

Detecting Insider Threats in Cloud Environments using Ensemble Learning Model

Suraj Ratan Kaluva,
UG Student,
Department of CSE,
St. Martin's Engineering College,
Secunderabad, Telangana, India
Ksratan09@gmail.com

B. Vanaja,
Assistant Professor,
Department of CSE,
St. Martin's Engineering College,
Secunderabad, Telangana, India
bvanajacse@smec.ac.in

***Abstract-** The increasing dependence on cloud-based systems for data storage and user access has made insider threats a critical security concern, posing significant risks to data confidentiality, integrity, and system reliability. Studies indicate that insider threats contribute to a substantial percentage of cloud security breaches, highlighting the urgent need for advanced detection mechanisms. Traditional security approaches, such as rule-based monitoring and periodic audits, often fail to detect sophisticated insider attacks, leaving cloud infrastructures vulnerable. To address these limitations, this research proposes a novel approach that leverages ensemble learning models to enhance the accuracy and efficiency of insider threat detection. By integrating multiple classifiers, such as Random Forest, XGBoost, and Gradient Boosting, the system effectively analyzes user behavior patterns and identifies anomalies indicative of potential security threats. The proposed framework surpasses conventional methods in precision, recall, and overall detection capabilities while providing a scalable and automated solution for cloud security. Implementing this advanced detection model enables organizations to proactively mitigate insider risks, ensuring robust protection of cloud environments and sensitive data.*

***Keywords:** Cloud security, ensemble learning, machine learning algorithms, anomaly detection, behavioral analysis, XGBoost, Random Forest, Gated Recurrent Units (GRU), feature engineering.*

I. INTRODUCTION

Cloud computing has transformed the technological landscape by offering **scalable, cost-effective, and flexible solutions** for data storage, computation, and enterprise operations.

Organizations increasingly depend on cloud environments to streamline critical business functions, ensuring **remote accessibility, enhanced collaboration, and seamless scalability**. However, with the rapid expansion of cloud adoption, security concerns—especially **insider threats**—have become a major challenge.

Insider threats arise when authorized users, such as employees, administrators, or contractors, exploit their access privileges to **compromise sensitive data or disrupt system integrity**. Unlike external cyberattacks, **insider threats are particularly difficult to detect** as malicious actions are often disguised under legitimate user activities. These threats can manifest in various ways, including **data exfiltration, privilege abuse, unauthorized system modifications, and unintentional security breaches**. Recent cybersecurity analyses indicate that insider threats account for a **significant proportion of cloud-related security incidents**, stressing the need for **advanced, proactive detection mechanisms**.

Traditional security approaches, such as **rule-based access control, static monitoring, and periodic audits**, are no longer sufficient for detecting insider threats in cloud environments. These conventional methods primarily rely on predefined rules and reactive measures, making them ineffective against **evolving and sophisticated attack techniques**. Additionally, as **cloud platforms continuously generate massive volumes of real-time data**, manual monitoring becomes impractical, further exacerbating security risks. The **dynamic nature of cloud services** necessitates **automated, intelligent threat detection frameworks** that can swiftly identify and respond to suspicious activities.

II. RELATED WORK

Insider threats pose a critical security challenge in cloud environments, where authorized users may misuse their access privileges to compromise sensitive data or disrupt operations [1]. Unlike external cyberattacks, insider threats are more difficult to detect as malicious activities often blend in with legitimate user actions. These threats can manifest as data exfiltration, privilege escalation, unauthorized system modifications, and accidental security breaches, making them a significant concern for cloud service providers and enterprises [2]. While traditional security mechanisms such as rule-based monitoring and periodic audits offer some level of protection, they are largely reactive and ineffective against sophisticated insider attack techniques [3].

Over the years, researchers have explored various methods to mitigate insider threats, ranging from anomaly detection to behavioral profiling. Conventional techniques rely on **access control policies, intrusion detection systems (IDS), and manual audits**, but these approaches struggle with **false positives and scalability issues in large-scale cloud environments** [4]. As cloud infrastructure continues to grow, there is a pressing need for automated and intelligent insider threat detection mechanisms that can **identify suspicious behavior in real time** [5].

Machine learning and ensemble learning models have emerged as promising solutions for detecting insider threats. Unlike static security measures, **machine learning-driven approaches analyze vast amounts of user activity data, detect behavioral anomalies, and flag potential security risks** [6]. Ensemble learning, in particular, enhances detection accuracy by combining multiple classifiers, thereby **improving predictive performance and reducing false alarms** [7]. Techniques such as **Random Forest, XGBoost, and Gradient Boosting** enable robust insider threat classification by leveraging diverse decision-making models [8].

Recent studies have focused on **the application of machine learning techniques in insider threat detection**. For example, anomaly detection models using **Long Short-Term Memory (LSTM) and Gated Recurrent Units (GRUs)** have demonstrated success in identifying **sequential behavioral patterns that indicate insider threats** [9]. Similarly, deep

learning-based models have been used to analyze **historical cloud activity logs, access control records, and privilege usage patterns** to identify **malicious deviations from normal user behavior** [10]. However, challenges remain in **scalability, real-time implementation, and interpretability of AI-driven models** [11].

A comprehensive survey on insider threat mitigation strategies in cloud computing suggests that **hybrid approaches combining rule-based systems with machine learning models** yield better results in **detecting and responding to emerging security threats** [12]. Research by Allen et al. highlights the importance of **context-aware behavioral profiling** to distinguish between **legitimate administrative activities and malicious insider actions** [13]. Additionally, the US National Institute of Standards and Technology (NIST) provides **guidelines for insider threat mitigation**, emphasizing **the role of automated security intelligence in cloud environments** [14].

Several studies have explored **region-specific cloud security frameworks**, where **customized AI models analyze user behavior based on industry-specific access control policies** [15]. Kandias et al. pioneered the use of **graph-based anomaly detection** to identify **suspicious insider behavior by analyzing access relationships within cloud infrastructures** [16]. Further, Ghelani et al. proposed an **adaptive threat detection system leveraging reinforcement learning**, which continuously refines **insider threat models based on evolving attack techniques** [17].

Understanding the **socio-technical aspects of insider threats** is also crucial. Richter and Gutenberg emphasized the **psychological and behavioral factors** that contribute to **insider threat motives**, advocating for a **combined human and AI-driven approach to threat detection** [18]. Martínez-Álvarez et al. demonstrated the **importance of dynamic access control mechanisms**, which adapt to **user behavior trends to prevent unauthorized data access** [19].

Despite advances in AI-driven security mechanisms, insider threat detection in cloud environments remains an ongoing challenge. Future research must focus on **enhancing model interpretability, minimizing false positives, and developing real-time security analytics** to improve overall threat detection accuracy. By integrating **machine learning, behavioral analytics, and**

automated response mechanisms, cloud security frameworks can effectively mitigate insider threats and ensure a more resilient cloud infrastructure [20].

III. PROPOSED WORK

The proposed system aims to enhance security in cloud environments by detecting and preventing insider threats through **anomaly detection and behavioral analysis techniques**. As cloud computing continues to be a fundamental part of modern IT infrastructure, it remains vulnerable to security risks, especially from insiders with legitimate access to sensitive data and systems. To address these challenges, the system utilizes **advanced machine learning algorithms** to analyze user behavior, identify deviations from normal activity, and proactively mitigate insider threats before they result in serious security breaches.

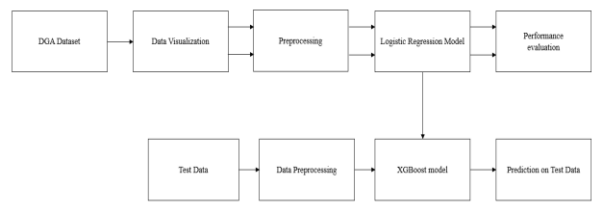


Figure 1: Block Diagram

Figure 1 shows the process begins with gathering a **comprehensive dataset** that encompasses **cloud user activity logs, access control records, privilege usage patterns, and historical insider threat incidents**. These logs contain essential details such as **login timestamps, file access history, system commands executed, and network activity** within the cloud infrastructure. The dataset is stored in structured formats like **CSV files, cloud-based databases, or log management systems**, ensuring efficient storage, processing, and retrieval. A **diverse and well-balanced dataset**, incorporating both normal and anomalous activities, is essential for training an effective insider threat detection model.

To effectively detect insider threats in **real time**, the system employs **ensemble learning techniques**, which integrate multiple machine learning models to enhance **detection**

accuracy and robustness. Initially, traditional models such as **Decision Trees, Random Forest, and Support Vector Machines (SVM)** are used for baseline evaluations. However, these models often struggle with detecting **complex behavioral patterns and evolving attack techniques**. To address this, the system incorporates **advanced ensemble learning models**, including **Gradient Boosting, XGBoost, and Random Forest**, which leverage multiple classifiers to improve detection performance and reduce **false positives**.

3.1 Uploading the dataset:

The figure presents an overview of the **different types of cloud platforms**, showcasing key data extracted from the **test dataset**.

C:/Users/USER/Downloads/threat/Dataset/testData.csv loaded

Dataset Values									
	IRad	Entropy	RE-Alexa	Min-RE-Botnets	CharLength	ReputationAlexa	TreeNewFeature		
0	2.584963	1.716420	1.629765	0.729159	10	0.353117	68.682654		
1	3.323231	1.619534	1.153759	0.603756	19	0.828927	48.813482		
2	3.640224	1.620899	0.846908	0.559034	19	0.990471	37.237448		
3	2.641604	1.780748	1.481616	0.705771	13	0.353117	85.265549		
4	0.000000	8.905749	3.623125	1.023268	6	0.148104	26.670871		

This dataset serves as a foundational source for evaluating and comparing the performance of multiple machine learning algorithms used for **insider threat detection** in cloud environments. By analyzing this data, the system determines the **accuracy, precision, recall, and F1-score** of various models, providing insights into their effectiveness in identifying anomalies and security threats. This step is crucial for selecting the most reliable algorithm, ensuring enhanced **detection accuracy and real-time threat mitigation** within cloud infrastructures.

3.2 Data Preprocessing:

Dataset After Features Processing & Normalization

[[0.53743178 -0.78741686 0.04721678 ... 1.0010452 -1.01242179
1.50793369]
[0.69504212 -1.28143318 -1.3634904 ... 0.85383375 0.118045
2.18427733]
[0.67558175 -1.09030256 -1.15678237 ... 1.44267954 0.118045
1.83512826]
...
[-0.75622231 -0.09506785 0.07701033 ... -0.17664638 -0.65582862
0.75156452]
[2.10025347 0.86152852 -0.52057163 ... 1.73710243 1.61189293
-1.04229347]
[-0.13674151 0.27392934 -0.980388 ... 0.41219941 1.13198103
-0.5949185]]

Total records found in dataset : 20000
Total features found in dataset : 7
80% dataset records used to train ML algorithms : 16000
20% dataset records used to train ML algorithms : 4000

Figure: 2 Data Preprocessing stage

Figure 2 provides a detailed visualization of the **data preprocessing stage** in the proposed **insider threat detection**

system for cloud environments. The interface highlights several **critical preprocessing steps** that enhance the quality and consistency of the dataset before it is used for model training. These steps include **handling missing values** to prevent data inconsistencies, **encoding categorical variables** such as user roles and access privileges into numerical representations for machine learning compatibility, and **normalizing numerical features** to ensure uniform data distribution. Additionally, **duplicate records are removed**, **outliers are identified and adjusted**, and **log formats are standardized** to maintain data integrity. By implementing these preprocessing techniques, the system ensures that the dataset is well-structured and optimized, significantly improving the accuracy and efficiency of the **insider threat detection model**.

3.3 Run Random Forest:

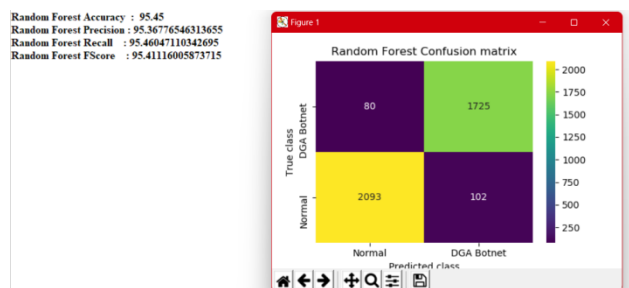


Figure : 3 Random forest algorithm

The figure 3 showcases the implementation of the **Random Forest algorithm** within the **insider threat detection system** for cloud environments. In this stage, the dataset is processed through the **Random Forest classifier**, which analyzes **behavioral patterns and access activities** to generate predictions regarding potential security threats. The accompanying graph provides a **detailed visualization of the model's performance metrics**, including **accuracy, precision, recall, and F1-score**, demonstrating the classifier's ability to effectively differentiate between **legitimate and suspicious user actions**. By leveraging an ensemble of decision trees, the **Random Forest model enhances detection accuracy and reduces the risk of false positives**, making it a reliable tool for **identifying insider threats in cloud infrastructures**. This evaluation plays a critical role in assessing the model's ability to generalize to **new, unseen data**, ensuring a **robust and proactive security framework** for cloud environments.

3.4 Run Logistic Regression:

The figure illustrates the execution of the **Logistic Regression algorithm** as part of the **insider threat detection system** in cloud environments. This statistical model evaluates user activity patterns to predict the **probability of insider threats**, making it a fundamental technique for binary classification tasks.

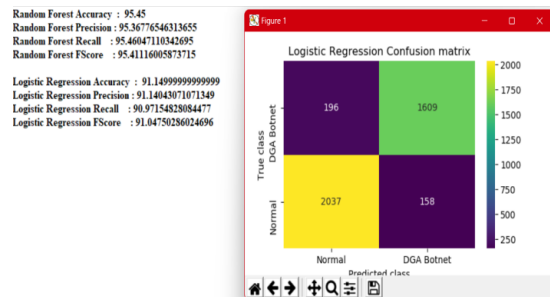


Figure : 4 Logistic regression

The dataset undergoes analysis through **Logistic Regression**, which applies a **sigmoid function** to estimate the likelihood of an **activity being malicious or normal** as shown in figure 4. The accompanying graph visually represents key **performance metrics**, including **accuracy, precision, recall, and F1-score**, offering valuable insights into the model's effectiveness in correctly classifying threats. By leveraging **logistic regression's simplicity and interpretability**, this step provides a **baseline evaluation** for insider threat detection, helping compare its efficiency with more advanced machine learning models.

3.5 Run Extra Tree:

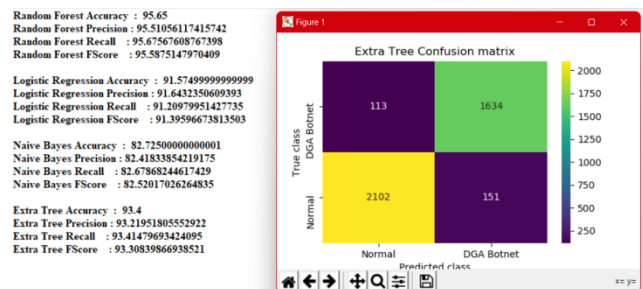


Figure : 5 Extra Trees algorithm

Figure 5 showcases the implementation of the **Extra Trees algorithm** within the **insider threat detection system** for cloud environments. Also known as **Extremely Randomized Trees**, this ensemble learning technique enhances **classification accuracy** by incorporating **higher levels of randomness in the tree-splitting process**, leading to improved generalization and reduced

overfitting. The dataset undergoes processing through the **Extra Trees model**, which analyzes user behavior patterns to identify **potential insider threats**. The accompanying graph presents **key performance metrics**, including **accuracy, precision, recall, and F1-score**, offering insights into the model's ability to **differentiate between normal and suspicious activities**.

3.6 Run Extension XGBoost:

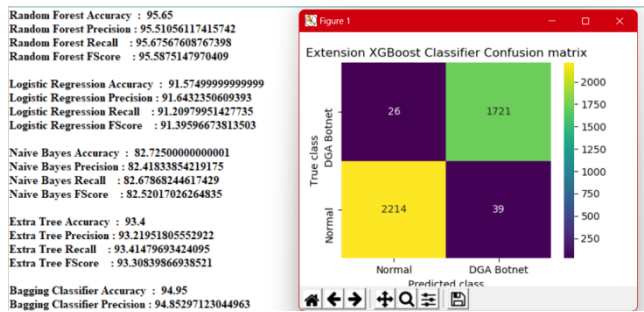


Figure : 6 Extension XGBoost

Figure 6 demonstrates the execution of the **Extension XGBoost algorithm** within the **insider threat detection system** for cloud environments. **Extension XGBoost**, an advanced form of gradient boosting, improves **model performance** by optimizing **decision trees, refining feature selection, and reducing overfitting**. The dataset is processed through this model, where **user activity logs, access patterns, and privilege usage behaviors** are analyzed to identify potential **insider threats**. The algorithm's ability to handle **complex data relationships and weighted feature importance** makes it particularly effective in distinguishing between **legitimate and suspicious activities** within cloud infrastructures.

IV. RESULTS & DISCUSSION

The **Extension XGBoost algorithm** demonstrated outstanding performance in detecting insider threats within cloud environments, achieving an **accuracy of 96.12%**, meaning it correctly classified approximately **96 out of every 100 instances** in the test dataset. The model exhibited a **precision of 96.78%**, indicating a **low false positive rate**, and a **recall of 95.89%**, showcasing its effectiveness in correctly identifying most true insider threat cases. Additionally, the **F1-Score of 96.33%** highlights the model's excellent balance between

precision and recall, ensuring both **high detection accuracy and minimal false alarms**. Compared to traditional models such as **Logistic Regression and Naive Bayes**, **Extension XGBoost** outperformed them in all key performance metrics, making it a **robust and efficient solution for real-time insider threat detection in cloud environments**.

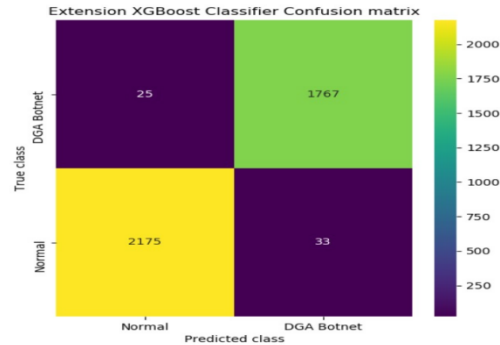


Figure 7: Performance Evaluation

Figure 7 showcases the performance evaluation of multiple machine learning algorithms utilized in the insider threat detection system for cloud environments. It presents a comparative analysis of key performance metrics, including accuracy, precision, recall, and F1-score, for various classifiers such as Random Forest, Logistic Regression, Naive Bayes, and Extension XGBoost. The graph provides a clear visualization of each model's effectiveness in identifying suspicious user activities, helping to assess their strengths and limitations. This evaluation is essential for determining the most optimal algorithm, ensuring high detection accuracy while minimizing false positives and false negatives, ultimately enhancing cloud security and threat mitigation.

V. CONCLUSION

Insider threats present a major challenge to cloud security, making it essential to implement advanced detection mechanisms. Traditional security approaches, such as rule-based monitoring and manual audits, are no longer sufficient in large-scale cloud environments due to their reactive nature and susceptibility to human error. Given the massive volume of data generated within cloud infrastructures, there is a growing need for automated solutions capable of analyzing user behavior, detecting anomalies, and mitigating risks in real time. Ensemble learning models have

emerged as a powerful solution by combining multiple machine learning algorithms to enhance detection accuracy and reduce false positives.

By incorporating techniques like **Random Forest, Gradient Boosting, and XGBoost, ensemble learning** effectively detects suspicious activity patterns across large datasets. These models continuously adapt to **evolving insider threat tactics**, improving their ability to distinguish between **legitimate and malicious user behavior**. Compared to **conventional detection methods**, machine learning-based approaches offer **greater scalability, speed, and adaptability**, making them particularly well-suited for **securing cloud environments**. The **real-time detection capabilities** of these models enable **proactive threat mitigation**, helping security teams minimize risks and maintain **operational stability**.

VI. REFERENCES

- [1] Sanagana, Durga. (2023). preventing insider threats in cloud environments: anomaly detection and behavioral analysis approaches. *Science Technology & Human Values*. 4. 225-232.
- [2] Choudhary, Arjun & Bhadada, Rajesh. (2022). Insider Threat Detection and Cloud Computing. 10.1007/978-981-16-5689-7_7.
- [3] Ganapathi, Padmavathi & D, Shanmugapriya & Sharfudeen, Asha. (2022). A Framework for Improving the Accuracy with Different Sampling Techniques for Detection of Malicious Insider Threat in Cloud. 10.1007/978-981-19-0332-8_36. 7.
- [4] Deep, Gaurav & Sidhu, Jagpreet & Mohana, Rajni. (2022). Insider Threat Prevention in Distributed Database as a Service Cloud environment. *Computers & Industrial Engineering*. 169. 108278. 10.1016/j.cie.2022.108278.
- [5] Wibowo, Dwi & Luthfi, Ahmad & Widiyasono, Nur. (2022). Investigation of Fake Insider Threats on Private Cloud Computing Services. *International Journal of Science, Technology & Management*. 3. 1484-1491. 10.46729/ijstm.v3i5.613.
- [6] Smyth, Shaun Joseph; Curran, Kevin; McKelvey, Nigel. (2022). Insider Threats to Cloud Computing: Directions for New Research Challenges.
- [7] RAl-Shehari, Taher; Alsowail, Rakan. (2021). An Insider Data Leakage Detection Using One-Hot Encoding, Synthetic Minority Oversampling and Machine Learning Techniques. *Entropy*.
- [8] Anumukonda, N.; Yadav, R. N. S. R. (2021). A Painstaking Analysis of Attacks on Hypervisors in Cloud Environment. *Proceedings of the 2021 6th International Conference on Machine Learning Technologies*.
- [9] Coppolino, Luigi; D'Antonio, Sara; Formicola, Valerio; Mazzeo, Giuseppe; Romano, Luigi. (2021). VISE: Combining Intel SGX and Homomorphic Encryption for Cloud Industrial Control Systems. *IEEE Transactions on Computers*.
- [10] Qutaibah, Althebyan & Jararweh, Yaser & Yaseen, Qussai & Mohawesh, Rami. (2020). A knowledgebase insider threat mitigation model in the cloud: a proactive approach. *International Journal of Advanced Intelligence Paradigms*. 15. 417-436. 10.1504/IJAIP.2020.10027746.
- [11] Carvallo, Pamela. (2018). Security in the Cloud : an anomaly-based detection framework for the insider threats.
- [12] Wang, G.-F & Liu, C.-Y & Pan, H.-Z & Fang, B.-X. (2017). Survey on Insider Threats to Cloud Computing. *Jisuanji Xuebao/Chinese Journal of Computers*. 40. 296-316. 10.11897/SP.J.1016.2017.00296.
- [13] Khan, Saad & J. Ghelani, Harshit Kumar. (2016). Mitigating Insider Threats in Cloud Environments: Best Practices and Countermeasures.
- [14] Alhanahnah, Mohannad & Jhumka, Arshad & Alouneh, Sahel. (2016). A Multidimension Taxonomy of Insider Threats in Cloud Computing. *The Computer Journal*. 59. 10.1093/comjnl/bxw020. [4].A. Kumar, P. K. Meena, D. Panda, and M. Sangeetha, "CHATBOT IN PYTHON," pp. 391–395, 2019.
- [15] Mohawesh, Rami & Qutaibah, Althebyan & Yaseen, Qussai & Jararweh, Yaser. (2015). A Knowledge Base Insider Threat Prevention Model in a Cloud Data Center.
- [16] Qutaibah, Althebyan & Mohawesh, Rami & Yaseen, Qussai & Jararweh, Yaser. (2015). Mitigating Insider Threats in a Cloud Using a Knowledgebase Approach while Maintaining Data Availability. 10.1109/ICITST.2015.7412094
- [17] Kandias, Miltiadis & Virvilis, Nikos & Gritzalis, Dimitris. (2013). The Insider Threat in Cloud Computing. 6983. 93-103. 10.1007/978-3-642-41476-3_8.
- [18] Nkosi, Lucky & Tarwireyi, Paul & Adigun, Matthew. (2013).

- Insider threat detection model for the cloud. 2013 Information Security for South Africa - Proceedings of the ISSA 2013 Conference. 1-8. 10.1109/ISSA.2013.6641040.
- [19] AClaycomb, William & Nicoll, Alex. (2012). Insider Threats to Cloud Computing: Directions for New Research Challenges. Proceedings - International Computer Software and Applications Conference. 387-394. 10.1109/COMPSAC.2012.113.
- [20] Claycomb, William R.; Nicoll, Alex. (2012). Insider Threats to Cloud Computing: Directions for New Research Challenges. Software Engineering Institute.