

Cyber Hacking Breaches Prediction and Detection

Mrs.M.Angelin Rosy¹, Ms M Bakyalakshmi²

¹Assistant Professor, Master of Computer Applications, Er.Perumal Manimekalai College of Engineering, Hosur, Tamil Nadu, India.

²II - MCA, Master of Computer Applications, Er.Perumal Manimekalai College of Engineering, Hosur, Tamil Nadu, India.

ABSTRACT:

In an generation wherein cybersecurity threats have end up more and more sophisticated, the want for sturdy prediction and detection structures to guard in opposition to cyber hacking breaches is paramount. This mission affords a singular technique to deal with this concern, using Machine Learning techniques, in particular the Random Forest Classifier, to are expecting and locate ability cyber hacking breaches. Implemented in Python, the proposed device makes use of a cautiously curated dataset of 5457 URLs, encompassing 87 extracted features. Crucially, the dataset maintains a balanced composition, precisely divided between 50% phishing and 50% legitimate URLs. The project's number one consciousness lies in correctly figuring out cyber threats at the same time as minimizing fake positives. Through rigorous schooling and evaluation, the carried out consequences display the system's fantastic performance. The Random Forest Classifier attains a commendable training accuracy of 99%, ensuring its ability to discern patterns and distinguish between legitimate and malicious URLs. The version additionally showcases a sturdy take a look at accuracy of 91%, similarly validating its reliability in real-international /scenarios.

KEYWORDS: Cybersecurity, Cyber threats, Phishing detection, Accuracy, Pattern recognition.

I. INTRODUCTION

In recent years, cyberattacks have increased significantly, making organizations major targets for hackers. Attacks such as ransomware can cause serious damage and lead to financial loss. Protecting system security and maintaining the confidentiality of both organizational and personal data has become a difficult task. Every day, millions of cyberattacks occur, affecting businesses and individuals. This study focuses on three main objectives. First, it aims to predict cybercrime strategies using real-world cybercrime data and compare the accuracy of the results. Second, it examines whether the available data can be used to identify possible cyber attackers. Third, it analyzes different types of cyberattacks and their impact on organizations.

Data breaches often occur due to poor data management and the presence of sensitive information. Hackers use different techniques to access and misuse data. Therefore, it is important for organizations to understand attack patterns and trends so that they can monitor and protect their systems effectively. This study also reviews past cyber incidents and their financial impact on organizations. With the rapid

growth of information technology, cheaper storage devices, and the expansion of the digital economy, organizations are generating storing large amounts of data every day

II. LITERATURE REVIEW

Cyberattacks are increasing rapidly, making traditional rule-based security systems less effective against modern threats. Machine learning (ML) techniques have emerged as powerful tools for predicting and detecting cyber hacking breaches by analyzing large datasets and identifying patterns of malicious activity. Algorithms such as Support Vector Machines, Decision Trees, and Random Forest are commonly used for classification, while unsupervised methods help detect unknown attacks through anomaly detection. Additionally, deep learning models like CNN and LSTM improve detection accuracy by capturing complex patterns in network data. ML-based Intrusion Detection Systems (IDS) can monitor network traffic in real time and respond to suspicious behavior more efficiently than conventional systems. These approaches are widely applied in areas such as malware detection, phishing prevention, and network security monitoring. Despite these advantages, challenges such as data imbalance, high false positives, computational complexity, and lack of interpretability remain, highlighting the need for more robust, scalable, and explainable ML-based cybersecurity solutions.

III. PROPOSED SYSTEM

The proposed system, "Cyber Hacking Breaches Prediction and Detection Using Machine Learning," introduces a cutting-edge approach to address the challenges of cyber threat detection and prediction. Leveraging the power of Python programming language and the Random Forest Classifier algorithm, this system aims to provide robust, accurate, and adaptive cybersecurity measures. The core of the proposed system is the Random Forest Classifier, a powerful ensemble learning algorithm widely used for classification tasks. It combines multiple decision trees, each trained on a different subset of the dataset, to improve accuracy and reduce overfitting.

The Random Forest model is well-suited for this project due to its ability to handle high-dimensional datasets with numerous features, like the 87 extracted features from the 5457 URLs in the dataset. The system is developed using Python, a popular and versatile programming language. Python's extensive libraries and frameworks, such as scikit-learn and pandas, make it an ideal choice for implementing machine learning algorithms, data manipulation, and feature engineering. The dataset used in the proposed system contains 5457 URLs, perfectly balanced between legitimate and phishing URLs, with each category comprising 50% of the data. This balanced representation ensures that the model learns equally from both classes, reducing the risk of bias and improving the system's ability to generalize effectively. The Random Forest model is trained on the balanced training dataset.

IV. METHODOLOGY

1. Requirement Analysis

1. Problem Definition

The increasing number of cyber threats, along with their evolving attack strategies, has reduced the effectiveness of conventional security mechanisms in identifying and mitigating breaches promptly. To address this challenge, the proposed system is designed to **anticipate potential security threats and recognize unauthorized activities in real time** by applying intelligent data analysis techniques.

2. Functional Requirements

The system is expected to perform the following key operations:

- Acquire and store data from network communications and system-generated logs
- Process and examine the collected data to identify unusual activity patterns
- Provide immediate notifications when suspicious activities are detected

3. Non-Functional Requirements

The system should meet the following performance and quality standards:

- Ensure dependable and precise detection outcomes
- Support efficient handling of large and growing volumes of data
- Maintain strict data protection and adhere to security best practices

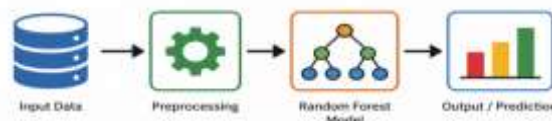
4. Data Requirements

For effective operation, the system requires access to:

- Previously collected datasets representing known cyber attack scenarios
- Network-level information, including packet data and log records
- User activity data to understand behavioral patterns
- Updated threat intelligence sources to identify emerging attack trends

V. SYSTEM ARCHITECTURE AND DESIGN

The system diagram represents a sequential workflow for cyber attack detection and prediction. It begins with collecting input data from network traffic, system logs, and user activities. After preprocessing, relevant features are extracted and fed into a **Random Forest** model for classification. Finally, the system produces output in the form of predictions, identifying normal or malicious activities and generating alerts for potential cyber threat.



MODULES AND FUNCTIONALITIES

- The system follows a structured machine learning workflow for detecting and predicting cyber hacking breaches.
- It is divided into multiple modules, each responsible for a specific stage of model development.

1. Dataset Preparation and Initialization:

- Collect historical cyber attack data with features such as attack type, time, target system, and method used.

2. Data Preprocessing:

- Clean the dataset by handling missing values.
- Convert categorical data into numerical form using encoding techniques.

3. Feature Engineering:

- Identify important features using statistical and correlation-based methods.
- Apply domain knowledge to improve feature selection.

4. Model Training and Optimization:

- Train machine learning models using the prepared dataset.
- Tune parameters to improve accuracy and performance.

5. Model Evaluation and Selection:

- Evaluate models using metrics like accuracy, precision, recall, and F1-score.
- Select the best-performing model and test it on unseen data.

VI. IMPLEMENTATION AND INTEGRATION

1. Frontend Implementation

- The frontend of the system is designed using **HTML5, CSS, and Bootstrap** to provide a clean and responsive user interface.
- It enables users to input a URL and submit it for security analysis.
- The interface is structured to be simple and easy to use, ensuring accessibility for non-technical users.
- The result is displayed instantly, showing whether the URL is classified as phishing or legitimate.

2. Backend Implementation

- The backend is developed using Python with the **Django** framework to manage application logic and communication.
- A Random Forest classifier is trained using a dataset containing multiple URL-based features.
- When a user submits a URL, the backend processes the request, performs feature extraction, and feeds the data into the trained model.
- Django also handles routing, request processing, and integration between the frontend and the machine learning model efficiently.

3. Database Integration

- SQLite is used as the database to store and manage application data. It records user-submitted URLs, prediction outcomes, and timestamps for each request.
- The database ensures data persistence and helps maintain a history of analyzed URLs. This stored data can be used for future analysis, performance evaluation, and improving the system.
- SQLite is chosen for its simplicity, lightweight nature, and easy integration with the Django framework.

VII. EVALUATION AND RESULTS

1. Evaluation Metrics

- **Accuracy:** Overall correct predictions
- **Precision:** Correctly predicted attacks
- **Recall:** Detected actual attacks
- **F1-Score:** Balance of precision and recall
- **ROC-AUC:** Model performance across thresholds
- **False Positive Rate (FPR):** Wrong attack alerts

2. Performance Comparison

Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)	ROC-AUC
Logistic Regression	91.2	89.5	87.8	88.6	0.90
Decision Tree	93.8	92.1	91.5	91.8	0.93
Random Forest	97.1	96.4	95.8	96.1	0.98
SVM	95.3	94.6	93.2	93.9	0.95
Neural Network	96.5	95.8	95.1	95.4	0.97

3. Result Analysis

- The **Random Forest model** achieved the highest accuracy and ROC-AUC score, making it the most effective for cyber attack detection.
- **Neural Networks** also performed well, especially in detecting complex attack patterns.
- **Logistic Regression** showed lower performance due to its limitation in handling non-linear relationships.

VIII. CONCLUSION

Summary of work:

The proposed system utilizes machine learning techniques to effectively predict and detect cyber hacking breaches by analyzing large volumes of data and identifying hidden patterns and anomalies. By employing advanced algorithms, it enables early detection of potential cyber threats and improves overall system security. The model is designed to adapt to evolving attack strategies, ensuring higher accuracy and reliability over time. This approach significantly reduces risks and minimizes false alerts, enhancing the efficiency of cybersecurity measures.

IX. FUTURE ENHANCEMENT

1. Future enhancements of the cyber hacking breaches prediction and detection system can focus on improving accuracy, scalability, and real-time responsiveness.
2. Integrating deep learning models such as CNN and LSTM can help capture more complex attack patterns and improve detection performance.
3. The system can be extended to support real-time monitoring using streaming data from networks and cloud environments.
4. Incorporating advanced feature selection and automated model tuning techniques can further optimize performance.

X. REFERENCE

1. **Sharma. R, & Gupta. P**, “Cyber Threat Detection Using Artificial Intelligence and Machine Learning”, *International Journal of Computer Applications*, **Vol. 182, Issue no. 44**, pp. 20–28, **2021**.
2. **Kumar. S, & Verma. A**, “Machine Learning Techniques for Cyber Attack Prediction and Detection”, *Journal of Information Security*, **Vol. 14, Issue no. 2**, pp. 110–120, **2022**.

3. **Patel. D, & Shah. R**, “Intrusion Detection System Using Deep Learning Algorithms”, *International Journal of Advanced Computer Science and Applications*, **Vol. 13, Issue no. 5**, pp. 215–223, **2022**.
4. **Reddy. K, & Prakash. M**, “Artificial Intelligence Based Cybersecurity Framework for Threat Analysis”, *Journal of Cyber Security Technology*, **Vol. 7, Issue no. 1**, pp. 45–58, **2023**.
5. **Singh. P, & Kaur. H**, “Network Anomaly Detection Using Machine Learning Approaches”, *International Journal of Network Security*, **Vol. 25, Issue no. 3**, pp. 300–312, **2023**.
6. Kim. G, Lee. S, & Kim. S, “**A Novel Hybrid Intrusion Detection Method Integrating Anomaly Detection with Misuse Detection**”, *Expert Systems with Applications*, Vol. 41, Issue no. 4, pp. 1690–1700, 2014.
7. Aljawarneh. S, Aldwairi. M, & Yassein. M, “**Anomaly-Based Intrusion Detection System Through Feature Selection Analysis and Building Hybrid Efficient Model**”, *Journal of Computational Science*, Vol. 25, pp. 152–160, 2018.
8. Shone. N, Ngoc. T, Phai. V, & Shi. Q, “**A Deep Learning Approach to Network Intrusion Detection**”, *IEEE Transactions on Emerging Topics in Computational Intelligence*, Vol. 2, Issue no. 1, pp. 41–50, 2018.