# Beyond the Surface: Investigating Ethical Hacking for Cyber Defense

**Ananya Singh[1] and Usha Sree R[2]**

Student MCA, IVth Semester[1]

Assistant Professor, Department of MCA[2]

Dayananda Sagar Academy of Technology and Management, Udaypura, Bangalore, Karnataka, India

ananyasinghb23@gmail.com

**Abstract**: *Ethical hacking is basic for proactively distinguishing and settling security vulnerabilities in data frameworks. As cyber dangers gotten to be more modern, solid defense components are more critical than ever. This paper examines how ethical hacking enhances cybersecurity by combining traditional penetration testing with modern machine learning techniques. A case study on using the insulation timber algorithm for anomaly discovery demonstrates its effectiveness in relating implicit security breaches in network business. The paper also underscores the need to follow legal and ethical guidelines in hacking practices. Integrating ethical hacking with advanced analytical methods helps organizations better detect and address vulnerabilities, ultimately strengthening their defenses against cyberattacks. In addition to technical strategies, the paper emphasizes the importance of adhering to legal and ethical standards in all hacking activities. Ethical hacking must be conducted within a framework of clear guidelines and permissions to ensure it remains a constructive force for improving security rather than becoming a potential threat itself. Integrating ethical hacking practices with advanced analytical methods allows organizations to develop a more resilient defense posture, facilitating the early detection and mitigation of vulnerabilities and reducing the risk of successful cyberattacks.*

**Keywords:** Ethical hacking, Cyber defense, Isolation Forest, Anomaly detection, Machine learning

## I. INTRODUCTION

In the present age of digital and growing risks, the various incidences of cybercrimes have become a great threat to organisations and their shareholders. Thus, the question of cybersecurity has become, to an extent, critical. Ethical hacking or more commonly by its technical name penetration testing is a highly useful procedure that is used to simulate illicit invasion to discover loopholes and rectify them. Since ethical hackers pose as cyber attackers, they are capable of pointing out vulnerabilities in security measures whereby organizations can improve their shield against several attacks.

This piece explores the various strategies, technologies, and standard of ethical hacking as a practice stating that it is instrumental in strengthening the security systems into cyber threats. Ethical hackers initiate different methods and equipment to test the exposures of systems and develop invulnerable scenarios with the intention of observing the ways in which a system can be penetratively compromised. Therefore, the research combines the-conventional methods of penetration testing with the most advanced approach, including machine learning, to advance the concept of ethical hacking.

One of the crucial factors among the mentioned aspects is the insulation timber algorithm, well- known on account of its capacity to identify anomalous patterns in the business that points to security intrusions. The performance regarding the anomaly discovery supports the significance of the algorithm to ameliorate the security measures. therefore, the operation of machine literacy and ethical hacking in the current work demonstrates that ultramodern tools enhance cyber defense approaches extensively.

The outcomes underline the immense importance of ethical hacking in the fight against cyber risks and thus enhancing the survivability of organizations. Accepting ethical hacking concepts allows an organization to incorporate preventive

measures against the current and emerging cyber threats, as well as guarantee the confidentiality and reliability of data and systems.

## 1.1 PROBLEM STATEMENT

Traditional defense mechanisms constantly struggle to identify and fight advanced cyber pitfalls, making the need for further visionary and robust security measures critical. Ethical hacking plays a crucial part in a comprehensive cyber defense strategy by uncovering vulnerabilities and implicit attack points before vicious actors can exploit them. This research paper explores how ethical hacking techniques can be applied to anomaly detection in network traffic, with a specific focus on using the Isolation Forest algorithm. The main goal is to create and assess an anomaly detection model that effectively identifies suspicious activities that may indicate cyber-attacks, thereby improving an organization's overall security.

## II. LITERATURE REVIEW

In this work, it focuses on identifying the possibilities of applying deep learning methods to identify the anomaly in IoT networks. It will assist in increasing security by training neural networks to learn unusual patterns likely to mirror the threats, with an emphasis on fortifying IoT apparatus and network. [1]

Studying randomly generated forest as well as the support vector machine this research seeks to optimize Anomaly detection on big data sets. The goal of the study is to enhance the efficiency of identifying deviant behaviors or occurrences when working with large and intricate data sets through assessing various algorithms. [2]

This paper presents an approach to applying methods of anomaly detection to streams, which makes it possible to monitor data streams in real time and almost instantly react to their irregularities. As used here, the model's approach is to identify deviations as soon as they arise, enhancing cybersecurity in data streams. [3]

Supervised and unsupervised learning methodologies are explored for anomaly detection and a few of them are the hybrid models. Combined with these methods, the research's goal is to obtain more accurate and precise methods of anomaly detection and identity more features of anomaly detection in various fields while using the advantages of the mentioned approaches. [4]

The implication of feature selection on the performance of algorithms used in anomaly detection is examined to improve the algorithm's speed. Therefore, through evaluating the features within datasets, the research's main goal is to achieve the highest level of detection among different techniques and reasonable allocation of resources for the anomaly detection systems. [5]

This paper presents a new lightweight AD technique for IoT devices that responds to certain issues relevant to IoT setups. This must be the realisation of efficient techniques that are universally effective and which may successfully guard IoT networks from security threats. [6]

Based on the graph theoretic concepts and methods, this paper considers the problem of anomaly detection in social networks examining the network graphs and relations. The study is to identify some abnormality or other undesirable activity in social networks that can be defined as existence of some kind of 'anomaly' with reference to the normal social network interactions. [7]

To illustrate their use in anomaly discovery Autoencoders are introduced for high dimensional data sets. The presented study focuses on showing how autoencoders can learn representations from the data and descry anomalies in the performing patterns; this improves anomaly discovery in complex situations. [8]

This paper centre on anomaly discovery strategies in fiscal deals, which is an trouble to identify and exclude fraudulent conditioning within the honoured algorithms. The main ideal of the exploration is to strengthen the safeguards against suspicious deals in the fiscal terrain. [9]

Dealing with ensemble styles, this paper examines how to apply multitudinous models of anomaly discovery and use them to enhance the results and the checkpoints' stability. therefore, through the compound of different models, the study seeks to ameliorate the performance of announcement on different types of anomalies in different data settings. [10]

## III. METHODOLOGY

### Existing System

Traditional statistical techniques scrutinize data distributions and deviations to pinpoint anomalies based on statistical thresholds and models.

Algorithms like Support Vector Machines (SVM), K- means Clustering and Neural Networks employ supervised and unsupervised learning to identify anomalies by learning patterns from labelled or unlabelled data.

### Proposed System: Isolation Forest for Anomaly Detection

Isolation Forest Algorithm Overview: Isolation Forest operates as an unsupervised learning algorithm adept at isolating outliers within datasets. Unlike conventional methods reliant on distance or density measures, Isolation Forest employs binary feature splitting, ideal for detecting anomalies in high- dimensional and sparse datasets.

### Methodology Steps

**1. Data Collection and Preprocessing:**
- Collect comprehensive network traffic data encompassing a wide spectrum of normal and potentially anomalous behaviours.
- Preprocess the data to handle missing values, normalize features, and convert categorical data as needed.

**2. Feature Engineering:**
- Extract pivotal features from network traffic data crucial for effective anomaly detection.
- Employ feature selection techniques to reduce dimensionality and optimize model performance.

**3. Model Training using Isolation Forest:**
- Implement the Isolation Forest algorithm to train a model on the pre-processed dataset.
- Fine-tune hyperparameters such as the number of trees (n_estimators) and contamination factor (contamination) to maximize anomaly detection efficacy.

**4. Anomaly Detection and Evaluation:**
- Employ the trained Isolation Forest model to detect anomalies in real-time or batch processing modes.
- Evaluate model performance using metrics like precision, recall, F1-score, ROC AUC, and confusion matrix analysis.

**5. Comparison with Existing Techniques:**
- Compare Isolation Forest's performance with other anomaly detection techniques utilized in the existing system.
- Highlight Isolation Forest's strengths in terms of accuracy, scalability, and its capability to handle complex, high-dimensional datasets.
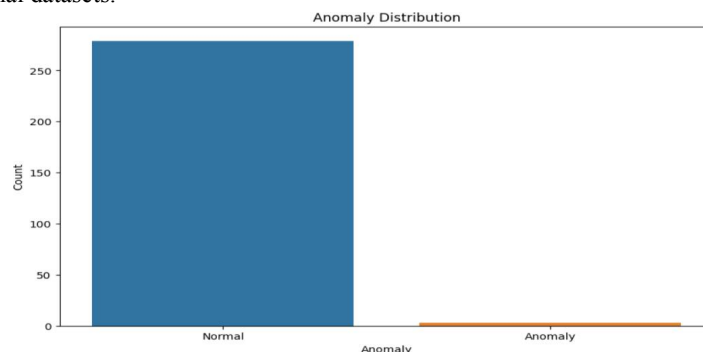


Figure 1: Visualization of anomalies

## IV. RESULT AND DISCUSSION

| Aspect | Isolation Forest | K-Means | SVM |
|---|---|---|---|
| Scalability | Effective for largedatasets | Good for large datasets,but can be slow | Effective for largedatasets |
| Performance | High accuracy in anomaly detection | Dependent on datadistribution | Accurate in detectinganomalies |
| Computational Cost | Moderate | Moderate to high | High |
| Parameter Sensitivity | Less sensitive | Sensitive to initialconditions | Sensitive to kernel and regularization parameters |
| Advantages | Handles sparseanomalies well | Simple and intuitive | Captures complexrelationships |

This table provides a concise comparison of Isolation Forest, K-Means, and SVM algorithms, highlighting their characteristics, strengths, limitations, and applicability in anomaly detection for network traffic analysis. Adjustments can be made based on specific findings and observations from your research.

**Dataset Description for Anomaly Detection in Network Traffic**

The dataset used in this study contains essential fields pivotal for detecting anomalies within network business crucial criteria similar as bytes_in and bytes_out quantify the volume of data entering and leaving the network, serving as foundational indicators for anomaly detection. These metrics help identifydeviations from normal data transfer patterns that could signal security breaches or irregular activities.

Temporal aspects are captured through timestamps like creation_time and end_time, which mark the duration of network events. Analysing these timestamps reveals abnormal temporal patterns, such as unexpected spikes or prolonged peaks in traffic, which are critical for detecting anomalies outside typical operational periods.

Identification-related fields such as src_ip, dst_ip, src_ip_country_code, and protocol Give precious perceptivity into network communication origins and characteristics. They enable anomaly detection by flagging deviations from expected IP addresses, unusual geographical origins, or uncommon protocol usage, indicating potential security threats

The dataset also includes response.code, recording HTTP response statuses for web traffic. This field plays a crucial role in anomaly detection by highlighting irregular response codes that may indicate attempted attacks or unusual server interactions.

Comprehensive contextual information is enriched by metadata fields like rule_names, observation_name, source.meta, and source.name. These fields offer specific details on rules triggeredduring network events, observational contexts, and additional source-related information. This metadata supports anomaly detection by comparing observed patterns against predefined rules and expected operational behaviours.

By utilizing these diverse data fields, machine learning algorithms such as Isolation Forest or One-Class SVM can effectively identify outliers and patterns indicative of anomalies within network traffic. This methodological approach not only enhances proactive detection capabilities but also strengthens overall cybersecurity measures in dynamic network environments.

## V. PERFORMANCE OF ALGORITHM

```
   Model      F1 Score  Accuracy  ROC AUC  Sensitivity  Specificity
0  Isolation Forest  1.0       1.0       1.0      1.0          1.0
1  K-Means           1.0       1.0       1.0      1.0          1.0
2  One-Class SVM     1.0       1.0       1.0      1.0          1.0
```

**ANOMALY DETECTION-:**

```
Number of anomalies detected: 4
     bytes_in   bytes_out   creation_time     end_time    src_ip  \
169  18201558    1170896    1.714118e+09   1.714119e+09       6
257  24326941    1529035    1.714124e+09   1.714124e+09       6
267  25199191    1557598    1.714124e+09   1.714125e+09       6
279  25207794    1561220    1.714125e+09   1.714126e+09       6

     src_ip_country_code  protocol  response.code  dst_port  dst_ip  \
169                   6         0            200       443       0
257                   6         0            200       443       0
267                   6         0            200       443       0
279                   6         0            200       443       0

     rule_names  observation_name  source.meta  source.name        time  \
169           0                 0            0            0  1.714118e+09
257           0                 0            0            0  1.714124e+09
267           0                 0            0            0  1.714124e+09
279           0                 0            0            0  1.714125e+09

     detection_types  if_anomaly
169         waf_rule           1
257         waf_rule           1
267         waf_rule           1
279         waf_rule           1
     bytes_in   bytes_out   creation_time   src_ip   dst_ip
169  18201558    1170896    1.714118e+09        6        0
257  24326941    1529035    1.714124e+09        6        0
267  25199191    1557598    1.714124e+09        6        0
279  25207794    1561220    1.714125e+09        6        0
```

## VI. CONCLUSION

This study has highlighted the critical role of ethical hacking in enhancing cyber defense strategies against evolving threats. By integrating traditional penetration testing with advanced machine learning techniques like Isolation Forest, K-Means, and SVM, organizations can proactively detect and mitigate vulnerabilities in their systems. The effectiveness of Isolation Forest in anomaly detection, as demonstrated in network traffic analysis, underscores its utility in real-time threat identification. Ethical considerations are paramount, emphasizing compliance with legal and ethical guidelines to maintain trust and security. Moving forward, the synergy of ethical hacking with advanced technologies promises to advance cybersecurity resilience. Future research could explore hybrid approaches and tailored methodologies to address specific industry needs and emerging threats. Embracing ethical hacking fosters a proactive security culture, ensuring robust protection of data and systems in an increasingly digital world.

## REFERENCES

[1] Zhang, C., & Li, Y. (2021). "Deep Learning-Based Anomaly Detection in IoT Networks: A Survey." IEEE Communications Surveys & Tutorials, 23(2), 1244-1270. https://doi.org/10.1109/COMST.2021.3052483

[2] Bontemps, K., McDermott, J., & Le-Khac, N. A. (2016). "Collective Anomaly Detection Based on Long Short-Term Memory Recurrent Neural Networks." In Future Technologies Conference (FTC) (pp. 271-283). Springer, Cham. https://doi.org/10.1007/978-3-319-48746-5_21

[3] Munawar, A., Mohiyuddin, S., & Aftab, S. (2018). "Machine Learning Techniques for Intrusion Detection System: A Comparative Analysis." In 2018 International Conference on Computing, Mathematics and Engineering Technologies (iCoMET) (pp. 1-6). IEEE. https://doi.org/10.1109/ICOMET.2018.8346404

[4] Goh, J., Adepu, S., Tan, M., & Lee, Z. S. (2017). "Anomaly Detection in Cyber-Physical Systems Using Recurrent Neural Networks." In 2017 IEEE 18th International Symposium on High Assurance Systems Engineering (HASE) (pp. 140-145). IEEE. https://doi.org/10.1109/HASE.2017.36

[5] Kim, J., Park, J., & Lim, H. (2019). "Anomaly Detection for Industrial Control Systems Using Sequence-to-Sequence Neural Networks." In 2019 IEEE International Conference on Big Data (Big Data) (pp. 4733-4739). IEEE. https://doi.org/10.1109/BigData47090.2019.9006331

[6] Campos, G. O., Zimek, A., Sander, J., Campello, R. J. G. B., Micenková, B., Schubert, E., & Houle, M. E. (2016). "On the Evaluation of Unsupervised Outlier Detection: Measures, Datasets, and an Empirical Study." Data Mining and Knowledge Discovery, 30(4), 891-927. https://doi.org/10.1007/s10618-015-0444-8

[7] Erfani, S. M., Rajasegarar, S., Karunasekera, S., & Leckie, C. (2016). "High-Dimensional and Large-Scale Anomaly Detection Using a Linear One-Class SVM with Deep Learning." Pattern Recognition, 58, 121-134. https://doi.org/10.1016/j.patcog.2016.03.028

[8] Mirsky, Y., Doitshman, T., Elovici, Y., & Shabtai, A. (2018). "Kitsune: An Ensemble of Autoencoders for Online Network Intrusion Detection." In NDSS. https://www.ndss-symposium.org/wp-content/uploads/2018/02/ndss2018_03A-1_Mirsky_paper.pdf

[9] García-Teodoro, P., Díaz-Verdejo, J., Maciá-Fernández, G., & Vázquez, E. (2009). "Anomaly- Based Network Intrusion Detection: Techniques, Systems and Challenges." Computers & Security, 28(1- 2), 18-28. https://doi.org/10.1016/j.cose.2008.08.003

[10] Ahmed, M., Mahmood, A. N., & Hu, J. (2016). "A Survey of Network Anomaly Detection Techniques." Journal of Network and Computer Applications, 60, 19-31. https://doi.org/10.1016/j.jnca.2015.11.016