

Novartis Server Hack Dataset: Model Training & Framework

Nawfal Ahmed Khan, *IBA-KHI* Ahmed Tariq, *IBA-KHI*
Zeehan Rashid, *IBA-KHI*

Abstract

'Nova Threatix Model' employs advanced machine learning models on the Novartis dataset for intelligent intrusion detection. Utilizing boosting models and ensemble methods, the project aims to discern and predict server hacking incidents with optimal accuracy. The initiative integrates a continuous monitoring framework and personalized Threat Intelligence to fortify network security. Meticulous data preprocessing, feature extraction, and model training contribute innovative insights, elevating the project's efficacy in safeguarding network integrity.

1. Work Break-Down: Project Scope

The core aim of this project is to craft and deploy advanced machine learning models utilizing the logged values from the Novartis dataset's Intrusion Detection System (IDS) and Intrusion Prevention System (IPS). The primary focus is on empowering the identification and prediction of server hacking incidents by unraveling intricate patterns and correlations within the network traffic data.

In addition to model development, the project endeavors to introduce a resilient framework for continuous monitoring. This framework is designed to facilitate real-time evaluation of new data, enabling the swift detection of potential security threats. Moreover, the project aims to elevate its proactive cybersecurity measures through the incorporation of threat intelligence.

An innovative aspect of this project is the potential to implement personalized Threat Intelligence (TI). Rather than relying solely on international breach data, the project scope explores the development of a local Threat Intelligence system within the organization. This localized TI system is envisioned to be more tailored and intelligent, drawing insights from the specific threats encountered by the organization. This strategic approach aims to enhance the effectiveness of cybersecurity measures by aligning threat intelligence with the organization's unique risk landscape.

1.1. Task For Project

Data Exploration, Understanding & Preprocessing:

Identify key features from IDS/IPS logs indicating malicious activities.

Clean and preprocess dataset, addressing missing values, outliers, and inconsistencies.

Engage in feature engineering: Extract relevant IDS/IPS log features for model input.

Model Selection and Training:

Evaluate and select suitable machine learning models for network security (e.g., decision trees, random forests).

Train multiple models on labeled data to learn normal and malicious behavior patterns.

Performance Evaluation:

Assess model performance using metrics (accuracy, precision, recall, F1-score).

Implement cross-validation for model robustness.

Hyperparameter Tuning:

Fine-tune model hyperparameters for optimal performance.

Use techniques like grid search for hyperparameter optimization.

Model Comparison:

Compare machine learning models' performance for effective network security.

Consider interpretability, efficiency, and false positive/negative rates.

Security Implications:

Analyze security implications in real-world network environments.

Address challenges like adversarial attacks and model interpretability.

Documentation and Reporting:

Documentation of methodology, including data preprocessing, models, & results.

Summarizing key findings & recommending machine learning model implementation in network security.

1.2. Project Tasks Assignment

Nawfal Ahmed Khan [Team - Lead]

- Data Exploration, Understanding & Preprocessing.
- Documentation and Reporting.

Ahmed tariq [Team-Member 01]

- Model Selection and Training.
- Hyperparameter Tuning.

Zeehan Rashid [Team-Member 02]

- Performance Evaluation.
- Model Comparison.

2. Inclusion: Dataset Description

The foundation of this project is rooted in the comprehensive Novartis dataset, which presents a unique opportunity to glean valuable insights from real-world network traffic data. With an expansive scope, the dataset encapsulates approximately 18 columns and 40,000 rows, allowing for a detailed exploration. Covering a diverse array of network events, the dataset provides a holistic perspective on various interactions within the network. Particularly relevant to this project are the specific logged values originating from the Intrusion Detection System (IDS) and Intrusion Prevention System (IPS). These values serve as focal points, capturing key features crucial for discerning and predicting server hacking incidents. The depth and breadth of the Novartis dataset serve as a rich source for the development and implementation of robust machine learning models aimed at bolstering network security.

Logged Values in Dataset Column:

`Attack_Indicator`, `Source_Port_Risk`, `Destination_Port_Risk`, `Protocol_Risk`, `Severity_Score`, `Signature_Frequency`, `Signature_Variety`, `Traffic_Flow_Score`, `Anomaly_Index`, `Confidence_Score`, `IP_Risk_Score`, `IP_Risk_Level`, `Network_Layer_Score`, `Event_Risk_Score`, `Response_Effectiveness`

2.1. Machine Learning Models

The project employs a sophisticated approach to machine learning models, focusing on classification techniques to distinguish between normal and malicious network behavior. The algorithmic selection is meticulous, with a strategic emphasis on decision trees, random forests, and support vector machines chosen for their particular suitability in the context of network security applications. Furthermore, the project incorporates boosting models, ensemble methods,

and stacking techniques to create a robust model that ensures the highest accuracy in predicting network security incidents. This comprehensive training approach involves deploying multiple models on the labeled dataset, facilitating the identification of intricate patterns associated with both normal and malicious network activities. Through the strategic integration of advanced algorithms and ensemble methods, the project aims to achieve an optimized and highly accurate predictive model for enhanced network security.

3. Exclusion

Out-of-scope elements for this project encompass specific types of network security threats and attacks that fall outside the purview of the analysis. Notably, the project does not delve into the detailed examination of certain attack vectors or data features that may not directly contribute to the identification and prediction of server hacking incidents. Additionally, the analysis does not extend to cover network security elements beyond the predefined scope, ensuring a focused and effective exploration of relevant factors.

4. Functional Requirements

For data preprocessing, the initial step involves a meticulous examination and cleaning of the Novartis dataset. This encompasses addressing missing values, handling outliers, and ensuring data consistency. Feature extraction involves the careful selection and extraction of pertinent features from the Intrusion Detection System (IDS) and Intrusion Prevention System (IPS) logs. These selected features serve as the input for the subsequent machine learning models. In the model training and evaluation phase, the project adopts a systematic approach. Multiple machine learning models, including boosting models, are trained on the labeled dataset, enabling them to discern patterns associated with normal and malicious network activities. Performance evaluation utilizes appropriate metrics such as accuracy, precision, recall, and F1-score, complemented by cross-validation techniques to ensure the robustness of the models. This comprehensive process ensures the development of accurate and reliable models for the identification and prediction of server hacking incidents.

5. Hardware Requirements

The hardware requirements for this project involve a computing environment capable of handling the computational demands associated with model training. Adequate processing power, memory, and storage are essential to ensure

efficient execution. High-performance CPUs or GPUs are recommended for expedited model training processes.

6. Software Requirements

The project will utilize a combination of open-source software tools and libraries. Key components include Python as the primary programming language for implementation, scikit-learn for machine learning algorithms, and TensorFlow for deep learning applications. Additional libraries such as NumPy, Pandas, and Matplotlib will be employed for data manipulation and visualization. The use of Google Colab, Jupyter Notebooks or similar platforms will facilitate a collaborative and interactive development environment.

7. Vulnerability of Data Analysis Method: Threat Modeling

Identifying potential threats to the data analysis methodology is crucial. Threats may arise from biases present in the dataset, overfitting of models to specific patterns, or ethical considerations in handling sensitive information. Moreover, there may be risks associated with the interpretability and generalization of the models in real-world scenarios.

7.1. Risk Mitigation

To mitigate potential threats, a proactive approach is adopted. This involves using a diverse dataset to minimize biases, implementing robust cross-validation techniques to ensure model generalization, and adhering to ethical guidelines in data handling and analysis. Regular model evaluations and performance monitoring are integrated to identify and rectify overfitting issues. Furthermore, transparency in the analysis process and documentation aids in addressing ethical considerations and ensures the reliability and integrity of the project's outcomes.