Web Services, the well-known storage and hosting platform used by a wide range of businesses, was affected (AWS, 2021).

Corresponding author: Hauwa Kulu Haruna.

✉ hauwakuluharuna@gmail.com

Department of mathematical Sciences, Abubakar Tafawa Balewa University, Bauchi.

© 2022. Faculty of Technology Education. ATBU Bauchi. All rights reserved



The failure led to a significant financial loss. It's because service providers for cloud computing often establish a Service Level Agreement (SLA) with clients. Customers who demand 99.99 percent availability, for instance, are requesting that the service be available 365 days a year, 99.99 percent of the time (Google Cloud, 2021). Compensation is necessary if the service is disrupted by more than 0.01 percent. The three prominent categories of Cloud anomaly detection have some setback. The statistical based approach suffered from its relied-on assumption, which may not hold, especially for the high-dimensional real data set (Shahzad, et, al., 2022). SVM is computational expensive as it requires to solve quadratic optimization problem (GeeksforGeeks, (2020), Unsuitable for large dataset, it has large training time, performs badly on high noise and have more features and more complexities (Patil, 2022). Further experiments are recommended to generalize the findings, including larger datasets, different Machine Learning algorithms, and exploring additional hyper parameters and parameter values (Aldallal & Alisa, 2021)).

Aim and Objectives of the Study

The aim of the research is to build a k-medoids clustering techniques to detect cloud security related anomalies.

The Objectives of the studies are;

- i. Design a Cloud-based Machine Learning anomaly detection Model using Recursive Feature Elimination (RFE) and k-medoid clustering algorithm
- ii. Implementation of the Model using the Standard Cloud Anomaly Dataset, a benchmark datasets.

LITERATURE REVIEW

Cloud Computing

In order to provide quicker innovation, adaptable resources, and scale economies, cloud computing is the distribution of computer services over the Internet ("the cloud"), including servers, storage, databases, networking, software, analytics, and intelligence. Typically, you only pay for the cloud services you actually use, which lowers your operational expenses,

makes it easier to manage your infrastructure, and enables you to scale as your company needs evolve. Some of the benefits of cloud computing are; Cost, which eliminate the cost of buying software, Speed, as it provides safe service in demand for the clients, Global Scale, Productivity, which removes the need of many tasks like, Hardware setup, software patching and other time-consuming IT management chores. Performance, Reliability and Security are also parts of the major benefits gain from cloud computing system. Apart from the benefits Cloud as a system is categorizes into three; which are public cloud, Private cloud and Hybrid Cloud deployment types. Finally, cloud computing can be used to create services like; Creating a cloud-native application, Test and Build application, Store, Backup and recovery data, analyze data, Stream audio and Video, embed intelligence and Deliver software on demand (Azure, 2022)

Cloud Security

Cloud security is a set of practices and tools created to address both internal and external security threats to businesses. As they implement their digital transformation strategy and integrate cloud-based tools and services into their infrastructure, organizations need cloud security. With the continued development of the digital environment, security threats have advanced. Due to an organization's general lack of visibility in data access and movement, these threats specifically target providers of cloud computing. Organizations may encounter serious governance and compliance issues when handling client information, regardless of where it is housed, if they don't take proactive measures to increase their cloud security. Despite your company size, cloud security needs to be a major talking point. Almost every aspect of contemporary computing is supported by cloud infrastructure, which spans several verticals and all sectors. Alouffi, et al. (2021), categorizes security research issues into four taxonomies. The categorize are; Security and Privacy, Attack Detection, Security Taxonomy and Resource scheduling. Research is basically on an Attack Detection in Cloud. The attack survey will literarily start with the five types of most

Corresponding author: Hauwa Kulu Haruna.

✉ hauwakuluharuna@gmail.com

Department of mathematical Sciences, Abubakar Tafawa Balewa University, Bauchi.

© 2022. Faculty of Technology Education. ATBU Bauchi. All rights reserved



discussed attack on a cloud computing. Khorshed, et al., (2012), described the five (5), most popular attacks on cloud system. We highlighted them as follows;

- i. Denial of Service Attack (DOS), which is described as an attack in which the attacker prevents the genuine users from accessing computer resources.
- ii. Cross VM side channel (CVMSC) attack, is a sophisticated in Cloud computing which exploit the information of the hardware resources usage such as Central processing unit.
- iii. Malicious Insider (MI) attacks, mostly hard to discuss, since it involves trusted insider, to carry out the attack.
- iv. Attacks targeting shared memory (ATSM); the attacker takes advantages of shared memory, physical/cache memory of a physical/virtual system.
- v. Phishing Attacks: It involves the use of social engineering techniques to harvest an information from group or individual information.

Threats and Vulnerabilities in Cloud

The risks and weaknesses associated with moving to the cloud are constantly changing. It is crucial to take into account additional difficulties and dangers related to cloud adoption that are unique to their objectives, systems, and data.

A Threat is anything that could gain access to, harmed, or destroyed an asset by intentionally or unintentionally exploiting vulnerability.

Vulnerability is an opening in a security program that threats could use to access an asset without authorization. Vulnerability is a flaw or opening in our defenses. Khorshed, et al., (2012), highlighted some of the top threats and the existing gaps that are yet to be addressed in the threat's vulnerability. AlertLogic, (2021), expressed some of the latest vulnerabilities in cloud that need urgent attention, which includes;

Shared Technology vulnerabilities – Increased resource leverage provides attackers with a single point of attack, which can cause harm that is out of proportion to its importance. A

hypervisor or cloud orchestration is examples of sharing technology.

Data Breach – The potential of an accidental, malicious, or intentional data breach is greater as data protection moves from the cloud user to the cloud service provider.

Account of Service traffic hijacking – Access to the cloud over the Internet is one of its main benefits, but it also comes with the possibility of account compromise. Losing access to a privileged account could result in service interruption. Ways to Hijack account among many are; Phishing, Keyloggers, Buffer overflow attacks, Cross site scripting attacks, Brute force attacks.

Denial of Service (DoS) – Any denial-of-service assault on the cloud provider has the potential to disrupt all of the tenets.

Malicious Insider – In a cloud scenario, a motivated insider can find more ways to attack and cover the trace.

Internet Protocol – IP spoofing, ARP spoofing, and DNS poisoning are only a few examples of IP vulnerabilities.

Injection Vulnerabilities – Vulnerabilities in the administration layer, such as SQL injection, OS injection, and LDAP injection, can pose serious problems for numerous cloud users.

API & Browser Vulnerabilities – in a cloud provider's API or interface poses a considerable hazard when paired with social engineering or browser-based attacks; the harm can be significant.

Changes to Business Model – A cloud consumer's business model may be significantly altered by cloud computing. The IT department, as well as the business, must adapt or risk being exposed to danger.

Abusive use – Certain cloud computing characteristics, such as the usage of a trial period of use to conduct zombie or DDoS attacks, can be used for malicious offensive goals.

Availability – the likelihood that a system will perform as expected and when expected.

Cloud Anomaly Detection

Corresponding author: Hauwa Kulu Haruna.

✉ hauwakuluharuna@gmail.com

Department of mathematical Sciences, Abubakar Tafawa Balewa University, Bauchi.

© 2022. Faculty of Technology Education. ATBU Bauchi. All rights reserved

An Anomaly is abnormal activity or deviation from the normal behavior. While Anomaly detection is the process of removing these abnormal or anomalous behaviors from the data or services. The services delivered to users by cloud service providers must have normal behavior. Anomaly detection aims to identify the presence of aberrant patterns in network traffic, and routine detection of such patterns can give network administrators access to additional

information sources to understand network behavior or track down and identify the source of network issues.

The Anomalies Classification

Anomalies are categorized into three, a point anomalies, Contextual anomalies and Collective anomalies. (Sari, 2015) explained the three taxonomy, which we represent it in a chart in figure 1.

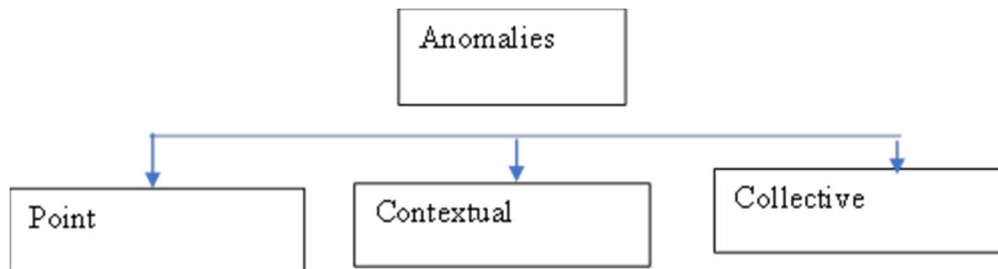


Figure 1: Anomalies Taxonomy

Intrusion Detection System

Is a program that recognizes network intrusions based on various telltale signals, such as known attack methods, IDS is typically installed at the network's gateway to monitor incoming and outgoing network traffic and find malicious intrusions. To put it another way, the IDS can be described as a tuple of (S, R), where S denotes the desired behavior (or the reference model). R, on the other hand, describes the similarity measurement that identifies the behavior that deviates from S (Kabla, et al., 2022). Two basic types of intrusion detection are the anomaly and the Signature base intrusion detection system.

Anomaly Based Intrusion Detection System

It is the technique of locating odd deviations in huge datasets that provide information that can be used by business users. A dataset may contain a single, erroneous data point or a collection of data points that have a distinct distribution within the dataset's overall distribution. Applications that are relevant to most enterprises, such as fraud or cyber security for finance teams, and equipment failure for operations teams, all depend on the integration of anomaly detection (Lee, 2021) Some of the benefits of anomaly detection mentioned in (Solana Networks, 2016).

Signature Based Detection

A signature-based detection method uses specified attack patterns as signatures, which are then utilized to identify network attacks. They typically use predetermined signatures to analyze network traffic each time the database is updated, an illustration of an intrusion detection system using signatures (Kumar & Sangwan, 2012).

Anomaly Detection in Cloud computing

Anomaly detection is the process of identifying unusual events, items, or observations that are unusual in comparison to standard behaviors or patterns. Standard deviations, outliers, noise, novelty, and exceptions are all terms used to describe data anomalies. Interesting incidents are not uncommon in the area of network anomaly detection/network intrusion and abuse detection. Unexpected

Corresponding author: Hauwa Kulu Haruna.

✉ hauwakuluharuna@gmail.com

Department of mathematical Sciences, Abubakar Tafawa Balewa University, Bauchi.

© 2022. Faculty of Technology Education. ATBU Bauchi. All rights reserved



spikes in activity, for example, are usually noticeable, even if they fall outside the scope of many classic statistical anomaly detection tools. Many outlier detection algorithms, particularly unsupervised techniques, miss this type of abrupt increase in activity as an outlier or rare object. A cluster analysis technique, on the other hand, may typically detect these types of micro clusters more easily. (Arif, 2015), made a comparison between four methods and techniques used in Cloud Based network.

METHODS AND TECHNIQUES OF ANOMALY DETECTION IN CLOUD

Different approaches or methodologies, such as threshold detection, statistical analysis, rule-based measurements, neural networks, genetic algorithms, data mining, and machine learning, have been employed in cloud networks to identify abnormal behaviors. This section presents a comparison of the many approaches to anomaly detection in cloud networks. Research would be done on statistical, data mining, and machine learning approaches and techniques in comparison to others.

Statistical Anomaly Detection

This technique for finding anomalies in cloud-based networks looks at network computations to find anomalies, and then builds a profile that retains or stores the created value to understand their behavior. Using this method, two profiles are formed; the first one stores the usual or anomalous rules or signatures while the second one updates on a regular basis. Scores for anomalies are calculated throughout the update. It is recognized to be abnormal and discovered if the threshold value is less than the most recent anomaly profile generated (Fuse & Kamiya, 2017).

Machine Learning Based Anomaly Detection Techniques

An essential strategy in the detection of anomalies is the ability of programs or software to get better over time at doing their job through learning. When an anomaly happens or is recognized, the machine learns its behavior and

saves the new sequence or rule. Verified values or the typical behavior of the data are recorded. Through the use of the past outcomes, this technique develops a system that can enhance the performance of the program. The intriguing aspect of this technique is that when performance is improved from prior results, new information is retrieved, and if better performance calls for a change in the execution strategy, the decision is made using the new knowledge from the prior results. Machine learning-based anomaly detection falls under a number of categories, including Bayesian Network, Genetic Algorithm, and Neural Network.

The historical knowledge or signature as well as the data used in anomaly detection can both be incorporated into the Bayesian Network's processing. This method combines with a statistical process that is quite helpful for detecting anomalies (Mi, et al., 2013). Neural networks have the ability to enhance incomplete data to develop the potential to recognize and comprehend hidden patterns. The neural network may detect previously undetected behavior or patterns as well as recent attacks (Wang, et al., 2010). The evolutionary algorithm techniques such as mutation, selection, etc. are used by genetic algorithms. This unique method is based on rules gathered from data from network analysis performed by the IDS (Safi, 2015).

ADAPTIVE ANOMALY DETECTION TECHNIQUES

For adaptive failure detection, the Adaptive Anomaly Detection Systems (AAD) uses data description using a hyper sphere. The AAD in cloud networks uses performance data from the cloud service to identify potential failures or anomalies that are discovered by cloud operators. The log of the discovered failure records handed in by the cloud operators is used or capitalized upon by the AAD detection systems to later identify new sorts of failures. The AAD detection algorithm adapts to fresh verified results or detections from the cloud provider by repeatedly learning from them in order to be ready for detections in the future. It was observed that the prototype was portable and that startup times for the detector, set adaption, and failure

Corresponding author: Hauwa Kulu Haruna.

✉ hauwakuluharuna@gmail.com

Department of mathematical Sciences, Abubakar Tafawa Balewa University, Bauchi.

© 2022. Faculty of Technology Education. ATBU Bauchi. All rights reserved

detection each required a few seconds to complete. A typical cloud server's Central Processing Unit consumption, job switching procedures, memory and swap space utilization, paging and page faults, input and output data transfer, interruptions, and other statistics were all profiled on 518 metrics per minute. The ADD algorithm was compared using failure detectors such subspace regularization (Safi, 2015).

THE BENCHMARK DATASETS USE IN THE RESEARCH

The dataset used to assess any IDS is a key factor in determining how effective it is (Aldallal & Alisa, 2021). The absence of real-world datasets containing actual anomalies is a major barrier to the advancement of anomaly detection (Pang, 2020). As a result, the KDD CUP 99 dataset has been used in a plethora of research projects to evaluate their system. Apart from KDD CUP 99, NSL KDD can also be used. But due to their age and lack of features, these datasets are unreliable for replicating the systems and environments of the present (Aldallal & Alisa, 2021). This work is meant to be used for real cloud data. But since it is impossible to get the actual cloud data needed to construct the system intended for this study, CICIDS2017, a pre-defined dataset that closely resembles the

existent cloud data will be utilized in its place, which is one of the most recent datasets used in machine learning applications (Aldallal & Alisa, 2021). This system will work well for cloud computing since it is anticipated to offer a high level of information security that will shield data from many forms of hostile attacks.

PERFORMANCE METRICS

A business's behavior, operations, and performance are all tracked using performance metrics. This should take the form of measurements of the necessary data within a range, enabling a foundation to be built in support of the accomplishment of overarching business objectives. False Positive Rates and False Negative Rates are the most crucial performance indicators for an anomaly detection program. Point- and range-based metrics can be used to evaluate anomaly detection performance. Only the abnormal and normal answer labels are considered valid in point-based techniques, and if the answer is the same, points are awarded. In contrast, range-based techniques assign various ratings depending on performance measures in addition to labeling things as abnormal or normal. This enables the evaluator to assign various weights and get various scores (Kim et al., 2022).

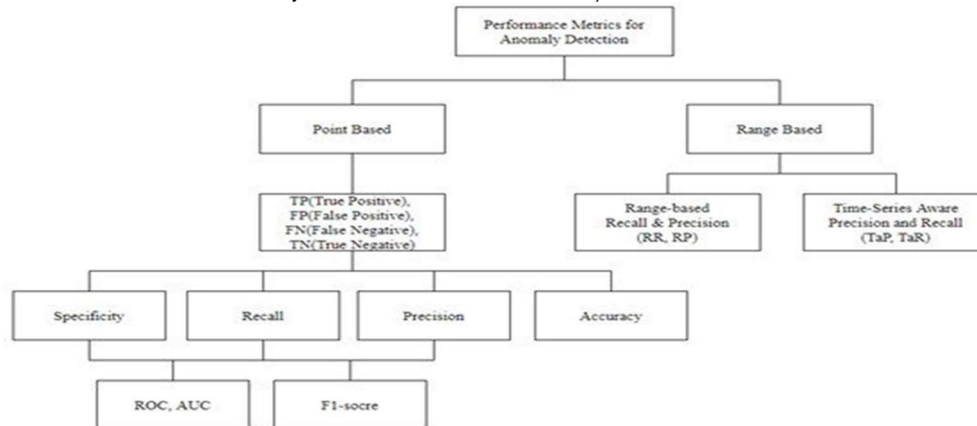


Figure 2: Performance metrics tree for anomaly detection evaluation (Kim et al., 2022)

The standard measure, a point-based performance index, was built utilizing the components of a confusion matrix. In addition to

accuracy, the function frequently takes precision and recall into account. The receiver operating characteristic curve (ROC) using true and false

Corresponding author: Hauwa Kulu Haruna.
 ✉ hauwakuluharuna@gmail.com
 Department of mathematical Sciences, Abubakar Tafawa Balewa University, Bauchi.
 © 2022. Faculty of Technology Education. ATBU Bauchi. All rights reserved

positive rates, the F1-score using precision and recall, and the AUC (area under the curve) representing the bottom region of the ROC curve are all available. A point-based performance evaluation index is incorrect since anomaly detection is defined as a range-based anomaly brought on by an unusual occurrence. As a result, the range-based performance indices recall (RR) and precision (RP) were applied. If the irregularity is sufficiently diagnosed, recall obtains a point; otherwise, it receives a penalty. However, recall is an improper performance evaluation indicator because time series data can be expressed as a range (Kim et. al., 2022).

RESEARCH METHODOLOGY

The proposed research used the modified parameter wise approach to solved the lingering machine learning cloud anomaly detection issues. Recursive Feature Elimination (RFE) techniques, is use to select the most useful features from the datasets. While the K-medoids clustering is use for the actual classification of the datasets as normal an anomaly. The proposed approach is seen in figure 3.

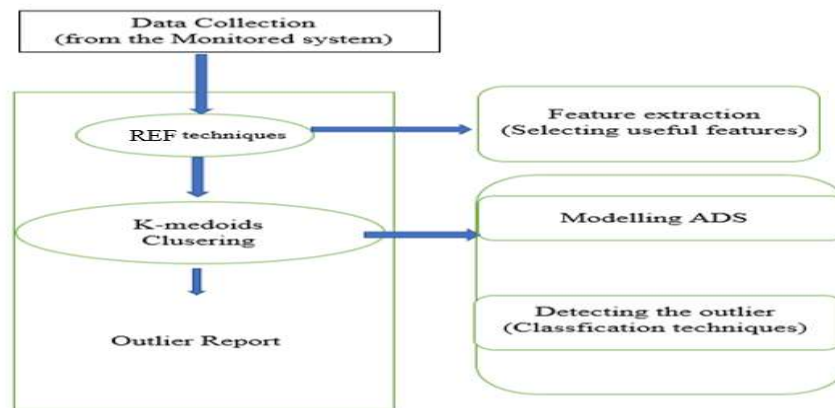


Figure 3: The proposed approach to Anomaly Detection in Cloud

GENERAL PROCESS OF THE PROPOSED MODEL

The most crucial security mechanisms against the complex and expanding network threats are intrusion detection systems (IDSs) and intrusion prevention systems (IPSs). Performance evolutions for anomaly-based intrusion detection techniques are inconsistent and inaccurate since there aren't enough trustworthy test and validation datasets.

The Proposed Model Dataset

CICIDS2017 dataset will be used for this research, together with the benchmark's datasets, KDD-CUP 99 and DARPA 98. The CICIDS2017 dataset is used because of its coverage of all the eleven requirements that must be met in order to create a trustworthy benchmark dataset. In which known of the

previous IDS cover (CICIDS2017, 2020). The most attractive feature of the CICIDS2017 Datasets in the proposed study, is it Attack diversity. Based on the 2016 McAfee study, the most frequent attacks were covered in this dataset, including Web-based, Brute-force, DoS, DDoS, infiltration, Heart-bleed, Bot, and Scan assaults.

First Phase of the Proposed Model

Generally, the proposed model starts from Data preparation using the Recursive Feature Elimination techniques to Clean and select the most important features of the Datasets. A more modern optimization technique based on swarm intelligence is the Recursive Feature Elimination (RFE) algorithm. Euclidean distances between potential solutions serve as the foundation for the update equations. The

Corresponding author: Hauwa Kulu Haruna.

✉ hauwakuluharuna@gmail.com

Department of mathematical Sciences, Abubakar Tafawa Balewa University, Bauchi.

© 2022. Faculty of Technology Education. ATBU Bauchi. All rights reserved



technique has been widely used to address challenging optimization issues (Sharma, et al., 2020). In the Data preprocessing stage, the RFE techniques is use to change data types to integer, replace the missing values with zero, to remove some of unwanted columns in the records and to performed scaling for data normalization.

The Second Phase of the Proposed Model

The second phase of the proposed model is Modeling and detection of the anomaly using K-medoids techniques, a Clustering technique. K-Medoids is a clustering approach that is similar to the K-Means method. It belongs to the class of machine learning without supervision. In terms of how it chooses the centers of the clusters, it differs significantly from the K-Means algorithm. While the latter always chooses the actual data points from the clusters as their centers (also known as "exemplars" or "medoids"), the former chooses the average of a cluster's points as its center (which may or may not be one of the data points). The K-Medians technique, which is identical to K-means except that it chooses the medians (instead of means) of the clusters as centers, is another way in which K-Medoids varies from K-Medians (Nikita, 2021). The Process of modeling with K-medoids clustering techniques is as follows;

- i. Pick k randomly from the input data; k is the number of clusters that will be produced. Utilizing techniques like the silhouette approach, one can evaluate the accuracy of the k's value selection.
- ii. Each piece of data is assigned to the cluster that contains the closest medoid.
- iii. The distance from each data point in cluster I to all other data points is calculated and added. Assigned as the medoid for that

cluster is the point in the cluster for which the estimated sum of distances from other points is minimal.

- iv. Steps (2) and (3) are repeated until convergence is reached i.e., the medoids stop moving.

The K-Medoids algorithm has an $O(N^2CT)$ complexity, where N, C, and T stand for the number of data points, clusters, and iterations, respectively. The complexity of the K-Means algorithm can be expressed as O using similar notations (NCT).

Benefits of the Method using K-Medoids Clustering

A measure that is significantly impacted by the extreme points is the mean of the data points. Therefore, if the data contains outliers, the centroid in the K-Means algorithm may be pushed to the erroneous place, which will lead to improper clustering. As opposed to this, a medoid in the K-Medoids algorithm is the cluster's center component, with the shortest distance to other locations. The K-Medoids algorithm is more resistant to outliers and noise than the K-Means algorithm since medoids are unaffected by extremities. The positions of the mean and medoid can change in the presence of an outlier, as shown in the following picture (Nikita, 2021).

PROPOSED MODEL ARCHITECTURE

As discussed, the architecture of the proposed Cloud Anomaly Detection using K-medoids and REF techniques of machine learning can be represented in a hybrid combination of the two classical approaches. The Model architecture can be presented in figure 4.

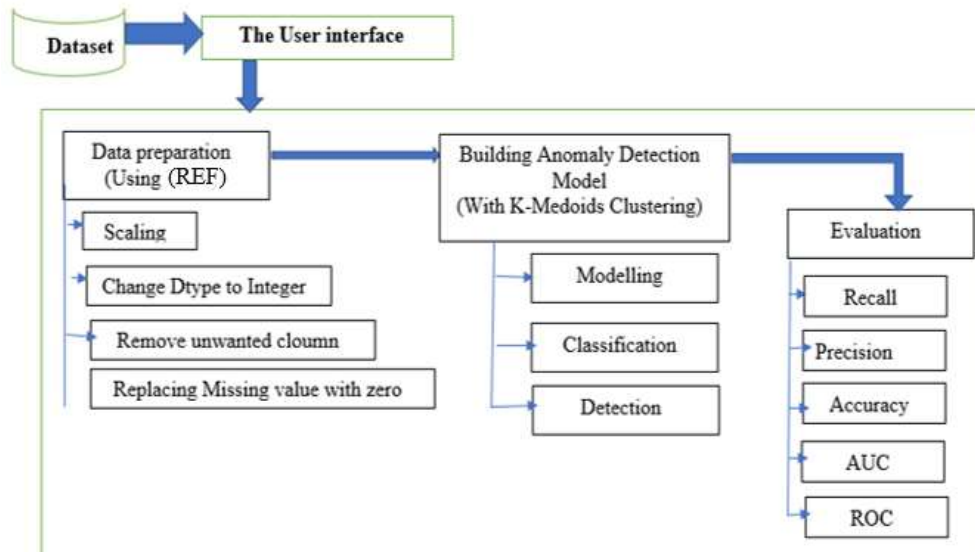


Figure 4: Architecture of the proposed ADS Model

The proposed System Algorithm

The algorithm is a set of clear directions about how to solve a specific problem. It receives one or more inputs and generates the required output. Our proposed research, take on two algorithms from the machine learning techniques. The Recursive Feature Elimination, mainly for feature selection and the K-Medoids Clustering for Classification and Detection. Hence the General Algorithm for the proposed modeling will take in two distinct sections of the algorithms.

The Feature Selection Process

The Feature Selection Algorithm use Recursive Feature Elimination (RFE), a machine learning approach, which is a group process that iterates collectively based on trial and error. Six phases make up the RFE process: Local Leader, Global Leader, Local Leader Learning, Global Leader Learning, Local Leader Decision, and Global Leader Decision.

The Algorithm: RFE

Step 1: initialize population, Limit for Global and Local leader, pr.

Step 2: Calculate fitness for all population.

Step 3: Find and select global and local leaders using greedy selection.

While the number of iterations **do**

- i. All group members' positions are updated based on personal, local leader, and member experience.
- ii. Apply a ruthless selection process based on fitness to all group members;
- iii. Use the formula to determine the probability $prob_i$ for each group member.
 $prob_i = fitness_i / \sum_{i=1}^N fitness_i$
- iv. Updating one's position throughout the global leader phase.
- v. By employing the greedy selection method on all the groups using the local and global leader learning phase, the position of local and global leaders is updated.
- vi. All members of that specific group should be redirected for foraging by the local leader choice phase if any local group leaders do not update their positions after a predetermined number of times (Localleaderlimit).
- vii. After a predetermined number of updates (Globalleaderlimit), if the global leader does not update her stance, she divides the group into smaller groups according to the Global leader choice phase

Corresponding author: Hauwa Kulu Haruna.

✉ hauwakuluharuna@gmail.com

Department of mathematical Sciences, Abubakar Tafawa Balewa University, Bauchi.

© 2022. Faculty of Technology Education. ATBU Bauchi. All rights reserved

Step 4: Stop

THE CLASSIFICATION AND DETECTION ALGORITHM

A clustering algorithm similar to the K-means algorithm is the K-medoids algorithm. Both techniques are partition and aim to reduce squared error, which is the separation between a location identified as the cluster center and a point labeled as being in the cluster. K-medoids, as opposed to the K-means method, selects data points as centers. The data set of n objects is clustered into k clusters using the traditional clustering partitioning approach known as K-medoid. In comparison to k-means, it is more resistant to noise and outliers. A medoid is an object in a cluster that is the most centrally situated in the supplied data set and has an average dissimilarity to all other objects in the cluster that is low.

Step 1: Out of n data points, the program randomly chooses k objects to serve as medoid points ($n > k$).

Step 2: Out of n data points, the program randomly chooses k objects to serve as medoid points ($n > k$).

Step 3: Assign each data object in the given data collection to the medoid that is the closest match after choosing the k medoid points. The distance metric used to define the similarity in this instance is either the Minkowski distance, the Manhattan distance, or the Euclidean distance.

Step 4: Pick a non-medoid object at random O'

Step 5: If $S > 0$, replace the original medoid with the new one (if $S > 0$, a new set of medoids will exist).

Step 6: Till there is no change in the medoid, repeat stages 2 through 5.

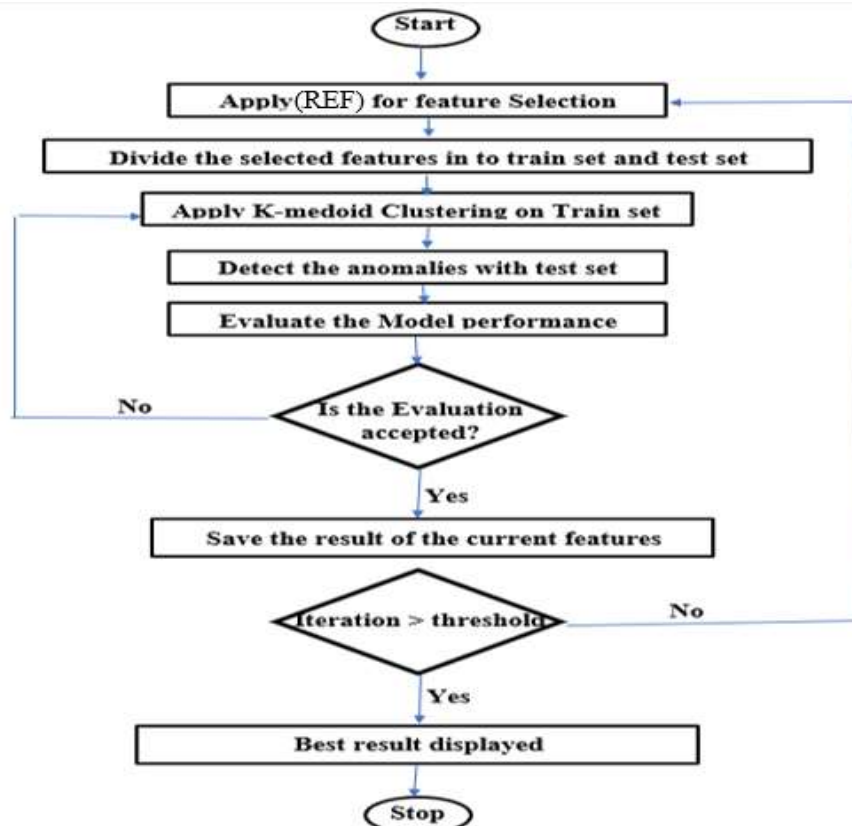


Figure 5: The flowchart of the proposed system

Corresponding author: Hauwa Kulu Haruna.

✉ hauwakuluharuna@gmail.com

Department of mathematical Sciences, Abubakar Tafawa Balewa University, Bauchi.

© 2022. Faculty of Technology Education. ATBU Bauchi. All rights reserved

Table1: The result of the model when implemented with the two prominent datasets

Dataset	Original Features	Best feature	Original Dataset	Sample Dataset	Match datapoints	Accuracy (%)
CICIDS2017	79	33	225745	2257	1228	54.4
KDD CUP 99	42	28	487877	2439	2157	88.4

The complete research results are seen in table 4.1 which represent the datasets of both the two permitted in the cloud anomaly detection model. In the table the original characteristics of the datasets were presented, which takes the headings as the original features of CICIDS2017 and KDDCUP 99 as 79 and 42 respectively. The total number of the original datasets were both 225745 and 487877 for both CICIDS2017 and KDDCUP 99 respectively. On

this table best features selected using the Recursive feature elimination technique were 33 and 28 for both CICIDS2017 and KDD CUP99 respectively. Also, the sample datasets were 2257 and 2439 for both CICIDS2017 and KDD CUP99 respectively with their matches points as 1228 with 54.4% accuracy and 2157 with 88.4% accuracy for CICIDS2017 and KDD CUP99 respectively. This vital information can be graphically explained in figure 4.18.

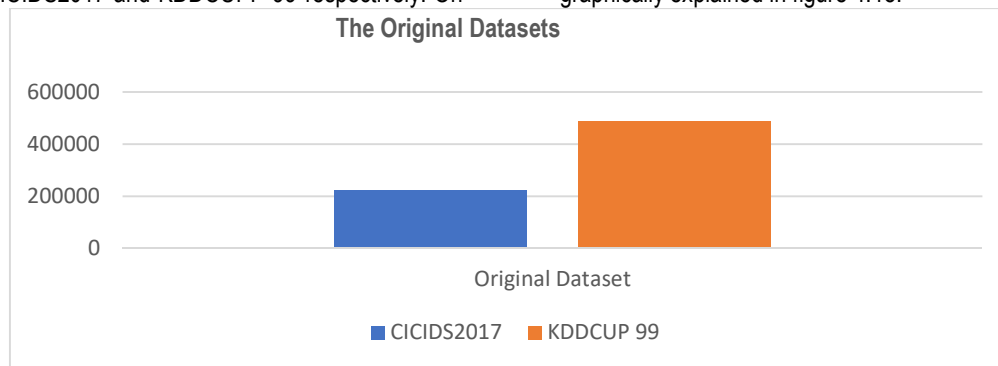


Figure 6: The Original Datasets used in the model

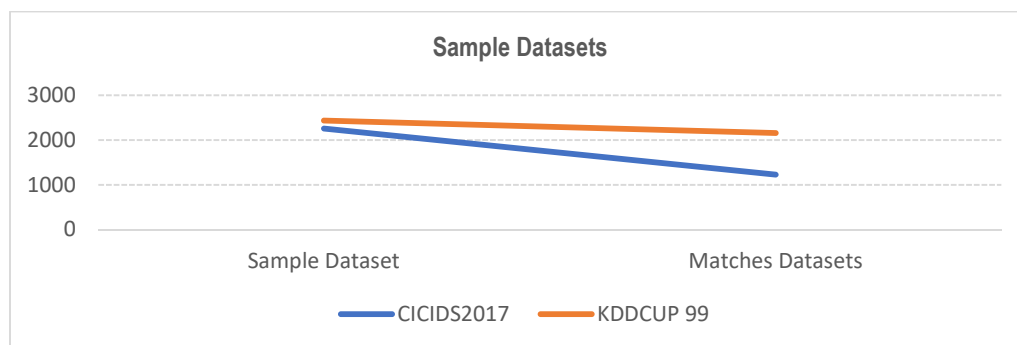


Figure 7: The Sample Datasets used in the model

Figure 7 shows the sample datasets considered in the implementation which is justified best on the number of the dataset. The graph shows how close is the datasets points to

balance the implementation, Despite that the matches show a wide gap between the KDDCUP 99 and CICIDS2017. Indeed KDDCUP 99 plays

Corresponding author: Hauwa Kulu Haruna.

✉ hauwakuluharuna@gmail.com

Department of mathematical Sciences, Abubakar Tafawa Balewa University, Bauchi.

© 2022. Faculty of Technology Education. ATBU Bauchi. All rights reserved

above CICIDS2017 dataset in clustering anomalies.

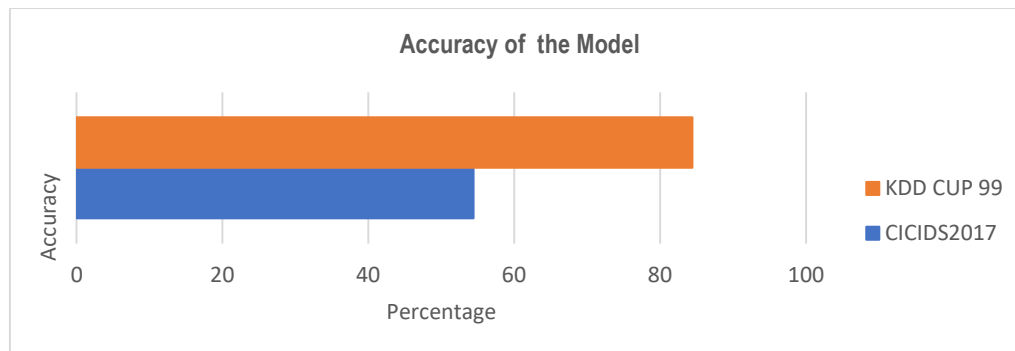


Figure 8: The accuracy of the model on the two different datasets

The expression of figure 8 is an indication that KDDCUP 99 is more accurate in the proposed model than the proposed CICIDS2017 dataset.

Comparing the Proposed System and the Existing System

Table2: The Comparative study of the two models with two different datasets

Author (s)	Accurarcy rate	Dataset	Achieved Improvement
Proposed ADS	88.4	KDD-CUP 99	-
Existing ADS	87.3		1.1
Proposed ADS	54	CICIDS2017	-
Existing ADS	32.5		21.5

Table 2 shows an achievement made in the system implementation. It shows that with KDD-CUP 99 dataset the proposed system slightly outperformed the existing model with 1.1%, which is a remarkable improvement in research. When comparing the two models on CICIDS2017 dataset also the proposed system achieved 21.5 accuracy rates on the existing model, this shows that the model can be clear alternative to the existing model.

RESULT DISCUSSION

The research is aimed at building a k-medoids clustering techniques to detect cloud security related anomalies. Basically, the proposed model was design using two main Machine Learning Algorithm, Two distinct datasets were used in the model, which are; the CICIDS2017 and KDDCUP99 datasets implemented differently in the model. Recursive

Feature Elimination (RFE) method was used to select the best features that were used in the model. In both the two datasets, about 28 out of the original 42 features were selected when KDDCUP99 was used. Then 33 out of the 79 original features were selected in CICIDS2017 datasets. The datasets were original 225745 and 487877 in CICIDS2017 and KDDCUP 99 respectively, 2257 and 2439 were sample from the CICIDS2017 and KDDCUP 99 datasets respectively. This is done for easy implementation considering the large number of the original datasets. After all the preprocessing activities like Data Encoding, Scaling Data type conversion and Data labeling, K-Medoid was built using those selected best features. K-Medoid or Partition around medoids (PAM), divides data points around clusters best on their agreement. The result shows both the two datasets used that KDDCUP 99 have more dissimilarities in terms of

Corresponding author: Hauwa Kulu Haruna.

✉ hauwakuluharuna@gmail.com

Department of mathematical Sciences, Abubakar Tafawa Balewa University, Bauchi.

© 2022. Faculty of Technology Education. ATBU Bauchi. All rights reserved



the anomalies detected in the dataset than the CICIDS2017 dataset. That means KDDCUP 99 has more attacks categories than the CICIDS2017 dataset. From the table 4., it was indicated that out of 2257 data points in CICIDS2017 it has only 1228 matches while in KDDCUP 99, out of 2439 data points it has 2157 matches which is up to 84.4% of the dataset, 34% more than the matches in CICIDS2017 data points. Garg, et al., (2019) designed model leverages Grey Wolf Optimization (GWO) and Convolutional Neural Network (CNN) in two distinct phases for efficient network anomaly detection in cloud setups. In the first phase, GWO is used for feature selection from the network traffic stream repository.

GWO is a meta-heuristic optimization algorithm that helps extract the most relevant features from the streaming data. The model includes an improved version of GWO called Improved-GWO (ImGWO) which enhances the exploration, exploitation, and initial population generation abilities of the standard GWO. The initial population generation in ImGWO is achieved using a uniform distribution approach to ensure a diverse and equally distributed population. The evaluation metrics used include Detection Rate (DR), False Positive Rate (FPR), precision, accuracy, and F-score. The results show that the proposed model outperforms other models in terms of accuracy, detection rate, false positive rate, and F-score. The model achieves high DR in detecting both normal and anomalous classes, low FPR, and high precision in achieving the desired results. The model performs better on KDD'99 dataset compared to DARPA'98 dataset. It also performs well on a synthetic dataset, achieving high DR, precision, accuracy, and F-score with low FPR. The result of the average performance shows 8.25% detection rate, 4.08% false positive rate improvement and 3.62% accuracy improvement.

The result in Table 4.2 indicates that the proposed system performed better than the existing system (Aldallal & Alisa, 2021), When considering the quality of the Algorithms used in the model building. The RFE were mainly applied to improved model performance, reduced

overfitting, most compatible with different algorithm and it help to improve the Robustness of the model. RFE also simplify the model and makes it a bit interpretable by selecting a smaller subset of features. This not only improves interpretability but also reduces computational requirements and training time. With fewer features, it becomes easier to understand and communicate the factors driving the model's predictions. The K-Medoids algorithm used also increased performance of a model by some attributes like, Robustness to outliers, Scalability, Interpretability and Non-parametric distance measures, which add up to the excellent performance of the proposed model, Combining those algorithms in a model harvest a great improvement in the model performance.

CONCLUSION

The main objective of the project is to develop a model for detecting anomalies in cloud security use a combination of RFE and k-medoids clustering techniques. RFE is a feature selection technique commonly used in machine learning to identify the most relevant features for a given task. In this project, RFE is applied to the dataset to extract a subset of features that are most important for detecting cloud security-related anomalies. K-medoids is a clustering algorithm that groups similar data points into clusters, where each cluster is represented by a central data point known as the "medoid." In this project, k-medoids clustering is used to regroup the dataset into subcategories of anomalies, helping to identify different patterns or types of anomalies in the cloud security context.

The implementation process involves several steps. Datasets relevant to cloud security are collected. These datasets contain various features related to cloud activities, configurations, and behaviors. The RFE technique is applied to the collected datasets. This involves iteratively selecting and removing less relevant features from the dataset, aiming to retain the most informative ones. Once the dataset is refined using RFE, the k-medoids clustering technique is applied to group data points into clusters based on their similarity. Each cluster represents a

Corresponding author: Hauwa Kulu Haruna.

✉ hauwakuluharuna@gmail.com

Department of mathematical Sciences, Abubakar Tafawa Balewa University, Bauchi.

© 2022. Faculty of Technology Education. ATBU Bauchi. All rights reserved



subcategory of cloud security-related anomalies. The model's accuracy in detecting anomalies is measured and evaluated. This likely involves comparing the detected anomalies with ground truth labels (if available) and assessing metrics such as precision, recall, F1-score, etc.

The research developed a k-medoids clustering-based approach for detecting cloud security anomalies using Recursive Feature Elimination (RFE) for feature selection. Two distinct datasets, CICIDS2017 and KDDCUP99, are employed, and RFE is utilized to extract the most relevant features. After preprocessing steps like data encoding, scaling, datatype conversion, and labeling, K-medoids clustering is applied using the selected features. The results show that KDDCUP99 exhibits more diverse anomaly categories than CICIDS2017. Specifically, KDDCUP99 achieves an 84.4% match rate compared to CICIDS2017's 54.4%, indicating higher anomaly detection in the former. Another model by Garg et al. (2019) utilizes Grey Wolf Optimization (GWO) and Convolutional Neural Network (CNN) for network anomaly detection in clouds. ImGWO, an enhanced GWO variant, is employed for feature selection, resulting in improved accuracy, detection rate, precision, and F-score. The proposed model surpasses existing ones, particularly excelling on KDD'99 and synthetic datasets. The model's efficiency stems from RFE's feature selection, enhancing model robustness and interpretability, and K-medoids' advantages, including outlier resilience and scalability, which collectively contribute to the model's outstanding performance.

RECOMMENDATIONS

Based on the results and analysis presented in the research, these are three essential recommendations for further study.

- i. **Enhanced Feature Selection Techniques:** While Recursive Feature Elimination (RFE) was effective in selecting relevant features for anomaly detection, exploring other advanced feature selection techniques could be beneficial. Techniques such as genetic algorithms, wrapper methods, and embedded methods like LASSO could be

investigated to determine if they provide even better feature subsets for improved anomaly detection accuracy. Additionally, evaluating the impact of different feature selection approaches on diverse datasets and anomaly types would contribute to a more comprehensive understanding of their effectiveness.

- ii. **Hybrid Model Integration:** Given the positive outcomes of the individual models (k-medoids clustering and GWO-CNN) in the research, a recommendation is to explore the potential benefits of integrating these models into a hybrid system. Combining clustering-based techniques with deep learning approaches could provide a more comprehensive anomaly detection solution, leveraging the strengths of both methods. Investigating the synergy between these techniques and optimizing their integration could lead to enhanced detection capabilities, particularly for complex and evolving cloud security threats.
- iii. **Real-time Anomaly Detection and Adaptation:** The research primarily focuses on offline analysis of cloud security anomalies. Further research could delve into real-time or near-real-time anomaly detection scenarios, which are crucial for prompt threat mitigation in cloud environments. Developing models capable of adapting to dynamic changes in cloud configurations, traffic patterns, and emerging attack vectors would be valuable. This could involve exploring online learning techniques, adaptive clustering algorithms, and methods to update feature selection and model parameters dynamically to ensure effective real-time anomaly detection and response.

REFERENCES

- Al Nafea, R., & Almaiah, M. A. (2021, July). Cyber security threats in the cloud: Literature review. In *2021 International Conference on Information Technology (ICIT)* (pp. 779-786). IEEE.

Corresponding author: Hauwa Kulu Haruna.

✉ hauwakuluharuna@gmail.com

Department of mathematical Sciences, Abubakar Tafawa Balewa University, Bauchi.

© 2022. Faculty of Technology Education. ATBU Bauchi. All rights reserved



- Aldallal, A., & Alisa, F. (2021). Effective Intrusion Detection System to Secure Data in Cloud Using Machine Learning. *Symmetry*, 13(12), 2306.
- Alghofaili, Y., Albattah, A., & Rassam, M. A. (2020). A financial fraud detection model based on LSTM deep learning technique. *Journal of Applied Security Research*, 15(4), 498-516.
- Alhenaki, L., Alwatban, A., Alamri, B., & Alarifi, N. (2019, May). A survey on the security of cloud computing. In 2019 2nd international conference on computer applications & information security (ICCAIS) (pp. 1-7). IEEE.
- Alrashdi, I., Alqazzaz, A., Aloufi, E., Alharthi, R., Zohdy, M., & Ming, H. (2019, January). Ad-iot: Anomaly detection of iot cyberattacks in smart city using machine learning. In 2019 IEEE 9th Annual Computing and Communication Workshop and Conference (CCWC) (pp. 0305-0310). IEEE.
- Assessing Threat-Specific Security Risks in the Cloud. In (pp. 1-6):
- AWS, (2021). Summary of the Amazon S3 Service Disruption in the Northern Virginia (US-EAST-1) Region. <https://aws.amazon.com/message/41926/>
- Binod, A., (2023). Top 15 Cloud Computing Challenges [with Solution]. <https://www.knowledgehut.com/blog/cloud-computing/cloud-computing-challenges>
- Cisco Intrusion Prevention System. Technical Report, Cisco. [Citation Time(s):1]
- Cloud Computing: Solutions and Future Directions. *J. Comput. Sci. Eng.* 2015, 9, 119–133.
- D., & Khan, K. M. (2019). ThreatRiskEvaluator: A Tool for
- David Puzas (2023, January 26). 12 Cloud Security Issues: Risks, Threats & Challenges. <https://www.crowdstrike.com/cybersecurity-101/cloud-security/cloud-security-risks-threats-challenges>
- Fuse, T., & Kamiya, K. (2017). Statistical anomaly detection in human dynamics monitoring using a hierarchical dirichlet process hidden markov model. *IEEE Transactions on Intelligent Transportation Systems*, 18(11), 3083-3092.
- GeeksforGeeks, (2020, December 20). Support vector machine in Machine Learning. GeeksforGeeks. <https://www.geeksforgeeks.org/support-vector-machine-in-machine-learning/>
- GoogleCloud, (2021). Compute Engine Service Level Agreement (SLA) | Google Cloud. <https://cloud.google.com/compute/sla>
- IBM (2021). What is Data Security? Data Security Definition and Overview. IBM. <https://www.ibm.com/topics/data-security>
- IBM Security Network Intrusion Prevention System. Technical Report. <http://www-01.ibm.com/software/tivoli/products/security-network-intrusion-prevention/> [Citation Time(s):1] IEEE.
- Islam, M. S., Pourmajidi, W., Zhang, L., Steinbacher, J., Erwin, T., & Miransky, A. (2021, May). Anomaly detection in a large-scale cloud platform. In 2021 IEEE/ACM 43rd International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP) (pp. 150-159). IEEE.
- John & Mello (2022, JUL 4). 11 Top Cloud Security Threats. CSO Online. <https://www.csoonline.com/article/3043030/top-cloud-security-threats.html>
- Kabla, A. H. H., Anbar, M., Manickam, S., Alamedy, T. A., Cruspe, P. B., Al-Ani, A. K., & Karupayah, S. (2022). Applicability of Intrusion Detection System on Ethereum Attacks: A Comprehensive Review. *IEEE Access*.
- Kazim, M., & Zhu, S. Y. (2015). A survey on top security threats in cloud computing. *International Journal of Advanced*

Corresponding author: Hauwa Kulu Haruna.

✉ hauwakuluharuna@gmail.com

Department of mathematical Sciences, Abubakar Tafawa Balewa University, Bauchi.

© 2022. Faculty of Technology Education. ATBU Bauchi. All rights reserved



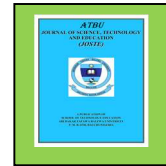
- Computer Science and Applications, 6(3), 109-113.
- Khorshed, M. T., Ali, A. S., & Wasimi, S. A. (2012). A survey on gaps, threat remediation challenges and some thoughts for proactive attack detection in cloud computing. *Future Generation computer systems*, 28(6), 833-851.
- Khraisat, A., Gondal, I., Vamplew, P., & Kamruzzaman, J. (2019). Survey of intrusion detection systems: techniques, datasets and challenges. *Cybersecurity*, 2(1), 1-22.
- Kim, G. Y., Lim, S. M., & Euom, I. C. (2022). A Study on Performance Metrics for Anomaly Detection Based on Industrial Control System Operation Data. *Electronics*, 11(8), 1213.
- Kumar, R., & Goyal, R. (2019). On cloud security requirements, threats, vulnerabilities and countermeasures: A survey. *Computer Science Review*, 33, 1-48.
- Lee, R. (2021). Why anomaly Detection is important, SISU Data, retrieved from <https://sisudata.com/blog/why-anomaly-detection-is-important>, on 8/7/2022
- Lee, T. J., Gottschlich, J., Tatbul, N., Metcalf, E., & Zdonik, S. (2018). Precision and recall for range-based anomaly detection. *arXiv preprint arXiv:1801.03175*.
- Mi, H.B., Wang, H.M., Zhou, Y.F., Lyu, M.R.T. and Cai, H. (2013) Toward Fine-Grained, Unsupervised, Scalable Performance Diagnosis for Production Cloud Computing Systems. *IEEE Transactions on Parallel and Distributed Systems*, 24, 1245-1255. <http://dx.doi.org/10.1109/TPDS.2013.21> [Citation Time(s):2]
- Mondal, A., Paul, S., Goswami, R. T., & Nath, S. (2020, January). Cloud computing security issues & challenges: A Review. In *2020 International Conference on Computer Communication and Informatics (ICCCI)* (pp. 1-5). IEEE.
- Nikita, S. (2021). *Comprehensive Guide To K-Medoids Clustering Algorithm*. Retrieved from: <https://analyticsindiamag.com/comprehensive-guide-to-k-medoids-clustering-algorithm/>. Accessed date: 8/July/2022
- Pang, G., (2020). *ADRepository: Anomaly Detection Datasets with Real Anomalies, A continuously updated anomaly dataset collection*. <https://towardsdatascience.com/adrepository-anomaly-detection-datasets-with-real-anomalies-2ee218f76292>
- Patil, A. (2022, October 31). *Disadvantages of SVM*. <https://iq.opengenus.org/disadvantages-of-svm/>
- Ramachandra, G., Iftikhar, M., & Aslam Khan. (2017). *A Comprehensive Survey on Security in Cloud Computing*. Elsevier, 110, 8. <https://doi.org/10.1016/j.procs.2017.06.124>
- Rittinghouse, J. W., & Ransome, J. F. (2017). *Cloud computing: implementation, management, and security*. CRC press.
- Sehgal, S. (2013, June 30). *Road Towards Cloud Computing – What are the issues? - Part - I. Major Concerns and Problem Workaround*. <https://www.simplilearn.com/cloud-computing-issues-part-i-article>
- Shahzad, F., Mannan, A., Javed, A. R., Almadhor, A. S., Baker, T., & Al-Jumeily OBE, D. (2022). Cloud-based multiclass anomaly detection and categorization using ensemble learning. *Journal of Cloud Computing*, 11(1), 1-12.
- Useni, D. E., Gital, A. Y. U., Lawal, M. A., Emmanuel, O. C., Job, G. K., & Ahmad, A. *A Review of Machine Learning-based Algorithms for Intrusion Detection System*.
- Veritis. (n.d.). *6 Emerging Technologies That Make IT 'Cloud' Stand-out!* Retrieved June 9, 2022, from

Corresponding author: Hauwa Kulu Haruna.

✉ hauwakuluharuna@gmail.com

Department of Mathematical Sciences, Abubakar Tafawa Balewa University, Bauchi.

© 2022. Faculty of Technology Education. ATBU Bauchi. All rights reserved



- <https://www.veritis.com/blog/6-emerging-technologies-that-make-cloud-stand-out/>
- Wang, C.W., Talwar, V., Schwan, K. and Ranganathan, P. (2010) Online Detection of Utility Cloud Anomalies Using Metric Distributions. IEEE Network Operations and Management Symposium (NOMS), Osaka, 19-23
- April 2010, 96- 103. [Citation Time(s):3]
- Wesley, C., & Stephen, B. (2021, December 16). What is Cloud Computing? Everything You Need to Know. SearchCloudComputing. <https://www.techtarget.com/searchcloudcomputing/definition/cloud-computing>.

Corresponding author: Hauwa Kulu Haruna.

✉ hauwakuluharuna@gmail.com

Department of mathematical Sciences, Abubakar Tafawa Balewa University, Bauchi.

© 2022. Faculty of Technology Education. ATBU Bauchi. All rights reserved