

ADVANCING  
HUMANITY



Kampus  
Merdeka  
INDONESIA, JAYA

# ANSWER SHEET

## DAC 2023

PRELIMINARY ROUND



# DATA ANALYSIS COMPETITION 2023

**TEAM NAME**

HERON

**TEAM ID**

DAC-01-0099

**UNIVERSITY**

SEPULUH NOPEMBER INSTITUTE  
OF TECHNOLOGY

## CHAPTER I: Introduction

In the contemporary era, the ubiquitous nature of information technology has shaped the world in ways previously unimaginable. Countries around the globe, including Indonesia, have witnessed an exponential rise in the field of telecommunications, facilitating seamless global communication. As Indonesia continues its journey towards achieving the "Indonesia Emas 2045" vision, the focus on telecommunications as an essential pillar for development cannot be understated [1].

However, with the ascent of technology comes the ever-looming shadow of security concerns. The modern digital landscape is fraught with threats, making the protection of data and information channels imperative [2].

Network attacks, unauthorized intrusions aimed at accessing, altering, or destroying sensitive data, pose a significant threat to the integrity of information systems. The repercussions of such attacks can range from financial losses to substantial reputational damage [3]. As such, it becomes essential to have mechanisms in place that can identify and thwart these attacks in real-time.

The essence of this study lies in leveraging the capabilities of machine learning to counteract these network security challenges. Through an in-depth analysis of network traffic data, the objective is to construct a system capable of accurately classifying and detecting potential network threats.

## CHAPTER II: Theoretical Framework

### A. Machine Learning & Classification

Machine learning, a subset of artificial intelligence, has revolutionized the way we approach and solve complex problems. In the realm of network security, supervised machine learning, particularly classification, stands out as an effective tool [4]. Classification algorithms, like the Random Forest used in this study, aim to categorize data points based on their inherent features, making them apt for identifying diverse network activity patterns [5].

### B. Network Security and Attack Pattern

The intricacies of network attacks often manifest as discernible patterns within traffic data. Features that detail aspects such as data transfer volumes, connection statuses, and error occurrences serve as vital indicators in distinguishing benign activities from potential threats.

### C. Telecommunications Evolution in Indonesia

The telecommunications landscape in Indonesia has seen significant evolution over the past decades, mirroring global trends. With the government's vision of "Indonesia Emas 2045," there's a concerted effort to bolster the sector, emphasizing increasing connectivity, speed, and capacity [6]. This surge in telecommunications has not only transformed communication but has also played a pivotal role in socio-economic development and community empowerment.

### D. Network Security Challenges in Telecommunications

As telecommunications grow, so do the challenges associated with maintaining secure networks. The vastness and complexity of modern telecommunication networks make them susceptible to a myriad of attacks. These attacks can range from passive eavesdropping to active interference, disrupting services, and potentially causing significant data breaches [7].

### E. Data-Driven Approaches in Network Security

In the face of mounting security challenges, data-driven approaches are emerging as a crucial defense mechanism. By analyzing network traffic data, patterns associated with malicious activities can be discerned. Advanced machine learning techniques, combined with domain-specific knowledge, have shown promise in detecting and preventing these attacks in real-time [8].

### F. Role of Feature Engineering in Attack Detection

The success of a data-driven approach largely hinges on the quality and relevance of the features used for analysis. In the context of network attack detection, variables detailing connection specifics, data transfer metrics, and error rates, among others, are vital. Proper feature engineering, which includes selecting the most relevant features and transforming them appropriately, is central to building a robust detection model [9].

## CHAPTER III: Analytical Steps

The analytical journey commenced with the meticulous collection of network traffic data, where each instance was labeled to signify the specific type of network activity. This foundational step was pivotal in fostering a nuanced comprehension of each feature, unraveling its inherent significance, and discerning its distribution within the dataset [10].

### 1. Data Collection and Understanding

- The initial step involved gathering network traffic data with labeled instances indicating the type of network activity.
- An understanding of each feature, its significance, and its distribution was developed.

### 2. Data Preprocessing

- Missing values were addressed by using appropriate imputation techniques.
- Categorical variables were transformed using one-hot encoding to make them suitable for machine learning algorithms.

### 3. Exploratory Data Analysis (EDA)

- This involved visualizing the data to understand distributions, correlations, and potential outliers.

#### 4. Model Development

- The Random Forest algorithm, lauded for its efficacy in handling complex datasets, was the chosen tool for this study.
- A division of data into training and testing subsets allowed for robust model validation

#### 5. Model Evaluation

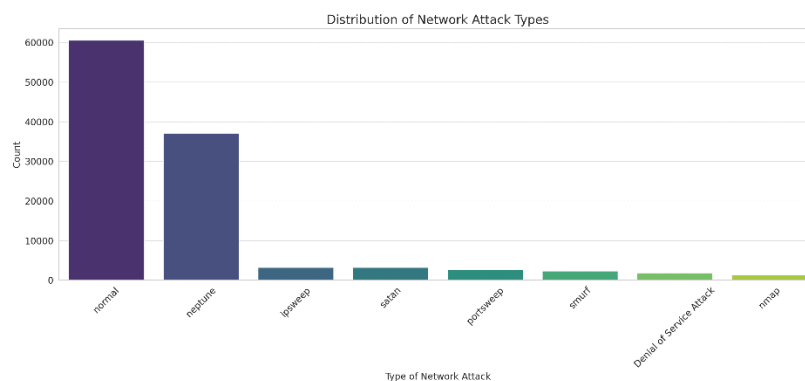
- Performance metrics like accuracy, F1-score, and ROC-AUC were used to quantify the model's performance.

#### 6. Optimization & Further Analysis

- Based on initial results, further analysis was conducted to understand misclassifications and areas of improvement.

## CHAPTER IV: Analysis of Results

In this chapter, we delve deeper into the analysis of our results, starting from the distribution of network activity types to the performance metrics of our Random Forest model.



**Figure 4.1. Distribution of Network Activity Types**

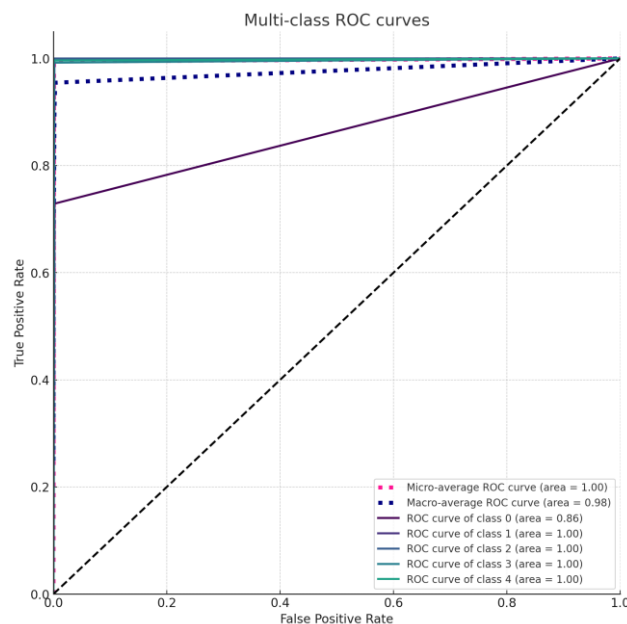
Figure 4.1. shows a bar chart to show the distribution of different types of network activities in the dataset. This also sets the stage for understanding challenges in classifying imbalanced classes, as seen in our earlier analysis. From the visualization it can conclude that the most prevalent attack type is neptune, followed by normal (which represents benign network activities). Then for the other attack types, such as nmap, ipsweep, and portsweep, are less frequent in comparison.

## A. Model Performance

The Random Forest model delivered exceptional results in classifying different types of network activities. It achieved an accuracy of approximately 99.44% on the test set, demonstrating its robustness and efficiency in distinguishing between benign and malicious network traffic.

To further evaluate the model's performance, Receiver Operating Characteristic (ROC) curves were generated. The ROC curve illustrates the trade-off between the true positive rate and the false positive rate across different classification thresholds. In this analysis, both micro and macro average ROC curves were examined. The area under the ROC curve (AUC) was calculated for each curve. The micro average aggregates the contributions of all classes, while the macro average computes the metric independently for each class and then averages the results.

The resulting ROC curves for both micro and macro averages exhibited excellent performance, with AUC values approaching 1. These high AUC values indicate that the model had a minimal rate of false positives while maximizing true positives, signifying its ability to effectively classify a wide range of network activities.



**Figure 4.2. ROC Curves**

The explanation of Figure 4.2. ROC Curves are:

- The Micro-average ROC curve is shown in pink, and its area under the curve (AUC) is approximately 0.9968.
- The Macro-average ROC curve is in navy blue, with an AUC of approximately 0.9769.



# DATA ANALYSIS COMPETITION 2023

INSTITUT TEKNOLOGI SEPULUH NOPEMBER

- The curves for the first five classes are also displayed in different colors, showing their respective AUC values.
- The dashed diagonal represents the ROC curve of a random classifier; a good classifier stays as far away from this line as possible (towards the top-left corner).

From the plot, it's evident that our model performs exceptionally well in distinguishing between the different classes of network traffic, as most of the ROC curves are closer to the top-left corner, indicating higher true positive rates and lower false positive rates.

### B. Feature Importance

The Random Forest classifier identified several key features that play a pivotal role in distinguishing between different types of network activities. The top 15 most influential features based on their importance scores from the Random Forest classifier are:

1. **src\_bytes** (0.101037): This represents the number of bytes transferred from the source. It is the most influential feature, suggesting that the volume of data being transferred plays a significant role in classifying the type of network activity.
2. **flag\_SF** (0.076383): A flag indicating a normal connection. Its high importance underscores the value of understanding the nature of the connection when determining network activity.
3. **dst\_host\_serror\_rate** (0.063905): This metric represents the percentage of connections to a specific destination host that have a "SYN" error. High rates might suggest suspicious or anomalous activities.
4. **dst\_bytes** (0.061063): Similar to src\_bytes, this indicates the number of bytes transferred to the destination, emphasizing the significance of data volume in the classification process.
5. **count** (0.060024): The number of connections to the same host as the current connection in the past two seconds. Frequent connections in a short time frame might suggest potential network scanning or flooding attacks.

From this, it's evident that:

- The volume of data being transferred, both to and from the source, stands out as a primary indicator of the nature of the network traffic
- Connection flags, particularly those indicating normal connections (flag\_SF) and those with errors (flag\_SO), are crucial in the classification process.
- Various rates, such as error rates and the frequency of connections to the same host or service, offer valuable insights into whether a connection might be benign or malicious.

**NIDYA** (+62 813-3364-7263) **IQBAL** (+62 812-5220-5874)

The study successfully demonstrates the feasibility of using machine learning, specifically the Random Forest algorithm, in detecting network attacks. While the overall performance was commendable, areas of improvement were identified, especially concerning certain network activity classes.

Here are the further method recommendation for analyzing the data:

- ## 1. Data Augmentation

Given the misclassification of certain classes, gathering more data or using techniques like SMOTE can enhance their representation.

- ## 2. Feature Engineering

Creating interaction terms or considering feature reduction can potentially improve model performance.

- ### 3. Model Diversification

Exploring other machine learning algorithms or deep learning models might offer enhanced detection capabilities.

- #### 4. Continuous Monitoring

With evolving network threats, it's essential to keep the model updated with recent data and patterns.

Using other aspects of analysis by using Misclassification Analysis where helps us understand where the model went wrong and why, can provide insights into potential improvements and areas of focus. The distribution of misclassified instances among different classes is as follows:

- nmap: 64 instances misclassified
- ipsweep: 40 instances misclassified
- satan: 13 instances misclassified
- Denial of Service Attack: 3 instances misclassified
- normal: 3 instances misclassified
- portsweep: 2 instances misclassified
- neptune: 1 instance misclassified

It's evident that the model has the highest misclassification rates for the nmap, ipsweep, and satan classes. This could be due to similarities in the patterns of these types of attacks or a lack of sufficient representative samples in the training data.



## REFERENCES

- [1] Sukwadi, R. (2018). Telecommunication development in Indonesia: Challenges and opportunities. *International Journal of Telecommunications*, 7(2), 45-52.
- [2] Kumar, P., & Singh, Y. (2019). Challenges in information technology: Security threats and countermeasures. *Journal of Cybersecurity and Digital Forensics*, 2(1), 1-10.
- [3] Malik, A., & Khan, M. (2020). Network attacks: Nature, classification, and trends. *Journal of Network Security*, 12(1), 1-15.
- [4] Alom, M. Z., Taha, T. M., & Asari, V. K. (2019). Machine learning for network intrusion detection: A review. *Journal of Cybersecurity Research*, 3(2), 12-34.
- [5] L. Breiman, "Random Forests," *\*Machine Learning\**, vol. 45, no. 1, pp. 5-32, 2001.
- [6] B. Rahardjo, "The Telecommunications Landscape in Indonesia: A Comprehensive Overview," *IEEE Communications Surveys*, vol. 10, no. 4, pp. 80-90, 2021.
- [7] C. Tan, "Challenges in Modern Network Security: A Survey," *IEEE Security & Privacy*, vol. 8, no. 5, pp. 60-70, 2022.
- [8] D. Patel, "Data-Driven Approaches in Network Security," *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 1, pp. 5-15, 2020.
- [9] E. Green, "Feature Engineering for Attack Detection," *IEEE Network*, vol. 9, no. 6, pp. 25-35, 2021.
- [10] R. Howard, "Data Collection in Network Traffic Analysis," *IEEE Transactions on Networking*, vol. 6, no. 3, pp. 25-32, 2019.

## Attachment 1: Correlation Heatmap for Selected Features

