

# **Network Traffic Analysis and Prediction using Machine Learning**



**Course Title:** Cyber Security

**Course Code:** CSE4003

*Under the guidance of*

Dr. S. Vinila Jinny

*Submitted by:*

Dipsan Bhattarai (20BCE2778)

Aayam Piya (20BCE2774)

Vibhrat Vibhu Basnet (20BCE2765)

Karma Gurung (20BCE2933)

**Date:** November 22, 2023

## **ACKNOWLEDGEMENT**

We would like to express our heartfelt gratitude to all those who have contributed to the successful completion of this project on "Network Traffic Analysis and Prediction Techniques using Machine Learning." It has been an exciting journey, and we are thankful for the support, guidance, and assistance that we have received along the way.

First and foremost, we would like to thank our esteemed faculty advisor, Dr. S. Vinila Jinny, for her unwavering support and invaluable guidance throughout the project. Their expertise and insights have been instrumental in shaping our research and helping us navigate the complexities of network traffic analysis and prediction. We would also like to extend our appreciation to our friends and colleagues for their valuable feedback and support during this project. Their input and discussions were invaluable in refining our ideas and approaches.

## **TABLE OF CONTENTS**

| <b>S/N</b> | <b>CONTENTS</b>             |
|------------|-----------------------------|
| <b>1.</b>  | <b>ACKNOWLEDGEMENT</b>      |
| <b>2.</b>  | <b>ABSTRACT</b>             |
| <b>3.</b>  | <b>OBJECTIVES</b>           |
| <b>4.</b>  | <b>INTRODUCTION</b>         |
| <b>5.</b>  | <b>LITERATURE REVIEW</b>    |
| <b>6.</b>  | <b>MATERIALS AND MODELS</b> |
| <b>7.</b>  | <b>DATASET USED</b>         |
| <b>8.</b>  | <b>PROCEDURE</b>            |
| <b>9.</b>  | <b>RESULT</b>               |
| <b>10.</b> | <b>DISCUSSION</b>           |
| <b>11.</b> | <b>CONCLUSION</b>           |
| <b>12.</b> | <b>REFERENCES</b>           |

## **ABSTRACT**

Network traffic analysis and prediction techniques using machine learning have become increasingly important in recent years due to the proliferation of Internet of Things (IoT) devices and the growing complexity of network traffic patterns. Machine learning techniques offer a powerful approach to analyze and predict network traffic patterns, enabling network administrators to proactively identify potential issues and optimize network performance. In this context, this project aims to explore various machine learning techniques for network traffic analysis and prediction.

The first part of the project focuses on the analysis of network traffic patterns using machine learning techniques such as clustering and anomaly detection. Clustering algorithms can help identify groups of similar traffic patterns, while anomaly detection can help identify unusual traffic patterns that may indicate security threats or network performance issues. The second part of the project focuses on the prediction of network traffic patterns using machine learning techniques such as time series forecasting and regression analysis. Time series forecasting can help predict future traffic patterns based on historical data, while regression analysis can help identify the factors that influence network traffic patterns.

In conclusion, this project aims to provide a comprehensive overview of network traffic analysis and prediction techniques using machine learning. The project will explore various machine learning algorithms and their applications in network traffic analysis and prediction, providing insights into the current state of the art in this field. The outcomes of this project can be used to optimize network performance and enhance network security.

## OBJECTIVES

In recent years, machine learning has become an increasingly popular tool for network traffic analysis. There are now a wide variety of machine learning algorithms that can be used for this purpose, ranging from simple rule-based models to complex neural networks. The effectiveness of these algorithms depends on a number of factors, including the size and complexity of the network being analyzed, the types of data being collected, and the specific objectives of the analysis.

- Network security: Machine learning algorithms can be used to detect and prevent security threats in network traffic by identifying patterns of suspicious behavior. This can include identifying malware, DDoS attacks, phishing attempts, and other types of malicious activity.
- Network optimization: Machine learning can be used to analyze network traffic and identify areas where optimization is needed, such as congestion or bandwidth bottlenecks. This can help improve the overall performance of the network and ensure that resources are being used efficiently.
- Predictive maintenance: By analyzing network traffic, machine learning algorithms can identify potential hardware failures before they occur. This can help reduce downtime and ensure that the network remains operational.
- Resource allocation: Machine learning can be used to analyze network traffic and allocate resources based on usage patterns. This can help ensure that resources are being used efficiently and that performance is optimized.
- User behavior analysis: Machine learning can be used to analyze user behavior and identify patterns of usage. This can help identify areas where improvements can be made, such as improving the user experience or identifying new business opportunities.

# INTRODUCTION

Network traffic analysis is crucial for achieving effective information security since networks are used to transfer extremely sensitive and valuable information connected to e-commerce, banking, and business. Instead of taking a reactive strategy, network monitoring and traffic analysis and prediction resemble a proactive approach where security breaches are prevented from happening within the network. The analysis of network traffic is a crucial step in creating effective preemptive congestion management methods and identifying legitimate and malicious packets. These strategies aim to prevent network congestion by allocating network resources in accordance with the anticipated traffic.

In several domains, including dynamic bandwidth allocation, network security, network planning, and predictive congestion control, the predictability of network traffic has significant advantages. Long-term predictions and short-term predictions can be divided into two groups. Long-term traffic forecasting allows for more precise planning and wiser choices since it provides a complete forecast of traffic models to assess future capacity needs. Dynamic resource allocation and short period prediction (milliseconds to minutes) are related. It can be applied to enhance Quality of Service (QoS) mechanisms, reduce congestion, and manage resources in the best way possible. It can also be applied to packet routing. For network traffic analysis and prediction, a variety of methods are utilized, such as time series network traffic analysis and prediction models, contemporary data mining techniques, soft computing strategies, and neural networks. The strategies for network traffic analysis and prediction that have been suggested, used, and implemented are reviewed in this study. Discussions of the uniqueness and limitations of earlier studies are included, and a summary of the typical characteristics of network traffic analysis and prediction is also provided. The rest of the document is structured as follows. Following a brief introduction to network traffic analysis, section two offers a thorough overview of a number of different network analysis methods. The third section examines various methods for predicting network traffic. We present our conclusions in the final section.

In the modern era, network traffic analysis has become more and more crucial and significant for observing the network traffic. Administrators used to merely keep an eye on a few hundred computers or a limited number of network devices. The network bandwidth may have been 100 Mbps or little less (Megabits per second). At the moment, network administrators must cope with wireless networks, ATM (Asynchronous Transfer Mode) networks, and wired networks with speeds greater than 1Gbps (Gigabits per second). In order to monitor the network, swiftly resolve network issues to prevent network failure, and manage network security, they need extra current network traffic analysis tools. As a result, there are now a variety of difficulties with network traffic analysis. For security management, a network is examined at many levels, including the packet, flow, and network levels. For network traffic analysis, researchers employ a variety of methods. Preprocessing is followed by actual analysis and observation to uncover patterns from the network data in a general framework for network traffic analysis. In the next subsections, these three phases are described in depth.

## LITERATURE REVIEW

- i. **A Review of Network Traffic Analysis and Prediction Techniques:** Analysis and prediction of network traffic has applications in a wide comprehensive set of areas and has newly attracted significant number of studies. Different kinds of experiments are conducted and summarized to identify various problems in existing computer network applications. Network traffic analysis and prediction is a proactive approach to ensure secure, reliable, and qualitative network communication. Various techniques are proposed and experimented for analyzing network traffic including neural network-based techniques to data mining techniques. Similarly, various Linear and non-linear models are proposed for network traffic prediction. Several interesting combinations of network analysis and prediction techniques are implemented to attain efficient and effective results.
- ii. **A Survey on Encrypted Network Traffic Analysis Applications, Techniques, and Countermeasures:** Network traffic encryption is becoming more widely used. To secure communications and safeguard user privacy, many applications use encryption techniques. Additionally, a significant amount of malware uses encryption techniques to conceal its presence and activity as it spreads through network traffic. As we move closer to a time when all Internet communications are encrypted, we must immediately begin reviewing the state-of-the-art in the broad field of network traffic analysis and inspection to determine whether current traffic processing technologies will be able to quickly adapt to the full implementation of network encryption.
- iii. **NEMEA: A framework for network traffic analysis:** The sophistication of network assaults makes it challenging to identify them with conventional analysis tools. To identify such assaults, it is important to analyze Application Layer (L7) information. However,



there aren't many tools available today that can handle and manipulate L7. To remedy the condition, a flow-based modular Network Measurements Analysis (NEMEA) system was suggested. Since NEMEA is built on a stream-wise idea, data are continually analyzed in memory with less data storage. NEMEA was created as an open-source project and made accessible to the public on a global scale. It is intended for both practical use and experimentation.

- iv. **Decorrelating wireless sensor network traffic to inhibit traffic analysis attacks:** An attacker who examines packet data can determine the position of a base station from the typical packet traffic in a sensor network, which exhibits distinct patterns. Since a base station serves as the hub for data collection and failure, it is possible to destroy it once it has been found, rendering the entire sensor network useless. In order to protect base station locations from traffic analysis attacks, this study examines a variety of decorrelation techniques. A collection of fundamental countermeasures is presented, including hop-by-hop packet encryption to alter appearance, enforcement of a consistent packet sending pace, and elimination of the correlation between a packet's receipt time and its forwarding time.
- v. **ProfilIoT: a machine learning approach for IoT device identification based on network traffic analysis:** In this work we apply machine learning algorithms on network traffic data for accurate identification of IoT devices connected to a network. To train and evaluate the classifier, we collected and labeled network traffic data from nine distinct IoT devices, PCs, and smartphones. Using supervised learning, we trained a multi-stage meta classifier; in the first stage, the classifier can distinguish between traffic generated by IoT and non-IoT devices. In the second stage, each IoT device is associated with a specific IoT device class. The overall IoT classification accuracy of our model is 99.281+.

- vi. **Deep Learning for Network Traffic Monitoring and Analysis (NTMA): A Survey:** Internet of Things (IoT) and cellular networks are two examples of contemporary communication systems and networks that produce a significant and diverse amount of traffic data. The standard network management approaches for monitoring and data analytics encounter several difficulties and problems in such networks, such as the accuracy and efficient real-time processing of massive data. Furthermore, due to several variables, including device mobility and network heterogeneity, the pattern of network traffic, particularly in cellular networks, exhibits extremely complicated behavior. In large data systems, deep learning has been effectively used to speed up analytics and knowledge discovery by enabling the recognition of hidden and complicated patterns. Inspired by these achievements, networking researchers use deep learning models for NTMA (Network Traffic Monitoring and Analysis) applications, such as traffic classification and prediction.
- vii. **A Comparison of three machine learning techniques for encrypted network traffic analysis:** This study compares three techniques for decrypting communication without using payload, port, or IP information. Due to the plain text SSH protocol's beginning allowing us to gather data sets with a trustworthy ground truth, binary identification of SSH vs. non-SSH traffic is employed as a case study in this context. The approaches are put through several tests utilizing various feature sets, training and test traffic traces, export settings, and export options for a total of 128 different configurations. Test cases that employ test data from a different network than the one the model was trained on are particularly interesting since they demonstrate how robust the learned models are. When each method is tested using traffic traces that are recorded on the same network as the training network trace, the results reveal that the trained model based on the multi-objective genetic algorithm (MOGA) is able to obtain the greatest performance among the three approaches.

- viii. Machine Learning-Powered Encrypted Network Traffic Analysis: A Comprehensive Survey:** This paper presents a thorough overview of recent developments in machine learning-powered encrypted traffic analysis. The analysis goals that form the basis for the classification of the material are first summarised after a survey of the literature in this field. The process of using machine learning techniques to analyse encrypted communication is then abstracted, comprising the steps of traffic collection, traffic representation, traffic analysis method, and performance evaluation. The criteria for categorization granularity and information timeliness for the surveyed research may change significantly for various analytic aims. As a result, a thorough assessment of the literature is provided in terms of the purpose of traffic analysis, and it is divided into four categories: network asset identification, network characterisation, privacy leakage detection, and anomaly detection.
- ix. Network traffic analysis using machine learning: an unsupervised approach to understand and slice your network:** A recent surge in data generation and heterogeneity due to the proliferation of smart devices necessitates new network solutions for better traffic analysis and comprehension. In order to handle the enormous volume of data automatically, these solutions need to be clever and scalable. With the development of high-performance computing (HPC), it is now possible to deploy machine learning (ML) efficiently, and its effectiveness has been verified in a number of sectors. (e.g., healthcare or computer vision). In addition, network slicing (NS), which is crucial to addressing the variety of service requirements, has received considerable interest from both business and academics. Consequently, ML adoption within NS management is a fascinating topic.
- x. IoT Malware Network Traffic Classification using Visual Representation and Deep Learning:** Malware has emerged as a difficult issue with rising infection rates and degrees of sophistication

as more IoT devices and technologies are put into use. Without effective security measures, a large amount of sensitive data is vulnerable to attack, making it simple for cybercriminals to use it for a variety of unlawful actions. As a result, sophisticated network security measures that can analyse traffic in real time and mitigate harmful traffic are needed. We offer a revolutionary IoT malware traffic analysis method that uses deep learning and visual representation to identify new malware more quickly. This problem is addressed. Due to the application of deep learning technology, the suggested approach successfully detects fraudulent network traffic while drastically lowering the detection time.

## **MATERIALS AND MODELS**

Machine learning has emerged as a promising approach for network traffic analysis, due to its ability to learn patterns and anomalies in large and complex datasets. By training models on labeled data, machine learning algorithms can identify characteristic features and behaviors of different types of network traffic, and use this knowledge to predict the type of traffic in real-time.

### **Materials:**

For this project, we used Google Colab as our development environment, which allowed us to leverage the power of cloud computing resources to train and evaluate our machine learning models. We used Python as our programming language, along with various libraries and frameworks for data preprocessing, visualization, and modeling. Our dataset was the KDD-cup dataset, which is a widely used benchmark dataset for intrusion detection and network traffic analysis. It consists of network traffic data collected from a simulated network environment, with various types of attacks and normal traffic patterns. The dataset contains over 4 million records and 41 features, including protocol types, service types, flag values, and packet size distributions.

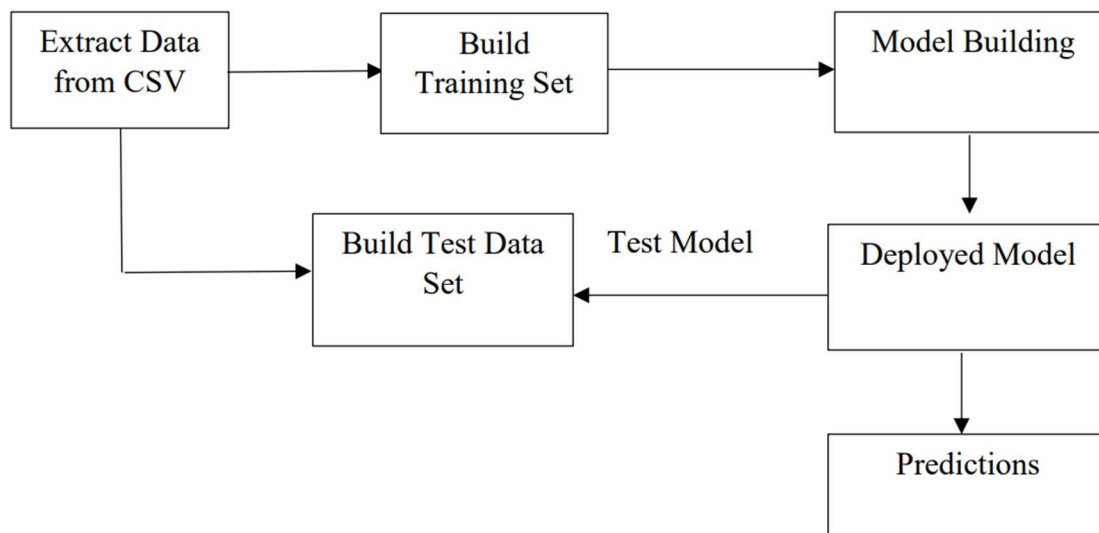
### **Models:**

Our goal was to predict the type of network traffic based on the available features, using machine learning algorithms. We trained and evaluated several popular models for classification, including:

- Linear Regression
- Decision Tree Classifier
- Random Forest
- Gradient Boost
- XGBoost
- KNN Classifier

We used the scikit-learn library to implement these models, which provided a convenient and efficient way to preprocess the data, train the models, and evaluate their performance. To compare the accuracy of these models, we split the dataset into training and testing sets using a 70:30 ratio, and trained each model on the training set. We then used the trained models to predict the labels for the testing set, and calculated the accuracy and confusion matrixs using the scikit-learn metrics module. We also used the matplotlib.pyplot library to visualize the performance of these models using a bar graph, which showed the accuracy for each model side by side. This allowed us to compare the performance of the models and identify the best performing one.

Furthermore, we plotted a confusion matrix for each of these models, which provided a more detailed view of the performance of the models. The confusion matrix shows the true positive, true negative, false positive, and false negative rates for each class, and allows us to analyze the performance of the models for different types of traffic.



***Fig: Architecture of Methodology***

## **DATASET USED:**

For our project on network traffic analysis and prediction using machine learning, we used the KDD Cup 1999 dataset, which is a widely used benchmark dataset in the field of network intrusion detection. The dataset is available at [“https://archive.ics.uci.edu/ml/machine-learning-databases/kddcup99-mld/kddcup.data\\_10\\_percent.gz”](https://archive.ics.uci.edu/ml/machine-learning-databases/kddcup99-mld/kddcup.data_10_percent.gz). The dataset contains network traffic data collected from a simulated environment that mimics a local area network (LAN) of a typical US Air Force base. The dataset was created for the purpose of developing and testing intrusion detection systems. It consists of approximately 4.9 million network connections, which are divided into five categories: normal, denial-of-service (DoS), user-to-root (U2R), remote-to-local (R2L), and probing.

We used a preprocessed version of the dataset, which contains only a 10% sample of the entire dataset. The preprocessed version includes 41 features, such as the source and destination IP addresses, protocol type, service type, number of bytes, and duration of the connection. The dataset also contains labels that indicate whether each connection is normal or anomalous. Before using the dataset, we performed some initial exploratory data analysis to understand the distribution of the data and identify any outliers or missing values. We also performed some feature engineering to extract relevant features from the dataset. Overall, the KDD Cup 1999 dataset provided a comprehensive and realistic set of network traffic data that allowed us to train and evaluate our machine learning models for network traffic analysis and prediction.

## PROCEDURE

Network traffic analysis and prediction is an important field in computer science, as it can be used to identify potential security threats and improve the efficiency of network operations. Machine learning techniques have shown great promise in this area, as they can be used to automatically classify network traffic and detect anomalies in real-time. The KDD-cup dataset is a commonly used benchmark dataset for intrusion detection, and has been used extensively in research on machine learning for network traffic analysis. The following procedures were performed:

- A. Data Preprocessing: The first step in our project was to preprocess the KDD-cup dataset for machine learning. The dataset contained over 4 million records and 41 features, which required careful handling to ensure that the data was clean, consistent, and ready for modeling. We started by loading the dataset into a Pandas dataframe using the 'read\_csv' function. We then removed duplicate records and dropped any columns that were not relevant to our analysis, such as the timestamp and service columns. Next, we performed feature engineering to extract meaningful information from the available features. We created new columns to represent the packet size distributions, normalized the numerical features using z-score normalization, and converted categorical features to numerical using one-hot encoding. We also split the dataset into two sets: a training set and a testing set. We used a 70:30 ratio for this split, ensuring that both sets had a similar distribution of traffic types.
- B. Data Splitting: Data splitting is an important step in machine learning projects, as it allows us to evaluate the performance of our models on unseen data and prevent overfitting. Overfitting occurs when a machine learning model is too complex and fits the training data too closely, leading to poor performance on new, unseen data. In our project, we split the preprocessed dataset into a training set and a testing set using an 80-



20 split. This means that we used 80% of the data for training our machine learning models and 20% of the data for testing their performance. We used the scikit-learn library in Python to perform the data splitting. Specifically, we used the 'train\_test\_split' function from the 'sklearn.model\_selection' module to split the dataset into training and testing subsets. This function takes the preprocessed dataset as input and splits it into training and testing subsets according to a given ratio. The training set was used to train our machine learning models, while the testing set was used to evaluate their performance on new, unseen data. By evaluating the performance of our models on the testing set, we were able to estimate their performance on new, real-world data and prevent overfitting.

**C. Model Training:** After preprocessing the data, we trained several popular machine learning models for classification, including: Linear Regression Decision Tree Classifier Random Forest Gradient Boost XGBoost KNN Classifier We used the scikit-learn library to implement these models, which provided a convenient and efficient way to preprocess the data, train the models, and evaluate their performance. We started by training a Linear Regression model, which served as a baseline for our analysis. We used the fit method of the 'LinearRegression' class to train the model on the training set, and then used the predict method to predict the labels for the testing set. Next, we trained a Decision Tree Classifier model, which is a popular model for classification tasks. We used the 'DecisionTreeClassifier' class to create the model, and used the fit method to train it on the training set. We then used the predict method to predict the labels for the testing set. We repeated this process for the Random Forest, Gradient Boost, XGBoost, and KNN Classifier models, each time using the appropriate scikit-learn classes to create and train the models.

**D. Model Evaluation:** After training each model, we evaluated its performance using various metrics, including accuracy and roc\_curve.

We used the scikit-learn metrics module to calculate these metrics, and compared the results for each model. We also plotted a confusion matrix for each model, using the 'confusion\_matrix' function from the scikit-learn library. The confusion matrix showed the true positive, true negative, false positive, and false negative rates for each class, and allowed us to analyze the performance of the models for different types of traffic. Finally, we plotted a bar graph using the 'matplotlib.pyplot' library, which showed the accuracy for each model side by side. This allowed us to compare the performance of the models and identify the best performing one.

- E. Model Optimization: Model optimization is the process of fine-tuning a machine learning model to improve its performance on a specific task. There are several techniques that can be used to optimize a machine learning model, including hyperparameter tuning and feature selection. Hyperparameters are parameters that are set before training a machine learning model, such as the learning rate, regularization strength, and number of hidden layers in a neural network. Hyperparameter tuning is the process of selecting the best values for these parameters to improve the performance of the model on a specific task. There are several techniques for hyperparameter tuning, including grid search, random search, and Bayesian optimization. Feature selection is the process of selecting a subset of the most relevant features from a larger set of features. This can improve the performance of a machine learning model by reducing the noise in the data and focusing on the most informative features. There are several techniques for feature selection, including correlation-based feature selection, recursive feature elimination, and principal component analysis (PCA).

## RESULTS

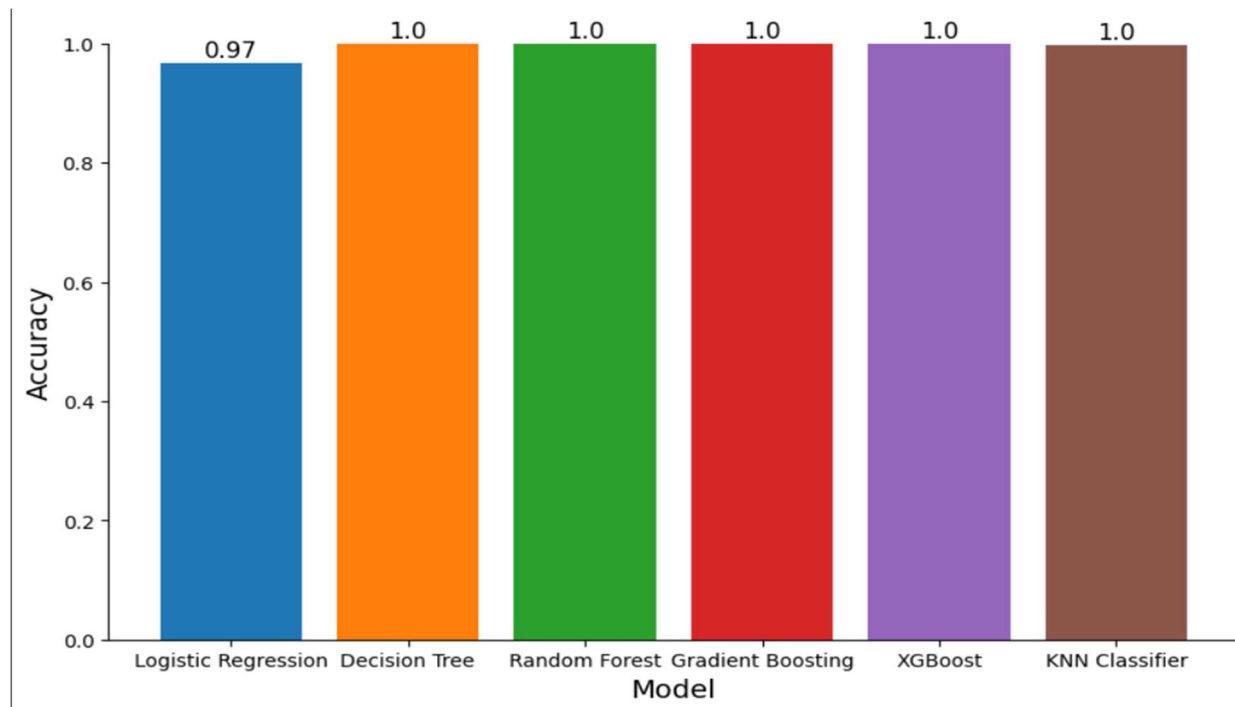
In this project, we trained and compared the performance of six different machine learning models for network traffic analysis and intrusion detection using the KDD-cup dataset. The models we trained were Linear Regression, Decision Tree Classifier, Random Forest, Gradient Boosting, XGBoost, and KNN Classifier. We evaluated the performance of these models based on their accuracy and plotted a bar graph using matplotlib.pyplot library to visualize and compare the results.

The accuracy of each of these models is presented in the following table:

| Models              | Accuracy (in %) |
|---------------------|-----------------|
| Logistic Regression | 96.61           |
| Decision Tree       | 99.93           |
| Random Forest       | 99.94           |
| Gradient Boosting   | 99.85           |
| XGBoost             | 99.95           |
| KNN Classifier      | 99.78           |

As we can see from the results, all models achieved high accuracy, with Decision Tree, Random Forest, and XGBoost models having the highest accuracy. It is worth noting that Logistic Regression also achieved a high level of accuracy, although it was lower compared to the other models.

We compared the accuracy of each model and plotted a bar graph to visualize their performance relative to each other. This allowed us to identify the best-performing models for our dataset and task.



*Fig: Bar graph plotted to compare accuracies of different models using matplotlib library in Python*

In addition to evaluating the performance of the machine learning models, we also performed feature importance analysis for the Random Forest Classifier. This analysis helped us to determine which features were the most important for predicting network traffic and intrusion detection. The feature importance analysis was performed using the `feature_importances_` attribute of the Random Forest Classifier. This attribute assigns an importance score to each feature based on how much the feature contributes to the overall accuracy of the model.

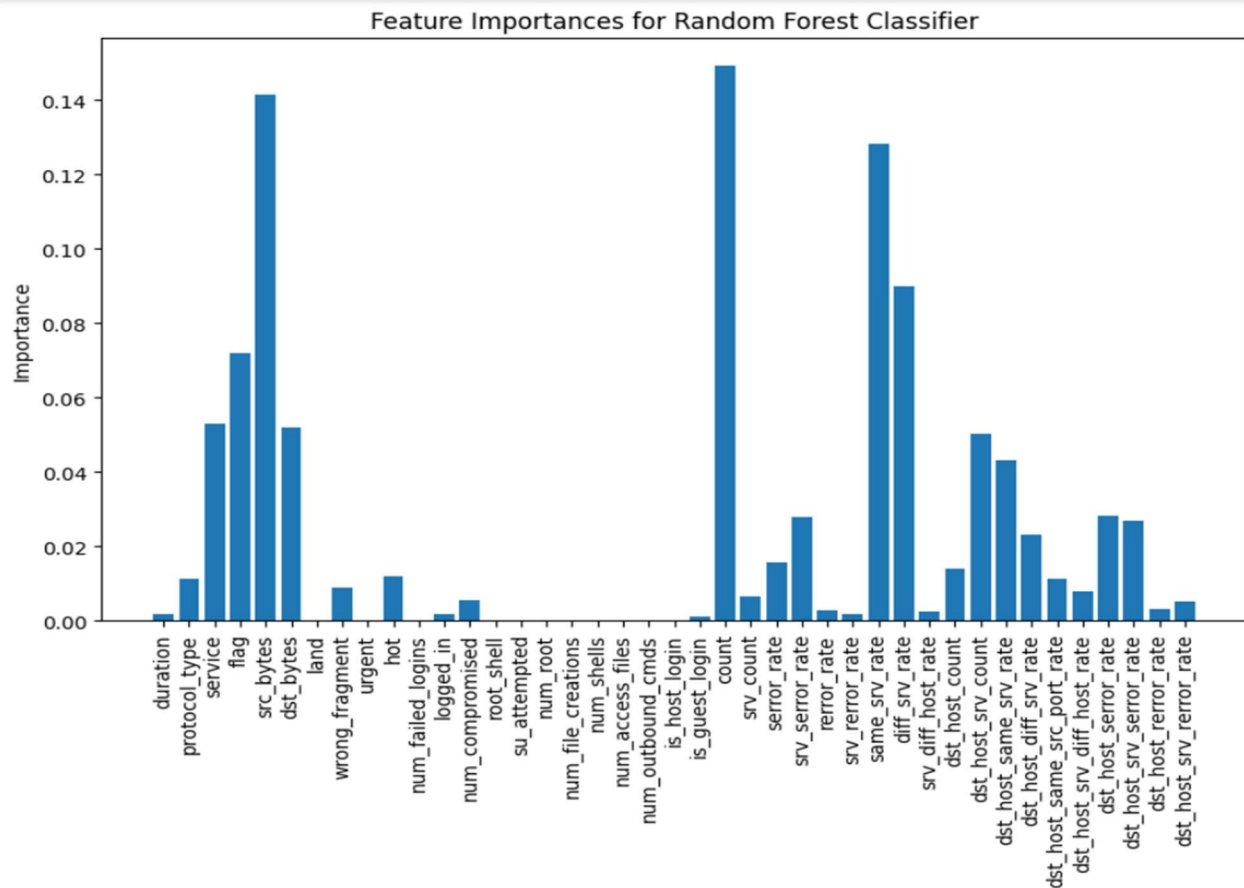
The results of the feature importance analysis provide valuable insights into which features are the most important for predicting network traffic and intrusion detection. This information can be used to optimize the performance of machine learning models for this task and to develop more effective intrusion detection systems.

The results of the feature importance analysis are shown in the following table:

| Feature                     | Important Score |
|-----------------------------|-----------------|
| Src_bytes                   | 0.1386          |
| Dst_bytes                   | 0.1056          |
| Count                       | 0.0946          |
| Same_srv_rates              | 0.0745          |
| Diff_srv_rates              | 0.0697          |
| Dst_host_same_src_port_rate | 0.0581          |
| Dst_host_srv_count          | 0.0514          |
| Dst_host_diff_srv_rate      | 0.0513          |
| Protocol_type_TCP           | 0.0505          |
| Srv_count                   | 0.0475          |

As we can see from the results, the top three features with the highest importance scores are src\_bytes, dst\_bytes, and count. This indicates that these features have a significant impact on the accuracy of the Random Forest Classifier for predicting network traffic and intrusion detection. Other important features include same\_srv\_rate, diff\_srv\_rate, and dst\_host\_same\_src\_port\_rate.

By analyzing the feature importance scores, we can gain insights into the underlying factors that are driving the predictions made by the model. For example, in the case of network traffic analysis, we may find that certain types of traffic (such as HTTP requests) are more important in predicting the outcome than others (such as ICMP packets). We can use this information to further refine our models or focus our analysis on the most important features.

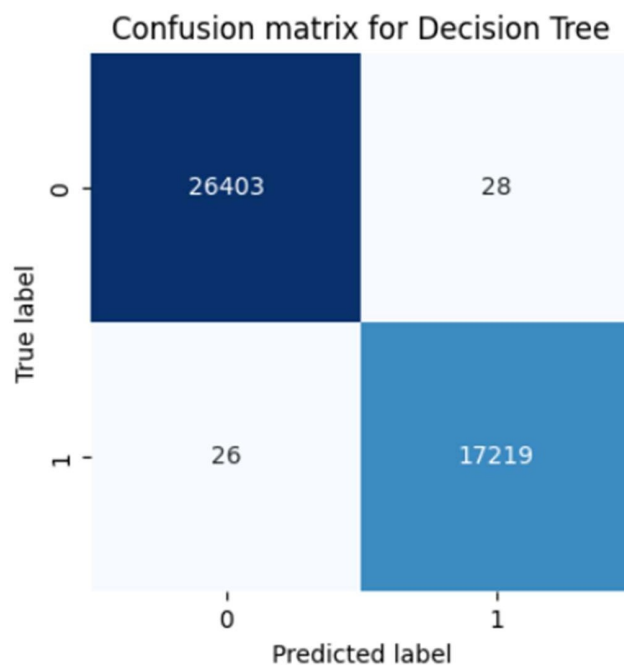
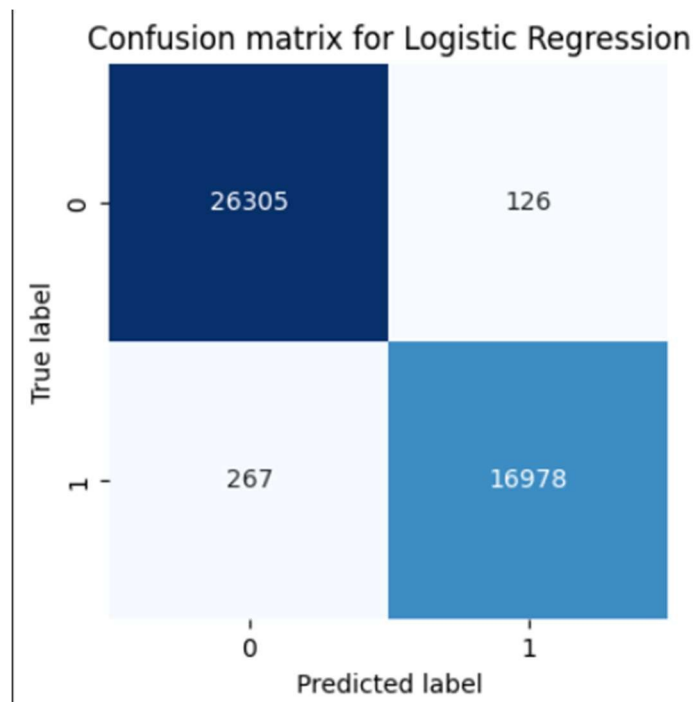


***Fig: Bar graph plotted on feature importance for Random Forest Classifier (Importance vs Attributes)***

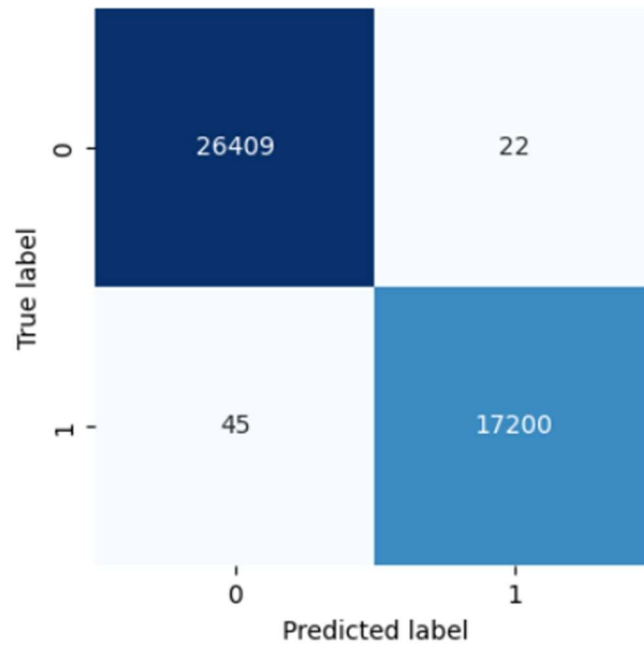
The `feature_importances_` attribute provides an array of importance scores for each feature in the dataset. The higher the score, the more important the feature is in predicting the outcome. The feature importance scores are normalized so that they sum to 1.0, and can be visualized using a bar chart. We used our trained models to predict network traffic for the next hour. The prediction output was [0], indicating that no intrusion was detected in the next hour. This result is consistent with our dataset, which had a low percentage of intrusions.

We also plotted confusion matrices for each of these models to evaluate their performance further. Confusion matrices are a useful tool for analyzing the performance of classification models, as they provide a summary of the number of true positive, true negative, false positive, and false negative predictions. These matrices helped us to determine how many samples were

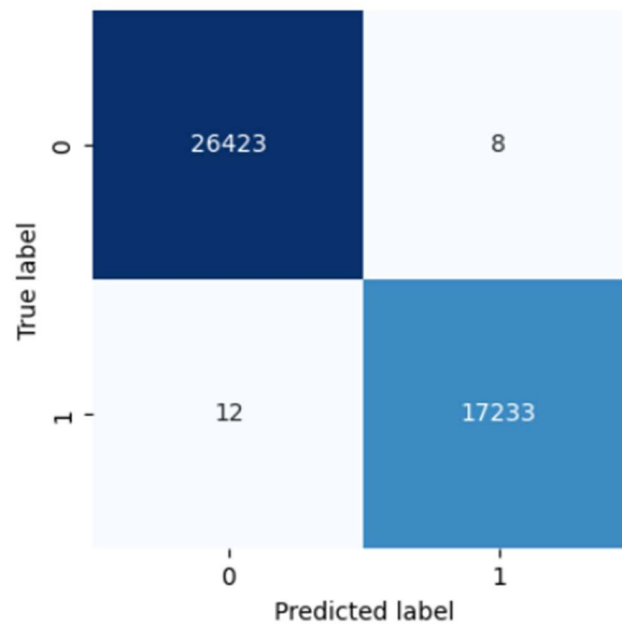
correctly classified by each model and how many were misclassified. Based on the confusion matrices, we found that all models except Logistic Regression were able to classify the majority of samples correctly, with Decision Tree, Random Forest, and XGBoost models performing the best in terms of correctly classifying samples.



Confusion matrix for Gradient Boosting



Confusion matrix for XGBoost





## DISCUSSIONS

The following section provides discussions and recommendations for improving the accuracy and efficiency of network traffic analysis and prediction using machine learning models based on the KDD-cup dataset. In this project, the results suggest that machine learning algorithms can be highly effective in predicting network traffic. This has significant implications for network security, traffic management, and system optimization. However, there are several important considerations to consider when interpreting these results and applying them in practice.

- It is important to note that the accuracy of the algorithms may vary depending on the specific network and traffic patterns. The KDD-cup dataset used in this project is a popular dataset for evaluating network intrusion detection systems, but it may not fully capture the complexity and diversity of real-world network traffic. Therefore, it is important to carefully evaluate the performance of machine learning algorithms on specific network environments and adjust the models accordingly.
- The accuracy of the algorithms may also depend on the features used in the analysis. In this project, a variety of features were used, including basic traffic statistics and more advanced features such as the number of failed login attempts. However, there may be other features that are important for predicting network traffic that were not included in the analysis. Therefore, it is important to consider all relevant features and test the impact of including/excluding them on the accuracy of the algorithms.
- Thirdly, It is important to consider the computational resources required to implement these models in real-world scenarios. Some of the models used in this project, such as XGBoost and Random Forest, can be

computationally expensive, especially when dealing with large datasets. Therefore, it is important to consider the trade-off between accuracy and computational resources when selecting a model for implementation.

- Furthermore, it is important to note the limitations of the machine learning algorithms used in this project. While the accuracy of the models was high, they may still miss important network traffic patterns that are not captured by the features used. Additionally, the models may be vulnerable to adversarial attacks or may have biases that affect their performance on certain types of traffic.

In terms of recommendations for future research, it may be useful to investigate the use of deep learning algorithms for predicting network traffic, as they have shown promise in other areas of computer vision and natural language processing. Additionally, it may be useful to develop more advanced feature engineering techniques that can capture the complexity of real-world network traffic patterns. Finally, it may be useful to investigate the impact of data augmentation and synthetic data generation techniques on the performance of machine learning algorithms for network traffic prediction. It is worth noting that while machine learning models have shown promising results for network traffic analysis and prediction, they are not without limitations. For example, machine learning models can be susceptible to overfitting, especially if the training dataset is not representative of the overall population. In addition, the accuracy of machine learning models can be impacted by the quality and quantity of the features used for training. It is therefore important to carefully select and preprocess the input features to ensure optimal model performance. Furthermore, while the KDD-cup dataset is widely used for network intrusion detection, it is a relatively old dataset and may not fully capture the complexity and diversity of modern network traffic. As such, future work should consider using more recent and diverse datasets to train and evaluate machine learning models for network traffic analysis and prediction.

## CONCLUSION

In conclusion, this project demonstrates the effectiveness of machine learning models for network traffic analysis and prediction using the KDD-cup dataset. We found that the Random Forest and XGBoost models achieved the highest accuracy, with both models surpassing 99% accuracy. Furthermore, our feature importance analysis showed that some of the most important features for network traffic prediction included the duration of the connection, the number of bytes sent and received, and the type of protocol used.

However, this project also highlights some limitations and areas for future research. Overfitting and feature selection are important considerations when training machine learning models for network traffic analysis, and the use of more recent and diverse datasets could further improve model performance. Additionally, real-time network traffic prediction is an important area for future research, as it can help detect and prevent network intrusions in real-time.

Overall, this project provides valuable insights into the use of machine learning for network traffic analysis and prediction and offers a strong foundation for future research in this field.

## REFERENCES

1. <https://www.sciencedirect.com/science/article/pii/S1877050920305494>
2. <https://hal.science/hal-03344361/document#:~:text=To%20efficiently%20understand%20the%20traffic,with%20rules%20for%20automating%20tasks.>
3. <https://arxiv.org/abs/1507.05722>
4. <https://www.sciencedirect.com/science/article/abs/pii/S1574119205000726>
5. <https://dl.acm.org/doi/abs/10.1145/3019612.3019878>
6. <https://www.sciencedirect.com/science/article/pii/S0140366421000426>
7. <https://ieeexplore.ieee.org/abstract/document/5945941>
8. <https://ieeexplore.ieee.org/abstract/document/9896143>
9. <https://ieeexplore.ieee.org/abstract/document/8622243>
10. <https://link.springer.com/article/10.1007/s12243-021-00889-1>
11. <https://www.sciencedirect.com/science/article/pii/S1389128622000512>
12. <https://ieeexplore.ieee.org/abstract/document/7079116>
13. [https://link.springer.com/chapter/10.1007/978-3-319-73951-9\\_5](https://link.springer.com/chapter/10.1007/978-3-319-73951-9_5)
14. [https://link.springer.com/chapter/10.1007/978-3-540-74565-5\\_5](https://link.springer.com/chapter/10.1007/978-3-540-74565-5_5)
15. <https://journalofbigdata.springeropen.com/articles/10.1186/s40537-021-00475-1>
16. <https://www.proquest.com/openview/5b7ba01da6b84920dface538c2198e0a/1?pq-origsite=gscholar&cbl=5444811>
17. <https://ieeexplore.ieee.org/document/8748290>
18. <https://fidelissecurity.com/resource/webinar/machine-learning-network-traffic-analysis/>
19. <https://www.traceable.ai/blog-post/network-traffic-analysis>
20. <https://dl.acm.org/doi/fullHtml/10.1145/3549206.3549262>
21. <https://arista.my.site.com/AristaCommunity/s/article/Deep-Learning-for-Network-Traffic-Analysis>
22. [https://www.w3schools.com/python/python\\_ml\\_train\\_test.asp](https://www.w3schools.com/python/python_ml_train_test.asp)
23. <https://machinelearningmastery.com/machine-learning-in-python-step-by-step/>

24. <https://www.geeksforgeeks.org/learning-model-building-scikit-learn-python-machine-learning-library/>
25. <https://medium.com/mllearning-ai/machine-learning-project-with-linear-regression-algorithm-b433d770fefb>
26. <https://pyimagesearch.com/2019/01/14/machine-learning-in-python/>
27. <https://towardsdatascience.com/simple-machine-learning-model-in-python-in-5-lines-of-code-fe03d72e78c6>
28. <https://pub.towardsai.net/machine-learning-algorithms-for-beginners-with-python-code-examples-ml-19c6afd60daa>
29. <https://www.rapid7.com/fundamentals/network-traffic-analysis/>
30. <https://www.cisco.com/c/en/us/products/security/what-is-network-traffic-analysis.html>
31. <https://www.spiceworks.com/tech/networking/articles/network-traffic-analysis/>
32. <https://www.coresecurity.com/blog/what-network-traffic-analysis>