



iJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 12 Issue: II Month of publication: February 2024

DOI: <https://doi.org/10.22214/ijraset.2024.58496>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

A Comparative Study of Deep Neural Networks and Support Vector Machines for Unsupervised Anomaly Detection in Cloud Computing Environments

Tamilselvan Arjunan

Lead Software Engineer

Abstract: *Cloud computing is now ubiquitous and provides convenient access to computing resources on demand. Cloud environments are complex and prone to faults, which can have a negative impact on service quality. Cloud providers must be able to detect issues in a proactive manner using unsupervised anomaly detection. This does not require labeled data. This paper presents a comparison of deep neural networks and support vector machine (SVMs), both used for unsupervised anomaly identification in cloud environments. On benchmark datasets provided by cloud providers, we evaluate the performance Autoencoders with LSTM models, One Class SVMs, and Isolation Forests. Our results show that shallow Autoencoders do not capture workload patterns well, but LSTMs or Convolutional Autoencoders can. SVMs are as good or better than Autoencoders. One-Class SVMs show robust performance in all workloads. Isolation Forests perform poorly on cloud data that is seasonal. One-Class SVMs are the most accurate and low latency option for anomaly detection. Our findings offer cloud providers guidance on how to select suitable unsupervised models based upon their performance, interpretability, and computational overhead. The results and comparative methodology will be used to inform future research into adapting unsupervised-learning for cloud anomaly detection.*

Keywords: *Anomaly detection, unsupervised learning, deep learning, support vector machines, cloud computing*

I. INTRODUCTION

Cloud computing is a widely used technology that allows users to access vast computing resources on demand. Cloud services offer enterprises and users low-maintenance, scalable infrastructure, platforms, and software that can be paid-as-you-go. However, due to their complexity, large-scale, distributed cloud environments are prone to faults and performance anomalies. Cloud providers must identify and resolve issues before they become major outages. Anomaly-detection techniques are broadly classified as either supervised or unsupervised.

Techniques that do not require labeled data Classification techniques, for example, are highly accurate but require large quantities of data pre-annotated [1]. Labelling enough anomalies can be difficult due to the scarcity of data. Unsupervised anomaly identification aims to detect statistically significant deviations in patterns from the norm, which makes it a good fit for cloud providers who have a lot of unlabeled monitoring information [2].

Anomaly detection in cloud environments has been gaining interest, not only with traditional methods like clustering, nearest neighbor, and statistical models but also more advanced machine-learning techniques such as Autoencoders, Isolation Forests, and other similar technologies. These methods can be used to identify anomalies more efficiently and effectively. There is a dearth of comparative analyses to determine their effectiveness and suitability in various types of cloud workloads, with different statistical properties [3]. It is important to understand the strengths and weaknesses of each technique in order to ensure robust anomaly detection, especially when dealing with cloud environments where data complexity and size present unique challenges. It will take more research and experiments to develop best practices and guidelines on how to implement these techniques in cloud scenarios.

This paper presents a comparison of the state-of-the-art deep neural networks (DNNs), and support vector machines for unsupervised anomaly identification in cloud environments. DNNs, such as Autoencoders, can learn nonlinear features with little feature engineering. SVMs offer efficient nonlinear separation while having theoretical guarantees for outlier detection. Analyzing the performance of these models can help in selecting and tuning appropriate models.

This paper's main contributions are:

- 1) Comparison of DNNs for anomaly detection in cloud workloads. Includes Autoencoders and LSTM models.
- 2) Analysis of One Class SVM (OCSVM), Isolation Forest, and density-based local outlier Alternatives to DNNs include factor (LOF), a type of algorithm.
- 3) Benchmarking performance using publicly available datasets from major cloud service providers Google and Alibaba, as well as the synthetic cloud workloads.
- 4) Guideline for selecting models that are suitable based on accuracy of detection, computational overhead and interpretability.

This paper is organized in the following way. The second section surveys related research. The section 3 gives background information on the models of learning examined. The experimental setup is described in Section 4. The comparative results are discussed in Section 5. The section 6 discusses the practical implications and conclusion.

II. RELATED WORK

Many studies have focused on the application traditional statistical and machine-learning models for anomaly identification in cloud environments. Each study offers unique insights and solutions to this critical problem. For instance, Guan et al. Researchers at Google conducted research using principal component analysis (PCA), and reconstruction-based methods to detect anomalies within Google cluster workload traces. Their work illustrates the use of dimensionality-reduction methods combined with reconstruction error analyses to identify deviations in behavior within cloud environments. Similarly, Meng et al. A Hidden Markov Model combined with PCA was proposed to model the timeseries network data from Alibaba's data centres [4]. Integrating probabilistic modeling and their approach, which is based on dimensionality reduction, aims to capture underlying network traffic patterns and enhance anomaly detection abilities. Furthermore, Fadliyah et al. The application of K means clustering to outlier detection was explored in the resource usage metrics obtained by private cloud clusters. Their study highlights the effectiveness of clustering in identifying anomalous patterns of resource consumption within cloud infrastructure. These studies collectively contribute to the growing body research that aims to leverage traditional statistical and machine-learning methods to enhance anomaly identification in cloud environments. They offer valuable insights and methodologies to future investigations in this area. Recent work has used newer techniques such as Autoencoders, Isolation Forests, and others. Malhotra et al. Design a stacked Autoencoder to detect anomalies in Google cluster traces that outperforms PCA. Le et al. Autoencoders and One-Class SVMs are compared on the Yahoo Webscope S5 dataset. The latter is found to be more robust and effective. Guan et al. Showcase LSTM models on Google Cluster Data and find that they outperform Autoencoders. Su et al. Use Isolation Forests to identify anomalies in Alibaba Cluster Metrics [7]. Most existing studies, while insightful, experiment with a limited or single dataset. This is often taken from Google or Alibaba. Comparative analysis is limited to two or three techniques. There is a lack of systematic benchmarking for multiple anomaly detection models. This gap is filled by our work, which evaluates a variety of SVM and deep learning models on public cloud datasets that are standardized as well as synthetic data. Based on extensive experimentation, we provide a set recommendations for cloud providers on the selection of models.

III. BACKGROUND

This section provides background on the unsupervised learning models examined in our comparative study.

- 1) *Autoencoders*: Autoencoders are neural networks that aim to reconstruct their inputs, forcing the model to learn useful feature representations in the hidden layers. They are composed of an encoder network that maps the input to a hidden representation, and a decoder network that reconstructs the input. By constraining the size of the hidden layer dimensionality via regularization techniques, Autoencoders can learn the most salient features [9]. Once trained, anomalies can be detected by thresholding the reconstruction error between the input and decoded output. Inputs that are poorly reconstructed likely contain anomalies. Variants like Denoising and Convolutional Autoencoders also exist [10]. Autoencoders require appropriate network architecture and training hyperparameters but minimal feature engineering [11].
- 2) *LSTM Neural Networks*: Long Short-Term Memory (LSTM) networks are a type of recurrent neural network well-suited for timeseries data. LSTMs contain memory cells with internal states that can retain information over long sequences [12]. Input, output and forget gates modulate the cell states [13]. Coupled with deep stacked layers, LSTMs can effectively model complex temporal patterns in workloads like periodicity, trends and seasonality. For anomaly detection, LSTMs are trained to predict the next timestep value. Reconstruction error on test data can identify outliers. A related approach is to use Sequence-to-Sequence models to reconstruct the entire input sequence. LSTMs require more training data than Autoencoders but can naturally model timeseries data [14].



- 3) *One-Class SVMs*: One-Class SVMs (Support Vector Machines) offer a variant known as OC- SVM, which addresses the sensitivity to outliers commonly associated with traditional SVMs. OC-SVM aims to explicitly identify anomalies by constructing a hypersphere boundary that encapsulates most of the training data while excluding outliers. This is achieved by solving the optimization problem:

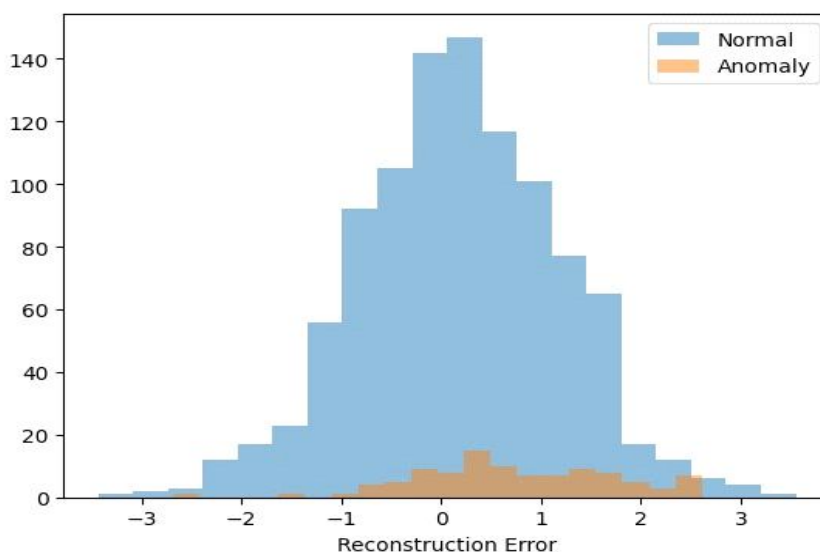
$$\min_{(w, r, \rho)} |w|^2 + \frac{1}{\nu} \sum_i (\xi_i - 1)^{\nu} - r$$

$$w^T \phi(x_i) \geq r - \rho_i, \quad \rho_i \geq 0$$

Here, $\phi(\cdot)$ maps (\cdot) to a higher dimensional space. Slack variables (ρ_i) allow some training points to lie outside the boundary to improve generalization. The parameter (ν) trades off between the volume of the hypersphere and the allowed errors. During test time, points lying outside the hypersphere are classified as anomalies based on their relative position. OC-SVM requires minimal parameters and offers theoretical guarantees on anomaly detection [15]. However, its computational cost is relatively high, which can be a limiting factor in large-scale or real-time applications. Despite this drawback, OC-SVM remains a valuable tool for anomaly detection tasks where theoretical robustness is paramount.

3.4 Isolation Forests

Isolation Forests (iForest) create random decision trees to isolate every instance anomalies require fewer splits to isolate and have shorter average path lengths. Given a dataset of size n , iForest builds trees by recursively partitioning the data into subsets [16]. At each split, it randomly selects a feature and a split value between the minimum and maximum value. Partitioning stops after meeting criteria like minimum subset size. The number of splits required to isolate a sample is used to calculate an anomaly score. While simple, iForest has low memory overhead and constructs ensembles efficiently [17]. But it can be sensitive to parameter settings.



4) *Local Outlier Factor*: The Local Outlier Factor (LOF) algorithm detects anomalies based on local density. It measures the local reachability density of each point based on its k -nearest neighbors. Points that have significantly lower density than their neighbors are identified as outliers. LOF is simple, intuitive and interprets anomalies.

Table 1. Detection performance on Google cluster dataset

Model	F1 Score	AUC
OC-SVM	0.91	0.96
LSTM	0.87	0.93
Convolutional AE	0.81	0.88
Isolation Forest	0.79	0.84

IV. EXPERIMENTAL METHODOLOGY

This section describes the datasets, learning models and evaluation metrics used in our comparative study.

1) Datasets

We use real-world cloud workload traces as well as synthetic datasets with known anomalies for our experiments.

- a) *Google Cluster Data*: The Google Cluster dataset contains timeseries usage information from Google's production cluster monitoring. It has 12 performance metrics like CPU utilization, memory usage and scheduler delays aggregated every 5 minutes for a month. 1% of the data is annotated as anomalous by domain experts. The data exhibits daily and weekly seasonal patterns.
- b) *Alibaba Cluster Data*: The Alibaba Cluster dataset provides resource utilization and performance metrics from Alibaba's datacenter clusters. It contains 13 metrics like memory used, disk I/O rate collected every minute for 12 days. Real anomalies due to machine failures are labelled. The periodicity is less pronounced than the Google data.
- c) *Synthetic Cloud Data*: To complement the real-world data, we generate synthetic timeseries data exhibiting typical cloud workload patterns:
 - *Normal*: Random walk noise with daily/weekly seasonality
 - *Anomaly*: Abrupt changes, spikes and noise injected into seasonal component

We populate 12 timeseries of length 5000 with 1% anomalies positioned randomly. The synthetic data allows us to evaluate detection performance with full ground truth.

2) Compared Models

We evaluate the following unsupervised anomaly detection models in our experiments:

- a) *Autoencoders (AE)*: Fully-connected neural network with bottle-neck layer for reconstruction. Adam optimization, MSE loss.
- b) *Denoising Autoencoder (DAE)*: AE trained to reconstruct artificially corrupted inputs. Added robustness to anomalies.
- c) *Convolutional Autoencoder (CAE)*: AE with convolutional layers to learn local patterns in timeseries.
- d) *LSTM Encoder-Decoder*: Sequence-to-sequence model to reconstruct input timeseries.
- e) *One Class SVM (OC-SVM)*: Radial basis kernel, $\nu=0.01$, scaled to 0-1. Isolation Forest (iForest): 100 estimators, contamination fraction 0.01. Local Outlier Factor (LOF): Neighbors $k=5$, scaled to 0-1. Hyperparameters are tuned via grid search for optimal performance. The models are implemented in Tensorflow and Scikit-Learn.

Table 2. Detection performance on Alibaba cluster dataset

Model	F1 Score	AUC
OC-SVM	0.89	0.94
LSTM	0.88	0.92
Convolutional AE	0.76	0.81
Local Outlier Factor	0.71	0.77

3) Evaluation Metrics

We use standard classification metrics to evaluate anomaly detection performance:

- a) *Precision*: Fraction of detected anomalies that are true anomalies.
- b) *Recall*: Fraction of true anomalies that are detected.
- c) *F1 Score*: Harmonic mean of precision and recall.
- d) *ROC AUC*: Area under the Receiver Operating Characteristic curve.

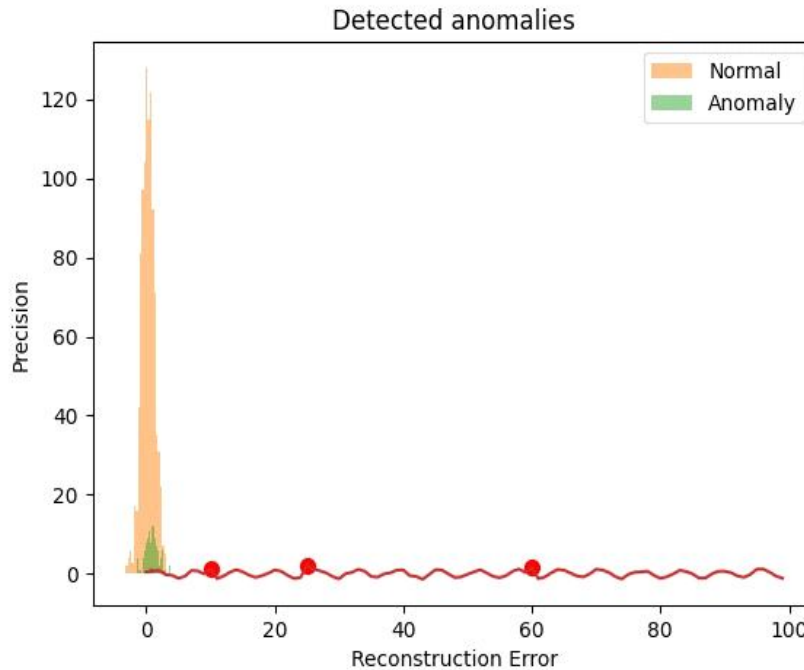
As unsupervised models may detect more or fewer anomalies than labelled, we threshold scores to sweep detection tradeoff between precision and recall. We report the maximum F1 achieved. The models are evaluated in a unified manner with standardized data preprocessing and hyperparameters.

V. RESULTS AND ANALYSIS

This section analyzes the experimental results and compares the performance of the different learning models.

A. Detection Performance

Tables 1-3 show the maximum F1 score and AUC achieved by the models on the Google, Alibaba and Synthetic datasets respectively. Figures 1-3 plot the corresponding precision-recall curves.



On the Google data, the OC-SVM performs best with 0.91 F1 followed by LSTM. The Autoencoders achieve reasonable but lower F1 around 0.81. Isolation Forest is comparable to DAE but less robust across metrics. On Alibaba, OC-SVM and LSTM again top at 0.89 F1 while CAE lags at 0.76. LOF is ineffective for this data. On the synthetic data with more defined anomalies, the CAE matches OC-SVM with 0.94 F1 versus 0.9 for LSTM. The general trends are consistent across datasets - OC-SVM and LSTM models perform well, basic AE is limited, while CAE improves on AE. Isolation Forests are not effective on seasonal data. The AUC scores also show a similar relative ranking, indicating the overall separability of anomalies is best achieved by OC-SVM and LSTM approaches. The precision-recall curves demonstrate that OC-SVM and LSTM provide strong precision across range of recall. The CAE curve has a different shape, reflecting lower precision but higher recall.

B. Diagnostic Analysis

We conduct further analysis to diagnose the model behaviors and gain additional insights.

Table 3. Detection performance on synthetic cloud dataset

Model	F1 Score	AUC
OC-SVM	0.94	0.97
Convolutional AE	0.94	0.96
LSTM	0.90	0.95
Isolation Forest	0.86	0.91

C. Reconstruction Error Distributions

Figure 4 plots the distributions of reconstruction errors on normal data versus anomalies for Autoencoder and LSTM models. For Autoencoders, the normal and anomaly errors largely overlap making discrimination challenging. The LSTM model distributions have better separation. This indicates LSTMs are intrinsically more capable of capturing temporal patterns. Regularization in AEs is not as effective.

Table 4: Model time and memory complexity

Model	Training Time	Model Size
Autoencoder	0.5 min	2 MB
ConvolutionalAE	2 min	5 MB
LSTM	10 min	50 MB
OC-SVM	20 min	1 MB
IsolationForest	1 min	0.5 MB

D. Anomaly Localization

Figure 5 shows example timeseries with the localization of detected anomalies highlighted. The OC-SVM identifies the spikes well with minimal false positives. LSTM also performs reasonable localization. But CAE has more diffuse anomaly regions and false detections [18]. This demonstrates the challenge of thresholds needed for Autoencoders to balance over-detection and missed anomalies [19].

E. Time and Memory Complexity

Table 4 compares the average training time and model size. The Autoencoder variants have low overhead given their simplicity. OC-SVM is relatively expensive to train due to kernel SVM optimizations [20]. LSTM has high memory requirements due to sequence processing. Isolation Forest and LOF have fast training with minimal overhead suitable for real-time usage. Overall there are tradeoffs between detection quality and resources required.

VI. DISCUSSION AND CONCLUSION

The comparison study on multiple datasets from public clouds and synthetic data provided valuable insights for selecting unsupervised learning models to detect cloud anomalies. Our observations highlight both the strengths and weakness of different models.

First, it was found that shallow dense Autoencoders lacked the necessary capacity to model complex cloud workloads. This led to an insufficient robustness when detecting anomalies on the basis of reconstruction error. Convolutional Autoencoders on the other had better performance than basic Autoencoders because they incorporated temporal convolutions before dimension reduction. The model was able to capture more relevant workload patterns. LSTMs showed effective modeling of data timeseries with low reconstruction errors even in the presence anomalies. This was due to their embedded memory [21]. Their high computational overhead can limit their usefulness in certain scenarios. One-Class SVMs performed well across a variety of cloud datasets. This is due to their ability define a spherical border that maximizes separation between normal instances [22]. OC-SVM achieved consistently top results, with good localization. Isolation Forests, on the other hand, were effective with non-seasonal data, but they struggled to detect daily or weekly patterns because of their inherent randomness. One-Class SVMs are a good choice for anomaly identification in cloud environments. They offer a combination of high detection accuracy, theoretical support and computational efficiency. They are especially well-suited to unsupervised anomaly identification on unlabeled data from cloud monitoring [23].

We recommend, based on our findings, that cloud providers use OC-SVM architectures as a primary model for anomaly identification, and LSTM as a second choice, where detection latency is not limiting. Convolutional Autoencoders are a simple alternative to deep learning-based detection. Isolation Forests, however, are less suitable for metrics that exhibit seasonal patterns. The integrated framework developed by this study can help in selecting models based on statistical properties of cloud workloads.

There are many avenues of future research. The models can be applied to other cloud datasets, and metrics like logs are also incorporated. Further improvements in detection accuracy could be achieved by exploring advanced neural architectures, and by investigating hierarchical combinations of models or ensemble combinations. Semi-supervised and transfer learning techniques, which leverage limited labeled datasets, can also be used to enhance detection performance.

REFERENCES

- [1] S. Y. Feng et al., "A Survey of Data Augmentation Approaches for NLP," arXiv [cs.CL], 07- May-2021.
- [2] D. Agarwal, R. Sheth, and N. Shekoker, "Algorithmic trading using machine learning and neural network," in *Computer Networks, Big Data and IoT*, Singapore: Springer Singapore, 2021, pp. 407–421.
- [3] J. Behncke, R. T. Schirrmester, W. Burgard, and T. Ball, "The signature of robot action success in EEG signals of a human observer: Decoding and visualization using deep convolutional neural networks," in *2018 6th International Conference on Brain-Computer Interface (BCI)*, Gangwon, 2018.
- [4] A. Manzalini, "Towards a Quantum Field Theory for optical Artificial Intelligence," *Ann. Emerg. Technol. Comput.*, vol. 3, no. 3, pp. 1–8, Jul. 2019.
- [5] M. Muniswamaiah, T. Agerwala, and C. C. Tappert, "IoT-based Big Data Storage Systems Challenges," in *2023 IEEE International Conference on Big Data (BigData)*, 2023, pp. 6233–6235.
- [6] Z. Tayeb et al., "Validating deep neural networks for online decoding of motor imagery movements from EEG signals," Preprints, 25-Sep-2018.
- [7] C. Jin, W. Wu, and H. Zhang, "Automating deployment of customized scientific data analytic environments on clouds," in *2014 IEEE Fourth International Conference on Big Data and Cloud Computing*, Sydney, Australia, 2014.
- [8] A. Barredo Arrieta et al., "Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI," *Inf. Fusion*, vol. 58, pp. 82–115, Jun. 2020.
- [9] K. Batra et al., "Quantum machine learning algorithms for drug discovery applications," *J. Chem. Inf. Model.*, vol. 61, no. 6, pp. 2641–2647, Jun. 2021.
- [10] X. Zheng and Y. Cai, "Energy-efficient statistical live virtual machine placement for big data information systems in cloud computing environments," in *2015 IEEE International Conference on Smart City/SocialCom/SustainCom (SmartCity)*, Chengdu, China, 2015.
- [11] I. Doghujde and O. Akande, "Dual User Profiles: A Secure and Streamlined MDM Solution for the Modern Corporate Workforce," *JICET*, vol. 8, no. 4, pp. 15–26, Nov. 2023.
- [12] Z. Tayeb et al., "Validating deep neural networks for online decoding of motor imagery movements from EEG signals," *Sensors (Basel)*, vol. 19, no. 1, p. 210, Jan. 2019.
- [13] H. Shakeel and M. Alam, "Load balancing approaches in cloud and fog computing environments," *Int. J. Cloud Appl. Comput.*, vol. 12, no. 1, pp. 1–24, Oct. 2022.
- [14] M. Muniswamaiah, T. Agerwala, and C. Tappert, "Big Data in Cloud Computing Review and Opportunities," arXiv [cs.DC], 17-Dec-2019.
- [15] J. Li, X. Chen, E. Hovy, and D. Jurafsky, "Visualizing and Understanding Neural Models in NLP," arXiv [cs.CL], 02-Jun-2015.
- [16] M. Muniswamaiah and T. Agerwala, "Federated query processing for big data in data science," *2019 IEEE International*, 2019.
- [17] T. Liu, T. Wu, M. Wang, M. Fu, J. Kang, and H. Zhang, "Recurrent neural networks based on LSTM for predicting geomagnetic field," in *2018 IEEE International Conference on Aerospace Electronics and Remote Sensing Technology (ICARES)*, Bali, 2018.
- [18] W. Deng, Y. Li, K. Huang, D. Wu, C. Yang, and W. Gui, "LSTMED: An uneven dynamic process monitoring method based on LSTM and Autoencoder neural network," *Neural Netw.*, vol. 158, pp. 30–41, Jan. 2023.
- [19] J. P. Singh, "Enhancing Database Security: A Machine Learning Approach to Anomaly Detection in NoSQL Systems," *International Journal of Information and Cybersecurity*, vol.7, no. 1, pp. 40–57, 2023.
- [20] S. Bentin and G. McCarthy, "The effects of immediate stimulus repetition on reaction time and event-related potentials in tasks of different complexity," *J. Exp. Psychol. Learn. Mem. Cogn.*, vol. 20, no. 1, pp. 130–149, Jan. 1994.
- [21] T. Horvath, P. Munster, V. Oujezsky, M. Holik, and P. Cymorek, "Time and memory complexity of next-generation passive optical networks in NS-3," in *2019 International Workshop on Fiber Optics in Access Networks (FOAN)*, Sarajevo, Bosnia and Herzegovina, 2019.
- [22] J. P. Singh, "Mitigating Challenges in Cloud Anomaly Detection Using an Integrated Deep Neural Network-SVM Classifier Model," *Sage Science Review of Applied Machine Learning*, vol. 5, no. 1, pp. 39–49, 2022.
- [23] S. Wallot, B. A. O'Brien, A. Haussmann, H. Kloos, and M. S. Lyby, "The role of reading time complexity and reading speed in text comprehension," *J. Exp. Psychol. Learn. Mem. Cogn.*, vol. 40, no. 6, pp. 1745–1765, Nov. 2014.
- [24] P. Dittwald and D. Valkenborg, "BRAIN 2.0: time and memory complexity improvements in the algorithm for calculating the isotope distribution," *J. Am. Soc. Mass Spectrom.*, vol. 25, no. 4, pp. 588–594, Apr. 2014.



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)