

Sahad Rafiuzzaman  
CYSE 635

Dr. MD Morshed Alam

October 3, 2024

### Supervised Learning based Attack Classifier Report

This assignment focused on developing a machine learning model to distinguish between attack and benign logs. The first step was to import the necessary libraries and upload the required CSV files. After that, several datasets, including benign, jamming, and spoofing logs, were combined into one for analysis. Preprocessing was done to convert features into numerical values using LabelEncoder, making the data ready for the machine learning models. Once the data was preprocessed, histograms and pie charts were used to visualize how the attack and benign logs were distributed, helping identify patterns for model comparison.

The exploratory data analysis used histograms to show the distribution of individual features, which helped identify patterns and any unusual data points. Pie charts were used to show the proportion of attack logs compared to benign logs. This helped visualize how the logs were distributed across the dataset, giving insights before training the models.

After preparing the data, two models were chosen: K-Nearest Neighbors and Logistic Regression. Both were used to classify data into two categories, which made them suitable for identifying attack logs and benign logs. A third model, the Support Vector Machine, was tested for comparison. The dataset was split into two parts: one part was used to train the models, and the other part was used to test them. This split helped evaluate how well the models would perform on new data.

K-Nearest Neighbors decides the log category by looking at the closest logs and choosing the category that appears the most. Logistic Regression predicts the category by creating a boundary that separates attack logs from benign logs. Each model was tested using

a confusion matrix, which provided insights into how well it classified the logs. The confusion matrix showed the number of true positives, true negatives, false positives, and false negatives. True positives represent correctly identified attack logs, while true negatives represent correctly identified benign logs. False positives refer to benign logs misclassified as attacks and false negatives refer to attack logs misclassified as benign.

The differences in how the models classified the logs were highlighted through the confusion matrices. Figure 1 shows that K-Nearest Neighbors performed well with true positives and true negatives. On the other hand, Figure 2 shows that Logistic Regression classified some of the true positives from the K-Nearest Neighbor model as false negatives. Meanwhile, Figure 3 shows that the Support Vector Machine classified some of the true negatives from the K-Nearest Neighbors model as false positives. The varying results are due to how each model approached classification and handled data patterns.

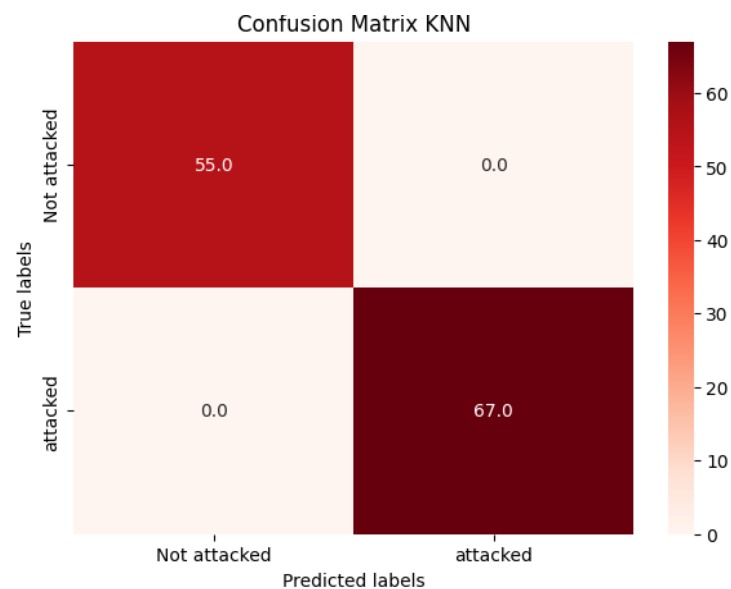
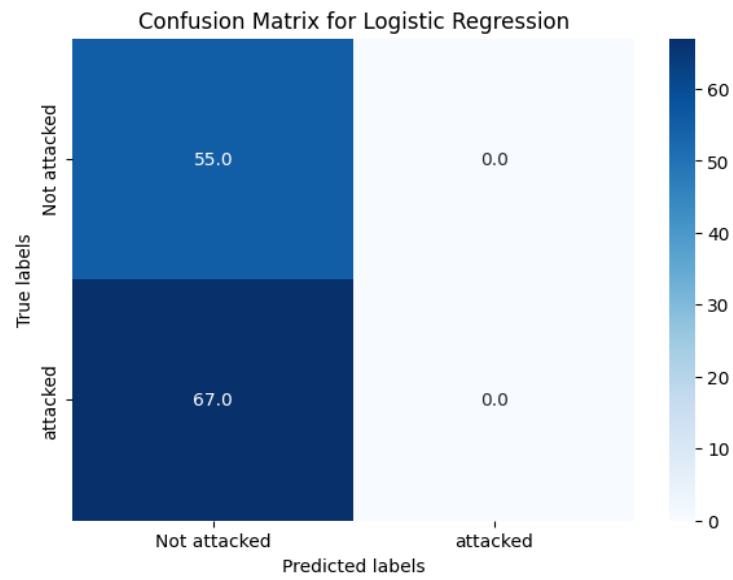
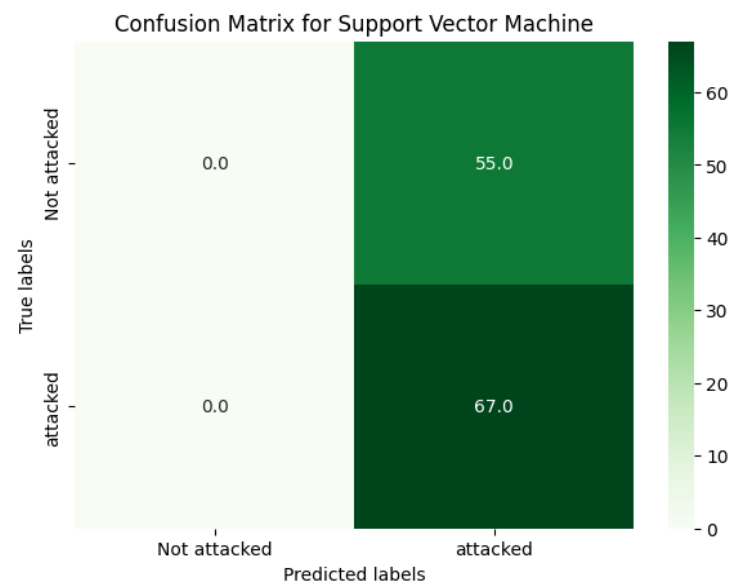


Figure 1: Confusion Matrix for K-Nearest Neighbors



**Figure 2: Confusion Matrix for Logistic Regression**



**Figure 3: Confusion Matrix for Support Vector Machine**