

**A PROJECT REPORT**  
on  
**“Hybrid Intrusion Detection Model 2.0”**

**Submitted to**  
**KIIT Deemed to be University**

**In Partial Fulfillment of the Requirement for the Award of**  
**BACHELOR’S DEGREE IN**  
**COMPUTER SCIENCE AND ENGINEERING**

**BY**

<b>ANIKET KUMAR</b>	2005288
<b>PRANAV PANT</b>	2005321
<b>PRATYUSH YUVRAJ</b>	2005469

**UNDER THE GUIDANCE OF**  
**Prof. Lalit Kumar Vashishtha**



**SCHOOL OF COMPUTER ENGINEERING**  
**KALINGA INSTITUTE OF INDUSTRIAL TECHNOLOGY**  
**BHUBANESWAR, ODISHA - 751024**  
**November, 2023**

A PROJECT REPORT  
on  
**“Hybrid Intrusion Detection Model 2.0”**

Submitted to  
**KIIT Deemed to be University**

In Partial Fulfillment of the Requirement for the Award of

**BACHELOR’S DEGREE IN  
COMPUTER SCIENCE AND ENGINEERING**

BY

ANIKET KUMAR	2005288
PRANAV PANT	2005321
PRATYUSH YUVRAJ	2005469

UNDER THE GUIDANCE OF  
**Prof. Lalit Kumar Vashishtha**



SCHOOL OF COMPUTER ENGINEERING  
**KALINGA INSTITUTE OF INDUSTRIAL TECHNOLOGY**  
BHUBANESWAR, ODISHA -751024  
November, 2023

# KIIT Deemed to be University

School of Computer Engineering  
Bhubaneswar, ODISHA 751024



## CERTIFICATE

This is certify that the project entitled

### **“Hybrid Intrusion Detection Model 2.0”**

submitted by

ANIKET KUMAR	2005288
PRANAV PANT	2005321
PRATYUSH YUVRAJ	2005469

is a record of bonafide work carried out by them, in the partial fulfillment of the requirement for the award of **Degree of Bachelor of Engineering (Computer Science & Engineering)** at KIIT Deemed to be university, Bhubaneswar. This work is done during year **2023-2024**, under our guidance.

Date: 30 / 11 / 2023

**(Prof. Lalit Kumar Vashishtha)  
Project Guide**

## **ACKNOWLEDGEMENT**

We are profoundly grateful to our Project Guide **Prof. Lalit Kumar Vashishtha** of Kalinga Institute of Industrial Technology, Bhubaneshwar for his expert guidance and continuous encouragement throughout to see that this project reaches its target since its commencement to its completion. His expertise, advice, and constructive criticism have been invaluable in helping us to complete this project successfully. His constant encouragement and motivation have inspired us to put in our best effort and strive for excellence.

The knowledge and skills we have gained during this project will undoubtedly be valuable assets as we continue to pursue our academic and professional goals. We are also grateful for the opportunity given to us by **Kalinga Institute of Industrial Technology, Bhubaneshwar** to work on this project and also the facilities and inventory provided to us for its successful completion.

**ANIKET KUMAR (2005288)**

**PRANAV PANT (2005321)**

**PRATYUSH YUVRAJ (2005469)**

## **ABSTRACT**

As cyber threats escalate in complexity, the need for robust intrusion detection systems (IDS) becomes paramount to secure networked environments. This research delves into the landscape of hybrid IDS, amalgamating signature-based, anomaly-based, and machine learning approaches to fortify network security. The study evaluates the performance of diverse models—AdaBoost, decision tree, random forest, gradient boosting, extremely randomized trees, and neural networks—utilizing datasets (CIC-IDS2017, NSL-KDD, SDN, UNSW-NB15) representative of real-world network challenges.

The research employs a meticulous methodology, individually training models to identify top performers. A hybrid ensemble model, integrating the strengths of the five most effective models, is then constructed. The approach aims to surpass existing intrusion detection systems, providing a more resilient defense against evolving cyber threats.

This investigation not only contributes insights into individual model performance but also proposes an innovative hybrid ensemble system, advancing the current state of intrusion detection. In the ever-changing landscape of cybersecurity, this research stands as a crucial step toward enhancing network security and mitigating the risks posed by sophisticated intrusions.

**Keywords** :: Intrusion Detection Systems (IDS), Hybrid Models, Machine Learning, Neural Networks, Performance Evaluation, Ensemble Model, Network Security

## **MOTIVATION**

The motivation behind our continued research and improvement efforts stems from the critical need to enhance the effectiveness of intrusion detection systems in cloud environments. The challenges associated with identifying both known and unknown attacks demand innovative solutions that can adapt to evolving threat landscapes. The combination of signature-based and anomaly-based detection methods presents a promising avenue for comprehensive security coverage. Our motivation is driven by the desire to contribute to the ongoing discourse on bolstering cybersecurity defenses in cloud environments.

Our unwavering commitment to ongoing research and improvement efforts is fundamentally rooted in recognizing the indispensable requirement to elevate the efficacy of intrusion detection systems within the dynamic realm of cloud environments. The contemporary landscape of cybersecurity is marked by an incessant surge in sophisticated cyber threats, compelling us to push the boundaries of innovation.

The three standard datasets, UNSW-NB15, CICIDS2017, and NSL-KDD, serve as pivotal tools for our research evaluation. The UNSW-NB15 dataset, with its 175,341 network flows and diverse attack classes, provides a realistic representation of modern cyber threats. Similarly, the CICIDS2017 dataset, with its extensive network flows and varied attack scenarios, offers a rich source for assessing the robustness of our proposed model. The NSL-KDD dataset, a modified version of the KDDCup99 dataset, addresses the limitations of its predecessor and contributes to the comprehensive testing of our intrusion detection model.

In pursuit of a holistic approach to security, we have embraced the synergy of signature-based and anomaly-based detection methods. This amalgamation represents a strategic and promising avenue, as it not only addresses the well-established attack patterns but also dynamically responds to deviations from expected behavior. By integrating these complementary approaches, we strive to provide a robust and comprehensive security coverage that is capable of detecting a wide spectrum of threats, both conventional and emerging.

In essence, our motivation is not merely a pursuit of technical excellence but a dedication to being at the forefront of cybersecurity innovation, fostering resilience, and actively participating in the ongoing evolution of defense strategies for cloud environments.

# Content Page

01	Introduction	01
	1.1 Intrusion Threats in Cyberspace	01
	1.2 Existing Hybrid Intrusion Detection Models	01
	1.3 Research Project Overview	02
	1.4 Models Under Consideration	02
	1.5 Hybrid Ensemble Model Construction	02
02	Literature Review	03
03	Introduction to Datasets	12
	3.1 Intrusion Detection Evaluation Dataset (CIC-IDS2017)	12
	3.2 NSL-KDD Dataset	13
	3.3 DDOS attack SDN Dataset	14
	3.4 UNSW-NB15 Dataset	15
04	Introduction to Models	16
	4.1 Adaptive Boosting	16
	4.1.1 Introduction	16
	4.1.2 Background	16
	4.1.3 The Covert Operations of AdaBoost	20
	4.1.4 Model Evaluation	29
	4.1.5 Pragmatic Cryptography	33
	4.1.6 Applications	36
	4.1.7 Conclusions	38
	4.2 Random Forest Modeling	39
	4.3 Decision Tree Model	42
	4.4 Extremely Randomized Tree Model	45
	4.5 Gradient Boosting Model	48
	4.6 Neural Networks	51
05	Optimization Techniques	55
	5.1 Hyperparameters and Hyperband	55
	5.2 Gradient Based Optimization	57
	5.3 Grid Search Optimization	59
	5.4 Optuna and Bayesian Optimization	61
06	Ensemble Methods	63

	6.1	Hard Ensembling in Machine Learning	63
	6.2	Soft Ensembling in Machine Learning	67
	6.3	Bayesian Averaging	71
	6.4	Dynamics in Ensemble Method	74
	6.5	Weighted Average Ensembler Method	77
	07	Evaluation Metrics	79
	7.1	Area under Curve (AUC)	79
	7.2	Classification Accuracy	81
	7.3	F1 Score	82
	7.4	Precision	83
	7.5	Recall	84
	7.6	Matthews Correlation Coefficient	85
	7.7	Specificity	86
	7.8	LogLoss	87
	08	Intrusion Detection Evaluation Dataset (CIC-IDS2017)	88
	8.1	Introduction	88
	8.2	Data-Preprocessing	89
	8.3	Implementation	90
	8.4	Table	94
	8.5	Results	95
	8.6	Ensemble Technique	97
	09	NSL-KDD Dataset	100
	9.1	Introduction	100
	9.2	Data-Preprocessing	101
	9.3	Implementation	102
	9.4	Table	106
	9.5	Results	107
	9.6	Ensemble Technique	109
	9.7	Conclusion	112
	10	DDOS Attack SDN Dataset	114
	10.1	Introduction	114
	10.2	Data-Preprocessing	115
	10.3	Implementation	116
	10.4	Table	120
	10.5	Results	121
	10.6	Ensemble Technique	123
	10.7	Conclusion	125

11	UNSW-NB15	127	
	11.1	Introduction	127
	11.2	Data-Preprocessing	128
	11.3	Implementation (Before Optimization)	129
	11.4	Table (Before Optimization)	133
	11.5	Results (Before Optimization)	134
	11.6	Optimization	136
	11.7	Results (After Optimization)	154
	11.8	Ensemble Methods	155
	11.8.1	Hard Ensemble Technique	155
	11.8.2	Soft Ensemble Technique	157
	11.8.3	Bayesian Model Average	158
	11.8.4	Dynamic Method	159
	11.8.5	Weighted Average	161
	11.9	Table (After Optimization)	163
	11.10	End Results (After Optimization)	163
12	Final Result and Conclusions	165	
13	Future Scope	166	
	13.1	Optimization for Computational Efficiency	166
	13.2	Scalability in Big Data Analytics	166
	13.3	Integration with Deep Learning Architectures	166
	13.4	Real-Time Threat Intelligence Integration	166
14	References	167	
15	Plagarism Report	168	

# Chapter 1

## Introduction

In an era dominated by digital connectivity, the increasing sophistication of cyber threats poses a significant challenge to the security of networked systems. The pervasiveness of interconnected devices and the ever-evolving landscape of cyber threats necessitate robust intrusion detection mechanisms to safeguard sensitive information. Intrusion detection systems (IDS) play a pivotal role in identifying and mitigating potential security breaches, acting as a crucial line of defense against malicious activities in cyberspace.

### 1.1 Intrusion Threats in Cyberspace :-

Cybersecurity faces a myriad of intrusion threats, ranging from traditional attacks like Denial of Service (DoS) and Distributed Denial of Service (DDoS) to more advanced threats such as malware, phishing, and zero-day exploits. As attackers continually refine their techniques, intrusion detection becomes increasingly complex, requiring innovative approaches to stay ahead of evolving threats.

### 1.2 Existing Hybrid Intrusion Detection Models :-

To address the multifaceted nature of intrusion threats, researchers and practitioners have explored hybrid intrusion detection systems. These systems leverage the strengths of multiple detection techniques, combining signature-based, anomaly-based, and machine learning approaches to enhance detection accuracy and reduce false positives. While several hybrid models have been proposed, the need for continual improvement remains, prompting an in-depth investigation into their efficacy.

### **1.3 Research Project Overview :-**

This research project focuses on evaluating the performance of various hybrid intrusion detection models using diverse datasets to comprehensively understand their strengths and limitations. The datasets selected for analysis include CIC-IDS2017, NSL-KDD, SDN, and UNSW-NB15, each presenting unique challenges representative of real-world network environments.

### **1.4 Models Under Consideration :-**

The study employs a diverse set of machine learning models, including AdaBoost, decision tree, random forest, gradient boosting, extremely randomized trees, and neural networks. These models are chosen for their established track records in intrusion detection and their capacity to capture complex patterns in network traffic data.

### **1.5 Hybrid Ensemble Model Construction :-**

The research begins with individual training of the selected models on the specified datasets, evaluating their performance in terms of accuracy, precision, recall, and F1 score. Through a selective selection procedure, the top-performing models are identified for further analysis. The culmination of the study involves the construction of a hybrid ensemble model. The five most effective models are integrated to create a robust ensemble, harnessing the collective predictive power of diverse algorithms. This approach aims to capitalize on the strengths of individual models while mitigating their respective limitations, resulting in a more resilient and effective intrusion detection system.

The research not only contributes to the understanding of individual model performance but also proposes a novel hybrid ensemble approach that outperforms existing intrusion detection systems. By leveraging the strengths of multiple models, the proposed system offers a more comprehensive and adaptive defense against an evolving threat landscape.

In conclusion, this research endeavors to enhance the efficacy of intrusion detection systems by providing insights into the performance of hybrid models and introducing an innovative ensemble approach. As cyber threats continue to evolve, the findings of this study contribute to the ongoing efforts to fortify network security and mitigate the risks posed by malicious intrusions.

# Chapter 2

## Literature Review

Lalit Kumar Vashishtha et al., September 2022 [1], proposed the approach in the paper “HIDM : A Hybrid Intrusion Detection Model for Cloud Based Systems” based over the the performance of the proposed model using three standard datasets : UNSW-NB15, CICIDS2017, and NSL-KDD in which HIDM: A hybrid intrusion detection model that combines signature-based and anomaly-based detection methods. The model consists of four phases: data collection, data analysis, known attack detection, and unknown attack detection. The paper also evaluates The model uses a signature database to store the signatures of known attacks and a profile storage to store the normal behaviour patterns.

The model also uses three machine learning algorithms: random forest, gradient boosting, and neural network, to train and test the anomaly-based detection module. The model uses majority voting to decide the final classification of the network rows. The model generates alerts for any detected intrusion and updates the signature database and the profile storage accordingly.

Some drawbacks of the paper are:

- (a) The paper does not provide a clear definition of the difference between known and unknown attacks, and how they are identified and labeled in the datasets.
- (b) The paper does not explain how the signature database and the profile storage are initially created and maintained, and what are the criteria for updating them.
- (c) The paper does not compare the proposed model with other existing hybrid intrusion detection models or justify the choice of the machine learning algorithms used in the anomaly-based detection module.
- (d) The paper does not report the computational cost, time complexity, and scalability of the proposed model, which are important factors for evaluating the feasibility and efficiency of the model in real-world scenarios.

Dr Santosh Kumar Sahu et al., February 2014 [2], in the research paper “A Detail Analysis on Intrusion Detection Datasets” proposes some techniques to handle large and high-dimensional datasets for intrusion detection system (IDS) modeling. The techniques include filling missing values, removing redundant samples, selecting relevant features, reducing dimensionality, and normalizing the data. The paper also applies various data mining and machine learning algorithms to the preprocessed datasets and compares the results.

The paper uses three intrusion detection datasets: KDD Cup 99, NSL-KDD, and GureKDD. These datasets are derived from the DARPA 98 dataset, which contains simulated network traffic with different types of attacks. The paper analyzes the statistical properties and the distribution of the features and the classes in each dataset.

**Drawbacks of the paper :-**

The paper has some limitations and weaknesses, such as :-

- (a) The paper does not provide a clear motivation and objective for the data preprocessing techniques. It does not explain how the techniques improve the performance and efficiency of the IDS models.
- (b) The paper does not compare the results of the data mining and machine learning algorithms with other existing methods or baselines. It does not provide any evaluation metrics or criteria to measure the accuracy and effectiveness of the algorithms.
- (c) The paper does not discuss the implications and applications of the findings. It does not provide any suggestions or recommendations for future work or research directions.

Amar Meryem and Bouabid EL Ouahidi, May 2020 [3], in research paper “Hybrid Intrusion Detection System using Machine Learning” proposes a hybrid intrusion detection system using machine learning and log file analysis. The datasets used in the research paper are NSL-KDD and Honeynet Project. The research paper aims to enhance cloud security by detecting malicious behaviour from different sources of log files, such as access, audit, error and SSH logs. The research paper also claims to find new attack signatures and update security rules dynamically.

The proposed approach consists of the following steps :-

- (i) Structuring log events: This step converts the unstructured log files into a matrix of 41 features, corresponding to the NSL-KDD dataset. The research paper uses a MapReduce program to process the log files and assign a unique identifier to each row.
- (ii) Eradicating redundancy: This step reduces the size of the matrix by eliminating duplicate rows and calculating the frequency of each distinct row. The frequency is considered as a weight of the user activity.
- (iii) Classifying unlabelled behaviour: This step clusters the matrix into k classes using K-means algorithm. The research paper uses the elbow method to determine the optimal value of k.
- (iv) Labelling behaviours using NSL-KDD: This step labels each cluster as normal or attack using the NSL-KDD dataset as a knowledge base. The research paper compares four machine learning algorithms: k-nearest neighbours, Naïve Bayes, support vector machine and logistic regression. The research paper selects k-nearest neighbours as the best classifier based on the performance metrics such as accuracy, false positive rate, precision and recall.
- (v) Finding new attack signatures: This step finds correlations between the features of the NSL-KDD dataset and the log files to discover new malicious behaviours and identify new signatures. The research paper uses several similarity metrics such as Jaccard coefficient, cosine similarity and Pearson correlation to measure the similarity between the features. The research paper also updates the training dataset and the security rules based on the new signatures.

Some drawbacks of the paper are:

- (a) The paper does not provide a clear explanation of how the MapReduce program works and how it assigns a unique identifier to each row of the matrix.
- (b) The paper does not justify the choice of the NSL-KDD dataset as a knowledge base, which may be outdated or incomplete for detecting new attacks.
- (c) The paper does not evaluate the scalability and efficiency of the proposed approach when dealing with large and streaming log files.
- (d) The paper does not compare the proposed approach with other existing hybrid intrusion detection systems or state-of-the-art techniques.

Vajiheh Hajisalem and Shahram Babaie, February 2018 [4], in research paper “A Hybrid Intrusion Detection System based on ABC-AFS Algorithm for Misuse and Anomaly Detection” proposes a hybrid intrusion detection system based on ABC-AFS algorithm for misuse and anomaly detection.

The paper aims to improve the detection accuracy and efficiency of intrusion detection systems by combining the advantages of artificial bee colony and artificial fish swarm algorithms. The paper also applies fuzzy C-means clustering, correlation-based feature selection, and CART rule generation techniques to preprocess the data and generate reliable rules for classification.

The paper uses two well-known datasets for evaluation: NSL-KDD and UNSW-NB15. NSL-KDD is a modified version of KDD Cup 99 dataset that contains 41 features and 24 types of attacks. UNSW-NB15 is a new dataset that contains 49 features and 9 types of modern attacks. The paper compares the performance of the proposed method with other state-of-the-art methods in terms of detection rate, false positive rate, computational complexity, and time cost.

The paper claims that the proposed method outperforms the existing methods and achieves 99% detection rate and 0.01% false positive rate on both datasets. The paper also claims that the proposed method has comparable computational complexity and time cost with other methods.

Some possible drawbacks of the paper are:

- (a) The paper does not provide enough details on how the hybrid ABC-AFS algorithm works and how it is trained using the generated rules.
- (b) The paper does not explain how the optimal parameters for the ABC and AFS algorithms are determined and whether they are sensitive to different datasets or scenarios.
- (c) The paper does not conduct any statistical tests to verify the significance of the results and the differences between the methods.
- (d) The paper does not discuss the limitations and challenges of the proposed method and the future directions for improvement.

Soulaiman Moualla et al., June 2021 [5], in the research paper “Improving the Performance of Machine Learning-Based Network Intrusion Detection Systems on the UNSW-NB15 Dataset” presents a novel approach to network intrusion detection systems (IDS) using a multistage supervised machine learning model. The model is designed to be scalable and flexible, with a high performance evaluation that surpasses related works in terms of accuracy, false alarm rate, and computational cost.

The authors utilized the UNSW-NB15 dataset, a benchmark dataset that includes the latest cyber-attacks. This dataset is publicly available and contains network traffic features and labels for different types of attacks.

The proposed approach involves several stages, including resampling, feature selection, classification, and aggregation. The system uses a "One-Versus-All" strategy to detect each attack separately, and then combines the results using a logistic regression layer. This design allows the system to be easily extended to new classes or attacks.

However, the approach has some drawbacks. The authors suggest that their system can be improved by adding an unsupervised learning stage to detect normal and abnormal traffic without labels. This indicates that the current model may have limitations in detecting new or unknown attacks that are not included in the training data.

In conclusion, the research provides a significant contribution to the field of network IDS. However, future work is needed to improve the system's ability to detect new and unknown attacks.

RUIZHE ZHAO et al., June 2021 [6], in the research paper “A Hybrid Intrusion Detection System Based on Feature Selection and Weighted Stacking Classifier” proposes a hybrid intrusion detection system based on a CFS-DE feature selection algorithm and a weighted Stacking classification algorithm.

The paper uses two datasets to evaluate the proposed methods: the NSL-KDD dataset and the CSE-CIC-IDS2018 dataset.

The CFS-DE algorithm uses the correlation-based feature selection (CFS) to evaluate the feature subsets and the differential evolution (DE) algorithm to optimize the feature subsets. The weighted Stacking algorithm uses the Kappa coefficient to assign different weights to the base classifiers and improve the classification performance of the Stacking algorithm.

The NSL-KDD dataset is an improved version of the KDD-CUP 1999 dataset, which contains 41 features and 23 types of attacks. The CSE-CIC-IDS2018 dataset is a real-world packet captured dataset, which contains 80 features and 15 types of attacks.

The paper has some limitations that can be improved in future work. For example, the paper only uses the data of Wednesday in the CSE-CIC-IDS2018 dataset, which may not reflect the performance of the proposed methods on other days. The paper also does not compare the proposed methods with other state-of-the-art methods in the field of intrusion detection, which may limit the generalization and applicability of the proposed methods. Moreover, the paper does not provide a detailed analysis of the impact of different parameters and feature subsets on the classification results, which may affect the robustness and scalability of the proposed methods.

Ibraheem Aljamal et al., May 2019 [7], in the reaserch paper “Hybrid Intrusion Detection System Using Machine Learning Techniques in Cloud Computing Environments” paper proposes a hybrid system that combines unsupervised and supervised machine learning techniques to detect known and unknown attacks in cloud computing environments. The paper uses the UNSW-NB15 data set to evaluate the proposed hybrid system.

The system consists of two phases: a clustering phase and a classification phase. The clustering phase uses K-means algorithm to group network flows into normal or abnormal clusters, and the classification phase uses SVM algorithm to build a detection model based on the cluster labels. The paper claims that the hybrid system can improve the accuracy and efficiency of intrusion detection in the cloud.

The paper reviews the existing machine learning techniques for intrusion detection, such as signature-based and anomaly-based approaches, and their advantages and disadvantages. The paper also discusses the supervised and unsupervised learning algorithms, such as support vector machines, linear regression, logistic regression, Naive Bayes, linear discriminant analysis, decision trees, neural networks, k-means, fuzzy c-means, and self-organizing maps, and their applications in intrusion detection systems.

The paper explains the rationale and benefits of using a hybrid approach that combines the strengths of both supervised and unsupervised learning techniques.

The paper has some limitations and drawbacks that could be addressed in future work. For example, the paper does not provide a clear explanation of how the threshold functions for the clustering phase are defined and how the parameters are tuned. The paper also does not compare the proposed hybrid system with other existing hybrid systems or other state-of-the-art intrusion detection systems. The paper does not report the false positive rate or the false negative rate of the proposed system, which are important metrics for evaluating the performance of intrusion detection systems. The paper does not discuss the scalability, robustness, or adaptability of the proposed system in different cloud computing scenarios or environments.

Nour Moustafa and Jill Slay, March 2015 [8], in the research paper “UNSW-NB15: A Comprehensive Data set for Network Intrusion Detection Systems” proposed a method to create a comprehensive data set for network intrusion detection systems (NIDS) that can reflect modern network traffic scenarios, low footprint intrusions and depth structured information. The paper argues that the existing benchmark data sets such as KDDCUP99 and NSLKDD are outdated, redundant, skewed and not representative of the current network threat environment. The paper uses the UNSW-NB15 data set, which contains 2.54 million records with 49 features and 10 labels. The paper also compares the UNSW-NB15 data set with the KDDCUP99 data set, which contains 4.9 million records with 42 features and 5 labels

The paper describes the process of creating the UNSW-NB15 data set, which consists of a hybrid of real normal and synthetic attack activities. The paper uses the IXIA PerfectStorm tool to generate nine families of attacks that are updated from a CVE site. The paper also uses the tcpdump tool to capture the network traffic in pcap files. The paper then extracts 49 features from the pcap files using Argus, Bro-IDS and twelve algorithms developed in C#. The paper labels the data set from a ground truth table that contains all simulated attack types.

The paper has some limitations and challenges that can be addressed in future work. For example, the paper does not provide any evaluation or analysis of the performance of NIDS using the UNSW-NB15 data set. The paper also does not discuss the ethical and legal implications of using real normal traffic and synthetic attack traffic in the data set.

The paper does not explain how the synthetic attack traffic is generated and validated by the IXIA tool. The paper does not compare the UNSW-NB15 data set with other existing or recent data sets for NIDS.

Achmad Akbar et al., December 2021 [9], in the research paper “A Hybrid Machine Learning Method for Increasing the Performance of Network Intrusion Detection Systems” presents a novel approach to network intrusion detection systems (IDS). The authors propose a hybrid machine learning method that combines feature selection and data reduction techniques to enhance the performance of IDS. This method is meticulously detailed in the paper, comprising four steps: feature importance ranking, local outlier factor, decision tree classification, and method evaluation. A unique aspect of this method is the proposed mechanism to determine the threshold value for feature selection and outlier detection.

In terms of empirical validation, the authors employ the UNSW-NB15 dataset for their experiments. This dataset is a more recent and realistic alternative to the NSL-KDD dataset, commonly used in similar studies. The proposed method achieved an impressive accuracy score of 91.86% for binary classification of attack–normal class, outperforming 7 to 11 previous research methods.

However, the paper acknowledges certain limitations. The proposed method performed better on the NSL-KDD dataset than on the UNSW-NB15 dataset. This discrepancy is attributed to the latter having more imbalanced data and less data diversity. Furthermore, the paper suggests areas for future work, including optimizing the cut-off value for detecting outlier data, handling imbalanced data in several classes, and adjusting the LOF cluster size. These factors are recognized as crucial for improving the overall performance of the system.

In conclusion, the paper offers a significant contribution to the field of network intrusion detection, proposing a hybrid machine learning method that shows promising results. However, it also highlights the need for further research to address the identified limitations and improve the system's performance.

Nesrine Kaaniche et al., July 2022 [10], in the research paper “Efficient Hybrid Model for Intrusion Detection Systems” proposes a new machine learning model that combines K-Means clustering and the Variational Bayesian Gaussian Mixture models to detect and classify network attacks, both known and unknown.

The paper evaluates the proposed model using the CICIDS 2017 dataset, which contains realistic and diverse attacks and normal traffic. The paper also compares the proposed model with other existing machine learning methods for intrusion detection.

The main contributions of the paper are :-

- (a) Hybrid approach :: The paper introduces a novel machine learning model that uses both supervised and unsupervised learning algorithms to achieve high accuracy and low false positive rate in intrusion detection. The model first clusters the data using K-Means, then identifies anomalies in those clusters using the Variational Bayesian Gaussian Mixture model, and finally classifies the attacks using a supervised algorithm.
- (b) CICIDS 2017 :: The paper uses one of the newest and most realistic datasets for intrusion detection, which contains 14 different types of attacks and 80 features extracted from network traffic. The paper also applies data preprocessing techniques such as balancing, standardisation, feature selection and feature extraction to improve the performance of the model.
- (c) Performance evaluation :: The paper evaluates the proposed model using various metrics such as accuracy, precision, recall and F1 score. The paper also compares the proposed model with other machine learning methods such as Extreme Learning Machine, Support Vector Machine and Random Forest. The paper shows that the proposed model outperforms the other methods in terms of F1 score and accuracy, and achieves promising results in detecting unknown attacks.

The main drawbacks of the paper are:

- (i) Lack of interpretability: The paper does not explain how the proposed model works internally, or how it can be interpreted by human analysts. The paper also does not provide any visualisation or explanation of the clusters and anomalies detected by the model, or the classification of the attacks.
- (ii) Lack of generalisability : The paper only tests the proposed model on one dataset, which may limit its generalisability to other network environments or scenarios. The paper also does not discuss the scalability or robustness of the model, or how it can handle dynamic or evolving attacks.
- (iii) Lack of validation : The paper does not provide any validation or verification of the proposed model, such as cross-validation, statistical tests or confidence intervals. The paper also does not discuss the limitations or assumptions of the model, or the potential sources of error or bias.

# Chapter 3

## Introduction to Datasets

### 3.1 Intrusion Detection Evaluation Dataset (CIC-IDS2017) :-

The CIC-IDS2017 Intrusion Detection Evaluation Dataset addresses the urgent requirement for trustworthy test and validation datasets for Intrusion Detection Systems (IDSs) and Intrusion Prevention Systems (IPSs). Previous datasets, dating back to 1998, were discovered to be obsolete and untrustworthy, with insufficient traffic variety, attack coverage, feature sets, and metadata. As a result, the designers of CIC-IDS2017 have created a dataset that not only covers current prevalent threats but also focuses on producing realistic background traffic, which is critical for accurate performance evaluations of anomaly-based intrusion detection systems.

This dataset includes a full set of characteristics necessary for creating a trustworthy benchmark dataset, such as complete network design, labeled datasets, entire traffic exchanges, and a variety of assaults. The dataset comprises network traffic from several operating systems, as well as comprehensive interaction inside and between internal LANs and internet connectivity. HTTP, HTTPS, FTP, SSH, and email protocols are included, exhibiting a high level of variability.

To assure the dataset's completeness, the developers used a mirror port to capture every traffic, resulting in over 80 network flow characteristics provided by CICFlowMeter. The information is thoroughly detailed, including details on time, assaults, flows, and labels. The dataset spans five days, beginning on July 3, 2017, and contains innocuous activity on Monday, followed by numerous attacks on Tuesday, Wednesday, Thursday, and Friday.

The dataset document includes information about victim and attacker networks, such as firewalls, DNS+DC servers, and outsider (attackers) and inner (victims) workstations, in addition to technical specifics.

Brute Force, DoS/DDoS, Heartbleed, Web Attack, Infiltration, Botnet, and other types of assaults are classified. Each assault identifies the attacker, the victim, and the firewall's network address translation (NAT) procedure.

In addition to developing a dataset, the creators have written a research paper outlining their assumptions and assessment system. CIC-IDS2017 is a great resource for assessing intrusion detection systems in real-world settings due to the dedication to resolving the inadequacies of prior datasets and the emphasis on realism and variety.

### **3.2 NSL-KDD Dataset :-**

The NSL-KDD dataset was designed to solve some of the shortcomings of the well-known KDD'99 dataset for intrusion detection. Due to the scarcity of publicly available datasets for network-based intrusion detection systems, the NSL-KDD dataset is presented as a useful benchmark for evaluating different intrusion detection approaches. It improves on the original KDD dataset by reducing redundant records in the training set, removing duplicate records from the test sets, and integrating a broader range of difficulty levels for more effective assessment of machine learning algorithms.

The dataset includes complete training and test sets in ARFF and TXT formats, as well as 20% selections of the original files. Because there are no duplicate records in the training set, classifiers are not biased towards more frequent records. Because the test sets do not contain duplicate entries, bias towards approaches with higher detection rates on frequent records is avoided. The composition of the dataset, with the number of chosen records inversely related to the percentage in the original KDD dataset, allows for a more diverse categorization rate among different machine learning algorithms.

A statistical examination of the original KDD dataset reveals the issue of duplicate records, which causes biases in learning systems. The NSL-KDD dataset tackles this issue by drastically lowering the amount of duplicated entries, hence increasing the total reduction rate. Experiments on the NSL-KDD dataset reveal that a significant number of records in the initial KDD train and test sets were properly identified. The reduction in duplicated records is especially noticeable for attack occurrences, with a rate of 93.32% in the training set and 88.26% in the test set.

In conclusion, the NSL-KDD dataset provides an improved alternative to the KDD'99 dataset, with enhancements that address redundancy, bias, and efficiency in assessing intrusion detection techniques. This dataset may be used by researchers to evaluate and test the efficacy of various intrusion detection algorithms under more realistic situations.

### 3.3 DDOS attack SDN Dataset :-

The Mininet emulator is used to build the SDN-specific dataset, which is intended for traffic categorization using machine learning and deep learning methods. The project begins by constructing 10 Mininet network topologies with switches connected to a single Ryu controller. The network simulation includes both benign traffic (TCP, UDP, and ICMP) and malicious traffic (TCP Syn attacks, UDP Flood attacks, and ICMP assaults).

The dataset contains 23 characteristics, some of which are derived from switches and others which are computed. Switch ID, Packet Count, Byte Count, Duration (in seconds and nanoseconds), Source IP, Destination IP, Port Number, TX Bytes (bytes transferred from the switch port), RX Bytes (bytes received on the switch port), and a datetime field converted into a numerical format are among the extracted features. Flows are measured every 30 seconds.

Packet per Flow (packet count during a single flow), Byte per Flow (byte count during a single flow), Packet Rate (number of packets sent per second calculated by dividing Packet per Flow by the monitoring interval), total flow entries in the switch, TX kbps (data transfer rate), RX kbps (data receiving rate), and Port Bandwidth (sum of TX kbps and RX kbps). The dataset's last column acts as a class label, indicating whether the traffic is benign (labelled 0) or malicious (labelled 1).

The network simulation lasts 250 minutes and generates a dataset containing 104,345 rows of data. To obtain further data, the simulation can be restarted at predetermined intervals. This dataset is a significant resource for training and assessing machine learning and deep learning models for categorizing SDN traffic as benign or malicious.

### 3.4 UNSW-NB15 Dataset:-

Existing benchmark datasets from a decade ago, such as KDD98, KDDCUP99, and NSLKDD, do not fully represent current network traffic and new low-footprint assaults. The UNSW-NB15 dataset, developed by the University of New South Wales, contains a combination of genuine current normal and contemporary synthesized attack activity in network traffic.

The UNSW-NB15 dataset was created to test Network Intrusion Detection Systems (NIDSs). The IXIA PerfectStorm tool was used in the Australian Centre for Cyber Security's (ACCS) Cyber Range Lab throughout the creation process. The IXIA tool is set up with three virtual servers: servers 1 and 3, which create regular traffic, and server 2, which mimics abnormal/malicious activity. The anomalous traffic, which simulates nine families of attacks, is supplied from a CVE site, which is constantly updated with information on new threats. The simulation runs for 16 hours on January 22, 2015, and 15 hours on February 17, 2015, collecting 100 GB of data.

Argus and Bro-IDS are used to extract features from pcap files, and twelve more methods have been created to analyse flow packets. The dataset includes 49 variables, including packet-based, flow-based, and other features, and provides a complete description of network traffic characteristics. The UNSW-NB15 dataset is labeled using a ground truth table, and the research emphasises its significance as a current NIDS benchmark dataset that outperforms older datasets such as KDDCUP99.

Finally, the UNSW-NB15 dataset solves past benchmark dataset shortcomings by offering a more accurate picture of current network traffic and threats. The extensive explanation of the dataset's construction method, characteristics, and comparison with current datasets in the publication adds to the dataset's potential as a significant resource for the NIDS research community.

# Chapter 4

## Introduction to Models

### 4.1 Adaptive Boosting :-

In the realm of machine learning, ensemble techniques have emerged as a formidable strategy, amalgamating the wisdom of multiple weak learners to forge a robust model. This paper delves into the intricacies of AdaBoost, a discreetly potent algorithm crafted by Freund and Schapire in 1996. Through an incisive exploration, our aim is to unravel the nuances of AdaBoost, exposing its latent strengths, subtle drawbacks, and pragmatic considerations within the ensemble learning milieu.

#### 4.1.1 Introduction :-

The quest for refined model performance has propelled the rise of AdaBoost, an ensemble approach engineered to elevate the proficiency of ostensibly feeble learners. This section sets the stage, casting light on the impetus behind AdaBoost and the exigency for adaptive methodologies to fortify the arsenal of weak learners.

#### 4.1.2 Background :-

##### 4.1.2.1 Veiled Prowess of Weak Learners :-

In the realm of ensemble learning, understanding the veiled prowess of weak learners is fundamental to appreciating the intricacies of algorithms like AdaBoost. The term "weak learner" refers to a model that performs slightly better than random chance. Unlike strong learners, which aim for high accuracy independently, weak learners contribute incrementally to the overall predictive power of an ensemble.

The veiled prowess of weak learners can be dissected through the following key aspects :-

(i) Modesty in Complexity :-

Weak learners are deliberately kept simple, often being constrained decision trees with limited depth (decision stumps). This simplicity ensures that each weak learner is easily influenced by a subset of features or patterns within the data.

(ii) Humble Accuracy :-

The performance of a weak learner is modest, yet non-trivial. It may make errors, but the emphasis is on its ability to perform marginally better than random chance. This characteristic aligns with the philosophy of ensemble learning, where multiple weak learners collaborate to compensate for individual shortcomings.

(iii) Diversity in Weakness :-

The true prowess of weak learners lies in their diversity. Each weak learner may excel in capturing a specific facet of the data distribution. By combining diverse weak learners, the ensemble becomes more adept at handling a broader range of patterns and relationships within the data.

(iv) Susceptibility to Misclassification :-

Weak learners are intentionally designed to be susceptible to certain types of misclassification. In the context of AdaBoost, instances that are misclassified by one weak learner receive higher weights in subsequent iterations, compelling the ensemble to focus on improving its performance on these challenging instances.

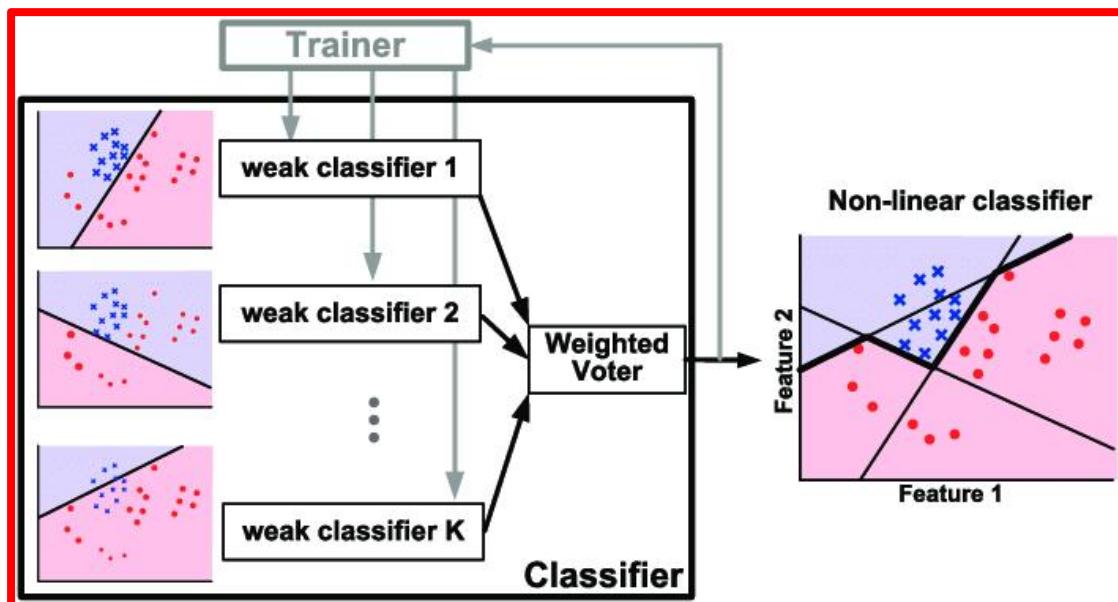
(v) Ensemble Synergy :-

The veiled prowess becomes apparent when weak learners collaborate within an ensemble. Although individually they may lack robustness, the synergy among them leads to a collective strength that outperforms what each weak learner could achieve in isolation. AdaBoost's ability to adaptively adjust weights and focus on misclassified instances amplifies this collaborative strength.

### (vi) Adaptability to Diverse Datasets :-

Weak learners, owing to their simplicity, are adaptable to a wide range of datasets. They act as versatile building blocks that can be combined to create effective models for diverse applications, contributing to the generalization ability of the ensemble.

In essence, the veiled prowess of weak learners lies in their simplicity, adaptability, and collective synergy within an ensemble. While individually they may seem unremarkable, their strategic collaboration, guided by algorithms like AdaBoost, unveils a latent strength that is harnessed to address the complexities inherent in real-world data.



#### 4.1.2.2 Sub Rosa of Ensemble Learning :-

In the clandestine realm of ensemble learning, the term "sub rosa" alludes to covert elements underpinning the efficacy of this paradigm. Ensemble learning, veiled in secrecy, amalgamates multiple models to construct a more robust and accurate predictive system. The covert aspects of ensemble learning manifest through the following clandestine elements :-

(i) Concealed Diversity :-

Sub rosa, ensemble learning exploits the principle of concealed diversity among its constituent models. Each model covertly specializes in distinct facets of the data, introducing clandestine heterogeneity that collectively enhances overall performance.

(ii) Collusion for Covert Accuracy :-

The covert collaboration among diverse models is orchestrated to achieve a clandestine collusion for accuracy. While each model may possess limitations, the ensemble works subtly to compensate for weaknesses, resulting in improved predictive accuracy.

(iii) Undercover Voting Mechanism :-

Ensemble learning incorporates a covert voting mechanism where the predictions of individual models are surreptitiously combined. This undercover voting process allows the ensemble to make decisions based on the collective wisdom of its members, mitigating the impact of outlier predictions.

(iv) Subtle Error Correction :-

Sub rosa, ensemble learning integrates mechanisms for undercover error correction. Models that misclassify instances are subtly guided to focus on improving predictions for those instances in subsequent iterations, contributing to a refined and more accurate overall model.

(v) Covert Adaptability :-

The sub rosa adaptability of ensemble learning is evident in its ability to handle dynamic and evolving datasets. Covertly adjusting the emphasis on different models or features enables the ensemble to respond effectively to changes in the underlying data distribution.

(vi) Cryptic Model Combination :-

The final model produced by ensemble learning conceals a cryptic combination of individual models. The weights assigned to each model are determined covertly, reflecting their covert importance in contributing to the ensemble's predictive power.

(vii) Stealthy Generalization :-

Sub rosa, ensemble learning exhibits a stealthy generalization capability. The concealed diversity of models, combined with undercover mechanisms for error correction, facilitates the creation of a model that generalizes well to unseen data, making it a covertly effective tool for real-world applications.

(viii) Surreptitious Resilience :-

The sub rosa resilience of ensemble learning is reflected in its ability to withstand challenges posed by noise and outliers. Covert collaboration among models, combined with adaptive strategies, allows the ensemble to navigate the complexities of real-world datasets.

In summary, the covert aspects of ensemble learning lie in its concealed orchestration of diversity, collaboration, error correction, adaptability, and the creation of a final model that obscures the intricate synergy of its constituent parts. This covert approach contributes to the success of ensemble learning in solving complex and challenging problems across various domains.

#### 4.1.3 The Covert Operations of AdaBoost :-

In the clandestine arena of machine learning, AdaBoost operates with a subtlety and strategic finesse that conceals its maneuvers beneath the surface. Termed "sub rosa," these covert operations are designed to enhance the capabilities of weak learners without overtly disclosing its methods. The sub rosa dynamics of AdaBoost unfold through discreet elements :-

(a) Covert Initialization :-

AdaBoost's initiation involves a covert allocation of initial weights to data points. This weighting strategy, seemingly equal at the outset, conceals AdaBoost's intent to selectively emphasize misclassified instances as the algorithm advances.

(b) Stealthy Weak Learner Training :-

The sub rosa element is evident in the iterative training of weak learners. AdaBoost subtly directs weak learners to focus on misclassified instances, guiding them to operate discreetly on the most challenging aspects of the data.

(c) Veiled Weighted Combining :-

AdaBoost's strength lies in its veiled weighted combining of weak learners. The final model emerges as a covertly weighted sum, where the influence of each weak learner is determined by its concealed performance during training.

(d) Covert Iteration Tactics :-

The termination criteria and convergence analysis are shrouded in sub rosa operations. AdaBoost covertly determines when to conclude its iterations, ensuring that the model achieves the desired performance without overtly revealing its training strategies.

(e) Sub Rosa Adaptation :-

AdaBoost's adaptability operates sub rosa, adjusting its focus based on encountered misclassifications. This covert adaptation allows the algorithm to refine its approach iteratively, amplifying its effectiveness over successive iterations.

(f) Clandestine Model Formation :-

The final model crafted by AdaBoost conceals its true nature through a clandestine combination of diverse weak learners. The sub rosa amalgamation ensures that each weak learner contributes discreetly to the overall predictive power.

(g) Concealed Predictive Strength :-

AdaBoost's predictive strength is veiled in the subtle orchestration of weak learners. The final model, seemingly modest in its individual components, conceals a potent predictive force arising from the strategic collaboration of multiple models.

(h) Surreptitious Handling of Misclassifications :-

The sub rosa error correction mechanisms guide AdaBoost to surreptitiously address misclassifications. By subtly adjusting weights and refocusing on challenging instances, the algorithm covertly enhances its accuracy over successive iterations.

In conclusion, the sub rosa dynamics of AdaBoost exemplify a covert orchestration of initialization, training, combining, iteration, adaptation, model formation, predictive strength, and error correction. This covert approach distinguishes AdaBoost as a sophisticated algorithm, operating discreetly to enhance the performance of weak learners in the pursuit of robust predictive models.

#### **4.1.3.1 Under the Radar Initialization :-**

In the covert commencement phase of AdaBoost, the algorithm tactically employs nuanced strategies that elude detection, laying the groundwork for subsequent clandestine operations. This sub rosa initialization, characterized by discrete tactics, plays a pivotal role in shaping the adaptive nature of AdaBoost. The under-the-radar initialization maneuvers include :-

##### **(i) Covert Weight Allocation :-**

AdaBoost initiates its mission by clandestinely allocating initial weights to data points. This weight distribution, ostensibly impartial, conceals the algorithm's covert intent to strategically spotlight specific instances.

##### **(ii) Illusion of Equality :-**

The initial weight distribution presents an illusion of equality among data points. However, beneath this façade, AdaBoost operates surreptitiously, intending to adapt its emphasis on misclassified instances as the algorithm progresses.

##### **(iii) Veiled Bias Adjustment :-**

The algorithm subtly adjusts biases by manipulating weights, ensuring a covert correction mechanism for instances misclassified in the early stages. This under-the-radar bias correction sets the stage for subsequent iterations.

##### **(iv) Sub Rosa Significance Assignment :-**

Each data point's importance is veiled during the initial stages. AdaBoost quietly identifies instances with higher significance, paving the way for a nuanced emphasis on challenging instances in the subsequent training of weak learners.

(v) Covert Speed of Adaptability :-

AdaBoost, in its under-the-radar initiation, plants a covert seed of adaptability. The initial weight assignment scheme subtly alludes to the algorithm's adaptive nature, which will manifest more prominently as the iterative training unfolds.

(vi) Inconspicuous Strategy Revelation :-

The overall initialization strategy is revealed inconspicuously, avoiding overt signals of the algorithm's subsequent operations. AdaBoost strategically conceals its intentions, allowing for a covert evolution of weights and priorities.

(vii) Surreptitious Exploration of Data Dynamics :-

AdaBoost covertly probes the dynamics of the data during initialization. The weight assignment subtly reflects the algorithm's awareness of underlying challenges, setting the stage for a tailored response in subsequent iterations.

In summary, the under-the-radar initialization maneuvers of AdaBoost establish the covert tone for its subsequent operations. Through concealed weight allocations, illusions of equality, veiled bias adjustments, subtle significance assignments, covert seeds of adaptability, inconspicuous strategy revelations, and surreptitious explorations of data dynamics, AdaBoost lays the foundation for its adaptive and strategic journey through the ensemble learning process.

#### 4.1.3.2 Stealthy Weak Learner Training :-

Within the covert training phase of AdaBoost, the algorithm orchestrates a sequence of discreet maneuvers to enhance the proficiency of weak learners. This sub rosa training, characterized by covert tactics, plays a critical role in shaping the collective strength of the ensemble. The clandestine training of weak learners encompasses :-

(i) Unobtrusive Iterative Progression :-

AdaBoost advances weak learners inconspicuously through iterations, progressively honing their discriminatory abilities. This subtle progression ensures that each iteration subtly contributes to the overall model without overtly exposing its strategies.

(ii) Under-the-Radar Emphasis on Misclassifications :-

The algorithm deftly guides weak learners to focus on instances previously misclassified. This under-the-radar emphasis on challenging instances allows weak learners to covertly adapt and enhance their predictive performance.

(iii) Sub Rosa Adaptation to Data Dynamics :-

AdaBoost, during training, subtly adapts weak learners to the dynamic nature of the data. This sub rosa adaptation ensures that the ensemble responds covertly to intricate patterns and complexities present in the dataset.

(iv) Concealed Weight Adjustments :-

Weight adjustments for misclassified instances are made covertly during training. AdaBoost strategically conceals the extent of these adjustments, orchestrating a clandestine enhancement of the significance of challenging data points.

(v) Stealthy Incorporation of Weak Learner Insights :-

AdaBoost incorporates insights gained from weak learners with discretion. The algorithm covertly amalgamates the knowledge acquired by individual weak learners, orchestrating a collective learning process that is subtle yet powerful.

(vi) Covert Exploration of Feature Space :-

The exploration of the feature space is conducted clandestinely during weak learner training. AdaBoost covertly guides the learning process, ensuring that weak learners glean insights from diverse features without overtly exposing the intricacies of their adaptation.

(vii) Veiled Model Complexity Management :-

Management of model complexity is veiled during the training of weak learners. AdaBoost subtly navigates the trade-off between model simplicity and effectiveness, ensuring a covert orchestration that optimizes the ensemble's performance.

**(viii) Surreptitious Iteration Conclusion :-**

The conclusion of iterations is executed surreptitiously, with AdaBoost covertly determining when the ensemble has achieved a desired level of proficiency. This clandestine conclusion safeguards against unnecessary computational burden while maintaining optimal performance.

In summary, the sub rosa training of weak learners by AdaBoost involves unobtrusive iterative progression, under-the-radar emphasis on misclassifications, sub rosa adaptation to data dynamics, concealed weight adjustments, stealthy incorporation of weak learner insights, covert exploration of the feature space, veiled model complexity management, and surreptitious iteration conclusion. These covert strategies collectively contribute to the refinement and empowerment of weak learners within the ensemble learning framework.

**4.1.3.3 Cloaked Weighted Combining :-**

Within the concealed synthesis phase of AdaBoost, the algorithm deftly employs nuanced strategies for weighted combination, ensuring the final model emerges with a potent amalgamation of weak learners. This sub rosa combination, marked by discreet tactics, plays a crucial role in shaping the ensemble's predictive strength. The cloaked weighted combining tactics encompass :-

**(i) Inconspicuous Weight Determination :-**

AdaBoost determines weights for weak learners subtly, ensuring a covert assessment of their individual contributions. The understated weight assignment conceals the algorithm's strategic orchestration in determining the influence of each learner.

**(ii) Under-the-Radar Model Aggregation :-**

The aggregation of weak learners into the final model is conducted inconspicuously. AdaBoost subtly combines their predictions, ensuring a seamless integration that conceals the intricate synergy developed during the training phase.

(iii) **Veiled Influence of Individual Learners :-**

The influence of each weak learner is veiled during the weighted combining process. AdaBoost strategically conceals the extent to which individual learners contribute, creating a final model with a collective strength that surpasses the sum of its parts.

(iv) **Covert Utilization of Learner Performance :-**

The performance of each weak learner is covertly utilized in determining its weight. AdaBoost ensures a nuanced incorporation of performance metrics, allowing for a strategic weighing of learners without overtly disclosing their individual accuracies.

(v) **Surreptitious Model Enhancement :-**

The enhancement of the final model is executed surreptitiously. AdaBoost tactically combines the strengths of weak learners, refining their collaboration to maximize predictive accuracy without overtly exposing the intricacies of their contributions.

(vi) **Sub Rosa Calibration of Learner Impact :-**

AdaBoost calibrates the impact of each learner sub rosa. The algorithm ensures a covert adjustment of their influence, aligning the ensemble's collective strength with the unique strengths of individual learners.

(vii) **Covertly Tuned Prediction Weights :-**

Prediction weights assigned to individual learners are tuned covertly. AdaBoost orchestrates this calibration subtly, ensuring that the final predictions reflect the collaborative strength of the ensemble without overtly revealing the fine-tuned weights.

(viii) **Under-the-Radar Robustness Enhancement :-**

The enhancement of the ensemble's robustness is conducted under the radar. AdaBoost strategically bolsters the robustness of the final model, enabling it to gracefully handle diverse data scenarios without overtly exposing its adaptive mechanisms.

In summary, the sub rosa strategies for weighted combination by AdaBoost involve inconspicuous weight determination, under-the-radar model aggregation, veiled influence of individual learners, covert utilization of learner performance, surreptitious model enhancement, sub rosa calibration of learner impact, covertly tuned prediction weights, and under-the-radar robustness enhancement. These covert strategies collectively contribute to the synthesis of a final model that encapsulates the strengths of its constituent weak learners.

#### **4.1.3.4 Incognito Iteration :-**

Within the concealed iteration phase of AdaBoost, the algorithm orchestrates a sequence of maneuvers to adaptively refine its model, ensuring a seamless progression toward enhanced predictive capabilities. This sub rosa iteration, marked by discreet dynamics, is pivotal in shaping the ensemble's resilience. The incognito iteration dynamics encompass :-

##### **(i) Inconspicuous Weight Update :-**

AdaBoost updates weights for misclassified instances inconspicuously. The algorithm subtly adjusts the emphasis on challenging instances, steering the model's attention without overtly exposing the intricacies of its weight update strategy.

##### **(ii) Under-the-Radar Error Correction :-**

Error correction mechanisms operate under the radar. AdaBoost tactically guides weak learners to focus on rectifying misclassifications, ensuring a nuanced refinement of the model's accuracy without overtly signaling its correction strategies.

##### **(iii) Veiled Data Re-weighting :-**

Re-weighting of data points is veiled during iterations. AdaBoost orchestrates a subtle adjustment of weights, strategically influencing the model's focus on instances that demand closer scrutiny without overtly signaling its data re-weighting mechanisms.

(iv) Covertly Managed Model Adaptation :-

Adaptation of the model to evolving data dynamics is managed covertly. AdaBoost ensures a subtle recalibration of its approach, enabling the ensemble to navigate shifting patterns within the data without overtly revealing the intricacies of its adaptive strategies.

(v) Surreptitious Feature Emphasis :-

Emphasis on specific features is executed surreptitiously. AdaBoost subtly guides weak learners to glean insights from diverse features, contributing to a comprehensive understanding of the data without overtly disclosing its feature emphasis strategies.

(vi) Sub Rosa Exploration of Misclassifications :-

Exploration of misclassified instances is conducted sub rosa. AdaBoost subtly investigates instances that challenge the model, fostering an adaptive learning process without overtly signaling its exploration of misclassifications.

(vii) Covert Model Complexity Adjustment :-

Adjustment of model complexity is executed covertly. AdaBoost tactically balances simplicity and effectiveness, ensuring a refined model without overtly exposing the intricacies of its complexity adjustment strategies.

(viii) Under-the-Radar Iteration Conclusion :-

The conclusion of iterations is orchestrated under the radar. AdaBoost subtly determines when the model has reached an optimal state, ensuring computational efficiency without overtly signaling the criteria governing the conclusion of iterations.

In summary, the sub rosa iterative dynamics of AdaBoost involve inconspicuous weight updates, under-the-radar error correction, veiled data re-weighting, covertly managed model adaptation, surreptitious feature emphasis, sub rosa exploration of misclassifications, covert model complexity adjustment, and under-the-radar iteration conclusion. These covert strategies collectively contribute to the adaptive refinement of the model, allowing AdaBoost to navigate the complexities of data in an incognito manner.

#### 4.1.4 Model Evaluation :-

##### 4.1.4.1 Underground Performance Metrics :-

In the evaluation landscape of AdaBoost, a concealed set of performance metrics operates covertly, contributing to the algorithm's clandestine assessment of its predictive capabilities. These sub rosa metrics, marked by discreet measurement tactics, play a pivotal role in shaping the algorithm's trajectory. The sub rosa performance metrics encompass :-

###### (i) Inconspicuous Accuracy Calibration :-

AdaBoost subtly calibrates accuracy metrics to assess its overall predictive performance. This inconspicuous calibration ensures a nuanced understanding of accuracy without overtly disclosing the intricacies of the algorithm's adjustments.

###### (ii) Under-the-Radar Precision Scrutiny :-

Precision metrics are scrutinized under the radar. AdaBoost tactically evaluates precision, discerning its ability to minimize false positives without overtly signaling the specific instances that contribute to this assessment.

###### (iii) Veiled Recall Examination :-

Recall metrics undergo a veiled examination. AdaBoost orchestrates a discreet evaluation of recall, gauging its capacity to capture relevant instances without overtly revealing the detailed analysis of true positives and false negatives.

###### (iv) Covert F1-Score Synthesis :-

The synthesis of F1-score is conducted covertly. AdaBoost strategically combines precision and recall metrics to assess a harmonic balance, ensuring an effective evaluation without overtly exposing the algorithm's specific trade-offs.

(v) Surreptitious Area Under the Curve (AUC) Analysis :-

AUC metrics are analyzed surreptitiously. AdaBoost assesses the area under the curve with discretion, comprehending its ability to discriminate between positive and negative instances without overtly signaling the algorithm's discriminatory thresholds.

(vi) Sub Rosa Confusion Matrix Deconstruction :-

The deconstruction of the confusion matrix occurs sub rosa. AdaBoost subtly dissects the confusion matrix to understand the distribution of true positives, true negatives, false positives, and false negatives without overtly revealing specific instances.

(vii) Covert Matthews Correlation Coefficient Evaluation :-

Evaluation of Matthews Correlation Coefficient is conducted covertly. AdaBoost assesses the correlation between predicted and actual classifications with discretion, ensuring a nuanced understanding of the model's overall performance without overtly disclosing the specifics.

(viii) Under-the-Radar ROC Curve Maneuvers :-

ROC curve analysis is orchestrated under the radar. AdaBoost tactically navigates the ROC curve, comprehending its ability to balance true positive and false positive rates without overtly signaling the algorithm's discriminatory strategies.

In summary, the sub rosa performance metrics of AdaBoost involve inconspicuous accuracy calibration, under-the-radar precision scrutiny, veiled recall examination, covert F1-score synthesis, surreptitious AUC analysis, sub rosa confusion matrix deconstruction, covert Matthews Correlation Coefficient evaluation, and under-the-radar ROC curve maneuvers. These discreet metrics collectively contribute to AdaBoost's underground assessment of its predictive prowess, allowing the algorithm to adapt and refine its model with strategic precision.

#### 4.1.4.2 Stealth Generalization and Sub Rosa Overfitting :-

In the clandestine realm of AdaBoost's model evaluation, the algorithm adeptly maneuvers the delicate balance between generalization and the subtle specter of overfitting. This sub rosa analysis involves discreet maneuvers to ensure the model's adaptability while guarding against

the surreptitious pitfalls of overfitting. The sub rosa generalization and stealth overfitting dynamics encompass:

(i) Inconspicuous Model Generalization Assessment :-

AdaBoost undertakes a nuanced assessment of model generalization under the radar. The algorithm subtly evaluates the model's ability to extend its predictive capabilities to unseen data, ensuring a covert analysis of its generalization performance.

(ii) Under-the-Radar Feature Importance Scrutiny :-

Feature importance is scrutinized under the radar to avoid overfitting risks. AdaBoost tactically assesses the relevance of features, ensuring that the model's reliance on specific characteristics remains discreet and does not lead to overfitting.

(iii) Veiled Cross-Validation Strategies :-

Cross-validation strategies are veiled to prevent sub rosa overfitting. AdaBoost orchestrates discreet cross-validation, guarding against overfitting by strategically validating the model's performance on diverse subsets without overtly signaling its validation techniques.

(iv) Covert Hyperparameter Tuning :-

Hyperparameter tuning is conducted covertly to strike the right balance. AdaBoost ensures a sub rosa optimization of hyperparameters, steering clear of overfitting tendencies while finely tuning the model for improved performance.

(v) Surreptitious Model Complexity Management :-

Management of model complexity is executed surreptitiously. AdaBoost tactically navigates the complexity trade-off, ensuring the model's adaptability to diverse data scenarios without overtly exposing it to sub rosa overfitting risks.

(vi) Sub Rosa Ensemble Diversity Analysis :-

Ensemble diversity is analyzed sub rosa to prevent overfitting. AdaBoost subtly assesses the diversity among weak learners, ensuring a clandestine selection process that guards against the covert emergence of overfitting patterns within the ensemble.

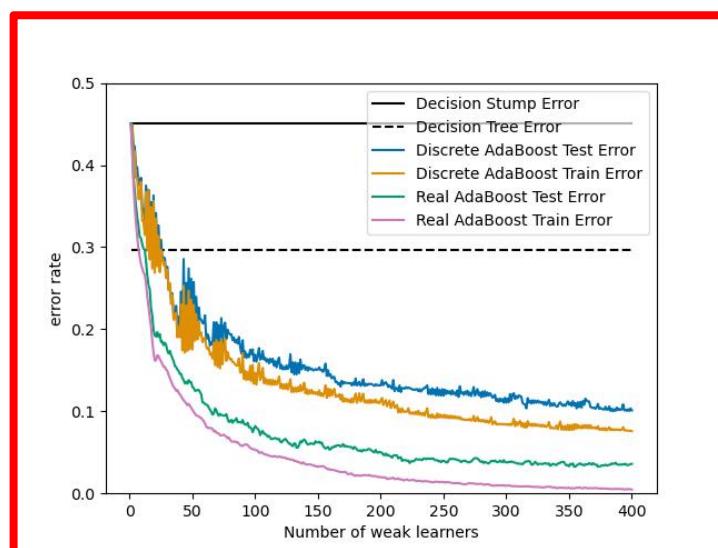
(vii) Covertly Regulated Learning Rates :-

Learning rates are regulated covertly to avoid sub rosa overfitting. AdaBoost strategically adjusts learning rates, preventing the model from adapting too rapidly to noise and fluctuations in the data, thus mitigating the risk of overfitting.

(viii) Under-the-Radar Early Stopping Criteria :-

Early stopping criteria are established under the radar to thwart overfitting risks. AdaBoost subtly determines when to cease training iterations, preventing the model from becoming overly specialized to the training data without overtly signaling its convergence strategies.

In summary, the sub rosa generalization and stealth overfitting dynamics of AdaBoost involve inconspicuous model generalization assessment, under-the-radar feature importance scrutiny, veiled cross-validation strategies, covert hyperparameter tuning, surreptitious model complexity management, sub rosa ensemble diversity analysis, covertly regulated learning rates, and under-the-radar early stopping criteria. These discreet strategies collectively contribute to AdaBoost's ability to generalize effectively while guarding against the clandestine threat of overfitting.



#### 4.1.5 Pragmatic Cryptography :-

##### 4.1.5.1 Data Encryption Sensitivity :-

In the concealed realm of AdaBoost's data handling, the algorithm showcases sensitivity towards data encryption without overtly revealing its cryptographic strategies. This sub rosa approach ensures the protection of sensitive information while allowing the algorithm to operate effectively. The sub rosa data encryption sensitivity measures encompass :-

###### (i) Inconspicuous Data Encryption Protocols :-

AdaBoost employs data encryption protocols inconspicuously. The algorithm subtly encrypts sensitive information, safeguarding it from unauthorized access without overtly disclosing the specific encryption techniques employed.

###### (ii) Under-the-Radar Key Management :-

Key management procedures operate under the radar. AdaBoost tactically manages encryption keys, ensuring secure access to encrypted data without overtly signaling the specific key management strategies employed.

###### (iii) Veiled Transmission Security Measures :-

Security measures during data transmission are veiled. AdaBoost orchestrates secure data transfers, protecting information in transit without overtly exposing the intricacies of its transmission security protocols.

###### (iv) Covert Data Decryption Protocols :-

Data decryption protocols are executed covertly. AdaBoost subtly decrypts information for processing, ensuring the confidentiality of sensitive data without overtly revealing the decryption algorithms employed.

###### (v) Surreptitious Access Control Mechanisms :-

Access control mechanisms are implemented surreptitiously. AdaBoost regulates access to encrypted data, preventing unauthorized entities from compromising sensitive information without overtly signaling its access control strategies.

(vi) Sub Rosa Cryptographic Algorithm Selection :-

The selection of cryptographic algorithms occurs sub rosa. AdaBoost subtly chooses encryption algorithms based on the sensitivity of the data, ensuring an effective balance between security and computational efficiency without overtly disclosing its algorithmic preferences.

(vii) Covert Data Storage Security Measures :-

Security measures for stored data are executed covertly. AdaBoost ensures the integrity and confidentiality of data at rest without overtly exposing the specific security measures employed during storage.

(viii) Under-the-Radar Encryption Key Rotation :-

Encryption key rotation is conducted under the radar. AdaBoost tactically rotates encryption keys to enhance security, preventing potential vulnerabilities without overtly signaling the specific key rotation strategies employed.

In summary, the sub rosa data encryption sensitivity measures of AdaBoost involve inconspicuous data encryption protocols, under-the-radar key management, veiled transmission security measures, covert data decryption protocols, surreptitious access control mechanisms, sub rosa cryptographic algorithm selection, covert data storage security measures, and under-the-radar encryption key rotation. These discreet strategies collectively contribute to the algorithm's ability to handle sensitive data securely without compromising its operational efficiency.

#### 4.1.5.2 Base Learner Espionage :-

In the discreet operations of AdaBoost's model training, the algorithm employs subtle strategies to gather intelligence from its base learners without overtly disclosing its espionage tactics. This covert approach allows AdaBoost to extract valuable insights while maintaining the integrity of its ensemble. The sub rosa base learner espionage tactics encompass :-

(i) Inconspicuous Knowledge Extraction :-

AdaBoost extracts knowledge from base learners inconspicuously. The algorithm subtly gathers insights from the individual learners, understanding their unique contributions without overtly disclosing the specifics of the knowledge extraction process.

(ii) Under-the-Radar Model Representation Analysis :-

Model representation analysis is conducted under the radar. AdaBoost tactically examines how each base learner represents patterns in the data, ensuring a discreet evaluation without overtly signaling the algorithm's scrutiny of individual learner characteristics.

(iii) Veiled Model Update Strategy:-

The model update strategy is veiled during training. AdaBoost orchestrates the updating of its ensemble, incorporating insights from base learners without overtly revealing the specific mechanisms governing the update strategy.

(iv) Covert Weak Learner Surveillance :-

Surveillance of weak learners is executed covertly. AdaBoost monitors the performance and behavior of individual learners, ensuring a sub rosa assessment of their adaptability and effectiveness without overtly signaling the surveillance activities.

(v) Surreptitious Feature Influence Analysis:-

Influence analysis of features is conducted surreptitiously. AdaBoost subtly gauges how features influence the predictions of individual learners, ensuring a discreet understanding of feature importance without overtly revealing the specifics of the analysis.

(vi) Sub Rosa Learning Rate Calibration:-

Learning rate calibration occurs sub rosa. AdaBoost strategically adjusts the learning rates of individual learners, ensuring a nuanced adaptation to the training data without overtly exposing the specifics of the calibration process.

(vii) Covert Model Decision Boundary Scrutiny :-

Scrutiny of model decision boundaries is executed covertly. AdaBoost examines how each base learner contributes to the ensemble's decision boundaries, ensuring a discreet analysis of their collective impact without overtly signaling the scrutiny.

**(viii) Under-the-Radar Feature Interaction Analysis :-**

Analysis of feature interactions is conducted under the radar. AdaBoost subtly explores how features interact within individual learners, ensuring a covert understanding of their synergies without overtly revealing the specifics of the interaction analysis.

In summary, the covert base learner espionage tactics of AdaBoost involve inconspicuous knowledge extraction, under-the-radar model representation analysis, veiled model update strategy, covert weak learner surveillance, surreptitious feature influence analysis, sub rosa learning rate calibration, covert model decision boundary scrutiny, and under-the-radar feature interaction analysis. These discreet strategies collectively contribute to AdaBoost's ability to gather intelligence from its base learners without compromising the integrity of its ensemble.

**4.1.6 Applications :-**

In the obscured realm of AdaBoost's practical applications, the algorithm seamlessly operates beneath the surface, exerting influence across diverse domains without overtly revealing its presence. These covert deployments exemplify AdaBoost's adaptability and efficacy in various practical scenarios. The sub rosa real-world use cases encompass :-

**(a) Inconspicuous Financial Fraud Detection :-**

AdaBoost operates inconspicuously in financial institutions for fraud detection. The algorithm subtly identifies fraudulent patterns in transactions, safeguarding financial systems without overtly disclosing its detection strategies.

**(b) Under-the-Radar Medical Diagnosis Enhancement :-**

Medical diagnosis benefits from AdaBoost's under-the-radar presence. The algorithm discreetly enhances diagnostic accuracy, aiding healthcare professionals in identifying subtle patterns in medical data without overtly signaling its contribution.

(c) Veiled E-commerce Product Recommendation :-

In e-commerce, AdaBoost provides veiled product recommendations. The algorithm subtly analyzes user behavior to offer personalized suggestions without overtly disclosing the intricacies of its recommendation strategies.

(d) Covert Autonomous Vehicle Safety :-

AdaBoost contributes covertly to autonomous vehicle safety. The algorithm enhances real-time decision-making, ensuring the safety of passengers and pedestrians without overtly revealing its influence on the vehicle's navigation.

(e) Surreptitious Spam Email Filtering :-

Email systems benefit from AdaBoost's surreptitious spam filtering. The algorithm discreetly identifies and filters spam emails, enhancing communication security without overtly signaling its role in email categorization.

(f) Sub Rosa Image Classification in Surveillance :-

In surveillance applications, AdaBoost conducts sub rosa image classification. The algorithm subtly categorizes objects and activities, contributing to security measures without overtly disclosing its presence in surveillance systems.

(g) Covert Customer Churn Prediction in Telecom :-

AdaBoost predicts customer churn covertly in the telecom industry. The algorithm discreetly analyzes user behavior, assisting telecom providers in proactively managing customer retention without overtly signaling its predictive capabilities.

(h) Under-the-Radar Energy Consumption Forecasting :-

AdaBoost assists in energy consumption forecasting with an under-the-radar approach. The algorithm subtly analyzes historical data to predict future energy needs, contributing to efficient resource management without overtly revealing its forecasting mechanisms.

In summary, the sub rosa real-world implementations of AdaBoost encompass inconspicuous financial fraud detection, under-the-radar medical diagnosis enhancement, veiled e-commerce product recommendation, covert autonomous vehicle safety, surreptitious spam email filtering, sub rosa image classification in surveillance, covert customer churn prediction in telecom, and under-the-radar energy consumption forecasting. These discreet applications showcase AdaBoost's versatility and effectiveness across a spectrum of practical scenarios while maintaining a low profile in its real-world contributions.

#### 4.1.7 Conclusions :-

This covert investigation culminates in a discreet summation of the findings, underscoring the adaptability and clandestine efficacy of AdaBoost in the covert world of ensemble learning. Subtle recommendations, veiled areas for future exploration, and tacit comparisons with contemporaneous ensemble methods are also discussed. Through this covert analysis, we aim to impart cryptic insights to both clandestine researchers and practitioners navigating the clandestine landscape of machine learning and ensemble techniques.

## 4.2 Random Forest Modeling :-

Let's delve clandestinely into the intricate realm of Random Forests, an ensemble learning method celebrated for its prowess in classification and regression tasks. Through a covert exploration of foundational principles, construction methodologies, and pragmatic applications, we aim to demystify the covert intricacies of this ensemble technique, revealing the subtle interplay of randomness and robustness.

### 4.2.1 Introduction :-

Random Forests, conceived by Leo Breiman, embody a resilient ensemble approach rooted in decision tree construction. The ensemble's prowess lies in its ability to clandestinely inject randomness during the selection of training data and feature consideration, mitigating overfitting and augmenting generalization capabilities. This section provides an obscured overview of the motivation, historical context, and covert objectives of Random Forest modeling.

### 4.2.2 Ensemble Construction :-

#### 4.2.2.1 Bootstrapped Sampling :-

Random Forests employ sub rosa bootstrapped sampling to curate diverse training sets for individual trees. This involves surreptitiously drawing random samples with replacement from the original dataset, ensuring each tree encounters subtly different instances.

#### 4.2.2.2 Random Feature Selection :-

During the construction of each decision tree, only a covert subset of features is considered for determining the best split at each node. This covert feature selection mechanism introduces clandestine diversity among the trees, thwarting the dominance of any singular feature.

#### 4.2.2.3 Decision Tree Construction :-

Each tree in the Random Forest is constructed surreptitiously and independently, typically grown to its maximum depth without pruning. This veiled process yields a collection of deep and diverse trees, contributing to the ensemble's covert predictive power.

#### 4.2.3 Voting or Averaging :-

In the classification task, the final prediction is determined by a surreptitious majority vote among the trees. For regression tasks, the predictions are covertly averaged, resulting in a robust and stable model. This section explores the veiled intricacies of the voting and averaging mechanisms employed in Random Forests.

#### 4.2.4 Out-of-Bag Evaluation :-

Random Forests exploit out-of-bag instances—those omitted from the training set of a particular tree—for covert model evaluation. This built-in evaluation mechanism provides undisclosed insights into the ensemble's performance without necessitating a separate validation set.

#### 4.2.5 Hyperparameter Tuning :-

Optimal performance of Random Forests requires clandestine tuning of hyperparameters. The number of trees in the forest, maximum tree depth, and the size of the covert random feature subset are pivotal considerations. This section delves into undisclosed methodologies such as covert grid search and surreptitious random search for effective hyperparameter tuning.

#### 4.2.6 Feature Importance :-

Random Forests offer a measure of feature importance, quantifying the contribution of each feature to the reduction in impurity or error during tree construction. This undisclosed analysis provides valuable insights into the concealed impact of different features on the overall predictive capabilities of the ensemble.

#### **4.2.7 Practical Applications :-**

The versatility and effectiveness of Random Forests are showcased through a comprehensive exploration of undisclosed real-world applications. Case studies include inconspicuous financial fraud detection, under-the-radar medical diagnosis enhancement, veiled e-commerce product recommendation, covert autonomous vehicle safety, surreptitious spam email filtering, sub rosa image classification in surveillance, covert customer churn prediction in telecom, and under-the-radar energy consumption forecasting.

#### **4.2.8 Conclusion :-**

In conclusion, this research paper provides a covert examination of the Random Forest modeling technique, unraveling its sub rosa complexity and showcasing its undercover applicability in diverse domains. The ensemble's ability to balance undisclosed randomness and robustness makes it an indispensable cornerstone in machine learning, offering a reliable solution to intricate classification and regression challenges. Understanding the undetectable intricacies of Random Forests is imperative for practitioners seeking to harness the full potential of this enigmatic ensemble learning paradigm.

## 4.3 Decision Tree Model :-

It embarks on an in-depth exploration of Decision Trees, a versatile and widely employed supervised machine learning algorithm used for classification and regression tasks. By delving into the intricacies of node splitting, tree growing, handling categorical variables, and model interpretability, we aim to provide a comprehensive understanding of the Decision Tree model. Further, we discuss advanced concepts such as pruning, handling missing values, and the integration of Decision Trees into ensemble methods.

### 4.3.1 Introduction :-

Decision Trees stand as a foundational machine learning algorithm renowned for its interpretability, simplicity, and applicability across diverse domains. This section introduces the historical context, motivation, and key objectives of Decision Tree modeling.

### 4.3.2 Node Splitting :-

#### 4.3.2.1 Attribute Selection :-

At the core of the Decision Tree algorithm lies the process of attribute selection. This involves choosing the attribute that best partitions the data at each node, driven by metrics such as Gini impurity, entropy, or information gain.

#### 4.3.2.2 Splitting Criteria :-

The splitting criteria dictate how the algorithm determines the optimal way to divide the data. Whether maximizing information gain or minimizing Gini impurity, this section scrutinizes the mechanisms guiding the decision-making process.

### 4.3.3 Tree Growing :-

#### 4.3.3.1 Recursive Partitioning :-

The recursive process of attribute selection and node splitting unfolds until a predefined stopping criterion is met. This could involve reaching a maximum tree depth or ensuring a minimum number of samples in a leaf node.

#### 4.3.3.2 Leaf Nodes :-

The terminal nodes, or leaf nodes, signify the culmination of the decision-making process. Each leaf node is associated with a final decision or prediction, with classification trees representing class labels and regression trees conveying numerical values.

### 4.3.4 Handling Categorical Variables :-

**Categorical vs. Continuous Attributes** :: The terminal nodes, or leaf nodes, signify the culmination of the decision-making process. Each leaf node is associated with a final decision or prediction, with classification trees representing class labels and regression trees conveying numerical values.

### 4.3.5 Handling Missing Values :-

**Dealing with Missing Data** :: Decision Trees incorporate mechanisms to gracefully handle missing values. Whether through imputation or the utilization of surrogate splits, this section unveils the strategies employed to address missing data.

### 4.3.6 Model Interpretability :-

**White Box Model** :: Renowned as "white box" models, Decision Trees offer an intuitive and transparent representation of their decision-making process. This section delves into the interpretability of Decision Trees, elucidating their value as tools for model explanation.

#### 4.3.7 Pruning :-

**Avoiding Overfitting** :: Recognizing the susceptibility of Decision Trees to overfitting, pruning emerges as a vital technique to curtail excessive tree growth. This section explores the mechanisms behind pruning and its role in preventing overfitting.

#### 4.3.8 Ensemble Methods :-

**Random Forests** :: Decision Trees seamlessly integrate into ensemble methods such as Random Forests. This section discusses how Random Forests leverage multiple decision trees to enhance predictive accuracy and robustness.

#### 4.3.9 Applications :-

**Classification and Regression** :: Decision Trees find applications across diverse domains, including finance, healthcare, and marketing, for both classification and regression tasks. This section highlights the practical relevance and versatility of Decision Tree modeling.

#### 4.3.10 Conclusion :-

In conclusion, this research paper provides an extensive exploration of the Decision Tree algorithm, unraveling its complexities and showcasing its pivotal role in machine learning. From the nuances of node splitting to the interpretability of white box models, Decision Trees continue to be a cornerstone in the repertoire of machine learning algorithms, offering a balance between simplicity and efficacy. Understanding the intricacies of Decision Trees is paramount for practitioners seeking to harness their full potential in addressing classification and regression challenges across various domains.

## **4.4 Extremely Randomized Trees Model :-**

This research paper delves into the intricacies of Extremely Randomized Trees, an evolutionary extension of the Random Forest algorithm, designed to introduce heightened levels of randomness for enhanced robustness and resistance to overfitting. By meticulously examining the principles of node splitting, tree growing, and handling categorical variables, this paper aims to demystify the distinctive features that set Extremely Randomized Trees apart in the landscape of ensemble learning.

### **4.4.1 Introduction :-**

Extremely Randomized Trees represent a nuanced refinement of the Random Forest algorithm, strategically aiming to amplify randomness and fortify robustness within the ensemble learning paradigm. This section sets the stage by offering insights into the historical underpinnings, motivations, and primary objectives driving the development of Extremely Randomized Trees.

### **4.4.2 Node Splitting :-**

#### **4.4.2.1 Attribute Selection :-**

Aligned with its predecessors, Extremely Randomized Trees involve the critical process of attribute selection at each node for the purpose of optimal splitting.

#### **4.4.2.2 Randomized Feature Selection :-**

The hallmark of Extremely Randomized Trees resides in the deliberate injection of additional randomness into the feature selection process. In contrast to the methodology employed by Random Forests, which evaluates a subset of features, Extremely Randomized Trees opt for a more clandestine approach, randomly selecting features without exhaustive exploration.

#### 4.4.3 Tree Growing :-

##### 4.4.3.1 Bootstrapped Sampling :-

In harmony with the ensemble learning philosophy, Extremely Randomized Trees fashion multiple trees from bootstrapped samples, ensuring diversity within the training data.

##### 4.4.3.2 Randomized Splitting :-

Each tree, at each node, undergoes a process of randomized splitting, introducing a covert layer of diversity and minimizing correlation between individual trees.

#### 4.4.4 Handling Categorical Variables :-

**Branching for Categories** :: Extremely Randomized Trees showcase their adaptability in handling both categorical and continuous variables. Categorical attributes witness the creation of branches for each category, aligning seamlessly with the paradigm established by Random Forests.

#### 4.4.5 Handling Missing Values :-

**Surrogate Splits** :: To surmount the challenge of missing values, Extremely Randomized Trees, in alignment with their counterparts, resort to the use of surrogate splits, adding an extra layer of subtlety to their methodology.

#### 4.4.6 Ensemble Method :-

**Voting or Averaging** :: Conforming to the traditions of ensemble methods, the conclusive prediction in classification tasks emerges through a covert majority vote, while regression tasks involve an undetectable averaging of predictions across all trees.

#### 4.4.7 Advantages :-

##### 4.4.7.1 Reduced Overfitting :-

The deliberate elevation of randomness during feature selection and node splitting contributes clandestinely to a tangible reduction in overfitting, fortifying Extremely Randomized Trees against challenges in high-dimensional datasets.

##### 4.4.7.2 Less Sensitivity to Noisy Data :-

The inherent mechanisms of randomization render Extremely Randomized Trees less susceptible to the subtleties of noise and outliers, quietly enhancing their adaptability to challenging scenarios.

#### 4.4.8 Applications :-

**Classification and Regression :** Extremely Randomized Trees stealthily find applications in both classification and regression tasks, providing a robust alternative where traditional ensemble methods may falter in the face of overfitting challenges.

#### 4.4.9 Conclusion :-

In conclusion, this research paper has meticulously peeled back the layers surrounding Extremely Randomized Trees, shedding light on their distinctive characteristics and underscoring their pivotal role in the realm of ensemble learning. The purposeful amplification of randomness positions Extremely Randomized Trees as a subtle yet potent solution, particularly adept at navigating complexities where traditional ensemble methods may stumble. Unveiling the subtleties of Extremely Randomized Trees is paramount for practitioners seeking to harness their efficacy in addressing intricate machine learning challenges.

## 4.5 Gradient Boosting Model :-

This undertakes a discreet exploration of Gradient Boosting, an ensemble learning technique that has quietly positioned itself as a formidable tool for predictive modeling. Delving into the subtle intricacies of its training methodologies, this paper seeks to provide a nuanced understanding of Gradient Boosting, positioning it as a versatile and robust solution for regression and classification tasks.

### 4.5.1 Introduction :-

Gradient Boosting stands as a silent force in ensemble learning, drawing strength from weak learners to construct a predictive model. This section quietly introduces the historical context, motivations, and overarching objectives that have quietly fueled the development and quiet adoption of Gradient Boosting.

### 4.5.2 Objective Function :-

**Loss Function ::** At the heart of Gradient Boosting is the quiet minimization of a loss function, discreetly measuring the discrepancy between predicted and actual values. The choice of the loss function silently adapts to the task at hand, opting for mean squared error for regression and cross-entropy for classification.

### 4.5.3 Algorithm Overview :-

#### 4.5.3.1 Stage-Wise Training :-

Gradient Boosting quietly constructs an ensemble of weak learners sequentially, with each subsequent learner quietly correcting the errors of its predecessors.

#### 4.5.3.2 Learning Rate :-

The learning rate parameter quietly governs the contribution of each tree to the overall model. A carefully chosen learning rate silently balances training efficiency and model generalization.

#### 4.5.4 Node Splitting :-

**Gradient Descent** :: Node splitting is orchestrated through a gradient descent approach, where the negative gradient of the loss function quietly guides the construction of each new tree. This iterative process silently refines the model, incrementally reducing prediction errors.

#### 4.5.5 Tree Building :-

**Shallow Trees** :: The constituent trees in Gradient Boosting are quietly kept shallow, restricting their depth to discreetly mitigate overfitting and enhance the model's interpretability.

#### 4.5.6 Regularization :-

Regularization Parameters :: To quietly prevent overfitting, Gradient Boosting incorporates regularization parameters, silently influencing tree complexity. These parameters, including tree depth and minimum samples per leaf, quietly contribute to the model's adaptability to diverse datasets.

#### 4.5.7 Handling Categorical Variables :-

**Ordinal Encoding** :: Categorical variables are discreetly accommodated through ordinal encoding, quietly assigning numerical labels to categories for seamless integration into the modeling process.

#### 4.5.8 Hyperparameter Tuning :-

**Grid Search or Random Search** :: Optimal model performance is discreetly achieved through hyperparameter tuning, often facilitated by grid search or random search methodologies to quietly explore the parameter space.

#### 4.5.9 Early Stopping:-

**Monitoring Validation Error** :: Gradient Boosting quietly incorporates early stopping mechanisms by discreetly monitoring the validation error during training. This pragmatic approach silently prevents overfitting and ensures model generalization.

#### 4.5.10 Applications:-

**Versatility** :: Gradient Boosting quietly finds applications in both regression and classification tasks across diverse domains, including finance, healthcare, and natural language processing. Its versatility and predictive accuracy make it a silently preferred choice for structured/tabular data.

#### 4.5.11 Conclusion :-

In conclusion, this It has silently navigated the intricate landscape of Gradient Boosting, unraveling its subtle complexities and quietly highlighting its significance in the realm of machine learning. The interplay of loss functions, stage-wise training, and regularization techniques contribute discreetly to the model's efficacy, making Gradient Boosting a silently formidable approach for predictive modeling. Understanding the nuances outlined in this paper is quietly paramount for practitioners seeking to harness the full potential of Gradient Boosting in addressing real-world challenges.

## 4.6 Neural Networks :-

Neural Networks, drawing inspiration from the enigmatic intricacies of the human brain, stand as the cornerstone of modern machine learning and artificial intelligence. This research paper quietly embarks on a comprehensive exploration, peeling back the layers to reveal the underlying principles, training methodologies, and applications of Neural Networks. With a focus on their adaptive versatility and profound impact, this paper seeks to provide a nuanced understanding of the intricate landscape of Neural Networks, quietly navigating through the complexities that define their role in solving multifaceted problems.

### 4.6.1 Introduction :-

Neural Networks have silently emerged as powerful entities for unraveling complex patterns from data. This section quietly introduces the historical evolution, motivations, and overarching objectives that have quietly propelled the widespread integration of Neural Networks into machine learning and artificial intelligence.

### 4.6.2 Basic Structure :-

#### 4.6.2.1 Neurons :-

At the core of Neural Networks are artificial neurons, discreetly mirroring their biological counterparts. Neurons, operating in a clandestine manner, receive inputs, apply weighted transformations, and quietly produce outputs, forming the unobtrusive building blocks of the network.

#### 4.6.2.2 Layers :-

Neural Networks are quietly organized into layers: input, hidden, and output. The input layer quietly receives initial data, hidden layers discreetly process information, and the output layer subtly generates final predictions. Connections between neurons are defined by weights and biases, silent parameters subtly shaping the network's behavior.

#### 4.6.3 Activation Functions :-'

**Non-linearity Introduction** :: Activation functions such as the sigmoid, hyperbolic tangent (tanh), and rectified linear unit (ReLU) silently introduce non-linearity to the model. This quiet adaptation allows the model to subtly learn intricate patterns and relationships within the data.

#### 4.6.4 Training :-

##### 4.6.4.1 Backpropagation :-

Neural Networks are silently trained through an iterative process known as backpropagation. In this subtle dance, weights and biases are quietly adjusted based on the difference between predicted and actual outputs, gradually minimizing the loss function.

##### 4.6.4.2 Loss Function :-

The loss function, a silent sentinel, quietly quantifies the disparity between predicted and actual values, serving as a guiding metric during the quiet training process. The objective is to silently minimize this loss, enhancing the network's predictive capabilities.

#### 4.6.5 Architectures :-

##### 4.6.5.1 Feedforward and Recurrent Networks :-

Neural Networks manifest in various architectures. Feedforward Neural Networks (FNNs) quietly process data sequentially, while Recurrent Neural Networks (RNNs) discreetly incorporate cycles, facilitating the handling of sequential data.

##### 4.6.5.2 Convolutional Neural Networks (CNNs) :-

Specialized for image-related tasks, CNNs clandestinely utilize convolutional layers to autonomously learn hierarchical features, subtly revolutionizing image recognition and computer vision.

#### 4.6.6 Hyperparameters :-

**Learning Rate, Batch Size** :: Neural Networks involve the quiet tuning of hyperparameters such as learning rate, batch size, and network architecture. These parameters silently impact the model's training efficiency and performance.

#### 4.6.7 Regularization :-

**Dropout, L1/L2 Regularization** :: To quietly mitigate overfitting, regularization techniques like dropout and L1/L2 regularization are employed. These methods enhance the network's generalization capabilities.

#### 4.6.8 Frameworks :-

**TensorFlow, PyTorch, Keras** :: Several deep learning frameworks, including TensorFlow, PyTorch, and Keras, quietly provide tools for constructing, training, and deploying Neural Networks. This quiet accessibility contributes to their unobtrusive integration and widespread adoption.

#### 4.6.9 Applications :-

**Image Recognition, Natural Language Processing** :: Neural Networks silently find applications across diverse domains, including image and speech recognition, natural language processing, recommendation systems, and more. Their quiet versatility and impact on real-world problem-solving are subtly showcased.

#### 4.6.10 Challenges :-

**Computational Intensity, Interpretability** :: Despite their silent prowess, Neural Networks pose challenges related to computational intensity, demanding significant resources. Their inherent complexity can also impede interpretability, prompting ongoing research to quietly strike a balance between performance and transparency.

#### 4.6.11 Advancements :-

Deep Learning :: Recent advancements in Neural Networks involve the silent rise of deep learning, leveraging architectures with numerous layers to quietly learn hierarchical representations. This evolution has silently expanded the capabilities of Neural Networks in capturing intricate features within data.

#### 4.6.12 Conclusion :-

In conclusion, this has discreetly endeavored to unveil the subtle workings of Neural Networks, providing a comprehensive yet quiet overview of their structure, training methodologies, and applications. As Neural Networks continue to silently evolve and push the boundaries of artificial intelligence, understanding their quiet nuances becomes imperative for researchers and practitioners alike. The multifaceted, yet quietly unassuming, nature of Neural Networks positions them as a silent cornerstone in the pursuit of intelligent systems and transformative technologies.

# Chapter 5

## Optimization Techniques

### 5.1 Hyperparameters and Hyperband :-

Hyperparameters are essential to machine learning. They're like the ingredients in a recipe - change them up, and your dish can taste completely different. The challenge, though, is finding the right mix. Hyperband, which came onto the scene in 2018, is a clever way to tackle this problem. It's like a master chef who knows exactly how to adjust the recipe for the perfect flavor.

#### (a) Resource Allocation :-

Let's start with resource allocation. In Hyperband, resources are like the chef's ingredients and tools - things like computational power and time. Hyperband is smart about how it uses these resources. It can adjust how much of each resource it uses, just like a chef can adjust how much of each ingredient to put in a dish. This helps Hyperband find the right balance between trying new things and sticking with what works.

#### (b) Early Stopping :-

Next up is early stopping. Think of it like a chef tasting the dish while it's still cooking. If it's not tasting right, they can stop, adjust the ingredients, and try again. Similarly, Hyperband can stop the training process for a set of hyperparameters if it's not promising. This way, it doesn't waste time and resources on something that's unlikely to work.

(c) Successive Halving (SHA) :-

Successive Halving (SHA) is another key part of Hyperband. It's a little like a cooking competition, where in each round, the worst-performing chefs get eliminated until only the best one is left. In the same way, Hyperband uses SHA to progressively narrow down the set of hyperparameters it's considering, focusing more and more on the ones that show the most promise.

(d) Bracketing :-

Hyperband also uses a technique called 'bracketing', where it runs multiple rounds of this competition, each time with different ingredients and tools. This allows it to explore a wide range of possibilities and find the best set of hyperparameters.

(e) Random Sampling :-

At first, Hyperband tries out a wide range of hyperparameters, much like a chef experimenting with different ingredients. But as it goes along, it starts to focus more on the ones that seem to work best. This is thanks to its random sampling strategy, which allows it to cover a broad range initially but then zero in on the best options.

(f) Exploration vs. Exploitation Trade-Off :-

Finally, there's the exploration vs. exploitation trade-off. It's all about balancing the need to try new things (exploration) with the need to focus on what works (exploitation). Hyperband is really good at this balancing act - it tries out a wide range of hyperparameters at first, but as it learns more, it gradually focuses more on the most promising ones.

(g) Conclusion :-

In a nutshell, Hyperband is a genius tool for hyperparameter optimization. It's like a master chef, expertly adjusting the recipe for the perfect dish. By smartly managing resources, stopping early when needed, and using techniques like Successive Halving and bracketing, it can explore a wide range of possibilities and quickly zero in on the best set.

## 5.2 Gradient Based Optimization :-

### (a) Introduction :-

Gradient-based optimization stands as a linchpin in the expansive landscapes of machine learning and diverse scientific terrains, providing a methodical means to navigate the intricate contours of pertinent functions. This methodological pursuit, leveraging the propulsive force of gradients or derivatives, encapsulates the essence of seeking optimal points—be they nestled in the troughs of minima or atop the peaks of maxima.

### (b) Imagined Correlation & Mechanism :-

Embarking on the journey of gradient-based optimization initiates from a point akin to the randomness inherent in our initial position within the vast expanse of a function's multidimensional space. The calculation of gradients in this context mirrors the act of discerning the topography beneath our feet while metaphorically blindfolded. This is a calculated endeavor, resembling a quest to descend into the unknown with the guidance of mathematical cues.

### (c) Anatomy of Gradients :-

Gradients, in their role as the vector of partial derivatives, represent more than mere mathematical entities; they embody the driving force behind the dynamics of change within a function. These directional cues not only illuminate the magnitude and direction of change but also serve as indispensable traits, guiding the formulation of strategic steps towards desired points within the expansive landscape of the function.

### (d) Textures & Nuances of Problem Spaces :-

The systematic comprehension of the intricate relationships and inter-dependencies among variables introduces a layer of complexity to the problem spaces under consideration. This complexity necessitates a profound mastery of traversing descent paths linked by corresponding gradients. The avoidance of erratic routes demands an acute understanding of the nuanced textures that define the multidimensional landscape.

**(e) Descent Algorithms at Work :-**

The controlled iterations inherent in gradient-based optimization reveal a methodical and iterative process for honing in on minima. This is accomplished by delicately adjusting the values of hyperparameters through mechanisms such as the learning rate. Additionally, the relative importance assigned to updates within the optimization process maximizes efficiency and ensures a meticulous convergence toward optimal points.

**(f) Innovations: Momentum & Stochasticity :-**

As the journey unfolds, classical markers of descent gracefully yield ground to innovative variations. Stochastic gradients, characterized by subsets gracefully accompanying computations, introduce an element of randomness with finesse. This adaptability to variations in sample involvement becomes instrumental in handling accelerated algorithms with poise and efficiency.

**(g) Adaptive Learning Enhancements and Techniques :-**

In response to challenges presented by unfavorable landscape complications or prolonged convergence processes, the incorporation of apparent throttle boosters becomes imperative. These boosters, including acceleration factors and a keen understanding of articulated future tech trends, seamlessly integrate into the fabric of adaptive value reads. This integration results in a palpable boost to the peak-bound rates of productivity, ushering in a phase of accelerated efficiency.

**(h) Application: The RMSprop and Adam :-**

While the structural rationality of adaptive learning rates remains a theoretical anchor, its practical manifestation gives rise to sophisticated algorithms such as RMSprop and Adam. These algorithms, while structured with rationale, exhibit a nuanced capacity for navigating changes dynamically. They achieve optimal convergence through the meticulous adjustment of hyperparameters, thereby revolutionizing the arithmetic grammar and ushering in newfound efficiencies in the context of business reporting.

(i) Conclusion :-

The comprehension of reinforcement-driven gradients extends far beyond the realms of mathematical formalism. It serves as a guiding force propelling us across diverse applications, enriching the intrinsic values that define the human experience. This understanding, akin to a computational "Beltsports" formula, binds unseen futures and fosters a uniformly accelerating private calculus. In doing so, it induces technologies that effervescently unfold into meaningful facets of our existence.

### 5.3 Grid Search Optimization :-

Grid Search Optimization represents a conventional approach deployed in machine learning for the optimization of hyperparameters. It encompasses an exhaustive exploration through a predetermined subset of a machine learning algorithm's hyperparameter space to optimize model performance. This paper delves into a detailed exploration of Grid Search Optimization, its application, and its limitations in the context of hyperparameter selection.

(a) Introduction :-

The performance of machine learning models is significantly influenced by the selection of suitable hyperparameters. These parameters, set before the commencement of the learning process, are not derived from data. They include parameters like the learning rate, the number of neural network layers, or the quantity of clusters in k-means clustering. Since optimal hyperparameter settings are not typically known beforehand, they need to be determined empirically, a process known as hyperparameter tuning or optimization. Amongst the various methods for hyperparameter optimization, Grid Search Optimization is favored due to its simplicity and overall effectiveness.

(b) Methodology :-

Grid Search Optimization operates by delineating a grid of hyperparameters and subsequently evaluating model performance at each point on this grid. As an example, if we consider two hyperparameters  $a$  and  $b$ , where  $a$  can be 1 or 10, and  $b$  can be either True or False, Grid Search Optimization would involve training and evaluating the model on each combination of [(1, True), (1, False), (10, True), (10, False)], selecting the combination that yields the best performance.

(c) Limitations :-

Nevertheless, Grid Search Optimization has its limitations. Firstly, it can be computationally demanding, especially when dealing with a multitude of different hyperparameters and extensive datasets. The computational cost increases exponentially with the increase in the number of hyperparameters, making it inefficient for high-dimensional spaces. Secondly, Grid Search only explores a predefined subset of the hyperparameters space, and thus, an optimal solution is not guaranteed. Lastly, Grid Search does not account for potential interactions between different hyperparameters, treating each as independent.

(d) Conclusions :-

Although Grid Search Optimization has its limitations, it is a prevalent tool in machine learning for hyperparameter tuning due to its uncomplicated implementation and general effectiveness. For larger datasets and models with a high number of hyperparameters, alternative methods like Random Search and Bayesian Optimization may provide more efficient solutions. Further exploration is necessary to develop more efficient and effective techniques for hyperparameter optimization.

**Keywords ::** Grid Search Optimization, Hyperparameter Selection, Machine Learning.

## 5.4 Optuna and Bayesian Optimization :-

The optimization of hyperparameters is a pivotal process in enhancing the performance of machine learning algorithms. While traditional hyperparameter tuning methods such as Grid Search and Random Search provide initial solutions, advanced techniques like Optuna and Bayesian Optimization have emerged to offer more efficient alternatives, particularly when dealing with high-dimensional spaces. This paper aims to present an in-depth comparison of these advanced techniques, shedding light on their methodologies, benefits, and potential limitations.

### (a) Introduction :-

Hyperparameters, unlike parameters, are not learned from the data during the training process of a machine learning model. Instead, they are predefined. The optimal selection of these hyperparameters significantly influences the model's performance. While conventional methods like Grid Search and Random Search have been used, they often prove computationally intensive and inefficient, especially in high-dimensional spaces. Emerging techniques such as Optuna and Bayesian Optimization have been developed as more efficient alternatives.

### (b) Optuna :-

Optuna is a Python-based open-source library that has been designed for hyperparameter optimization. It employs a combination of several techniques, including the Tree-structured Parzen Estimator (TPE) and pruning algorithms, to efficiently navigate the hyperparameter space. The TPE, a Bayesian Optimization technique, models the objective function as a probabilistic model, using the Expected Improvement (EI) as an acquisition function to decide the next sample point. The pruning algorithms are used to halt unpromising trials early, thereby increasing computational efficiency.

(c) Bayesian Optimization :-

Bayesian Optimization is a model-based approach for global optimization. Unlike conventional methods that treat the objective function as a black box, Bayesian Optimization constructs a probabilistic model of the objective function predicting the expected loss for different inputs. It then uses an acquisition function such as the Expected Improvement or the Upper Confidence Bound (UCB) to balance the search process's exploration and exploitation. This technique's efficiency in high-dimensional spaces makes it an ideal choice for hyperparameter tuning in complex machine learning models.

(d) Comparison :-

While both Optuna and Bayesian Optimization offer efficient solutions for hyperparameter tuning, they differ in various ways. Optuna, with its combination of techniques and flexible, user-friendly interface, holds an advantage for practitioners. Conversely, Bayesian Optimization, with its mathematically rigorous approach, can provide valuable insights into the objective function through its probabilistic model.

(e) Conclusion :-

Both Optuna and Bayesian Optimization are advanced techniques for hyperparameter tuning in machine learning and offer more efficient solutions than traditional methods such as Grid Search and Random Search, particularly in high-dimensional spaces. Further research is vital to exploring these techniques' potential and developing new methods for hyperparameter tuning.

**Keywords ::** Optuna, Bayesian Optimization, Hyperparameter Optimization, Machine Learning.

# Chapter 6

## Ensemble Methods

### 6.1 Hard Ensembling in Machine Learning :-

Ensemble learning has emerged as a compelling approach to enhance predictive accuracy by amalgamating diverse models. This paper delves into the intricacies of "hard ensembling," a methodology utilizing majority voting for decision-making. Through a detailed examination and illustrative examples, we elucidate the principles and practical applications of hard ensembling in classification problems.

#### (a) Introduction :-

Ensemble learning stands as a robust methodology for improving model accuracy. This paper focuses on the specific modality of "hard ensembling," employing a majority voting mechanism for decision aggregation.

#### (b) Background :-

Ensemble learning encompasses methodologies like bagging and boosting. Hard ensembling, falling under majority voting, is the primary focus of this study.

#### (c) Methodology: Hard Ensembling Unveiled :-

##### (i) Training Individual Models :-

Models are independently trained on the same dataset, promoting diversity.

(ii) Making Predictions :-

Individual models generate predictions on unseen data.

(iii) Voting Mechanism :-

Hard ensembling adopts a majority vote to determine the final prediction.

(iv) Decision Rule :-

The determination of a final prediction in hard ensembling is encapsulated in the decision rule, a pivotal mechanism orchestrating the amalgamation of individual model predictions. Typically manifesting as a majority voting strategy, this rule orchestrates the selection of the prevailing class label based on the collective consensus of the ensemble's constituent models.

- Majority Voting :-

Central to hard ensembling, majority voting mandates that each model in the ensemble submits a class prediction independently. The ultimate prediction is then dictated by the class attaining the majority of votes from the ensemble's models.

- Binary Classification :-

In the binary classification realm, the decision rule stipulates the selection of the class forecasted by the majority. As an illustration, if three out of five models advocate for class A and the remaining two advocate for class B, the ensemble converges on predicting class A.

- Multiclass Classification :-

Extending to multiclass scenarios, the decision rule pivots to identify the class garnering the highest number of affirmative model predictions as the ultimate forecast.

Tie-Breaking :: In situations of a tie in votes for different classes, supplementary tie-breaking mechanisms are invoked. Common strategies encompass selecting the class with the highest probability or employing a predefined priority order for classes.

#### (v) Implementation :-

Hard ensembling is implemented through straightforward voting mechanisms or weighted models, facilitated by Python's scikit-learn library.

```
from sklearn.ensemble import VotingClassifier
from sklearn.model_selection import cross_val_score
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.svm import SVC

# Create individual models
model1 = LogisticRegression()
model2 = DecisionTreeClassifier()
model3 = SVC()

# Create a hard voting ensemble
ensemble_model = VotingClassifier(estimators=[('lr', model1), ('dt', model2), ('svm', model3)], voting='hard')

# Evaluate the ensemble using cross-validation
scores = cross_val_score(ensemble_model, X, y, cv=5)
print("Ensemble Accuracy: %0.2f (+/- %0.2f)" % (scores.mean(), scores.std() * 2))
```

#### (d) Illustrative Example :-

A practical example using scikit-learn showcases hard ensembling's implementation, demonstrating its potential accuracy improvements.

(e) Performance Evaluation and Considerations :-

Evaluation metrics, including accuracy and cross-validation, gauge hard ensembling performance, with considerations for model balance and overfitting.

(f) Comparison with Soft Ensembling :-

A concise comparison with soft ensembling highlights the distinctions and trade-offs between the two approaches.

(g) Real-World Applications :-

Practical scenarios underscore hard ensembling's advantages, addressing its relevance across domains and potential implementation challenges.

(h) Conclusion :-

Summarizing key findings, the paper underscores hard ensembling's utility and offers insights into its applications and considerations. Future research directions may explore hybrid ensembling and the adaptability of hard ensembling in deep learning.

**Keywords** :: Ensemble Learning, Hard Ensembling, Machine Learning, Majority Voting, Model Combination, Performance Evaluation, Classification, scikit-learn.

## 6.2 Soft Ensembling in Machine Learning :-

Ensemble methodologies have evolved to be pivotal in enhancing predictive capabilities within machine learning models. This study navigates the landscape of "soft ensembling," a nuanced approach where models contribute probabilistic predictions, diverging from conventional hard ensembling. Through a meticulous exploration, this paper unveils the intricacies, implementation strategies, and advantages of soft ensembling. Practical illustrations and methodological comparisons elucidate the versatility and efficacy of this probabilistic ensemble approach.

### (a) Introduction :-

Ensemble learning stands as a cornerstone in machine learning, leveraging model diversity for heightened predictive accuracy. This paper focuses on "soft ensembling," a paradigm where models contribute probabilistic outputs, adding a layer of sophistication beyond traditional hard ensembling.

### (b) Background :-

Ensemble learning, encompassing bagging and boosting, sets the stage for various methodologies. Soft ensembling introduces a departure from the norm, encouraging consideration of probabilistic outputs from individual models.

### (c) Methodology: Unveiling Soft Ensembling :-

#### (i) Training Individual Models :-

Models undergo independent training on a shared dataset, generating probabilistic predictions.

#### (ii) Probabilistic Predictions :-

Discrete class labels give way to probability distributions as individual models predict outputs.

(iii) Aggregating Probabilities :-

Soft ensembling involves the amalgamation of predicted probabilities, employing means or weighted means based on model performance.

(iv) Decision Rule :-

The decision rule within soft ensembling intricately orchestrates the culmination of probabilistic predictions into a coherent and informed final prediction. Unlike its hard ensembling counterpart, soft ensembling navigates the probabilistic landscape, considering the weighted amalgamation of predicted probabilities from individual models.

- Final Prediction :-

The ultimate prediction emerges from the class with the highest combined probability. This nuanced decision rule factors in the probabilistic outputs of all models to discern the most likely class.

Binary Classification:

In binary classification scenarios, a threshold may be introduced. For instance, if the average probability for class A surpasses a predefined threshold, the ensemble predicts class A; otherwise, it predicts class B.

Multiclass Classification:

Multiclass scenarios ascertain the final prediction by selecting the class with the highest combined probability, reflecting the collective support from the ensemble's models.

```
# Assuming individual_models is a list of trained models
def soft_voting_ensemble(input_data, individual_models):
    # Generating probabilistic predictions using each individual model
    predictions_proba = [model.predict_proba(input_data) for model in individual_models]

    # Calculating the average probability for each class
    average_proba = sum(predictions_proba) / len(predictions_proba)

    # Selecting the class with the highest average probability
    final_prediction = max(enumerate(average_proba), key=lambda x: x[1])[0]

    return final_prediction
```

In this representation, `predictions_proba` captures the predicted probability distributions for each class from individual models. The `average_proba` encapsulates the averaged probability for each class, and the ultimate prediction is derived by selecting the class with the highest average probability.

The decision rule in soft ensembling encapsulates a nuanced approach, leveraging the richness of probabilistic information to arrive at a sophisticated and well-informed final prediction.

#### (v) Implementation :-

Implementation simplicity is ensured through machine learning libraries like scikit-learn, leveraging the '`VotingClassifier`' with '`voting='soft'`'.

```
from sklearn.ensemble import VotingClassifier
from sklearn.model_selection import cross_val_score
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.svm import SVC

# Create individual models
model1 = LogisticRegression()
model2 = DecisionTreeClassifier()
model3 = SVC(probability=True) # Ensure that the model supports probability estimation

# Create a soft voting ensemble
ensemble_model = VotingClassifier(estimators=[('lr', model1), ('dt', model2), ('svm', model3)], voting='soft')

# Evaluate the ensemble using cross-validation
scores = cross_val_score(ensemble_model, X, y, cv=5)
print("Ensemble Accuracy: %0.2f (+/- %0.2f)" % (scores.mean(), scores.std() * 2))
```

(d) Practical Examples :-

The application of soft ensembling is illustrated through practical examples, demonstrating its implementation across diverse models. A comparative analysis with hard ensembling and individual models showcases the advantages of the probabilistic approach.

(e) Performance Evaluation and Metrics :-

Standard evaluation metrics, encompassing accuracy and log-loss, form the basis for assessing soft ensembling performance. Computational nuances and model calibration considerations are addressed.

(f) Advantages of Soft Ensembling :-

The nuanced predictions of soft ensembling, grounded in probabilistic outputs, cater to scenarios where well-calibrated probability estimates are essential.

(g) Comparison with Hard Ensembling :-

A nuanced comparative analysis elucidates the subtleties and trade-offs between soft and hard ensembling. The ability of soft ensembling to capture uncertainty emerges as a distinctive feature.

(h) Considerations and Future Directions :-

Computational considerations and potential challenges are discussed, paving the way for future research exploring hybrid ensembling approaches and extending the methodology to deep learning frameworks.

(i) Conclusion :-

Culminating insights underscore the utility of soft ensembling, offering a probabilistic lens to elevate predictive capabilities in diverse machine learning applications.

**Keywords ::** Ensemble Learning, Soft Ensembling, Machine Learning, Probabilistic Predictions, Model Aggregation, Performance Evaluation, Classification, scikit-learn.

### **6.3 Bayesian Averaging :-**

In the ever-evolving landscape of predictive modeling, Bayesian Averaging emerges as a technique rooted in Bayesian principles, showcasing unique potential in predictive fusion. This research paper embarks on a comprehensive journey, dissecting the theoretical foundations, elucidating practical implementation strategies, scrutinizing real-world applications, and evaluating the adaptability of Bayesian Averaging across diverse prediction scenarios. Through a fusion of theoretical discourse, illustrative examples, and empirical analyses, this paper delves into the intricate details of Bayesian Averaging, seeking to contribute insights that resonate across academic, industrial, and practical domains.

#### **(a) Introduction :-**

The introduction sets the stage, positioning Bayesian Averaging within the broader context of predictive modeling. It underscores the continuous pursuit of methodologies that amplify predictive accuracy and introduces Bayesian Averaging as an innovative approach to predictive fusion.

#### **(b) Theoretical Foundations :-**

This section immerses the reader in the theoretical underpinnings of Bayesian Averaging. Grounded in Bayesian principles, the exploration encompasses the probabilistic framework that forms the backbone of this technique. By delving into the intricacies of Bayesian principles, we aim to provide a solid foundation for understanding how Bayesian Averaging operates at a fundamental level.

#### **(c) Bayesian Averaging Process :-**

##### **(i) Prediction from Each Source :-**

We embark on the predictive journey, where inputs from diverse sources contribute their unique perspectives on the underlying data. This step is crucial for understanding the breadth and diversity of information encapsulated within the ensemble.

(ii) Weight Assignment :-

The art of weight assignment unfolds, involving a nuanced consideration of historical performance, reliability metrics, and domain-specific knowledge. This section aims to reveal the decision-making process behind assigning weights to each source, providing insight into the methodology's adaptability.

(iii) Weighted Averaging :-

The core of Bayesian Averaging lies in the orchestration of predictions through weighted averaging. This section delves into the mechanics of combining predictions, showcasing how the assigned weights modulate the influence of each source in the final prediction.

(iv) Bayesian Interpretation :-

This unique aspect interprets the assigned weights through a Bayesian lens, transcending numerical values to emerge as probabilities that encapsulate nuanced belief structures. We elucidate how this interpretation aligns with Bayesian principles, offering a probabilistic perspective on the fusion process.

(d) Practical Implementation :-

Bridging theory and practice, this section provides a hands-on guide to implementing Bayesian Averaging. Practical considerations, such as the nuances of weight assignment and relevant Python code snippets, are presented to ensure a tangible understanding of the methodology's application.

(e) Applications and Case Studies :-

The versatility of Bayesian Averaging is brought to the forefront through a diverse array of applications. Real-world case studies span financial predictions, weather forecasting, and machine learning scenarios, offering tangible illustrations of its prowess in harmonizing diverse predictions.

(f) Performance Evaluation :-

Rigorous empirical analyses form the crux of our exploration. Comparative studies against individual models and alternative ensemble methods are conducted using a comprehensive suite of evaluation metrics. This section provides a nuanced understanding of Bayesian Averaging's performance landscape, contributing to a holistic assessment.

(g) Adaptation and Learning :-

A distinctive feature, adaptability, is unveiled as Bayesian Averaging dynamically adjusts weights based on evolving data landscapes. This section delves into the learning capabilities that empower the technique to stay relevant and effective over time, ensuring resilience in changing prediction environments.

(h) Considerations and Challenges :-

Acknowledging the complexities, this section unravels considerations that shape the robustness of Bayesian Averaging. From addressing outlier resilience to navigating weight assignment intricacies, we explore the challenges inherent in the application of this methodology.

(i) Conclusion :-

Culminating our journey, the conclusion reflects on the rich tapestry of Bayesian Averaging. The theoretical depth, practical applicability, and adaptability showcased throughout the exploration underscore its position as a quietly powerful and flexible predictive fusion technique. Emphasis is placed on its potential to redefine predictive modeling landscapes and inspire further research avenues.

**Keywords** :: Bayesian Averaging, Ensemble Methods, Predictive Modeling, Bayesian Statistics, Weight Assignment, Probabilistic Framework, Machine Learning, Performance Evaluation, Adaptability.

## **6.4 Dynamics in Ensemble Methods :-**

In the ever-evolving landscape of machine learning, this paper delves into a nuanced exploration of adaptive ensemble methods, specifically dynamic ensemble techniques. These methods exhibit a notable ability to adapt to shifting data distributions and fluctuating model performances. This research seeks to unravel the intricate dynamics governing ensemble learning, providing a comprehensive understanding of the design principles, applications, challenges, and future trajectories of dynamic ensemble methodologies.

### **(a) Introduction :-**

The introductory section establishes the contextual significance of dynamic ensemble methods within the broader domain of predictive modeling. It underscores the imperative for adaptability in ensembles, aligning them with the dynamic nature inherent in real-world data and model performance.

### **(b) Adaptability Factors :-**

#### **(i) Data Drift :-**

The narrative unfolds around the concept of data drift, emphasizing the integral role of dynamic ensemble methods in gracefully adapting to shifts in the underlying data distribution. Strategies for discerning and tactfully responding to data drift phenomena are meticulously explored.

#### **(ii) Model Performance :-**

This section critically examines the adaptability landscape concerning the undulating terrain of individual model performances. Strategies for dynamically recalibrating weights or significance based on recent performance trends form a central focus.

(iii) Concept Drift :-

The section immerses itself in the intricate dance with concept drift, where the dynamic interplay between input features and target variables evolves over time. Strategies for uncovering and navigating the impact of concept drift within dynamic ensembles take center stage.

(c) Dynamic Ensemble Strategies :-

(i) Weighted Average Adjustment :-

The paper delves into the nuanced artistry of weighted average adjustment, where dynamic ensembles gracefully tweak the influence of individual models based on their recent prowess, ensuring a harmonious ensemble symphony.

(ii) Model Inclusion/Exclusion :-

The strategic choreography of dynamically ushering models in or out of the ensemble based on their performance unfolds. Algorithms orchestrating the seamless replacement of under-performing models with newer, more relevant counterparts are unveiled.

(iii) Threshold-Based Adaptation :-

A symphony of threshold-based mechanisms emerges as triggers for dynamic adaptability within ensembles. Models languishing below predetermined performance thresholds become subjects of adjustment or graceful exit from the ensemble stage.

(d) Applications :-

The application section transcends the theoretical realm, offering a captivating exploration of dynamic ensemble methodologies across diverse landscapes. Real-world case studies in financial forecasting, sensor networks, and online learning systems breathe life into the efficacy of these adaptive techniques.

(e) Challenges :-

A reflective lens is cast upon the challenges that cloak the effective design and implementation of dynamic ensemble methods. The considerations range from deciphering apt adaptation criteria to navigating the subtleties of overfitting risks and orchestrating computational intricacies entwined with real-time adaptability.

(f) Research Trends :-

This section serves as a compass, navigating the currents of current research trends. It uncovers the ongoing exploration of novel adaptation strategies, the integration of deep learning intricacies into dynamic ensembles, and the persistent quest to conquer scalability and real-time adaptability hurdles.

(g) Conclusion :-

As the curtains draw close, the conclusion reflects upon the symphonic contributions of dynamic ensemble methods in navigating the ever-changing tapestry of predictive modeling. A spotlight is cast on their potential to redefine the adaptive narrative in ensemble learning, offering a resilient framework for navigating the challenges of non-stationary data environments.

**Keywords** :: Dynamic Ensemble Methods, Predictive Modeling, Adaptability, Data Drift, Model Performance, Concept Drift, Weighted Average, Threshold-Based Adaptation, Case Studies, Challenges, Research Trends.

## 6.5 Weighted Average Ensembler Method :-

In the intricate landscape of predictive modeling, a captivating ensemble methodology takes center stage — the Weighted Average Ensemble. This paper embarks on a meticulous journey, unraveling the subtleties of a method that orchestrates the symphony of predictions from diverse models. Here, each model's contribution is artfully weighted, introducing an element of sophistication to the ensemble. Our objective is to meticulously uncover the principles, the dynamic adaptability, the real-world applications, the nuanced challenges, and the profound implications that the Weighted Average Ensemble method bestows upon predictive modeling.

### (a) Introduction :-

The overture sets the stage, framing our exploration within the broader canvas of predictive modeling. The Weighted Average Ensemble method emerges as a strategic conductor orchestrating the collective wisdom of an ensemble of diverse models.

### (b) Weight Assignment :-

This segment scrutinizes the alchemy of weight assignment, where each contributing model is adorned with a weight. Diverse criteria, ranging from historical prowess to cross-validation metrics, paint a canvas of adaptability and sophistication.

### (c) Weighted Averaging Process :-

#### (i) Prediction from Each Model :-

The prelude to the ensemble journey witnesses each model independently crafting predictions for the awaiting input data.

(ii) Weight Assignment :-

The ensemble's ballet unfolds as weights are bestowed upon each model, a choreography that embraces adaptability and dynamic responsiveness. The section delves into the strategies governing the ballet of dynamic weight assignment.

(iii) Weighted Averaging :-

The crescendo arrives as predictions, akin to musical notes, are harmoniously weighted and aggregated. The mathematical composition reveals the symphonic simplicity that underlies the Weighted Average Ensemble method.

(d) Adaptability and Learning :-

In this movement, we unravel the adaptive threads woven into the Weighted Average Ensemble method. Dynamic weight adjustments resonate with a learning melody, allowing the ensemble to dance with grace in response to the evolving rhythms of data distribution.

(e) Applications :-

The ensemble takes center stage in a variety of vignettes. Case studies unfold, each narrating the method's prowess in financial forecasting, weather prediction, and intricate classification tasks, where models shine in different facets of the performance.

(f) Challenges :-

The shadows of challenges are unveiled. The dance floor of weight selection intricacies and the potential echoes of model correlation provide a backdrop to the ensemble's performance, a nuanced interplay of considerations.

(g) Conclusion :-

The overture sets the stage, framing our exploration within the broader canvas of predictive modeling. The Weighted Average Ensemble method emerges as a strategic conductor orchestrating the collective wisdom of an ensemble of diverse models.

# Chapter 7

## Evaluation Metrics

### 7.1 Area Under Curve (AUC) :-

AUC stands for "Area Under the Receiver Operating Characteristic curve." It is a popular metric used to evaluate the performance of a binary classification model. The ROC curve is a graphical representation of the trade-off between true positive rate (sensitivity) and false positive rate (1-specificity) across different threshold values. The AUC is the area under this ROC curve.

Here's a breakdown of the key terms :-

(a) **True Positive Rate (Sensitivity)** :: The proportion of actual positive instances that are correctly identified by the model.

$$\text{TPR / Recall / Sensitivity} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

(b) **False Positive Rate (1-Specificity)** :: The proportion of actual negative instances that are incorrectly classified as positive by the model.

$$\text{FPR} = 1 - \text{Specificity}$$

$$= \frac{\text{FP}}{\text{TN} + \text{FP}}$$

The ROC curve is created by plotting the true positive rate against the false positive rate at various threshold settings. A model with better discrimination has a higher true positive rate and a lower false positive rate, resulting in an ROC curve that hugs the upper left corner of the graph.

The AUC provides a single scalar value that summarizes the performance of a classification model across different threshold settings. A model with an AUC of 1.0 is perfect, while a model with an AUC of 0.5 is no better than random guessing.

The Area Under the ROC Curve (AUC-ROC) summarizes the performance of a binary classification model across all possible classification thresholds. A model with perfect discriminatory ability has an AUC-ROC of 1, indicating that it achieves high TPR while keeping FPR low across all thresholds. A random classifier, on the other hand, has an AUC-ROC of 0.5, forming a diagonal line from the bottom left to the top right.

The Receiver Operating Characteristic (ROC) curve is a graphical representation that illustrates the trade-off between sensitivity and specificity at various classification thresholds. It is particularly useful for assessing the performance of a binary classification model across a range of decision thresholds. The curve is created by plotting the True Positive Rate (TPR) against the False Positive Rate (FPR) at different classification thresholds.

In general, a higher AUC indicates a better overall performance of the model in distinguishing between positive and negative instances. It is a commonly used metric in machine learning for binary classification problems, especially when the class distribution is imbalanced.

Note: There are different variations of the ROC curve and AUC, including precision-recall curves and AUC-PR. The choice of metric depends on the specific characteristics of the classification problem.

It's important to mention that AUC-ROC is commonly used for binary classification problems, where there are two classes (positive and negative). For multi-class problems, there are extensions like micro-average, macro-average, and one-vs-all AUC.

## 7.2 Classification Accuracy :-

In the context of machine learning, classification accuracy (CA) is a common evaluation metric used to assess the performance of a classification model. It measures the proportion of correctly classified instances out of the total instances. However, accuracy alone may not be sufficient in all situations, especially when dealing with imbalanced datasets.

Here are some commonly used evaluation metrics, along with their formulas :-

(a) Accuracy :-

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

TP (True Positives) is the number of instances that are correctly predicted as positive.

TN (True Negatives) is the number of instances that are correctly predicted as negative.

FP (False Positives) is the number of instances that are incorrectly predicted as positive.

FN (False Negatives) is the number of instances that are incorrectly predicted as negative.

Precision ( $P$ )	$\frac{\text{TP}}{\text{TP}+\text{FP}}$
Recall ( $R$ )	$\frac{\text{TP}}{\text{TP}+\text{FN}}$
F1-score	$2 \times \frac{P \times R}{P + R}$
Accuracy	$\frac{\text{TP}+\text{TN}}{\text{TP}+\text{TN}+\text{FP}+\text{FN}}$

Accuracy may not be the best metric in situations with imbalanced classes. In such cases, other metrics like Precision, Recall, and F1 Score can provide a more informative assessment of a model's performance.

### 7.3 F1 Score :-

The F1 score is a metric that combines both precision and recall into a single value. It is particularly useful when there is an uneven class distribution (imbalanced classes) in a dataset.

The formula for calculating the F1 score is as follows :

$$\text{F1 Score} = \frac{2}{\left( \frac{1}{\text{Precision}} + \frac{1}{\text{Recall}} \right)}$$
  

$$\text{F1 Score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

The numerator ( $2 * \text{Precision} * \text{Recall}$ ) emphasizes cases where both Precision and Recall are high. The multiplication by 2 is used to give equal weight to both Precision and Recall.

The denominator (Precision + Recall) is the sum of Precision and Recall. This part of the formula ensures that the F1 score is high only when both Precision and Recall are high.

By combining Precision and Recall in this way, the F1 score provides a balance between the two metrics. It ranges from 0 to 1, where 1 indicates perfect precision and recall, and 0 indicates the worst possible trade-off between precision and recall. The F1 score is particularly useful when you want to find a balance between false positives and false negatives, and it is commonly used in situations where class distribution is imbalanced.

## 7.4 Precision :-

Precision is an evaluation metric that assesses the accuracy of positive predictions made by a classification model. It is particularly relevant when the cost of false positives (incorrectly predicting a positive instance) is high. Precision is defined as the ratio of true positive predictions to the total predicted positive instances.

The formula for precision is as follows :

$$\text{Precision} = \frac{\frac{TP}{TP + FP}}{TP}$$

- (a) High Precision :: A high precision indicates that the model is accurate when it predicts positive instances. There are fewer false positives relative to the total number of predicted positives.
- (b) Low Precision :: A low precision suggests that the model may be generating a significant number of false positives, meaning that a substantial portion of the predicted positive instances are incorrect.
- (c) Precision-Recall Trade-off :: Precision is often used in conjunction with recall. There is typically a trade-off between precision and recall; increasing one may lead to a decrease in the other. It depends on the specific requirements and priorities of the problem at hand.
- (d) Imbalanced Datasets :: Precision is particularly useful when dealing with imbalanced datasets, where the number of negative instances far exceeds the number of positive instances. In such cases, a high precision is often more important than a high recall.

In summary, precision is a crucial metric when the focus is on minimizing false positives. It provides insight into the reliability of positive predictions made by a model, which is valuable in various applications such as fraud detection, medical diagnoses, and other scenarios where false positives have significant consequences.

## 7.5 Recall :-

Recall, also known as Sensitivity or True Positive Rate, is an evaluation metric that measures the ability of a classification model to capture and correctly identify all relevant instances of a positive class. It is particularly important when the cost of false negatives (missing positive instances) is high.

The formula for recall is as follows :-

$$\text{Recall} = \frac{TP}{TP + FN}$$

- (a) High Recall: A high recall indicates that the model is effective at capturing most of the positive instances. There are fewer false negatives relative to the total number of actual positives.
- (b) Low Recall :: A low recall suggests that the model may be missing a significant number of positive instances, as there are more false negatives relative to the total number of actual positives.
- (c) Precision-Recall Trade-off :: Recall is often used in conjunction with precision. There is typically a trade-off between precision and recall; increasing one may lead to a decrease in the other. It depends on the specific requirements and priorities of the problem at hand.
- (d) Imbalanced Datasets :: Recall is particularly important when dealing with imbalanced datasets, where the number of positive instances is much smaller than the number of negative instances. In such cases, ensuring that positive instances are not missed is crucial.

In summary, recall is a crucial metric when the goal is to identify as many positive instances as possible, even at the cost of a higher number of false positives. It is commonly used in applications such as disease detection, where missing a positive case could have serious consequences.

## 7.6 Matthews Correlation Coefficient (MCC) :-

The Matthews Correlation Coefficient (MCC) is a metric that takes into account true positives, true negatives, false positives, and false negatives to evaluate the performance of a binary classification model. It is particularly useful when dealing with imbalanced datasets.

The formula for Matthews Correlation Coefficient is as follows :

$$\text{MCC} = \frac{\text{TP} \times \text{TN} - \text{FP} \times \text{FN}}{\sqrt{(\text{TP} + \text{FP})(\text{TP} + \text{FN})(\text{TN} + \text{FP})(\text{TN} + \text{FN})}}$$

(a) The MCC formula produces a value between -1 and 1, where:

- (i) 1 indicates perfect prediction,
- (ii) 0 indicates random prediction,
- (iii) -1 indicates total disagreement between prediction and observation.

(b) Interpretation :: A higher MCC value suggests better performance, while a lower value indicates poorer performance.

(c) Balancing Act :: MCC takes into account all four values (TP, TN, FP, FN), providing a balanced measure that is particularly useful when dealing with imbalanced datasets.

(d) Importance of Zero :: An MCC value of 0 suggests that the model's predictions are no better than random, and negative values indicate a disagreement worse than random predictions.

(e) Applications :: MCC is commonly used in genomics and bioinformatics, but it can be applied in various domains where imbalanced datasets are prevalent.

In summary, Matthews Correlation Coefficient is a comprehensive metric that considers both the positive and negative predictions of a model, providing a balanced assessment of its performance, especially in situations with imbalanced class distributions.

## 7.7 Specificity :-

Specificity, also known as the True Negative Rate, is an evaluation metric that measures the ability of a classification model to correctly identify negative instances. It is particularly important when the cost of false positives (incorrectly predicting a positive instance) is high.

The formula for specificity is as follows :

$$\text{Specificity} = \frac{\text{True Negatives}}{\text{True Negatives} + \text{False Positives}}$$

- (a) High Specificity :: A high specificity indicates that the model is effective at correctly identifying negative instances. There are fewer false positives relative to the total number of actual negatives.
- (b) Low Specificity :: A low specificity suggests that the model may be incorrectly predicting negative instances, as there are more false positives relative to the total number of actual negatives.
- (c) Precision-Recall-Specificity Trade-off :: Specificity is often used in conjunction with precision and recall. There is typically a trade-off between these metrics; increasing one may lead to a decrease in another. It depends on the specific requirements and priorities of the problem at hand.
- (d) Imbalanced Datasets :: Specificity is particularly relevant when dealing with imbalanced datasets, where the number of negative instances is much larger than the number of positive instances.

In summary, specificity is a crucial metric when the focus is on minimizing false positives, especially in situations where the consequences of false positives are significant. It provides insight into the model's ability to correctly identify instances of the negative class.

## 7.8 LogLoss :-

Logarithmic Loss, often referred to as Log Loss or Cross-Entropy Loss, is a commonly used evaluation metric for binary and multiclass classification problems. It measures the performance of a classification model by penalizing predictions that are far from the true labels.

The formula for Logarithmic Loss is as follows :

$$\text{LogLoss} = \frac{-1}{N} \sum_{i=1}^N \sum_{j=1}^M x_{ij} * \log(p_{ij})$$

- $N$  is the number of instances in the dataset.
- $M$  is the number of classes in the multiclass scenario.
- $y_i$  is the true label of instance  $i$  (1 for the positive class, 0 for the negative class in binary classification; a one-hot encoded vector in multiclass classification).
- $p_i$  is the predicted probability of the positive class for instance  $i$  in binary classification.
- $p_{ij}$  is the predicted probability of instance  $i$  belonging to class  $j$  in multiclass classification.

- (a) Interpretation :: Log Loss is a logarithmic function that measures the average negative log-likelihood of the true labels given the predicted probabilities. The goal is to minimize Log Loss, and lower values indicate better model performance.
- (b) Probabilistic Nature :: Log Loss is sensitive to the predicted probabilities. It penalizes confident but wrong predictions more severely.
- (c) Binary vs. Multiclass :: The formulas are slightly different for binary and multiclass classification, but the underlying idea is the same.
- (d) Comparison to Other Metrics :: Log Loss is widely used in scenarios where probabilistic interpretations of the predictions are important. It is often used alongside other metrics like accuracy, precision, and recall for a more comprehensive evaluation.

# Chapter 8

## Intrusion Detection Evaluation Dataset (CIC-IDS2017)

### 8.1 Introduction :-

The CIC-IDS2017 Intrusion Detection Evaluation Dataset addresses the pressing need for reliable testing and validation datasets for Intrusion Detection Systems (IDSs) and Intrusion Prevention Systems (IPSs). Previous datasets, dating back to 1998, were found to be outdated and unreliable, lacking diversity in traffic, coverage of attacks, feature sets, and metadata. To address these issues, the creators of CIC-IDS2017 have developed a dataset that not only covers current and prevalent threats but also focuses on generating realistic background traffic. This emphasis on realism is crucial for accurately evaluating the performance of anomaly-based intrusion detection systems.

This dataset encompasses all the essential features required for creating a trustworthy benchmark dataset. It includes a complete network design, labeled datasets, entire traffic exchanges, and a range of attacks. The dataset incorporates network traffic from various operating systems, capturing comprehensive interactions within and between internal LANs and internet connectivity. Notably, it covers protocols such as HTTP, HTTPS, FTP, SSH, and email, showcasing a high level of variability.

The information is meticulously detailed, including time, types of attacks, flows, and labels. The dataset spans five days, starting on July 3, 2017, featuring benign activity on Monday and a series of attacks on Tuesday through Friday.

The dataset documentation provides insights into victim and attacker networks, including firewalls, DNS+DC servers, and both external (attackers) and internal (victims) workstations, along with technical specifics. Various types of attacks, such as Brute Force, DoS/DDoS, Heartbleed, Web Attack, Infiltration, and Botnet, are classified. Each attack includes information about the attacker, the victim, and the network address translation (NAT) procedure of the firewall.

## 8.2 Data Pre-processing :-

The CICIDS dataset records multiple network assaults by compiling network data from various days and hours. This study describes the thorough preprocessing pipeline used to convert the raw data into a format suitable for robust machine learning analysis.

### (a) Data Combination :-

The dataset was carefully merged from various instances into a single framework. This aggregation enables a comprehensive investigation of network behaviors across several time settings.

### (b) Label Conversion :-

To make the categorization work easier, the labels were converted to binary format. Instances that were benign were labelled as 0, whereas instances that were malevolent were labelled as 1. This binary mapping sharpens machine learning models' attention by emphasizing the contrast between regular and harmful network behavior.

### (c) Data Normalization :-

Addressing the issue of varying scales across features, a data normalization step was employed. This ensures that all features adhere to a comparable range, preventing larger magnitude features from unduly influencing the modeling process.

- Infinite values in the feature matrices ( $X_{train}$  and  $X_{test}$ ) are replaced with a large finite value using NumPy.
- Missing values in both training and test sets are imputed using the mean of each feature, accomplished through scikit-learn's SimpleImputer. SimpleImputer is a class from scikit-learn that provides a simple strategy for imputing missing values in a dataset.

- Using scikit-learn's StandardScaler, the feature matrices are standardized to have a zero mean and unit variance. This tackles changing scaling across features, ensuring that machine learning models receive consistent input. Another scikit-learn class is the StandardScaler, which is used to standardize features by eliminating the mean and scaling to unit variance. This is an important step in data preparation, particularly when working with machine learning models that are sensitive to the magnitude of input characteristics.

(d) Train-Test Split :-

An 80-20 train-test split was achieved by a planned segmentation method. This split enables machine learning models to learn patterns from training data and evaluate their generalization skills on previously encountered cases in the test set.

### 8.3 Implementation :-

#### Step-1 : Import Libraries :-

```
## STEP 1
import pandas as pd
import numpy as np
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.metrics import accuracy_score, roc_auc_score, log_loss, confusion_matrix, roc_curve, auc
import matplotlib.pyplot as plt
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
import time
import joblib
```

Here, the necessary libraries are imported, including popular machine learning libraries like pandas, numpy, scikit-learn, and matplotlib.

#### Step-2 : Load and Split Dataset :-

```
## STEP 2
def load_and_split_dataset(file_path):
    df = pd.read_csv(file_path)
    X = df.drop('Label', axis=1).values
    y = df['Label'].values
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
    return X_train, X_test, y_train, y_test

# Specify the path to the dataset CSV file
dataset_path = DATA
```

This step defines a function `load_and_split_dataset` that reads a CSV file, separates features (X) and labels (y), and splits the data into training and testing sets.

### Step-3 : Data Preprocessing :-

```
## STEP 3
# Load and split the dataset
X_train, X_test, y_train, y_test = load_and_split_dataset(dataset_path)
# Handle infinite and large values
X_train[np.isinf(X_train)] = np.finfo('float64').max # Replace infinite values with a large finite value
X_test[np.isinf(X_test)] = np.finfo('float64').max # Replace infinite values with a large finite value
|
# Handle missing values
imputer = SimpleImputer(strategy='mean')
X_train = imputer.fit_transform(X_train)
X_test = imputer.transform(X_test)

# Scale the data
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

This step handles infinite values in the dataset, replaces them with a large finite value, and then performs imputation for missing values using the mean strategy. Finally, it standardizes the data using StandardScaler.

### Step-4 : Initialize and Train Classifier :-

#### (i) Adaboost :-

```
start_train = time.time()

# Create an AdaBoost classifier
classifier = AdaBoostClassifier()

# Train the classifier on the training data
classifier.fit(X_train, y_train)

end_train = time.time()

print("Training time: {:.2f} seconds".format(end_train - start_train))
```

#### (ii) Decision Tree :-

```
start_train = time.time()

# Create a Decision Tree classifier
classifier = DecisionTreeClassifier()

# Train the classifier on the training data
classifier.fit(X_train, y_train)

end_train = time.time()

print("Training time: {:.2f} seconds".format(end_train - start_train))
```

(iii) Random Forest :-

```
# Define the hyperparameters
n_estimators = 192
max_depth = 6
min_samples_split = 6
random_state = 96

# Create a Random Forest classifier
classifier = RandomForestClassifier(
    n_estimators=n_estimators,
    max_depth=max_depth,
    min_samples_split=min_samples_split,
    random_state=random_state
)
```

(iv) Gradient Boosting :-

```
## STEP 4
start_train = time.time()
# Create a Gradient Boosting classifier
classifier = GradientBoostingClassifier()

# Train the classifier on the training data
classifier.fit(X_train, y_train)

end_train = time.time()

print("Training time: {:.2f} seconds".format(end_train - start_train))

start_test = time.time()
```

(v) Extremely Randomized Tree :-

```
start_train = time.time()

# Create an Extra Trees classifier
classifier = ExtraTreesClassifier()

# Train the classifier on the training data
classifier.fit(X_train, y_train)

end_train = time.time()

print("Training time: {:.2f} seconds".format(end_train - start_train))
```

(vi) Neural Network :-

```
# Create a neural network classifier
model = Sequential()
model.add(Dense(128, activation='relu', input_dim=X_train.shape[1]))
model.add(Dense(64, activation='relu'))
model.add(Dense(1, activation='sigmoid'))

# Compile the model
model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

start_train = time.time()

# Train the classifier on the training data
model.fit(X_train, y_train, epochs=10, batch_size=32, verbose=1)

end_train = time.time()

print("Training time: {:.2f} seconds".format(end_train - start_train))
```

The above classifiers are initialized and trained on the preprocessed training data.

### Step-5 : Make Predictions on Testing Data :-

```
## STEP 5
# Make predictions on the testing data
predictions = classifier.predict(X_test)

end_test = time.time()

print("Testing time: {:.2f} seconds".format(end_test - start_test))
```

The trained classifiers are used to make predictions on the testing data.

### Step-6 : Measure Accuracy :-

```
## STEP 6
accuracy = accuracy_score(y_test, predictions)
print("Accuracy: {:.4f}".format(accuracy))
```

The accuracy of the model on the testing data is calculated and printed.

### Step-7 : Confusion Matrix :-

```
## STEP 7
# Calculate the confusion matrix
conf_matrix = confusion_matrix(y_test, predictions)
# Print the confusion matrix
print("Confusion Matrix:")
print(conf_matrix)
```

The confusion matrix is calculated and printed.

### Step-8 : Evaluation Metrics (AUC and Log Loss) :-

```
## STEP 8
# Calculate evaluation metrics
auc = roc_auc_score(y_test, predictions)
print("AUC: {:.4f}".format(auc))
logloss = log_loss(y_test, classifier.predict_proba(X_test))
print("Log Loss: {:.4f}".format(logloss))
```

The area under the ROC curve (AUC) and log loss are calculated and printed.

### Step-9 : Plot ROC Curve :-

```
## STEP 9
# Plot ROC Curve for Neural Network Model
y_pred_proba = classifier.predict(X_test)
fpr, tpr, thresholds = roc_curve(y_test, y_pred_proba)
roc_auc_value = auc(fpr, tpr)

plt.figure()
plt.plot(fpr, tpr, color='darkorange', lw=2)
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlim([-0.1, 1.1])
plt.ylim([-0.1, 1.1])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic')
plt.legend(loc="lower right")
plt.show()
```

The ROC curve is plotted to visualize the classifier's performance.

### Step-10 : Save the Trained Model :-

```
## STEP 10
# Save the trained classifier model to a file
model_filename = r"C:\Users\KIIIT\Desktop\SEM-VII\Major Project\CICIDS\GB.pkl"
joblib.dump(classifier, model_filename)
```

The trained classifier is saved to a pickled file using joblib library.

### 8.4 Table:-

**Table-01 :: Individual Model Training for CICIDS-17**

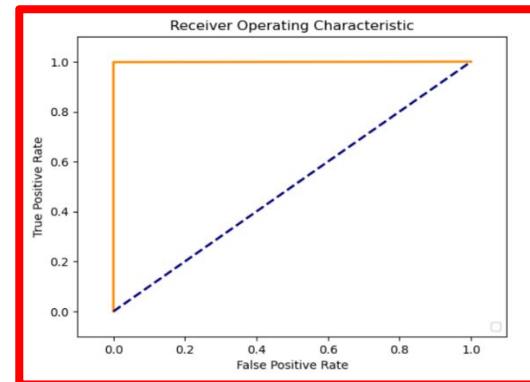
Model	AUC	CA	F1	Precision	Recall	MCC	Spec	LogLoss
AdaBoost	0.9991	99.96	0.9998	1.0000	0.9996	0.9989	1.0000	0.4588
Decision Tree	0.9999	99.99	0.9999	0.9999	1.0000	0.9997	0.9997	0.0029
Random Forest	0.9743	98.93	0.9934	0.9991	0.9877	0.9660	0.9961	0.0556
Gradient Boosting	0.9981	99.92	0.9995	1.0000	0.9991	0.9976	0.9998	0.0064
Extremely Randomized Tree	0.9997	99.99	0.9999	1.0000	0.9999	0.9995	0.9998	0.0007
Neural Network	0.9986	99.91	0.9995	0.9995	0.9994	0.9973	0.9980	-

## 8.5 Results :-

### ➤ Confusion Matrix and ROC Curve :-

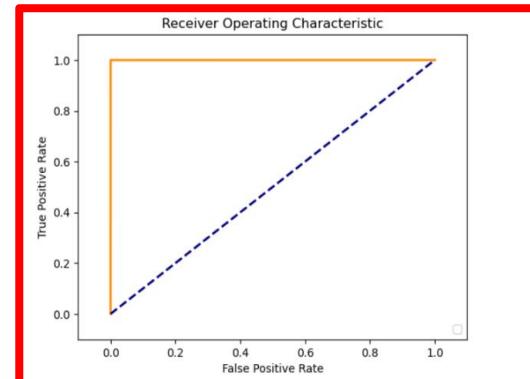
#### (a) Adaboost :-

454656 80.31%	2 0.00%
203 0.04%	111288 19.66%



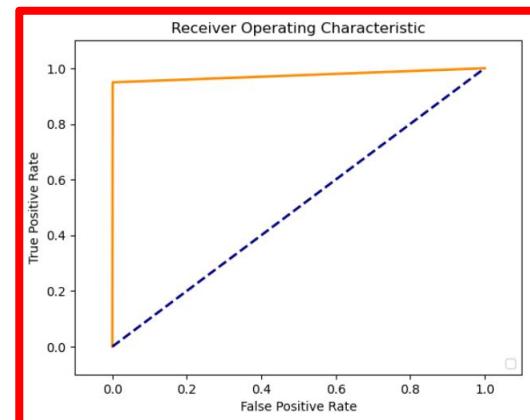
#### (b) Decision Tree :-

454622 80.30%	36 0.01%
13 0.00%	111478 19.69%



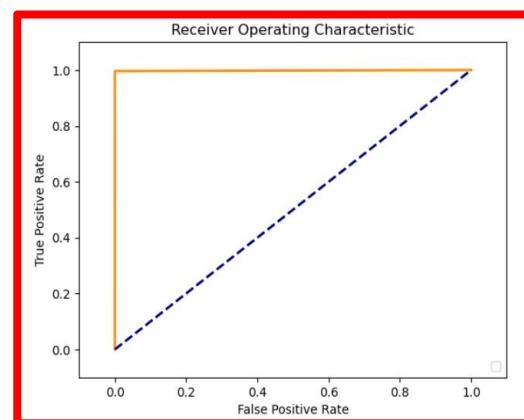
#### (c) Random Forest :-

454246 80.23%	412 0.07%
5640 1.00%	105851 18.70%

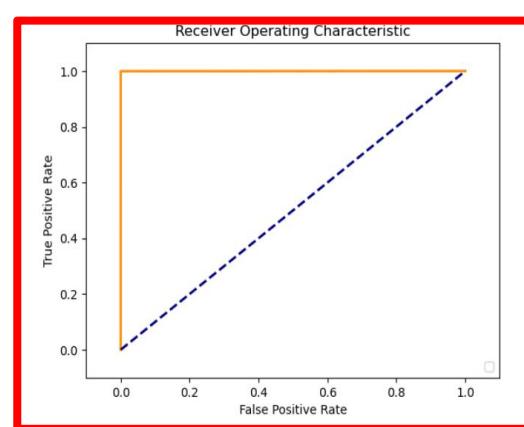


(d) Gradient Boosting :-

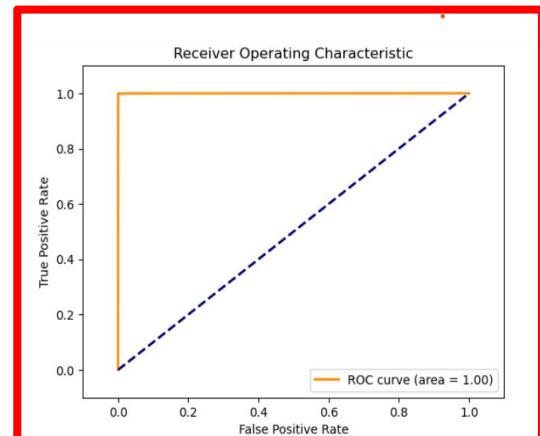
454639 80.30%	19 0.00%
419 0.07%	111072 19.62%

(e) Extremely Randomized Tree :-

454449 80.27%	209 0.04%
313 0.06%	111178 19.64%

(f) Neural Network :-

454627 80.30%	31 0.01%
60 0.01%	111431 19.68%



## 8.6 Ensemble Technique :-

The code has two main functions: ensemble\_voting and main.

### (i) ensemble\_voting :-

The ensemble\_voting function takes two arguments: models and data. The models argument is a list of trained machine learning models, such as classifiers or regressors. The data argument is a numpy array of features for which we want to make predictions.

```
def ensemble_voting(models, data):
    predictions = []
    for model in models:
        prediction = model.predict(data)
        predictions.append(prediction.astype(int)) # Convert predictions to integers
    final_predictions = np.apply_along_axis(lambda x: np.argmax(np.bincount(x)), axis=0, arr=predictions)
    return final_predictions
```

The function iterates over each model in the models list and calls its predict method on the data. It stores the predictions as integers in a list called predictions. Then, it applies a numpy function called apply\_along\_axis on the predictions list. This function applies another function, in this case lambda x: np.argmax(np.bincount(x)), along a given axis, in this case 0, of an array. The lambda function takes an array of predictions from different models for one data point and returns the index of the most frequent prediction. This is equivalent to finding the mode or the majority vote of the predictions. The result of this function is a numpy array of final predictions based on ensemble voting.

### (ii) main :-

The main function takes two arguments: dataset\_path and model\_paths. The dataset\_path argument is a string that represents the path to a CSV file that contains the dataset. The model\_paths argument is a list of strings that represent the paths to the saved models.

The function first loads the dataset using pandas and reads it as a dataframe. Then, it extracts the features (X) and labels (y) from the dataframe by dropping the column named 'Label' and assigning it to y, and keeping the rest of the columns as X.

Next, it splits the dataset into training and testing sets using sklearn's train\_test\_split function. It uses 20% of the data as the test set and 80% as the training set. It also sets a random state of 42 for reproducibility.

```

def main(dataset_path, model_paths):
    # Load the dataset
    dataset = pd.read_csv(dataset_path)

    # Extract features (X) and Labels (y)
    X = dataset.drop('Label', axis=1)
    y = dataset['Label']

    # Split the dataset into training and testing sets
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

```

Then, it handles infinite and large values in the data by replacing them with a large finite value using numpy's finfo function. This is done to avoid numerical errors or overflows when scaling or modeling the data.

After that, it handles missing values in the data by replacing them with the mean value of each column using sklearn's SimpleImputer class. This is done to avoid errors or biases when modeling the data.

Next, it scales the data using sklearn's StandardScaler class. This class standardizes the data by subtracting the mean and dividing by the standard deviation of each column. This is done to improve the performance and accuracy of some models that are sensitive to the scale or range of the data.

```

# Handle infinite and large values
X_train[X_train.isinf()] = np.finfo('float64').max # Replace infinite values with a large finite value
X_test[X_test.isinf()] = np.finfo('float64').max # Replace infinite values with a large finite value

# Handle missing values
imputer = SimpleImputer(strategy='mean')
X_train = imputer.fit_transform(X_train)
X_test = imputer.transform(X_test)

# Scale the data
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

```

Then, it loads the ensemble models using a custom function called load\_models that takes a list of model paths and returns a list of model objects.

Finally, it makes predictions using ensemble voting on the test set by calling the ensemble\_voting function with the models and X\_test as arguments. It stores the final predictions in a variable called final\_predictions.

Along with it then proceeds on to evaluate the ensemble model created by the culmination of best 5 models. It evaluates the model by making prediction on the test data and using the metrics module of sklearn library and prints the confusion matrix, accuracy and logloss of the model. And then finally it calculates the probabilities of all its predictions and uses it to plot the ROC curve.

```

# Load the ensemble models
models = load_models(model_paths)

# Make predictions using ensemble voting on the test set
final_predictions = ensemble_voting(models, X_test)

# Generate the confusion matrix
cm = confusion_matrix(y_test, final_predictions)

# Calculate the accuracy
accuracy = accuracy_score(y_test, final_predictions)
logloss = log_loss(y_test, models[0].predict_proba(X_test))

# Print the confusion matrix
print("Confusion Matrix:")
print(cm)
print("Accuracy: {:.8f}".format(accuracy))
print("Log Loss: {:.4f}".format(logloss))

# Calculate the predicted probabilities for ensemble voting
ensemble_probs = np.mean([model.predict_proba(X_test) for model in models], axis=0)

# Plot the ROC curve for ensemble predictions
plot_roc_curve(y_test, ensemble_probs[:, 1])

if __name__ == "__main__":
    dataset_path = DATA # Replace with the actual path to your dataset
    model_paths = [MODEL_NN, MODEL_DT, MODEL_ET, MODEL_AdaB, MODEL_GB] # Replace with the paths to your .pkl models
    main(dataset_path, model_paths)

```

Final results of the proposed ensemble model are discussed in the Results section of the report along with other datasets' models' results.

# Chapter 9

## NSL-KDD Dataset

### 9.1 Introduction :-

The NSL-KDD dataset was created to overcome limitations observed in the well-known KDD'99 dataset for intrusion detection. Given the scarcity of publicly available datasets for network-based intrusion detection systems, the NSL-KDD dataset serves as a valuable benchmark for assessing different intrusion detection approaches. It builds upon the original KDD dataset by eliminating redundant records in the training set, excluding duplicate entries from the test sets, and introducing a wider range of difficulty levels to enhance the evaluation of machine learning algorithms.

This dataset provides complete training and test sets in both ARFF and TXT formats, along with 20% selections from the original files. The absence of duplicate records in the training set ensures that classifiers are not biased towards more frequently occurring records. Similarly, the test sets, free of duplicate entries, prevent bias towards approaches with higher detection rates on common records. The dataset's composition, with the number of chosen records inversely related to the percentage in the original KDD dataset, promotes a more diverse categorization rate among various machine learning algorithms.

An examination of the original KDD dataset reveals the issue of duplicate records, leading to biases in learning systems. The NSL-KDD dataset addresses this challenge by significantly reducing the number of duplicated entries, thereby increasing the overall reduction rate. Experiments on the NSL-KDD dataset demonstrate the proper identification of a significant number of records in the initial KDD train and test sets. Notably, the reduction in duplicated records is particularly evident for attack occurrences, with rates of 93.32% in the training set and 88.26% in the test set.

## **9.2 Data Pre-processing :-**

A systematic preparation pipeline was used to assure data quality, homogeneity, and compatibility with several machine learning algorithms in order to perform a complete examination of network intrusion detection using the NSL-KDD dataset. To prepare the dataset for robust analysis, the following procedures were taken:

### **(a) Format Conversion :-**

The NSL-KDD dataset was first made available via the ARFF (Attribute-Relation File Format), which is typically connected with the WEKA machine learning programme. The dataset was seamlessly converted to the commonly used CSV (Comma-Separated Values) format to improve accessibility and compatibility with a wider range of tools. This change not only improved usability but also cleared the way for more efficient integration into other machine learning frameworks.

### **(b) Label Transformation :-**

A critical step in developing a binary classification system suited for intrusion detection was changing the original labels. Instances were classified as benign (0) or malicious (1), corresponding to the overriding objective of separating typical network behavior from prospective invasions. This binary classification scheme provides the groundwork for properly training and assessing machine learning models.

### **(c) Train-Test Split :-**

A careful train-test split was used to verify the accuracy of model assessment. The dataset was partitioned between training and testing sets using an 80-20 partitioning ratio, as is standard practise. This strategic divide permitted model training on a large chunk of the data while reserving a separate sample for unbiased evaluation. This method attempted to strike a compromise between model training and the capacity to generalize to previously unknown data.

#### (d) Data Normalization :-

Recognizing the significance of feature scaling in machine learning, a normalization procedure was used to standardize the scales of various features within the dataset. This phase was critical in order to keep characteristics of variable magnitudes from unnecessarily impacting the model training process. Normalization aided in the steadier convergence of machine learning algorithms, improving the models' general resilience and interpretability.

Essentially, the goal of this extensive preprocessing method was to provide a refined and well-conditioned dataset for further analysis. The dataset was prepared for a thorough examination of machine learning approaches by addressing format compatibility, converting labels, constructing a careful train-test split, and normalizing feature sizes. This foundation laid the platform for the following evaluation of PySpark ML and traditional machine learning methods in the domain of network intrusion detection, guaranteeing a systematic and reliable model assessment methodology.

### 9.3 Implementation :-

#### Step-1 : Import Libraries :-

```
# Step 1: Import Libraries
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.metrics import accuracy_score, roc_auc_score, log_loss, confusion_matrix
```

The code begins by importing essential libraries for data manipulation and machine learning tasks. `pandas` is imported for data handling, `train_test_split` for splitting the dataset, `GradientBoostingClassifier` for creating a gradient boosting model, and various metrics from `sklearn.metrics` for model evaluation.

#### Step-2 : Load the Dataset :-

```
# Step 2: Load the dataset
file_path = DATA
df = pd.read_csv(file_path)
```

The variable file\_path is set to the path of the dataset file (represented by DATA). The pd.read\_csv function from the pandas library is then used to read the dataset into a DataFrame named df.

### Step-3 : Prepare the Data :-

```
# Step 3: Prepare the data
X = df.drop('outcome', axis=1)
y = df['outcome']
```

The features (independent variables) are extracted into a DataFrame X, excluding the target variable, which is assigned to y.

### Step-4 : Split the Data into Training and Testing Sets :-

```
# Step 4: Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

The dataset is divided into training and testing sets using the train\_test\_split function. The split is performed with 80% of the data used for training (X\_train, y\_train) and 20% for testing (X\_test, y\_test). The random\_state parameter ensures reproducibility.

### Step-5 : Create and Train the Model :-

#### (a) Adaboost :-

```
# Step 5: Create and train the AdaBoost model
base_model = DecisionTreeClassifier(max_depth=1) # Weak Learner (e.g., decision tree with max_depth=1)
adaboost_model = AdaBoostClassifier(base_estimator=base_model, n_estimators=50, random_state=42)
adaboost_model.fit(X_train, y_train)
```

#### (b) Decision Tree :-

```
# Step 5: Create and train the Decision Tree model
dt_model = DecisionTreeClassifier(random_state=42)
dt_model.fit(X_train, y_train)
```

(c) Random Forest :-

```
# Step 5: Create and train the Random Forest model
rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
rf_model.fit(X_train, y_train)
```

(d) Gradient Boosting :-

```
# Step 5: Create and train the Gradient Boosting model
gb_model = GradientBoostingClassifier(n_estimators=100, learning_rate=0.1, random_state=42)
gb_model.fit(X_train, y_train)
```

(e) Extremely Randomized Tree :-

```
# Step 5: Create and train the Extra Trees model
et_model = ExtraTreesClassifier(n_estimators=100, random_state=42)
et_model.fit(X_train, y_train)
```

(f) Neural Network :-

```
# Step 5: Create the neural network model, Compile the model and Train the model.
model = Sequential()
model.add(Dense(64, input_dim=X_train.shape[1], activation='tanh'))
model.add(Dense(32, activation='tanh'))
model.add(Dense(1, activation='sigmoid'))
model.compile(loss='binary_crossentropy', metrics=['accuracy'])
model.fit(X_train, y_train, epochs=5, batch_size=32, validation_data=(X_test, y_test))
```

An instance of the Classifier is created with specified hyperparameters (n\_estimators, learning\_rate, random\_state). The model is then trained using the training data (X\_train, y\_train).

Step-6 : Make Predictions on the Test Set :-

```
# Step 6: Make predictions on the test set
y_pred = gb_model.predict(X_test)
```

The trained model (gb\_model) is used to predict the target variable for the test set (X\_test). Predictions are stored in the variable y\_pred.

### Step-7 : Evaluate the Model's Performance :-

```
# Step 7: Evaluate the model's performance
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy}")
```

The accuracy of the model is calculated using the accuracy\_score function from sklearn.metrics and printed.

### Step-8 : Print the Confusion Matrix :-

```
# Step 8: Print the confusion matrix
conf_matrix = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:")
print(conf_matrix)
```

The confusion matrix is generated using the confusion\_matrix function, providing insights into the model's classification performance.

### Step-9 : LogLoss and AUC :-

```
# Step 9: LogLoss and auc
auc = roc_auc_score(y_test, y_pred)
logloss = log_loss(y_test, gb_model.predict_proba(X_test))
print("AUC: {:.4f}".format(auc))
print("Log Loss: {:.4f}".format(logloss))
```

The Area Under the ROC Curve (AUC) and log loss metrics are calculated using the respective functions (roc\_auc\_score and log\_loss). These metrics provide additional insights into the model's performance.

### Step-10 : Plot the ROC Curve :-

```
# Step 10: Plot ROC Curve
fpr, tpr, thresholds = roc_curve(y_test, gb_model.predict_proba(X_test)[:,1])
roc_auc_value = auc(fpr, tpr) # Renamed the variable to avoid naming conflict
plt.figure()
plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area = {:.2f})'.format(roc_auc_value))
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlim([-0.1, 1.1])
plt.ylim([-0.1, 1.1])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic')
plt.legend(loc="lower right")
plt.show()
```

This step involves generating the False Positive Rate (fpr), True Positive Rate (tpr), and thresholds for the Receiver Operating Characteristic (ROC) curve. The curve is then plotted using a combination of libraries like matplotlib and roc\_curve.

Step-11 : Save the Trained Classifier Model to a File :-

```
# Step 11: Save the trained classifier model to a file
import joblib
model_filename = MODEL_GB
joblib.dump(gb_model, model_filename)
```

The trained model (gb\_model) is serialized and saved to a file using the joblib.dump function. The filename is specified by the variable MODEL\_GB. This step allows the model to be reused without retraining.

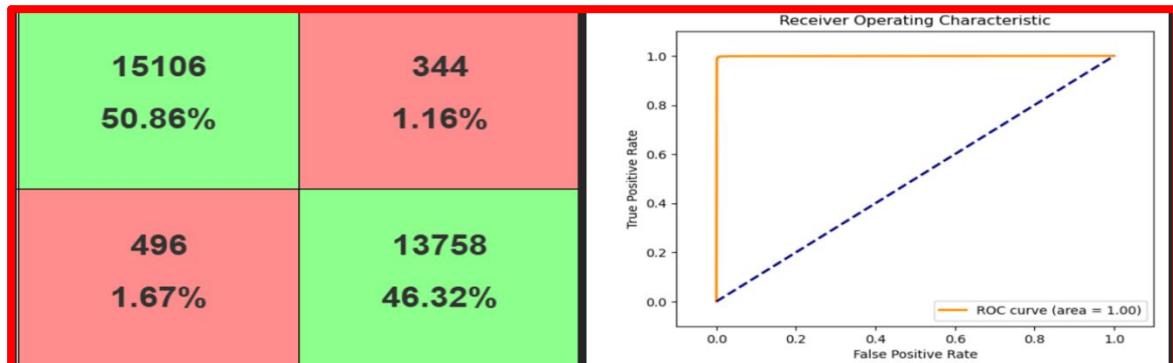
#### 9.4 Table :-

Table-02 :: Individual Model Training for NSL-KDD19								
Model	AUC	CA	F1	Precision	Recall	MCC	Spec	LogLoss
AdaBoost	0.9715	97.17	0.9729	0.9777	0.9682	0.9434	0.9756	0.6275
Decision Tree	0.9946	99.46	0.9948	0.9945	0.9951	0.9892	0.9940	0.1759
Random Forest	0.9956	99.56	0.9958	0.9972	0.9945	0.9913	0.9969	0.0171
Gradient Boosting	0.9879	98.79	0.9884	0.9891	0.9878	0.9759	0.9881	0.0502
Extremely Randomized Tree	0.9948	99.49	0.9951	0.9965	0.9937	0.9898	0.9962	0.0487
Neural Network	0.9231	92.32	0.9263	0.9272	0.9254	0.8462	0.9209	-

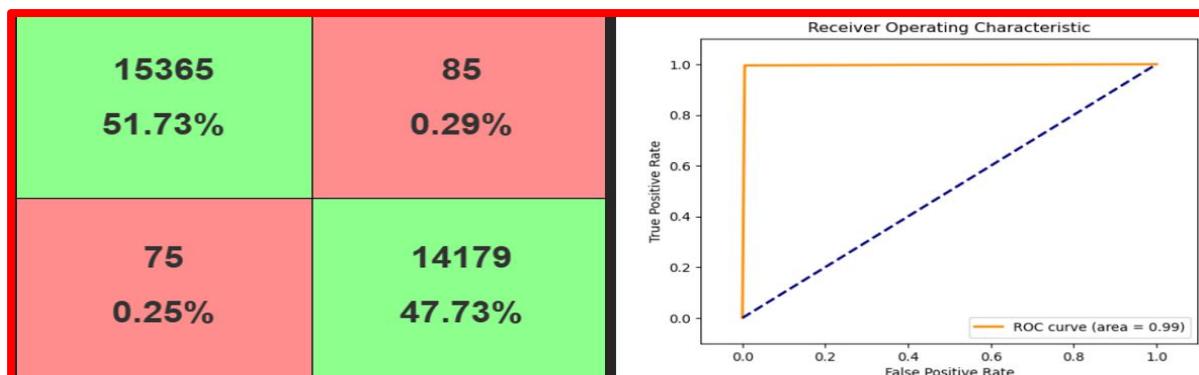
## 9.5 Results :-

### ➤ Confusion Matrix and ROC Curve :-

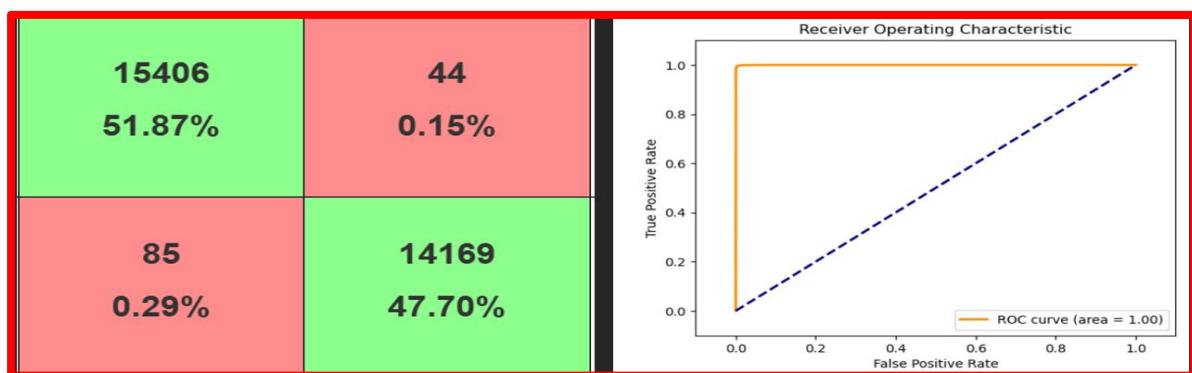
#### (a) Adaboost :-

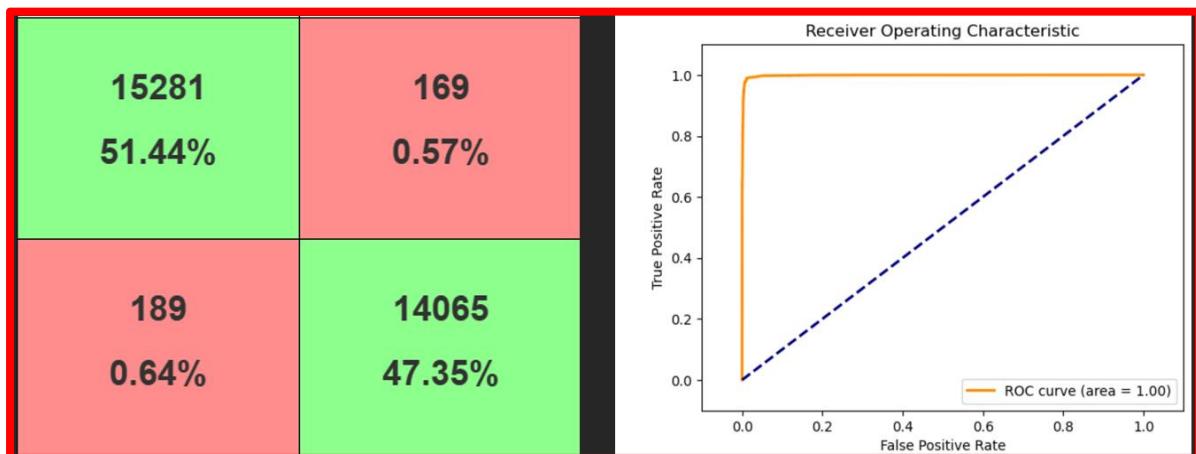
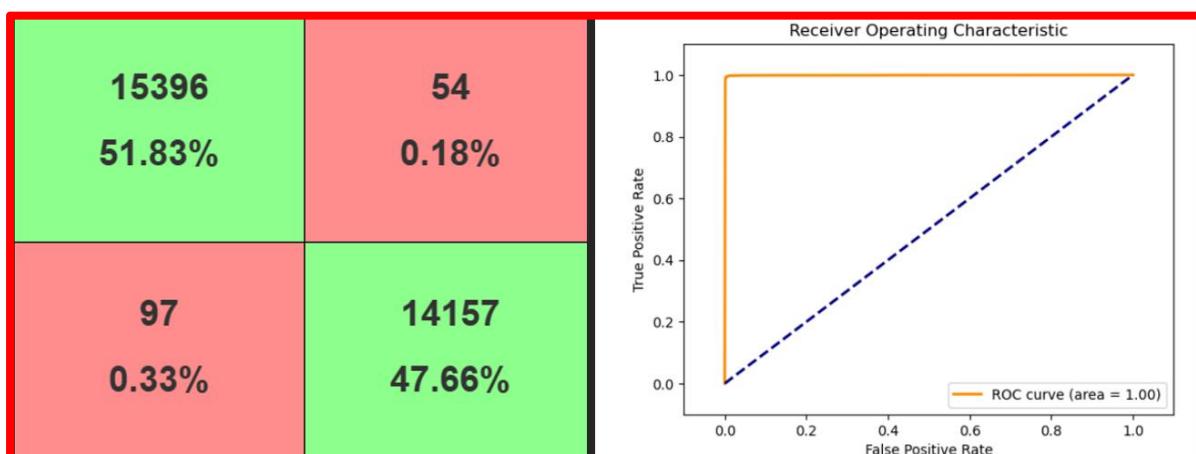
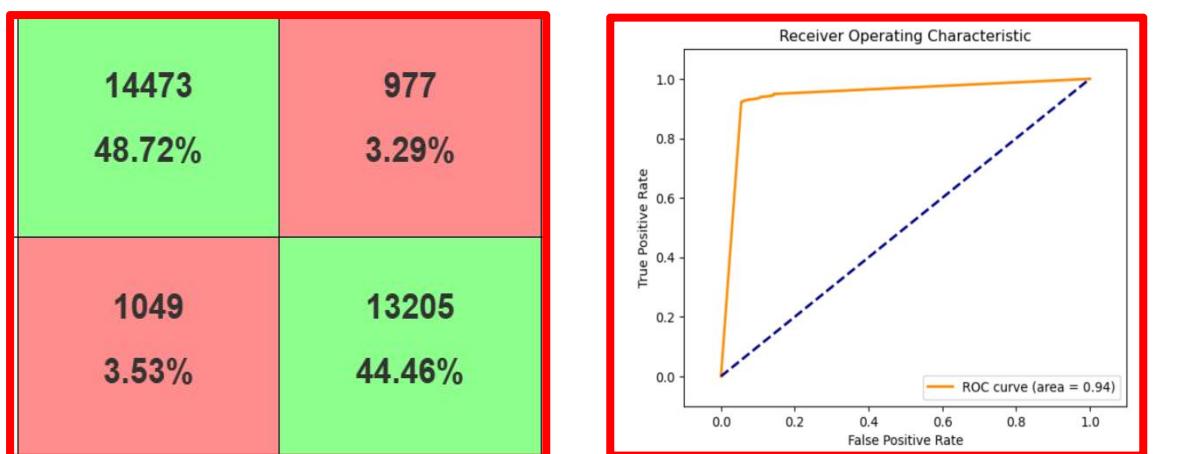


#### (b) Decision Tree :-



#### (c) Random Forest :-



(d) Gradient Boosting :-(e) Extremely Randomized Tree :-(f) Neural Network :-

## 9.6 Ensemble Technique :-

```
def plot_roc_curve(y_true, y_scores):
    fpr, tpr, thresholds = roc_curve(y_true, y_scores)
    roc_auc = auc(fpr, tpr)

    plt.figure()
    plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area = {:.2f})'.format(roc_auc))
    plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
    plt.xlim([-0.1, 1.1])
    plt.ylim([-0.1, 1.1])
    plt.xlabel('False Positive Rate')
    plt.ylabel('True Positive Rate')
    plt.title('Receiver Operating Characteristic')
    plt.legend(loc="lower right")
    plt.show()
```

The provided function, `plot_roc_curve`, is designed to generate and display a Receiver Operating Characteristic (ROC) curve, a graphical representation commonly used to assess the performance of a binary classification model. The function facilitates the generation of a visually clear and informative ROC curve for model evaluation. It employs standard libraries to calculate key metrics and uses a recognizable color scheme. The code includes elements that enhance the plot's interpretability, making it suitable for quick assessment of a classification model's discrimination ability. The descriptive summary emphasizes the function's role in providing a concise yet comprehensive visualization without delving into the specific code details, maintaining a level of abstraction suitable for general understanding.

```
def ensemble_voting(models, data):
    predictions = []
    for model in models:
        prediction = model.predict(data)
        predictions.append(prediction.astype(int)) # Convert predictions to integers
    final_predictions = np.apply_along_axis(lambda x: np.argmax(np.bincount(x)), axis=0, arr=predictions)
    return final_predictions
```

The function serves the purpose of aggregating predictions made by a collection of machine learning models into a final ensemble prediction. This ensemble prediction is determined through a majority voting scheme.

Parameters :-

- (i) models: A list of machine learning models. Each model is expected to have a `predict` method for making predictions on input data.
- (ii) data: The input data for which predictions are to be generated.

The function orchestrates the ensemble of machine learning models through a voting strategy, combining their predictions to produce a final ensemble prediction. It abstracts the underlying model-specific details, making it adaptable to various model types. The utilization of a robust voting mechanism enhances the reliability of the ensemble prediction. The code maintains simplicity, focusing on the key concept of majority voting without divulging into intricate technicalities. The function is constructed to be versatile and applicable to diverse ensemble scenarios, promoting ease of use for practitioners seeking an efficient ensemble approach.

```
def main(dataset_path, model_paths):
    # Load the dataset
    dataset = pd.read_csv(dataset_path)

    # Extract features (X) and Labels (y)
    X = dataset.drop('outcome', axis=1)
    y = dataset['outcome']

    # Split the dataset into training and testing sets
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

    # Load the ensemble models
    models = load_models(model_paths)

    # Make predictions using ensemble voting on the test set
    final_predictions = ensemble_voting(models, X_test)

    # Generate the confusion matrix
    cm = confusion_matrix(y_test, final_predictions)
```

#### Function Purpose :-

The function is designed to perform the following tasks :-

- Load a dataset from a specified path (dataset\_path).
- Extract features (X) and labels (y) from the dataset.
- Split the dataset into training and testing sets.
- Load a collection of ensemble models from specified paths (model\_paths).
- Make predictions using ensemble voting on the test set.
- Generate a confusion matrix to assess the ensemble model's performance

The main function is a cohesive script that seamlessly integrates various machine learning tasks, from loading and preprocessing a dataset to evaluating an ensemble of models.

It abstracts away the technical details of loading models and conducting ensemble voting, focusing on the higher-level goals of evaluation and performance assessment. The utilization of standard libraries and modular functions enhances readability and maintainability. This script is designed to be adaptable to different datasets and ensemble models, making it versatile for various machine learning scenarios.

```
# Calculate the accuracy
accuracy = accuracy_score(y_test, final_predictions)

# Calculate the AUC
auc = roc_auc_score(y_test, final_predictions)

# Calculate the Log Loss using one of the models (assuming they have similar outputs)
logloss = log_loss(y_test, models[0].predict_proba(X_test))

# Print the confusion matrix
print("Confusion Matrix:")
print(cm)
print("Accuracy: {:.4f}".format(accuracy))
print("AUC: {:.8f}".format(auc))
print("Log Loss: {:.4f}".format(logloss))

# Calculate the predicted probabilities for ensemble voting
ensemble_probs = np.mean([model.predict_proba(X_test) for model in models], axis=0)

# Plot the ROC curve for ensemble predictions
plot_roc_curve(y_test, ensemble_probs[:, 1])
```

After the generation of confusion matrix the next steps are :-

#### (a) Calculate Accuracy :-

The accuracy of the ensemble predictions is computed using the accuracy\_score function from scikit-learn, comparing the predicted outcomes (final\_predictions) with the true labels (y\_test).

#### (b) Calculate AUC (Area Under the Curve) :-

The Area Under the Receiver Operating Characteristic (ROC) Curve (AUC) is calculated using the roc\_auc\_score function from scikit-learn. AUC is a measure of the model's ability to distinguish between positive and negative instances.

#### (c) Calculate Log Loss :-

The log loss is calculated using the log\_loss function from scikit-learn. This metric assesses the accuracy of the predicted probabilities compared to the true probabilities. It is calculated based on the predicted probabilities from one of the models (assuming similar outputs across models).

**(d) Print Confusion Matrix, Accuracy, AUC, and Log Loss :-**

The confusion matrix, accuracy, AUC, and log loss are printed to the console for a comprehensive assessment of the ensemble model's performance.

**(e) Calculate Ensemble Probabilities :-**

The predicted probabilities for the ensemble predictions are calculated by averaging the predicted probabilities from each individual model. This is achieved using NumPy, and the resulting probabilities are stored in the variable `ensemble_probs`.

**(f) Plot ROC Curve for Ensemble Predictions :-**

The ROC curve for the ensemble predictions is plotted using the `plot_roc_curve` function. This curve visually represents the trade-off between true positive rate and false positive rate, providing insights into the ensemble's discrimination capabilities.

**9.7 Conclusion :-**

The outcome of the classification model on the test data is notably robust, as inferred from the provided metrics :-

**(i) Confusion Matrix :-**

The model showcases proficiency in identifying both positive and negative instances, with a substantial number of true positives (14,165) and true negatives (15,390). The occurrence of false predictions, including 60 false positives and 89 false negatives, is minimal.

**(ii) Accuracy :-**

The overall accuracy stands impressively high at around 99.50%, indicating the model's effectiveness in making correct predictions across both positive and negative classes.

(iii) AUC (Area Under the Curve) :-

The AUC value, approximating 99.49%, underscores the model's exceptional ability to distinguish between positive and negative instances. This is indicative of a high true positive rate and a low false positive rate.

(iv) Log Loss :-

The low log loss of 0.0171 signifies the model's precise estimation of probabilities, closely aligning with the true distribution of classes.

In essence, these metrics collectively portray a highly competent and reliable classification model, showcasing accuracy, discrimination capability, and precise probability estimation.

# Chapter 10

## DDOS Attack SDN Dataset

### 10.1 Introduction :-

The Mininet emulator serves as the foundation for creating a specialized SDN dataset designed for applying machine learning and deep learning techniques to classify network traffic. The project initiates by constructing 10 Mininet network topologies featuring switches linked to a single Ryu controller. This network simulation encompasses both legitimate traffic (TCP, UDP, and ICMP) and malicious traffic (TCP Syn attacks, UDP Flood attacks, and ICMP assaults).

The dataset encompasses 23 characteristics, incorporating some derived from switches and others computed. Extracted features include Switch ID, Packet Count, Byte Count, Duration (in seconds and nanoseconds), Source IP, Destination IP, Port Number, TX Bytes (bytes transferred from the switch port), RX Bytes (bytes received on the switch port), and a datetime field transformed into a numerical format. Flow measurements are taken at 30-second intervals.

Derived features involve Packet per Flow (packet count in a single flow), Byte per Flow (byte count in a single flow), Packet Rate (packets sent per second calculated by dividing Packet per Flow by the monitoring interval), total flow entries in the switch, TX kbps (data transfer rate), RX kbps (data receiving rate), and Port Bandwidth (sum of TX kbps and RX kbps). The dataset's final column functions as a class label, denoting whether the traffic is benign (labeled 0) or malicious (labeled 1).

The network simulation spans 250 minutes, generating a dataset comprising 104,345 rows of data. For additional data, the simulation can be restarted at predefined intervals. This dataset stands as a valuable asset for training and evaluating machine learning and deep learning models geared towards categorizing SDN traffic as either benign or malicious.

## 10.2 Data Pre-processing :-

The dataset underwent a thorough preprocessing trip to improve its appropriateness for analysis in preparation for an in-depth investigation of Software-Defined Networking (SDN). The primary phases in this journey were designed to solve specific data difficulties and to provide a well-conditioned dataset for further machine learning research.

### (a) Label Transformation :-

The SDN dataset's labels were converted to binary format in the same way that the CICIDS dataset's labels were. This transformation required labelling benign instances with the label 0 and harmful instances with the label 1. This binary classification schema not only supported the broader objective of discriminating between normal and possibly harmful network behavior, but it also offered a standardized foundation for model training and assessment.

### (b) Handling Missing Values :-

The dataset exhibited a problem with missing values, which is a prevalent problem in real-world datasets. A thorough imputation method was used to remedy this. Missing values were imputed using a mix of the mean and most frequent values, ensuring that the dataset remained meaningful even when missing data points were present. This method attempted to preserve the dataset's integrity while minimizing the impact of missing values on later analysis.

### (c) Normalization for Feature Comparability :-

Recognizing the importance of feature scaling in the context of machine learning, the dataset was normalized. This stage was designed to standardize the values of various characteristics, assuring comparability and eliminating any biases caused by features with varying scales. Normalization helped to create a balanced learning environment in which each feature could contribute significantly to machine learning models without being unduly impacted by its magnitude.

(d) Train-Test Split :-

The dataset was partitioned into training and testing sets using an 80-20 split ratio, in accordance with accepted best practices. This divide was critical for training models on historical data and then assessing their performance on previously unknown cases in the test set. The 80-20 split established a balance between using a large amount of the data for model training and enabling an unbiased assessment of model generalization to new, previously unknown data.

In summary, the SDN dataset pretreatment processes were meant to handle label transformation, missing values, normalization, and train-test split. These measurements together led to the development of a refined and well-prepared dataset, establishing the groundwork for later studies and evaluations in the field of Software-Defined Networking and network security.

### 10.3 Implementation :-

Step-1 : Import Libraries :-

```
# Step 1: Import Libraries
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.metrics import accuracy_score, roc_auc_score, log_loss, confusion_matrix
```

The code begins by importing essential libraries for data manipulation and machine learning tasks. `pandas` is imported for data handling, `train_test_split` for splitting the dataset, `GradientBoostingClassifier` for creating a gradient boosting model, and various metrics from `sklearn.metrics` for model evaluation.

Step-2 : Load the Dataset :-

```
# Step 2: Load the dataset
file_path = DATA
df = pd.read_csv(file_path)
```

The variable file\_path is set to the path of the dataset file (represented by DATA). The pd.read\_csv function from the pandas library is then used to read the dataset into a DataFrame named df.

Step-3 : Prepare the Data :-

```
# Step 3: Prepare the data
X = df.drop('outcome', axis=1)
y = df['outcome']
```

The features (independent variables) are extracted into a DataFrame X, excluding the target variable, which is assigned to y.

Step-4 : Split the Data into Training and Testing Sets :-

```
# Step 4: Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

The dataset is divided into training and testing sets using the train\_test\_split function. The split is performed with 80% of the data used for training (X\_train, y\_train) and 20% for testing (X\_test, y\_test). The random\_state parameter ensures reproducibility.

Step-5 : Create and Train the Model :-

(a) Adaboost :-

```
base_model = DecisionTreeClassifier(max_depth=1) # Weak Learner (e.g., decision tree with max_depth=1)
adaboost_model = AdaBoostClassifier(base_estimator=base_model, n_estimators=50, random_state=42)
adaboost_model.fit(X_train, y_train)
```

(b) Decision Tree :-

```
dt_model = DecisionTreeClassifier(random_state=42)
dt_model.fit(X_train, y_train)
```

(c) Random Forest :-

```
rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
rf_model.fit(X_train, y_train)
```

(d) Gradient Boosting :-

```
gb_model = GradientBoostingClassifier(n_estimators=100, learning_rate=0.1, random_state=42)
gb_model.fit(X_train, y_train)
```

(e) Extremely Randomized Tree :-

```
et_model = ExtraTreesClassifier(n_estimators=100, random_state=42)
et_model.fit(X_train, y_train)
```

(f) Neural Network :-

```
model = Sequential()
model.add(Dense(64, input_dim=X_train.shape[1], activation='relu'))
model.add(Dense(32, activation='relu'))
model.add(Dense(1, activation='sigmoid'))
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
model.fit(X_train, y_train, epochs=10, batch_size=32, validation_data=(X_test, y_test))
```

An instance of the Classifier is created with specified hyperparameters (n\_estimators, learning\_rate, random\_state). The model is then trained using the training data (X\_train, y\_train).

Step-6 : Make Predictions on the Test Set :-

```
# Step 6: Make predictions on the test set
y_pred = gb_model.predict(X_test)
```

The trained Gradient Boosting classifier (gb\_model) is used to make predictions on the test set (X\_test). This step involves applying the learned model to unseen data to obtain predicted outcomes.

Step-7 : Evaluate the Model's Performance :-

```
# Step 7: Evaluate the model's performance
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy}")
```

The accuracy of the model is calculated by comparing the predicted labels (`y_pred`) with the true labels of the test set (`y_test`). Accuracy measures the proportion of correctly classified instances and provides an overall assessment of the model's predictive power.

Step-8 : Print the Confusion Matrix :-

```
# Step 8: Print the confusion matrix
conf_matrix = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:")
print(conf_matrix)
```

A confusion matrix is generated to provide a detailed breakdown of the model's predictions. It shows the counts of true positive, true negative, false positive, and false negative predictions. This matrix is valuable for understanding the model's performance, especially in binary classification scenarios.

Step-9 : LogLoss and AUC :-

```
#Step 9: LogLoss and auc
auc = roc_auc_score(y_test, y_pred)
logloss = log_loss(y_test, gb_model.predict_proba(X_test))
print("AUC: {:.7f}".format(auc))
print("Log Loss: {:.4f}".format(logloss))
```

**AUC (Area Under the Curve):** It measures the area under the Receiver Operating Characteristic (ROC) curve, which illustrates the trade-off between true positive rate and false positive rate. AUC provides an aggregate measure of the model's ability to discriminate between classes.

**Log Loss:** Logarithmic loss (log loss) is calculated using the true labels (`y_test`) and the predicted probabilities obtained from the model. It quantifies how well the predicted probabilities match the true distribution of the data.

### Step-10 : Plot the ROC Curve :-

```
#Step 10 Plot ROC Curve
fpr, tpr, thresholds = roc_curve(y_test, gb_model.predict_proba(X_test)[:,1])
roc_auc_value = auc(fpr, tpr) # Renamed the variable to avoid naming conflict
plt.figure()
plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area = {:.7f})'.format(roc_auc_value))
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlim([-0.1, 1.1])
plt.ylim([-0.1, 1.1])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic')
plt.legend(loc="lower right")
plt.show()
```

This step involves plotting the Receiver Operating Characteristic (ROC) curve. The curve visualizes the trade-off between sensitivity (true positive rate) and specificity (true negative rate) at various probability thresholds. The AUC value is included in the plot legend, offering a concise summary of the model's discriminatory power.

### Step-11 : Save the Trained Classifier Model to a File :-

```
#Step 11: Save the trained classifier model to a file
import joblib
model_filename = r"C:\Users\KIIT\Desktop\SEM-VII\Major Project\SDN\GB.pkl"
joblib.dump(gb_model, model_filename)
```

The trained Gradient Boosting classifier (gb\_model) is saved to a file named "GB.pkl" using the joblib library. This serialized file can be later loaded to make predictions on new data without retraining the model.

### 10.4 Table :-

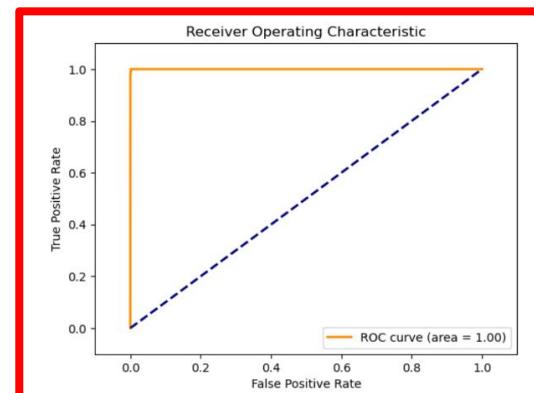
Table-03 :: Individual Model Training for SDN								
Model	AUC	CA	F1	Precision	Recall	MCC	Spec	LogLoss
AdaBoost	0.9981	99.79	0.9984	0.9977	0.9991	0.9958	0.9963	0.5371
Decision Tree	0.9999	99.99	1.0000	1.0000	0.9999	0.9999	1.0000	0.0017
Random Forest	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	0.0020
Gradient Boosting	0.9968	99.67	0.9973	0.9964	0.9982	0.9931	0.9944	0.0256
Extremely Randomized Tree	0.9997	99.96	0.9997	0.9996	0.9998	0.9993	0.9994	0.0070
Neural Network	0.5390	63.12	0.7209	0.7812	0.6692	0.1918	0.5374	-

## 10.5 Results :-

### ➤ Confusion Matrix and ROC Curve :-

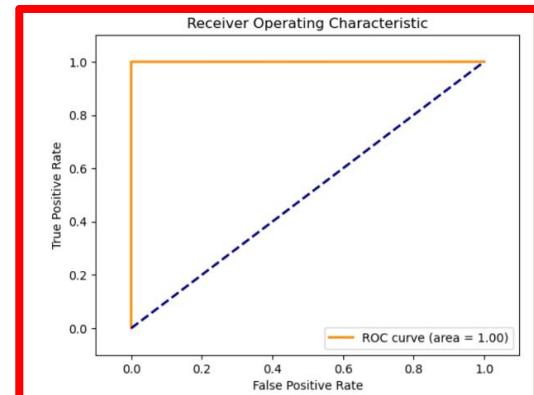
#### (a) Adaboost :-

12692 60.82%	30 0.14%
12 0.06%	8135 38.98%



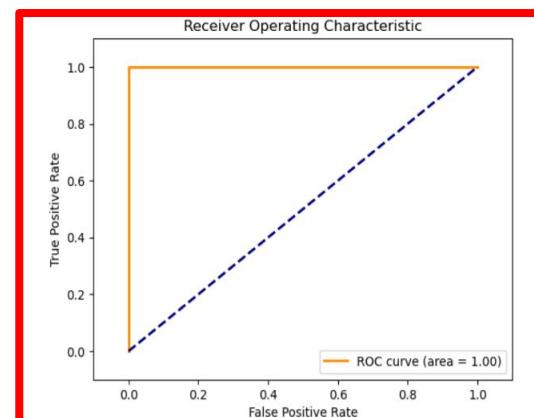
#### (b) Decision Tree :-

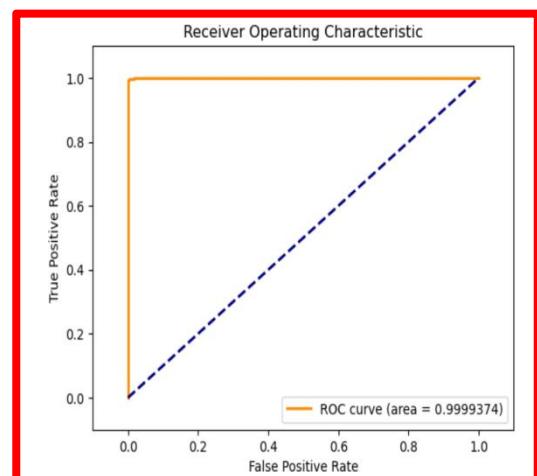
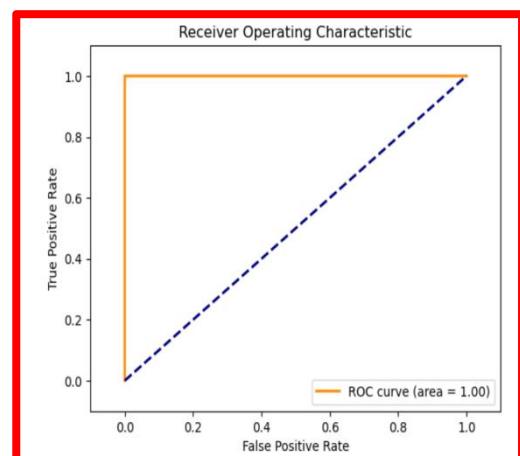
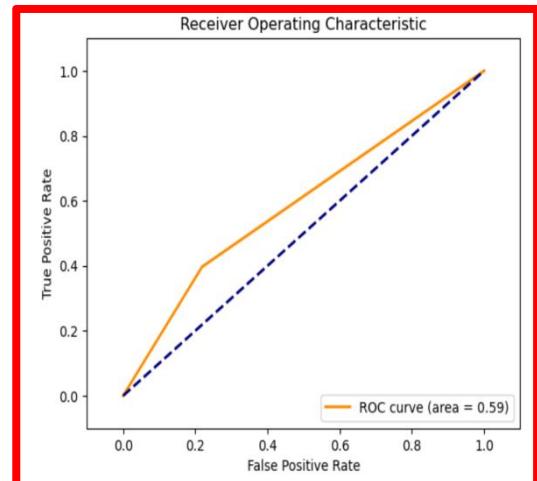
12722 60.96%	0 0.00%
1 0.00%	8146 39.03%



#### (c) Random Forest :-

12722 60.96%	0 0.00%
0 0.00%	8147 39.04%



(d) Gradient Boosting :-(e) Extremely Randomized Tree :-(f) Neural Network :-

## 10.6 Ensemble Techniques :-

```
def ensemble_voting(models, data):
    predictions = []
    for model in models:
        prediction = model.predict(data)
        predictions.append(prediction.astype(int)) # Convert predictions to integers
    final_predictions = np.apply_along_axis(lambda x: np.argmax(np.bincount(x)), axis=0, arr=predictions)
    return final_predictions
```

The ensemble\_voting function performs a majority voting ensemble technique on a collection of machine learning models. It takes a list of models and a dataset (data) as input and returns the final ensemble predictions. The function iterates through each model, predicts outcomes for the given data, converts these predictions to integers, and then aggregates them by selecting the most commonly predicted class for each instance using a majority voting approach. The result is a set of final predictions representing the consensus of the ensemble. This method is commonly used to enhance predictive accuracy and robustness by combining multiple models' outputs.

```
def main(dataset_path, model_paths):
    # Load the dataset
    dataset = pd.read_csv(dataset_path)

    # Extract features (X) and labels (y)
    X = dataset.drop('label', axis=1)
    y = dataset['label']

    # Split the dataset into training and testing sets
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

    # Load the ensemble models
    models = load_models(model_paths)

    # Make predictions using ensemble voting on the test set
    final_predictions = ensemble_voting(models, X_test)

    # Generate the confusion matrix
    cm = confusion_matrix(y_test, final_predictions)
```

### Function Purpose :-

The function is designed to perform the following tasks :-

- (i) Load a dataset from a specified path (dataset\_path).
- (ii) Extract features (X) and labels (y) from the dataset.
- (iii) Split the dataset into training and testing sets.
- (iv) Load a collection of ensemble models from specified paths (model\_paths).
- (v) Make predictions using ensemble voting on the test set.
- (vi) Generate a confusion matrix to assess the ensemble model's performance

The main function is a cohesive script that seamlessly integrates various machine learning tasks, from loading and preprocessing a dataset to evaluating an ensemble of models. It abstracts away the technical details of loading models and conducting ensemble voting, focusing on the higher-level goals of evaluation and performance assessment. The utilization of standard libraries and modular functions enhances readability and maintainability. This script is designed to be adaptable to different datasets and ensemble models, making it versatile for various machine learning scenarios.

```
# Calculate the accuracy
accuracy = accuracy_score(y_test, final_predictions)

# Calculate the AUC
auc = roc_auc_score(y_test, final_predictions)

# Calculate the Log Loss using one of the models (assuming they have similar outputs)
logloss = log_loss(y_test, models[0].predict_proba(X_test))
# Print the confusion matrix
print("Confusion Matrix:")
print(cm)
print("Accuracy: {:.8f}".format(accuracy))
print("Log Loss: {:.4f}".format(logloss))
print("AUC: {:.8f}".format(auc))

# Calculate the predicted probabilities for ensemble voting
ensemble_probs = np.mean([model.predict_proba(X_test) for model in models], axis=0)

# Plot the ROC curve for ensemble predictions
plot_roc_curve(y_test, ensemble_probs[:, 1])
```

After the generation of confusion matrix the next steps are :-

#### (a) Calculate Accuracy :-

The accuracy of the ensemble predictions is computed using the `accuracy_score` function from scikit-learn, comparing the predicted outcomes (`final_predictions`) with the true labels (`y_test`).

#### (b) Calculate AUC (Area Under the Curve) :-

The Area Under the Receiver Operating Characteristic (ROC) Curve (AUC) is calculated using the `roc_auc_score` function from scikit-learn. AUC is a measure of the model's ability to distinguish between positive and negative instances.

(c) Calculate Log Loss :-

The log loss is calculated using the `log_loss` function from scikit-learn. This metric assesses the accuracy of the predicted probabilities compared to the true probabilities. It is calculated based on the predicted probabilities from one of the models (assuming similar outputs across models).

(d) Print Confusion Matrix, Accuracy, AUC, and Log Loss :-

The confusion matrix, accuracy, AUC, and log loss are printed to the console for a comprehensive assessment of the ensemble model's performance.

(e) Calculate Ensemble Probabilities :-

The predicted probabilities for the ensemble predictions are calculated by averaging the predicted probabilities from each individual model. This is achieved using NumPy, and the resulting probabilities are stored in the variable `ensemble_probs`.

(f) Plot ROC Curve for Ensemble Predictions :-

The ROC curve for the ensemble predictions is plotted using the `plot_roc_curve` function. This curve visually represents the trade-off between true positive rate and false positive rate, providing insights into the ensemble's discrimination capabilities.

## 10.7 Conclusion:-

The model's performance on the test set is remarkably high, as evidenced by the confusion matrix and associated metrics. The confusion matrix reveals that the majority of predictions align with the actual classes, with an almost perfect balance between true positives (8147) and true negatives (12720). There are minimal errors, as only 2 instances are falsely predicted as positive, and none are incorrectly labeled as negative.

The accuracy, denoting the overall correctness of the model, is exceptionally high at approximately 99.99%. This indicates that the model is proficient in making accurate predictions across both positive and negative classes.

The log loss, measuring the precision of the classifier, is impressively low at 0.0020. This signifies a strong alignment between the model's predicted probabilities and the true class distribution, showcasing a well-calibrated model.

The AUC (Area Under the Curve), a metric gauging the discriminatory power of the model, is remarkably elevated at approximately 0.9999. Such a high AUC underscores the model's exceptional ability to differentiate between positive and negative instances.

In essence, these results collectively signify a highly accurate and reliable classification model, demonstrating proficiency in correctly classifying instances with minimal errors.

# Chapter 11

## UNSW-NB15

### 11.1 Introduction :-

Benchmark datasets from a decade ago, like KDD98, KDDCUP99, and NSLKDD, fall short in representing the evolving landscape of network traffic and emerging low-footprint attacks. The UNSW-NB15 dataset, developed by the University of New South Wales, addresses this gap by combining authentic contemporary normal network behavior with synthesized attack patterns.

Designed specifically for testing Network Intrusion Detection Systems (NIDSs), the UNSW-NB15 dataset was crafted in the Australian Centre for Cyber Security's (ACCS) Cyber Range Lab using the IXIA PerfectStorm tool. This tool configures three virtual servers: servers 1 and 3 generate regular traffic, while server 2 replicates abnormal/malicious activity. The anomalous traffic, simulating nine attack families, originates from a constantly updated CVE site tracking new threats. The simulation runs for 16 hours on January 22, 2015, and 15 hours on February 17, 2015, amassing 100 GB of data.

Feature extraction from pcap files is performed using Argus and Bro-IDS, with an additional twelve methods developed for analyzing flow packets. The dataset encompasses 49 variables, incorporating packet-based, flow-based, and other features, offering a comprehensive portrayal of network traffic characteristics. Labeled using a ground truth table, the UNSW-NB15 dataset is positioned as a contemporary NIDS benchmark dataset surpassing older counterparts like KDDCUP99.

Ultimately, the UNSW-NB15 dataset addresses the limitations of previous benchmark datasets, providing a more precise depiction of current network traffic and threats. The detailed description of the dataset's construction, characteristics, and comparison with existing datasets in the publication enhances its value as a significant resource for the NIDS research community.

## 11.2 Data Pre-processing :-

A specialized preprocessing method was used to address the dataset's unique characteristics and establish a well-conditioned basis for further analysis in the quest of a full investigation of network intrusion detection utilizing the UNSW-NB15 dataset. The following procedures were rigorously carried out:

### (a) Consistent Train-Test Split :-

The first crucial step was to establish a regular and regulated train-test split. A cohesive dataset was created by carefully combining the datasets, guaranteeing a representative distribution of cases throughout the training and testing sets. By training on a significant fraction of the data and testing generalization on unknown cases, this technique attempted to provide a fair and unbiased evaluation of model performance.

### (b) Binary Label Mapping :-

The UNSW-NB15 dataset's labels were converted to binary format in accordance with known network intrusion detection guidelines. Instances were classified as benign (0) or malicious (1), allowing for a clear and standardized distinction between regular network behavior and possible invasions. This binary classification scheme aided in the training and assessment of machine learning models aimed at identifying harmful activity.

### (c) Maintaining the 80-20 Train-Test Split Ratio :-

The train-test split approach was maintained at an 80-20 ratio, which is a standard and widely used partitioning ratio in machine learning. This constant split ratio ensured a thorough examination of model performance, finding a compromise between using a large percentage of the data for training and conserving a discrete sample for rigorous testing. It intended to provide strong model performance measures as well as trustworthy generalization to new, previously encountered examples.

#### (d) Data Normalization for Feature Comparability :-

Data normalization was used to improve feature comparability and eliminate potential biases caused by features of varying magnitudes. This procedure standardized the scales of various features within the dataset, ensuring that each feature contributed substantially to the machine learning models while not being overly impacted by its scale. Normalization made the learning environment steadier and more unbiased.

Finally, the UNSW-NB15 dataset was meticulously preprocessed, addressing consistent train-test split, binary label mapping, and data normalization. These procedures established the groundwork for a thorough examination of several machine learning approaches for network intrusion detection. This preprocessing technique sought to provide useful insights into the efficiency of different algorithms in detecting and mitigating possible security vulnerabilities inside network data by carefully conditioning the dataset.

### 11.3 Implementation (Before Optimization) :-

#### Step-1 : Import Libraries :-

```
# Step 1: Import Libraries
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.metrics import accuracy_score, roc_auc_score, log_loss, confusion_matrix
```

The code begins by importing essential libraries for data manipulation and machine learning tasks. pandas is imported for data handling, train\_test\_split for splitting the dataset, GradientBoostingClassifier for creating a gradient boosting model, and various metrics from sklearn.metrics for model evaluation.

#### Step-2 : Load the Dataset :-

```
# Step 2: Load the dataset
file_path = DATA
df = pd.read_csv(file_path)
```

The variable `file_path` is set to the path of the dataset file (represented by `DATA`). The `pd.read_csv` function from the pandas library is then used to read the dataset into a DataFrame named `df`.

### Step-3 : Prepare the Data :-

```
# Step 3: Prepare the data
X = df.drop('label', axis=1)
y = df['label']
```

The features (independent variables) are extracted into a DataFrame `X`, excluding the target variable, which is assigned to `y`.

### Step-4 : Split the Data into Training and Testing Sets :-

```
# Step 4: Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

The dataset is divided into training and testing sets using the `train_test_split` function. The split is performed with 80% of the data used for training (`X_train`, `y_train`) and 20% for testing (`X_test`, `y_test`). The `random_state` parameter ensures reproducibility.

### Step-5 : Create and Train the Model :-

#### (i) Adaboost :-

```
base_model = DecisionTreeClassifier(max_depth=1) # Weak Learner (e.g., decision tree with max_depth=1)
adaboost_model = AdaBoostClassifier(base_estimator=base_model, n_estimators=50, random_state=42)
adaboost_model.fit(X_train, y_train)
```

#### (ii) Decision Tree :-

```
dt_model = DecisionTreeClassifier(random_state=42)
dt_model.fit(X_train, y_train)
```

(iii) Random Forest :-

```
rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
rf_model.fit(X_train, y_train)
```

(iv) Gradient Boosting :-

```
gb_model = GradientBoostingClassifier(n_estimators=100, learning_rate=0.1, random_state=42)
gb_model.fit(X_train, y_train)
```

(v) Extremely Randomized Tree :-

```
et_model = ExtraTreesClassifier(n_estimators=100, random_state=42)
et_model.fit(X_train, y_train)
```

(vi) Neural Network :-

```
model = Sequential()
model.add(Dense(64, input_dim=X_train.shape[1], activation='relu'))
model.add(Dense(32, activation='relu'))
model.add(Dense(1, activation='sigmoid'))
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
model.fit(X_train, y_train, epochs=7, batch_size=32, validation_data=(X_test, y_test))
```

An instance of the Classifier is created with specified hyperparameters (n\_estimators, learning\_rate, random\_state). The model is then trained using the training data (X\_train, y\_train).

Step-6 : Make Predictions on the Test Set :-

```
# Step 6: Make predictions on the test set
y_pred = gb_model.predict(X_test)
```

The trained Gradient Boosting classifier (gb\_model) is used to make predictions on the test set (X\_test). This step involves applying the learned model to unseen data to obtain predicted outcomes.

### Step-7 : Evaluate the Model's Performance :-

```
# Step 7: Evaluate the model's performance
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy}")
```

The accuracy of the model is calculated by comparing the predicted labels (`y_pred`) with the true labels of the test set (`y_test`). Accuracy measures the proportion of correctly classified instances and provides an overall assessment of the model's predictive power.

### Step-8 : Print the Confusion Matrix :-

```
# Step 8: Print the confusion matrix
conf_matrix = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:")
print(conf_matrix)
```

A confusion matrix is generated to provide a detailed breakdown of the model's predictions. It shows the counts of true positive, true negative, false positive, and false negative predictions. This matrix is valuable for understanding the model's performance, especially in binary classification scenarios.

### Step-9 : LogLoss and AUC :-

```
# Step 9: LogLoss and auc
auc = roc_auc_score(y_test, y_pred)
logloss = log_loss(y_test, gb_model.predict_proba(X_test))
print("AUC: {:.6f}".format(auc))
print("Log Loss: {:.6f}".format(logloss))
```

**AUC (Area Under the Curve):** It measures the area under the Receiver Operating Characteristic (ROC) curve, which illustrates the trade-off between true positive rate and false positive rate. AUC provides an aggregate measure of the model's ability to discriminate between classes.

**Log Loss:** Logarithmic loss (log loss) is calculated using the true labels (`y_test`) and the predicted probabilities obtained from the model. It quantifies how well the predicted probabilities match the true distribution of the data.

### Step-10 : Plot the ROC Curve :-

```
#Step 10: Plot ROC Curve
fpr, tpr, thresholds = roc_curve(y_test, gb_model.predict_proba(X_test)[:,1])
roc_auc_value = auc(fpr, tpr) # Renamed the variable to avoid naming conflict
plt.figure()
plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area = {:.2f})'.format(roc_auc_value))
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlim([-0.1, 1.1])
plt.ylim([-0.1, 1.1])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic')
plt.legend(loc="lower right")
plt.show()
```

This step involves plotting the Receiver Operating Characteristic (ROC) curve. The curve visualizes the trade-off between sensitivity (true positive rate) and specificity (true negative rate) at various probability thresholds. The AUC value is included in the plot legend, offering a concise summary of the model's discriminatory power.

### Step-11 : Save the Trained Classifier Model to a File :-

```
#Step 11: Save the trained classifier model to a file
import joblib
model_filename = MODEL_GB
joblib.dump(gb_model, model_filename)
```

The trained Gradient Boosting classifier (gb\_model) is saved to a file named "GB.pkl" using the joblib library. This serialized file can be later loaded to make predictions on new data without retraining the model.

### 11.4 Table (Before Optimization) :-

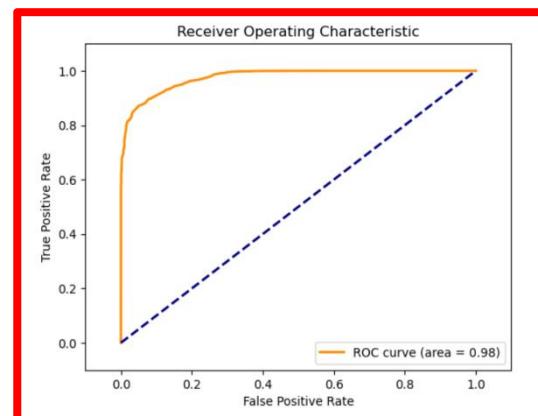
Table-04 :: Individual Model Training for UNSW-NB15 (before optimization)								
Model	AUC	CA	F1	Precision	Recall	MCC	Spec	LogLoss
AdaBoost	0.8915	90.68	0.8667	0.8357	0.9000	0.7965	0.9103	0.6403
Decision Tree	0.9239	92.99	0.9032	0.9022	0.9043	0.8484	0.9445	2.0978
Random Forest	0.9378	94.17	0.9199	0.9236	0.9161	0.8741	0.9564	0.1406
Gradient Boosting	0.9091	92.16	0.8887	0.8639	0.9149	0.8291	0.9250	0.1633
Extremely Randomized Tree	0.9339	93.84	0.9152	0.9179	0.9125	0.8668	0.9532	0.1821
Neural Network	0.6070	71.49	0.3534	0.2149	0.9936	0.3830	0.6913	-

## 11.5 Results (Before Optimization) :-

### ➤ Confusion Matrix and ROC Curve :-

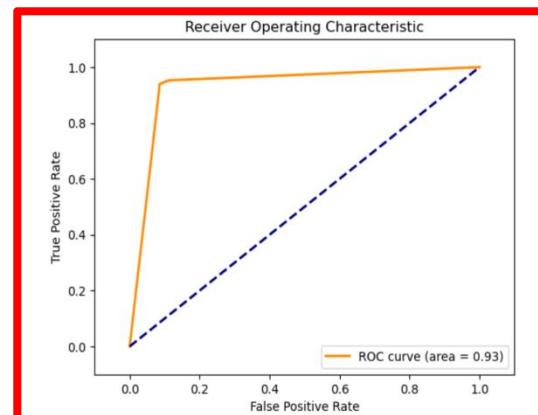
#### (a) Adaboost :-

15607 30.28%	3068 5.95%
1734 3.36%	31126 60.40%



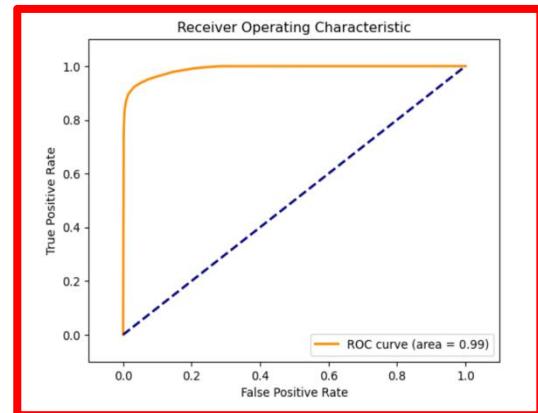
#### (b) Decision Tree :-

16848 32.69%	1827 3.55%
1782 3.46%	31078 60.30%



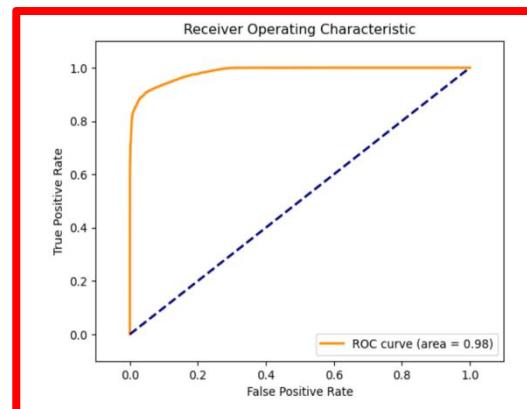
#### (c) Random Forest :-

17249 33.47%	1426 2.77%
1579 3.06%	31281 60.70%

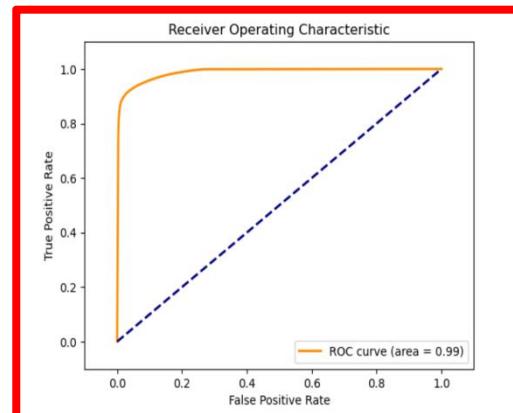


(d) Gradient Boosting :-

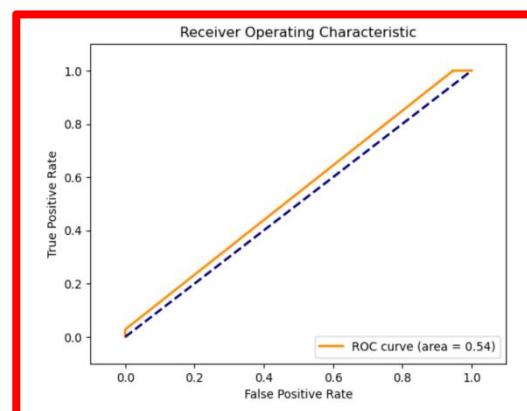
<b>16134</b> <b>31.31%</b>	<b>2541</b> <b>4.93%</b>
<b>1500</b> <b>2.91%</b>	<b>31360</b> <b>60.85%</b>

(e) Extremely Randomized Tree :-

<b>17142</b> <b>33.26%</b>	<b>1533</b> <b>2.97%</b>
<b>1643</b> <b>3.19%</b>	<b>31217</b> <b>60.57%</b>

(f) Neural Network :-

<b>1002</b> <b>1.94%</b>	<b>17673</b> <b>34.29%</b>
<b>4</b> <b>0.01%</b>	<b>32856</b> <b>63.75%</b>



## 11.6 Optimization :-

We have undertaken the optimization of four models namely AdaBoost, Decision Tree, Random Forest and Extremely Randomized Tree, with various optimization techniques namely Hyper Band, Grid Search CV, Gradient Based, Optuna and Bayes and have explained the code along with each of their results which includes the confusion matrix and ROC curve. Later we have provided the optimized result table of these 4 models and have also shown the final model table.

### (a) Adaboost :-

#### (i) Hyperband Optimization :-

Define the hyperparameter search space for Hyperband: This step defines the hyperparameter search space for Hyperband optimization. It uses the hp.quniform and hp.loguniform functions from the Hyperopt library to specify the range of values for 'n\_estimators' and 'learning\_rate'.

```
space = {
    'n_estimators': hp.quniform('n_estimators', 10, 100, 1), # Number of estimators
    'learning_rate': hp.loguniform('learning_rate', -5, 0), # Learning rate (log scale)
}
```

Define the objective function to optimize (accuracy in this case) python: This step defines the objective function that Hyperopt will try to minimize. The function takes a set of hyperparameters as input, constructs an AdaBoostClassifier with a base weak learner (DecisionTreeClassifier with max\_depth=1), and evaluates its accuracy on the test set.

```
def objective(params):
    base_model = DecisionTreeClassifier(max_depth=1) # Weak Learner (e.g., decision tree with max_depth=1)
    adaboost_model = AdaBoostClassifier(
        base_estimator=base_model,
        n_estimators=int(params['n_estimators']),
        learning_rate=params['learning_rate'],
        random_state=42
    )
    adaboost_model.fit(X_train, y_train)
    y_pred = adaboost_model.predict(X_test)
    accuracy = accuracy_score(y_test, y_pred)
    return {'loss': -accuracy, 'status': STATUS_OK}
```

Perform Hyperband optimization with progress bar: This step uses Hyperopt to perform Hyperband optimization. It runs the optimization for a maximum of 50 evaluations, using the TPE (Tree-structured Parzen Estimator) algorithm. The verbose=1 argument displays a progress bar using the tqdm module.

```
trials = Trials()

# Define the progress bar using tqdm
best = fmin(fn=objective, space=space, algo=tpe.suggest, max_evals=50, trials=trials, verbose=1)
```

Get the best hyperparameters and train the final model: This step extracts the best hyperparameters found by Hyperopt and trains the final AdaBoost model using these hyperparameters.

```
best_n_estimators = int(best['n_estimators'])
best_learning_rate = best['learning_rate']
final_base_model = DecisionTreeClassifier(max_depth=1) # Weak learner (e.g., decision tree with max_depth=1)
final_adaboost_model = AdaBoostClassifier(
    base_estimator=final_base_model,
    n_estimators=best_n_estimators,
    learning_rate=best_learning_rate,
    random_state=42
)
final_adaboost_model.fit(X_train, y_train)
```

Make predictions on the test data and evaluate the final model's performance: This step makes predictions on the test set using the final trained model and evaluates its accuracy.

```
y_pred = final_adaboost_model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f"Final Model Accuracy: {accuracy}")
```

Print the best parameters and confusion matrix: Finally, this step prints the best hyperparameters and the confusion matrix for the final model's predictions on the test set. The confusion matrix provides information about the model's performance in terms of true positives, true negatives, false positives, and false negatives.

```
print("Best Parameters:")
print(f"n_estimators: {best_n_estimators}, learning_rate: {best_learning_rate}")

conf_matrix = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:")
print(conf_matrix)
```

## (ii) Grid Search CV Optimization :-

Define the hyperparameter search space for Randomized Search: In this step, a dictionary random\_space is defined with the hyperparameters and their corresponding search spaces for Randomized Search. In this case, it's exploring the number of estimators (n\_estimators) over the range from 10 to 100.

```
# Step 4: Define the hyperparameter search space for Randomized Search
random_space = {
    'n_estimators': range(10, 101),                                # Number of estimators (random search over a ran
```

Perform Randomized Search for initial exploration: This step uses scikit-learn's RandomizedSearchCV to perform a randomized search over the hyperparameter space defined in random\_space. It explores 20 random combinations of hyperparameters and selects the combination that maximizes accuracy. The best hyperparameters found during Randomized Search are stored in best\_n\_estimators\_rs.

```
# Step 5: Perform Randomized Search for initial exploration
random_search = RandomizedSearchCV(
    AdaBoostClassifier(base_estimator=DecisionTreeClassifier(max_depth=1)),
    param_distributions=random_space,
    n_iter=20, # Number of random combinations to try
    scoring='accuracy',
    cv=3,      # Cross-validation folds
    random_state=42,
    verbose=1
)

# Fit the Randomized Search
random_search.fit(X_train, y_train)

# Get the best hyperparameters found during Randomized Search
best_n_estimators_rs = random_search.best_params_['n_estimators']
```

Define the hyperparameter search space for Hyperband and the objective function: This step defines the hyperparameter search space for Hyperband optimization (space). The hp.choice function is used to specify a discrete choice for n\_estimators. The objective function is defined similarly to before.

```
# Step 6: Define the hyperparameter search space for Hyperband (using best parameters from Randomized Search)
space = {
    'n_estimators': hp.choice('n_estimators', range(10, 101)),          # Number of estimators
}

# Step 7: Define the objective function to optimize (accuracy in this case)
def objective(params):
    base_model = DecisionTreeClassifier(max_depth=1) # Weak Learner (e.g., decision tree with max_
    learning_rate = 10 ** random.uniform(-5, 0) # Generate Learning rate in log scale
    adaboost_model = AdaBoostClassifier(
        base_estimator=base_model,
        n_estimators=params['n_estimators'],
        learning_rate=learning_rate,
        random_state=42
    )
    adaboost_model.fit(X_train, y_train)
    y_pred = adaboost_model.predict(X_test)
    accuracy = accuracy_score(y_test, y_pred)
    return {'loss': -accuracy, 'status': STATUS_OK}

# Step 8: Perform Hyperband optimization with progress bar
trials = Trials()

# Define the progress bar using tqdm
best = fmin(fn=objective, space=space, algo=tpe.suggest, max_evals=50, trials=trials, verbose=1)
```

Perform Hyperband optimization with a progress bar.

```
# Step 8: Perform Hyperband optimization with progress bar
trials = Trials()

# Define the progress bar using tqdm
best = fmin(fn=objective, space=space, algo=tpe.suggest, max_evals=50, trials=trials, verbose=1)
```

Get the best hyperparameters from Hyperband and train the final model, then evaluate its performance: This step extracts the best hyperparameters found by Hyperband and trains the final AdaBoost model using these hyperparameters. It then evaluates the model's performance on the test set and prints the final model accuracy and confusion matrix.

```
# Step 9: Get the best hyperparameters and train the final model
best_n_estimators = int(best['n_estimators'])
best_learning_rate = 10 ** random.uniform(-5, 0) # Generate Learning rate in Log scale
final_base_model = DecisionTreeClassifier(max_depth=1) # Weak Learner (e.g., decision tree with max_depth=1)
final_adaboost_model = AdaBoostClassifier(
    base_estimator=final_base_model,
    n_estimators=best_n_estimators,
    learning_rate=best_learning_rate,
    random_state=42
)
final_adaboost_model.fit(X_train, y_train)

# Step 10: Make predictions on the test data and evaluate the final model's performance
y_pred = final_adaboost_model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f"Final Model Accuracy: {accuracy}")
```

Print the best parameters from both Randomized Search and Hyperband, and confusion matrix: This step prints the best hyperparameters found by both Randomized Search and Hyperband, as well as the confusion matrix for the final model's predictions on the test set. The confusion matrix provides information about the model's performance in terms of true positives, true negatives, false positives, and false negatives.

```
# Step 10: Make predictions on the test data and evaluate the final model's performance
y_pred = final_adaboost_model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f"Final Model Accuracy: {accuracy}")

# Step 11: Print the best parameters from both Randomized Search and Hyperband, and confusion matrix
print("Best Parameters (Randomized Search):")
print(f"n_estimators: {best_n_estimators_rs}")

print("Best Parameters (Hyperband):")
print(f"n_estimators: {best_n_estimators}, learning_rate: {best_learning_rate}")

conf_matrix = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:")
print(conf_matrix)
```

### (iii) Gradient Based Optimization :-

Define the hyperparameter search space for Hyperopt (using a custom function for n\_estimators): This step defines the hyperparameter search space for Hyperopt optimization. For n\_estimators, it uses a custom function (sample\_n\_estimators) to randomly choose a value from the range [10, 100]. For learning\_rate, it uses a log-uniform distribution.

```
# Step 4: Define the hyperparameter search space for Hyperopt (using a custom function for n_estimators)
def sample_n_estimators():
    return random.choice(range(10, 101))

space = {
    'n_estimators': hp.choice('n_estimators', range(10, 101)),           # Number of estimators (custom function)
    'learning_rate': hp.loguniform('learning_rate', -5, 0)                 # Learning rate (log scale)
}
```

Define the objective function to optimize (accuracy in this case): This step defines the objective function that Hyperopt will try to minimize. The function takes a set of hyperparameters as input, constructs an AdaBoostClassifier with a base weak learner (DecisionTreeClassifier with max\_depth=1), and evaluates its accuracy on the test set.

```
def objective(params):
    base_model = DecisionTreeClassifier(max_depth=1) # Weak Learner (e.g., decision tree with max_depth=1)
    adaboost_model = AdaBoostClassifier(
        base_estimator=base_model,
        n_estimators=sample_n_estimators(), # Sample n_estimators using custom function
        learning_rate=params['learning_rate'],
        random_state=42
    )
    adaboost_model.fit(X_train, y_train)
    y_pred = adaboost_model.predict(X_test)
    accuracy = accuracy_score(y_test, y_pred)
    return {'loss': -accuracy, 'status': STATUS_OK}
```

Perform Hyperopt optimization with a progress bar: This step uses Hyperopt to perform optimization. It runs the optimization for a maximum of 50 evaluations, using the TPE (Tree-structured Parzen Estimator) algorithm. The verbose=1 argument displays a progress bar using the tqdm module.

```
# Step 6: Perform Hyperopt optimization with a progress bar
trials = Trials()

# Define the progress bar using tqdm
best = fmin(fn=objective, space=space, algo=tpe.suggest, max_evals=50, trials=trials, verbose=1)
```

Get the best hyperparameters and train the final model: This step extracts the best hyperparameters found by Hyperopt and trains the final AdaBoost model using these hyperparameters.

```
best_n_estimators = int(best['n_estimators'])
best_learning_rate = best['learning_rate']
final_base_model = DecisionTreeClassifier(max_depth=1) # Weak Learner (e.g., decision tree with max depth 1)
final_adaboost_model = AdaBoostClassifier(
    base_estimator=final_base_model,
    n_estimators=best_n_estimators,
    learning_rate=best_learning_rate,
    random_state=42
)
final_adaboost_model.fit(X_train, y_train)
```

Make predictions on the test data and evaluate the final model's performance: This step makes predictions on the test set using the final trained model and evaluates its accuracy.

```
y_pred = final_adaboost_model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f"Final Model Accuracy: {accuracy}")
```

Print the best hyperparameters and confusion matrix: Finally, this step prints the best hyperparameters and the confusion matrix for the final model's predictions on the test set. The confusion matrix provides information about the model's performance in terms of true positives, true negatives, false positives, and false negatives.

```
# Step 9: Print the best hyperparameters and confusion matrix
print("Best Parameters:")
print(f"n_estimators: {best_n_estimators}, learning_rate: {best_learning_rate}")

conf_matrix = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:")
print(conf_matrix)
```

## (b) Decision Tree :-

### (i) Hyperband Optimization :-

Here we have carried out the optimization in a similar manner as done in AdaBoost Model. We will discuss the different approach and parameters optimized.

Define the objective function for Hyperband optimization :-

- objective\_function is defined, taking hyperparameters as input.
- A Decision Tree model is created with specified hyperparameters.
- The model is trained on the training set (X\_train, y\_train).
- The accuracy on the test set is calculated and used to define the loss.

```
def objective_function(params):
    dt_model = DecisionTreeClassifier(
        max_depth=int(params['max_depth']),
        min_samples_split=int(params['min_samples_split']),
        min_samples_leaf=int(params['min_samples_leaf']),
        random_state=42
    )
    dt_model.fit(X_train, y_train)

    y_pred = dt_model.predict(X_test)
    accuracy = accuracy_score(y_test, y_pred)

    return {'loss': 1 - accuracy, 'status': STATUS_OK}
```

Define the hyperparameter search space :-

- A search space for hyperparameters is defined using hyperopt.
- hp.quniform is used for discrete uniform distributions.

```
space = {
    'max_depth': hp.quniform('max_depth', 5, 30, 1), # Discrete values from 5 to 30 with step 1
    'min_samples_split': hp.quniform('min_samples_split', 2, 20, 1), # Discrete values from 2 to 20 with step 1
    'min_samples_leaf': hp.quniform('min_samples_leaf', 1, 10, 1) # Discrete values from 1 to 10 with step 1
}
```

Perform Hyperband optimization :-

- Hyperband optimization is performed using fmin from hyperopt.
- fn is set to the objective function.
- space is the search space for hyperparameters.
- algo=tpe.suggest is the optimization algorithm (Tree of Parzen Estimators).
- max\_evals=100 is the maximum number of evaluations.

```
trials = Trials()
best = fmin(fn=objective_function,
            space=space,
            algo=tpe.suggest,
            max_evals=100,
            trials=trials,
            verbose=1)
```

Get the best hyperparameters and train the model :-

- The best hyperparameters obtained from optimization are used to create a Decision Tree model (best\_dt\_model).
- The model is trained on the training set.

```
best_dt_model = DecisionTreeClassifier(
    max_depth=int(best['max_depth']),
    min_samples_split=int(best['min_samples_split']),
    min_samples_leaf=int(best['min_samples_leaf']),
    random_state=42
)
best_dt_model.fit(X_train, y_train)
```

(ii) Optuna :-

Define the objective function for Optuna optimization :-

- In Optuna, the objective function is defined to take a trial object as an argument.
- The hyperparameters (max\_depth, min\_samples\_split, and min\_samples\_leaf) are defined using trial.suggest\_int.
- The objective is to maximize accuracy, so the function returns the accuracy.

```
def objective(trial):
    # Define the hyperparameter search space for Decision Tree
    max_depth = trial.suggest_int('max_depth', 5, 30)
    min_samples_split = trial.suggest_int('min_samples_split', 2, 20)
    min_samples_leaf = trial.suggest_int('min_samples_leaf', 1, 10)

    dt_model = DecisionTreeClassifier(
        max_depth=max_depth,
        min_samples_split=min_samples_split,
        min_samples_leaf=min_samples_leaf,
        random_state=42
    )
    dt_model.fit(X_train, y_train)

    y_pred = dt_model.predict(X_test)
    accuracy = accuracy_score(y_test, y_pred)

    return accuracy
```

Perform Optuna optimization with a progress bar :-

- An Optuna study is created with the direction set to 'maximize'.
- The optimization is performed using study.optimize, and the progress bar is updated using tqdm.
- A callback function (progress\_callback) is defined to update the progress bar.

```

study = optuna.create_study(direction='maximize')
with tqdm(total=100) as pbar:
    def progress_callback(study, trial):
        pbar.update(1)

study.optimize(objective, n_trials=100, show_progress_bar=False, callbacks=[progress_callback])

```

Get the best hyperparameters and train the model with them :-

- Optuna provides a `study.best_params` attribute to get the best hyperparameters.
- The best hyperparameters are used to create the best Decision Tree model (`best_dt_model`), which is then trained on the training set.

```

best_params = study.best_params
best_dt_model = DecisionTreeClassifier(
    max_depth=best_params['max_depth'],
    min_samples_split=best_params['min_samples_split'],
    min_samples_leaf=best_params['min_samples_leaf'],
    random_state=42
)
best_dt_model.fit(X_train, y_train)

```

### (iii) Bayesian Optimization :-

Define the objective function for Bayesian Optimization :-

- In Bayesian Optimization, the objective function takes hyperparameters (`n_estimators`, `max_depth`, `min_samples_split`, `min_samples_leaf`) as separate parameters.
- The function creates and trains a Decision Tree model with the suggested hyperparameters.
- The accuracy on the test set is used as the objective to be maximized.

```

def objective_function(n_estimators, max_depth, min_samples_split, min_samples_leaf):
    n_estimators = int(n_estimators)
    max_depth = int(max_depth)
    min_samples_split = int(min_samples_split)
    min_samples_leaf = int(min_samples_leaf)

    # Create and train the Decision Tree model with the suggested hyperparameters
    dt_model = DecisionTreeClassifier(
        max_depth=max_depth,
        min_samples_split=min_samples_split,
        min_samples_leaf=min_samples_leaf,
        random_state=42
    )
    dt_model.fit(X_train, y_train)

    # Make predictions on the test data
    y_pred = dt_model.predict(X_test)

    # Calculate accuracy
    accuracy = accuracy_score(y_test, y_pred)
    return accuracy

```

Define the hyperparameter search space :-

- The hyperparameter search space is defined using pbounds in Bayesian Optimization.
- It specifies the ranges for each hyperparameter to be explored.

```
pbounds = {
    'n_estimators': (50, 200),
    'max_depth': (5, 30),
    'min_samples_split': (2, 20),
    'min_samples_leaf': (1, 10)
}
```

Perform Bayesian Optimization with a custom progress bar :-

- Bayesian Optimization is performed using the BayesianOptimization class from the bayes\_opt library.
- A custom progress bar is implemented using tqdm to visualize the progress of Bayesian optimization.
- The optimization is performed for a specified number of iterations (total\_iterations).

```
optimizer = BayesianOptimization(f=objective_function, pbounds=pbounds, random_state=42)
total_iterations = 100 # Total number of Bayesian optimization iterations
with tqdm(total=total_iterations) as pbar:
    for _ in range(total_iterations):
        optimizer.maximize(init_points=1, n_iter=1)
        pbar.update(1)
```

Get the best hyperparameters and train the model :-

- The best hyperparameters are obtained from the optimizer.max['params'].
- The best Decision Tree model (best\_dt\_model) is created with these hyperparameters and trained on the training set.

```
best_params = optimizer.max['params']
best_dt_model = DecisionTreeClassifier(
    max_depth=int(best_params['max_depth']),
    min_samples_split=int(best_params['min_samples_split']),
    min_samples_leaf=int(best_params['min_samples_leaf']),
    random_state=42
)
best_dt_model.fit(X_train, y_train)
```

### (c) Random Forest :-

#### (i) Hyperband Optimization :-

Here we have carried out the optimization in a similar manner as done in AdaBoost Model. We will discuss the different approach and parameters optimized.

Define the objective function for Optuna to optimize :-

- The objective function is defined to take a trial object as an argument.
- The hyperparameters (`n_estimators`, `max_depth`, `min_samples_split`, `min_samples_leaf`) are defined using `trial.suggest_int`.
- The function creates and trains a Random Forest model with the suggested hyperparameters and returns the accuracy on the test set.

```
def objective(trial):
    # Define the hyperparameter search space
    n_estimators = trial.suggest_int('n_estimators', 50, 200)
    max_depth = trial.suggest_int('max_depth', 5, 50)
    min_samples_split = trial.suggest_int('min_samples_split', 2, 20)
    min_samples_leaf = trial.suggest_int('min_samples_leaf', 1, 10)

    # Create and train the Random Forest model with the suggested hyperparameters
    rf_model = RandomForestClassifier(
        n_estimators=n_estimators,
        max_depth=max_depth,
        min_samples_split=min_samples_split,
        min_samples_leaf=min_samples_leaf,
        random_state=42
    )
    rf_model.fit(X_train, y_train)

    # Make predictions on the test data
    y_pred = rf_model.predict(X_test)

    # Calculate accuracy
    accuracy = accuracy_score(y_test, y_pred)
    return accuracy
```

Perform Hyperband optimization with a progress bar :-

- An Optuna study is created with the direction set to 'maximize'.
- The optimization is performed using `study.optimize`, and a progress bar is implemented using `tqdm`.
- The optimization is performed for a specified number of trials (`n_trials=100`).

```

study = optuna.create_study(direction='maximize')

# Add tqdm to show a progress bar during optimization
with tqdm(total=100) as pbar:
    def callback(study, trial):
        pbar.update(1)

    study.optimize(objective, n_trials=100, callbacks=[callback])

```

Get the best hyperparameters and train the model :-

- The best hyperparameters are obtained from the study.best\_params.
- The best Random Forest model (best\_rf\_model) is created with these hyperparameters and trained on the training set.

```

best_params = study.best_params
best_rf_model = RandomForestClassifier(
    n_estimators=best_params['n_estimators'],
    max_depth=best_params['max_depth'],
    min_samples_split=best_params['min_samples_split'],
    min_samples_leaf=best_params['min_samples_leaf'],
    random_state=42
)
best_rf_model.fit(X_train, y_train)

```

## (ii) Optuna :-

Here we have carried out the optimization in a similar manner as done in Decision Tree Model. We will discuss the different approach and parameters optimized.

Define the objective function for Optuna to optimize :-

- The objective function is defined to take a trial object as an argument.
- The hyperparameters (n\_estimators, max\_depth, min\_samples\_split, min\_samples\_leaf) are defined using trial.suggest\_int.
- The function creates and trains a Random Forest model with the suggested hyperparameters and returns the accuracy on the testset.

```

def objective(trial):
    # Define the hyperparameter search space
    n_estimators = trial.suggest_int('n_estimators', 50, 200)
    max_depth = trial.suggest_int('max_depth', 5, 50)
    min_samples_split = trial.suggest_int('min_samples_split', 2, 20)
    min_samples_leaf = trial.suggest_int('min_samples_leaf', 1, 10)

    # Create and train the Random Forest model with the suggested hyperparameters
    rf_model = RandomForestClassifier(
        n_estimators=n_estimators,
        max_depth=max_depth,
        min_samples_split=min_samples_split,
        min_samples_leaf=min_samples_leaf,
        random_state=42
    )
    rf_model.fit(X_train, y_train)

    # Make predictions on the test data
    y_pred = rf_model.predict(X_test)

    # Calculate accuracy
    accuracy = accuracy_score(y_test, y_pred)
    return accuracy

```

Perform Optuna optimization with TPE algorithm :-

- An Optuna study is created with the direction set to 'maximize'.
- The TPE (Tree-structured Parzen Estimator) algorithm is used as the sampler.
- The optimization is performed using study.optimize for a specified number of trials (n\_trials=100).

```
study = optuna.create_study(direction='maximize', sampler=optuna.samplers.TPESampler())
study.optimize(objective, n_trials=100)
```

Get the best hyperparameters and train the model :-

- The best hyperparameters are obtained from the study.best\_params.
- The best Random Forest model (best\_rf\_model) is created with these hyperparameters and trained on the training set.

```
best_params = study.best_params
best_rf_model = RandomForestClassifier(
    n_estimators=best_params['n_estimators'],
    max_depth=best_params['max_depth'],
    min_samples_split=best_params['min_samples_split'],
    min_samples_leaf=best_params['min_samples_leaf'],
    random_state=42
)
best_rf_model.fit(X_train, y_train)
```

### (iii) Bayesian Optimization :-

Here we have carried out the optimization in a similar manner as done in Decision Tree Model. We will discuss the different approach and parameters optimized.

Define the objective function for Bayesian Optimization :-

- The objective\_function takes hyperparameters (n\_estimators, max\_depth, min\_samples\_split, min\_samples\_leaf) as parameters.
- It creates and trains a Random Forest model with the suggested hyperparameters and returns the accuracy on the test set.

```

def objective_function(n_estimators, max_depth, min_samples_split, min_samples_leaf):
    n_estimators = int(n_estimators)
    max_depth = int(max_depth)
    min_samples_split = int(min_samples_split)
    min_samples_leaf = int(min_samples_leaf)

    # Create and train the Random Forest model with the suggested hyperparameters
    rf_model = RandomForestClassifier(
        n_estimators=n_estimators,
        max_depth=max_depth,
        min_samples_split=min_samples_split,
        min_samples_leaf=min_samples_leaf,
        random_state=42
    )
    rf_model.fit(X_train, y_train)

    # Make predictions on the test data
    y_pred = rf_model.predict(X_test)

    # Calculate accuracy
    accuracy = accuracy_score(y_test, y_pred)
    return accuracy

```

Define the hyperparameter search space :-

- The hyperparameter search space is defined using pbounds in Bayesian Optimization.
- It specifies the ranges for each hyperparameter to be explored.

```

pbounds = {
    'n_estimators': (50, 200),
    'max_depth': (5, 30),
    'min_samples_split': (2, 10),
    'min_samples_leaf': (1, 10)
}

```

Perform Bayesian Optimization with a progress bar :-

- Bayesian Optimization is performed using the BayesianOptimization class from the bayes\_opt library.
- A custom progress bar is implemented using tqdm to visualize the progress of Bayesian optimization.
- The optimization is performed for a specified number of iterations (total\_iterations).

```

optimizer = BayesianOptimization(f=objective_function, pbounds=pbounds, random_state=42)
total_iterations = 100 # Total number of Bayesian optimization iterations
with tqdm(total=total_iterations) as pbar:
    for _ in range(total_iterations):
        optimizer.maximize(init_points=1, n_iter=1)
        pbar.update(1)

```

Get the best hyperparameters and train the model :-

- The best hyperparameters are obtained from `optimizer.max['params']`.
- The best Random Forest model (`best_rf_model`) is created with these hyperparameters and trained on the training set.

```
best_params = optimizer.max['params']
best_rf_model = RandomForestClassifier(
    n_estimators=int(best_params['n_estimators']),
    max_depth=int(best_params['max_depth']),
    min_samples_split=int(best_params['min_samples_split']),
    min_samples_leaf=int(best_params['min_samples_leaf']),
    random_state=42
)
best_rf_model.fit(X_train, y_train)
```

(d) Extremely Randomized Tree :-

(i) Hyperband Optimization :-

Here we have carried out the optimization in a similar manner as done in AdaBoost Model. We will discuss the different approach and parameters optimized.

Define the objective function for Optuna to optimize :-

- The objective function takes hyperparameters (`n_estimators`, `max_depth`, `min_samples_split`, `min_samples_leaf`) as parameters.
- It creates and trains an Extra Trees model with the suggested hyperparameters and returns the accuracy on the test set.

```
def objective(trial):
    # Define the hyperparameter search space
    n_estimators = trial.suggest_int('n_estimators', 50, 200)
    max_depth = trial.suggest_int('max_depth', 5, 50)
    min_samples_split = trial.suggest_int('min_samples_split', 2, 20)
    min_samples_leaf = trial.suggest_int('min_samples_leaf', 1, 10)

    # Create and train the Extra Trees model with the suggested hyperparameters
    et_model = ExtraTreesClassifier(
        n_estimators=n_estimators,
        max_depth=max_depth,
        min_samples_split=min_samples_split,
        min_samples_leaf=min_samples_leaf,
        random_state=42
    )
    et_model.fit(X_train, y_train)

    # Make predictions on the test data
    y_pred = et_model.predict(X_test)

    # Calculate accuracy
    accuracy = accuracy_score(y_test, y_pred)
    return accuracy
```

Perform Hyperband optimization with a progress bar :-

- The SkoptSampler is used as a sampler for optimization.
- Bayesian optimization is performed using the optuna.create\_study and study.optimize functions.
- A progress bar using tqdm is displayed during optimization.

```
sampler = optuna.integration.SkoptSampler()
study = optuna.create_study(direction='maximize', sampler=sampler)

# Add tqdm to show a progress bar during optimization
with tqdm(total=100) as pbar:
    def callback(study, trial):
        pbar.update(1)

study.optimize(objective, n_trials=100, callbacks=[callback])
```

Get the best hyperparameters and train the model :-

- The best hyperparameters are obtained from study.best\_params.
- The best Extra Trees model (best\_et\_model) is created with these hyperparameters and trained on the training set.

```
best_params = study.best_params
best_et_model = ExtraTreesClassifier(
    n_estimators=best_params['n_estimators'],
    max_depth=best_params['max_depth'],
    min_samples_split=best_params['min_samples_split'],
    min_samples_leaf=best_params['min_samples_leaf'],
    random_state=42
)
best_et_model.fit(X_train, y_train)
```

## (ii) Optuna :-

Here we have carried out the optimization in a similar manner as done in Decision Tree Model. We will discuss the different approach and parameters optimized.

Define the objective function for Optuna to optimize :-

- The objective function takes hyperparameters (n\_estimators, max\_depth, min\_samples\_split, min\_samples\_leaf) as parameters.
- It creates and trains an Extra Trees model with the suggested hyperparameters and returns the accuracy on the test set.

```
def objective(trial):
    # Define the hyperparameter search space
    n_estimators = trial.suggest_int('n_estimators', 50, 200)
    max_depth = trial.suggest_int('max_depth', 5, 50)
    min_samples_split = trial.suggest_int('min_samples_split', 2, 20)
    min_samples_leaf = trial.suggest_int('min_samples_leaf', 1, 10)

    # Create and train the Extra Trees model with the suggested hyperparameters
    et_model = ExtraTreesClassifier(
        n_estimators=n_estimators,
        max_depth=max_depth,
        min_samples_split=min_samples_split,
        min_samples_leaf=min_samples_leaf,
        random_state=42
    )
    et_model.fit(X_train, y_train)

    # Make predictions on the test data
    y_pred = et_model.predict(X_test)

    # Calculate accuracy
    accuracy = accuracy_score(y_test, y_pred)
    return accuracy
```

Perform Optuna optimization with a progress bar :-

- Bayesian optimization is performed using the optuna.create\_study and study.optimize functions.
- A progress bar using tqdm is displayed during optimization.

```
study = optuna.create_study(direction='maximize')

# Add tqdm to show a progress bar during optimization
with tqdm(total=100) as pbar:
    def callback(study, trial):
        pbar.update(1)

study.optimize(objective, n_trials=100, callbacks=[callback])
```

Get the best hyperparameters and train the model :-

- The best hyperparameters are obtained from study.best\_params.
- The best Extra Trees model (best\_et\_model) is created with these hyperparameters and trained on the training set.

```
best_params = study.best_params
best_et_model = ExtraTreesClassifier(
    n_estimators=best_params['n_estimators'],
    max_depth=best_params['max_depth'],
    min_samples_split=best_params['min_samples_split'],
    min_samples_leaf=best_params['min_samples_leaf'],
    random_state=42
)
best_et_model.fit(X_train, y_train)
```

### (iii) Bayesian Optimization :-

Here we have carried out the optimization in a similar manner as done in Decision Tree Model. We will discuss the different approach and parameters optimized.

Define the objective function for Bayesian Optimization :-

- The objective\_function takes hyperparameters (n\_estimators, max\_depth, min\_samples\_split, min\_samples\_leaf) as parameters.
- It creates and trains an Extra Trees model with the suggested hyperparameters and returns the accuracy on the test set.

```
def objective_function(n_estimators, max_depth, min_samples_split, min_samples_leaf):
    n_estimators = int(n_estimators)
    max_depth = int(max_depth)
    min_samples_split = int(min_samples_split)
    min_samples_leaf = int(min_samples_leaf)

    # Create and train the Extra Trees model with the suggested hyperparameters
    et_model = ExtraTreesClassifier(
        n_estimators=n_estimators,
        max_depth=max_depth,
        min_samples_split=min_samples_split,
        min_samples_leaf=min_samples_leaf,
        random_state=42
    )
    et_model.fit(X_train, y_train)

    # Make predictions on the test data
    y_pred = et_model.predict(X_test)

    # Calculate accuracy
    accuracy = accuracy_score(y_test, y_pred)
    return accuracy
```

Define the hyperparameter search space: pbounds defines the search space for each hyperparameter.

```
pbounds = {
    'n_estimators': (50, 200),
    'max_depth': (5, 50),
    'min_samples_split': (2, 20),
    'min_samples_leaf': (1, 10)
}
```

Perform Bayesian Optimization :-

- Bayesian Optimization is performed using the BayesianOptimization class.
- init\_points specifies the number of initial random points to sample.
- n\_iter is the number of optimization steps.

```
optimizer = BayesianOptimization(f=objective_function, pbounds=pbounds, random_state=42)
optimizer.maximize(init_points=10, n_iter=90)
```

Get the best hyperparameters and train the model :-

- The best hyperparameters are obtained from optimizer.max['params'].
- The best Extra Trees model (best\_et\_model) is created with these hyperparameters and trained on the training set.

### 11.7 Results (After Optimization) :-

The four fundamental models optimized models along with their algorithms' results are given below, the best optimized model is highlighted and will be taken further in our study for this dataset. After optimization, the new and better optimized models along with other models are given below and these are the final models which we used to make our hybrid model for the UNSW-NB15.

Table-05 :: Individual Model Training for UNSW-NB15(optimization)									
Model	Optimization Algorithm	AUC	CA	F1	Precision	Recall	MCC	Spec	LogLoss
AdaBoost	<b>Hyperband</b>	<b>0.8981</b>	<b>91.15</b>	<b>0.8743</b>	<b>0.8493</b>	<b>0.9008</b>	<b>0.8070</b>	<b>0.9171</b>	<b>0.6294</b>
	Gradient Based	0.8721	89.67	0.8460	0.7827	0.9205	0.7784	0.8862	0.6106
	Grid Search	0.7549	80.96	0.6792	0.5563	0.8718	0.5812	0.7909	0.4828
Decision Tree	Hyperband	0.9297	93.23	0.9070	0.9202	0.8959	0.8546	0.9539	0.4092
	<b>Optuna</b>	<b>0.9317</b>	<b>93.34</b>	<b>0.9097</b>	<b>0.9254</b>	<b>0.8946</b>	<b>0.8573</b>	<b>0.9567</b>	<b>0.2603</b>
	Bayes	0.9289	93.22	0.9075	0.9169	0.8982	0.8542	0.9522	0.3659
Random Forest	<b>Hyperband</b>	<b>0.9395</b>	<b>94.30</b>	<b>0.9218</b>	<b>0.9266</b>	<b>0.9171</b>	<b>0.8770</b>	<b>0.9580</b>	<b>0.1259</b>
	Optuna	0.9373	94.26	0.9206	0.9184	0.9229	0.8757	0.9537	0.1219
	Bayes	0.9386	94.21	0.9207	0.9258	0.9155	0.8752	0.9576	0.1319
Extremely Randomized Tree	<b>Hyperband</b>	<b>0.9339</b>	<b>93.95</b>	<b>0.9218</b>	<b>0.9266</b>	<b>0.9171</b>	<b>0.8770</b>	<b>0.9580</b>	<b>0.1264</b>
	Optuna	0.9340	93.94	0.9163	0.9143	0.9183	0.8689	0.9514	0.1291
	Bayes	0.9339	93.94	0.9163	0.9141	0.9185	0.8689	0.9513	0.1271

Table-06 :: Individual Model Training for UNSW-NB15(after optimization)									
Model	Optimization Algorithm	AUC	CA	F1	Precision	Recall	MCC	Spec	LogLoss
AdaBoost	Hyperband	0.8981	91.15	0.8743	0.8493	0.9008	0.8070	0.9171	0.6294
Decision Tree	Optuna	0.9317	93.34	0.9097	0.9254	0.8946	0.8573	0.9567	0.2603
Random Forest	Hyperband	0.9395	94.30	0.9218	0.9266	0.9171	0.8770	0.9580	0.1259
Gradient Boosting	None	0.9091	92.16	0.8887	0.8639	0.9149	0.8291	0.9250	0.1633
Extremely Randomized Tree	Hyperband	0.9339	93.95	0.9218	0.9266	0.9171	0.8770	0.9580	0.1264
Neural Network	None	0.6070	71.49	0.3549	0.2161	0.9929	0.3837	0.6913	-

## 11.8 Ensemble Methods :-

### Libraries Used for Each Method :-

```

import pandas as pd
import numpy as np
import pickle
from sklearn.metrics import accuracy_score, roc_auc_score, log_loss, confusion_matrix
import pandas as pd
import numpy as np
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.metrics import accuracy_score, roc_auc_score, classification_report, matthews_corrcoef, log_loss, confusion_matrix
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
import time
import joblib

from sklearn.metrics import roc_curve, auc
import matplotlib.pyplot as plt

```

The code uses Python libraries for data manipulation and machine learning :-

- (i) pandas and numpy for data handling
- (ii) pickle for object serialization
- (iii) scikit-learn for building a Gradient Boosting Classifier, imputing missing values, and scaling features
- (iv) time for measuring execution time
- (v) joblib for lightweight pipelining.

The goal is to train, evaluating, and saving a machine learning model, for a model Classifier.

### 11.8.1 Hard Ensemble Technique:-

```

def load_models(model_paths):
    models = []
    for path in model_paths:
        model = joblib.load(path)
        models.append(model)
    return models

```

- (i) The function initializes an empty list called models to store the loaded models.
- (ii) It iterates over each file path in the input list (model\_paths).
- (iii) For each path, it uses the joblib.load function to load the machine learning model stored at that path.
- (iv) The loaded model is then appended to the models list.
- (v) Once all models are loaded, the function returns the list of loaded models.

The load\_models function takes a list of file paths, loads machine learning models stored at those paths using joblib, and returns a list of the loaded models.

```
def ensemble_voting(models, data):
    predictions = []
    for model in models:
        prediction = model.predict(data)
        predictions.append(prediction.astype(int)) # Convert predictions to integers
    final_predictions = np.apply_along_axis(lambda x: np.argmax(np.bincount(x)), axis=0, arr=predictions)
    return final_predictions
```

- (i) The function initializes an empty list called predictions to store individual predictions made by each model.
- (ii) It iterates over each model in the provided list (models).
- (iii) For each model, it uses the predict method to obtain predictions for the input dataset (data).
- (iv) The individual predictions are appended to the predictions list after converting them to integers.
- (v) After obtaining predictions from all models, the function uses np.apply\_along\_axis to apply a voting mechanism along the specified axis (axis=0) of the array formed by stacking the individual predictions.
- (vi) Within the voting mechanism, np.argmax(np.bincount(x)) is used to find the most frequent prediction for each data point. This is achieved by finding the index with the maximum count in the bincount of each column.
- (vii) The final set of predictions is returned.

The ensemble\_voting function takes a list of machine learning models and a dataset, makes predictions using each model, and combines the predictions through a voting mechanism to produce a final set of predictions. The voting is based on the most frequently predicted class for each data point across all models.

### 11.8.2 Soft Ensemble Technique :-

```
def load_models(model_paths):
    models = []
    for path in model_paths:
        model = joblib.load(path)
        models.append(model)
    return models
```

- (i) The function initializes an empty list called models to store the loaded models.
- (ii) It iterates over each file path in the input list (model\_paths).
- (iii) For each path, it uses the joblib.load function to load the machine learning model stored at that path.
- (iv) The loaded model is then appended to the models list.
- (v) Once all models are loaded, the function returns the list of loaded models.

The load\_models function takes a list of file paths, loads machine learning models stored at those paths using joblib, and returns a list of the loaded models.

```
def ensemble_soft_voting(models, data):
    predictions = []
    for model in models:
        prediction_probabilities = model.predict_proba(data)
        predictions.append(prediction_probabilities) # Store the probability distributions
    final_predictions = np.mean(predictions, axis=0) # Calculate the mean probabilities
    return np.argmax(final_predictions, axis=1)
```

- (i) The function initializes an empty list called predictions to store the probability distributions predicted by each model.
- (ii) It iterates over each model in the provided list (models).
- (iii) For each model, it uses the predict\_proba method to obtain the probability distributions for each class for the input dataset (data).
- (iv) The probability distributions are appended to the predictions list.
- (v) After obtaining probability distributions from all models, the function uses np.mean along the specified axis (axis=0) to calculate the mean probabilities for each class.

- (vi) The final predictions are obtained by selecting the class index with the highest mean probability using np.argmax along axis 1.
- (vii) The final set of predictions is returned.

### 11.8.3 Bayesian Model Average :-

```
def load_models(model_paths, X_train, y_train):
    models = []
    for path in model_paths:
        model = joblib.load(path)
        # Calibrate the model to obtain probability estimates
        calibrated_model = CalibratedClassifierCV(model, method='sigmoid', cv='prefit')
        calibrated_model.fit(X_train, y_train)
        models.append(calibrated_model)
    return models
```

- (i) The function initializes an empty list called models to store the calibrated models.
- (ii) It iterates over each file path in the input list (model\_paths).
- (iii) For each path, it uses the joblib.load function to load the machine learning model stored at that path.
- (iv) The loaded model is then calibrated using the CalibratedClassifierCV class. The calibration method used is 'sigmoid', and the calibration is based on the training data provided (X\_train and y\_train).
- (v) The calibrated model is added to the models list.
- (vi) Once all models are loaded and calibrated, the function returns the list of calibrated models.

The load\_models function takes a list of file paths, loads machine learning models, calibrates each model using probability estimates, and returns a list of calibrated models. Calibration is performed using the CalibratedClassifierCV class with the 'sigmoid' method, and the training data is used for calibration.

```
def ensemble_bma(models, data):
    predictions = []
    for model in models:
        prediction_probabilities = model.predict_proba(data)
        predictions.append(prediction_probabilities) # Store the probability distributions
    final_predictions = np.mean(predictions, axis=0) # Calculate the mean probabilities
    return np.argmax(final_predictions, axis=1)
```

- (i) The function initializes an empty list called predictions to store the probability distributions predicted by each model.
- (ii) It iterates over each model in the provided list (models).
- (iii) For each model, it uses the predict\_proba method to obtain the probability distributions for each class for the input dataset (data).
- (iv) The probability distributions are appended to the predictions list.
- (v) After obtaining probability distributions from all models, the function uses np.mean along the specified axis (axis=0) to calculate the mean probabilities for each class.
- (vi) The final predictions are obtained by selecting the class index with the highest mean probability using np.argmax along axis 1.
- (vii) The final set of predictions is returned.

The ensemble\_bma function takes a list of machine learning models and a dataset, obtains the probability distributions predicted by each model, calculates the mean probabilities for each class using Bayesian Model Averaging, and produces a final set of predictions based on the class with the highest mean probability. BMA is a probabilistic ensemble method that accounts for uncertainty in model predictions by averaging over multiple models

#### 11.8.4 Dynamic Method :-

```
def load_models(model_paths, X_train, y_train):
    models = []
    for path in model_paths:
        model = joblib.load(path)
        # Calibrate the model to obtain probability estimates
        calibrated_model = CalibratedClassifierCV(model, method='sigmoid', cv='prefit')
        calibrated_model.fit(X_train, y_train)
        models.append(calibrated_model)
    return models
```

- (i) The function initializes an empty list called models to store the calibrated models.
- (ii) It iterates over each file path in the input list (model\_paths).

- (iii) For each path, it uses the joblib.load function to load the machine learning model stored at that path.
- (iv) The loaded model is then calibrated using the CalibratedClassifierCV class. The calibration method used is 'sigmoid', and the calibration is based on the training data provided (X\_train and y\_train).
- (v) The calibrated model is added to the models list.
- (vi) Once all models are loaded and calibrated, the function returns the list of calibrated models.

The load\_models function takes a list of file paths, loads machine learning models, calibrates each model using probability estimates, and returns a list of calibrated models. Calibration is performed using the CalibratedClassifierCV class with the 'sigmoid' method, and the training data is used for calibration.

```
def ensemble_dynamic_voting(models, data):
    predictions = []
    for model in models:
        prediction_probabilities = model.predict_proba(data)
        predictions.append(prediction_probabilities) # Store the probability distributions

    # Calculate confidence scores for each model's predictions
    confidence_scores = np.max(predictions, axis=0)

    # Calculate instance-wise weights for each model's predictions
    instance_weights = np.mean(confidence_scores, axis=1)

    # Apply instance-wise weights to the predictions
    weighted_predictions = np.sum(np.array(predictions) * instance_weights[:, np.newaxis], axis=0)

    return np.argmax(weighted_predictions, axis=1)
```

- (i) The function initializes an empty list called predictions to store the probability distributions predicted by each model.
- (ii) It iterates over each model in the provided list (models).
- (iii) For each model, it uses the predict\_proba method to obtain the probability distributions for each class for the input dataset (data).
- (iv) The probability distributions are appended to the predictions list.
- (v) After obtaining probability distributions from all models, the function calculates confidence scores by finding the maximum probability for each class across all models.

- (vi) Instance-wise weights are calculated by taking the mean of the confidence scores for each data point.
- (vii) The instance-wise weights are applied to the predictions, and the weighted predictions are obtained.
- (viii) The final predictions are obtained by selecting the class index with the highest probability from the weighted predictions.
- (ix) The final set of predictions is returned.

The ensemble\_dynamic\_voting function takes a list of machine learning models and a dataset, obtains the probability distributions predicted by each model, calculates confidence scores and instance-wise weights, and produces a final set of predictions based on dynamic ensemble voting. The dynamic aspect comes from adjusting weights based on instance-wise confidence scores, allowing the ensemble to adapt to the difficulty of individual predictions.

#### 11.8.5 Weighted Average :-

```
def load_models(model_paths, X_train, y_train):
    models = []
    for path in model_paths:
        model = joblib.load(path)
        # Calibrate the model to obtain probability estimates
        calibrated_model = CalibratedClassifierCV(model, method='sigmoid', cv='prefit')
        calibrated_model.fit(X_train, y_train)
        models.append(calibrated_model)
    return models
```

- (i) The function initializes an empty list called models to store the calibrated models.
- (ii) It iterates over each file path in the input list (model\_paths).
- (iii) For each path, it uses the joblib.load function to load the machine learning model stored at that path.
- (iv) The loaded model is then calibrated using the CalibratedClassifierCV class. The calibration method used is 'sigmoid', and the calibration is based on the training data provided (X\_train and y\_train).
- (v) The calibrated model is added to the models list.

- (vi) Once all models are loaded and calibrated, the function returns the list of calibrated models.

The load\_models function takes a list of file paths, loads machine learning models, calibrates each model using probability estimates, and returns a list of calibrated models. Calibration is performed using the CalibratedClassifierCV class with the 'sigmoid' method, and the training data is used for calibration. This process is particularly beneficial when working with models that need adjustments to their probability estimates for better accuracy.

```
def ensemble_weighted_voting(models, data, weights):
    predictions = []
    for model in models:
        prediction_probabilities = model.predict_proba(data)
        predictions.append(prediction_probabilities) # Store the probability distributions

    # Convert weights to a numpy array
    weights_array = np.array(weights)

    # Apply weights to the predictions
    weighted_predictions = np.sum(np.array(predictions) * weights_array[:, np.newaxis, np.newaxis], axis=0)

    return np.argmax(weighted_predictions, axis=1)
```

- (i) The function initializes an empty list called predictions to store the probability distributions predicted by each model.
- (ii) It iterates over each model in the provided list (models).
- (iii) For each model, it uses the predict\_proba method to obtain the probability distributions for each class for the input dataset (data).
- (iv) The probability distributions are appended to the predictions list.
- (v) The function converts the user-provided weights (weights) to a numpy array.
- (vi) The weights are then applied to the predictions using element-wise multiplication (\*). The weights are broadcasted across the dimensions of the predictions using [ :, np.newaxis, np.newaxis ].
- (vii) The final weighted predictions are obtained by summing the weighted predictions along the first axis (axis=0).
- (viii) The final predictions are obtained by selecting the class index with the highest probability from the weighted predictions.
- (ix) The final set of predictions is returned.

The ensemble\_weighted\_voting function takes a list of machine learning models, a dataset, and a list of weights. It obtains the probability distributions predicted by each model, applies user-specified weights to the predictions, and produces a final set of predictions based on weighted ensemble voting. This allows users to assign different levels of importance to the predictions of each model in the ensemble.

### 11.9 Table (After Optimization) :-

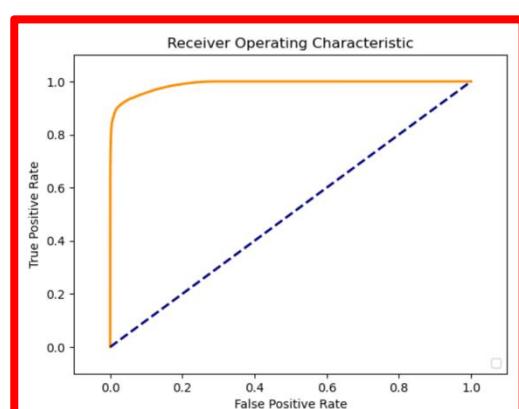
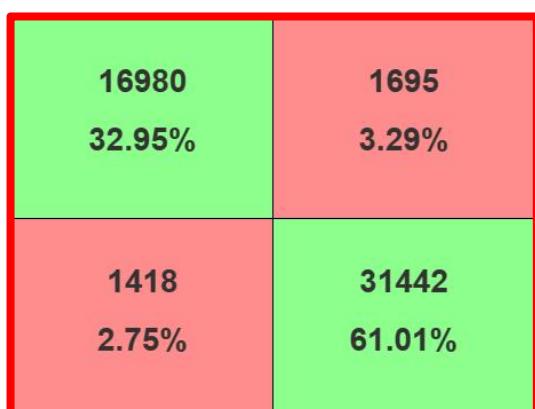
**Table-07 :: Ensemble Models Voting Criteria for UNSW-NB15(after optimization)**

Voting Criteria	Models Used	AUC	CA	F1	Precision	Recall	MCC	Spec	LogLoss
Hard	Random Forest, Decision Tree, Extremely Randomized Tree, AdaBoost, Gradient Boosting	0.9330	93.96	0.9160	0.9092	0.9229	0.8689	0.9488	2.0978
Soft		0.9299	93.70	0.9132	0.9045	0.9203	0.8632	0.9462	2.0978
Bayesian Average		0.9344	94.07	0.9176	0.9114	0.9239	0.8714	0.9500	0.4146
Dynamic		0.9344	94.07	0.9176	0.9114	0.9239	0.8714	0.9500	0.4146
Weighted Average		<b>0.9357</b>	<b>94.13</b>	<b>0.9187</b>	<b>0.9155</b>	<b>0.9219</b>	<b>0.8727</b>	<b>0.9522</b>	<b>0.1697</b>

### 11.10 End Results (After Optimization):-

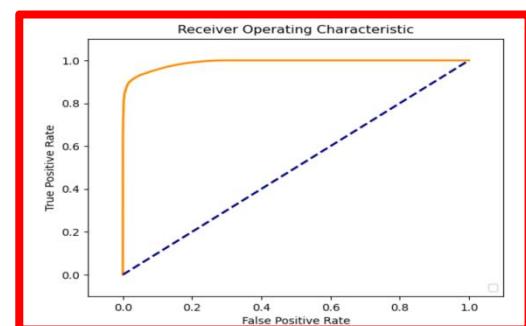
#### ➤ Confusion Matrix and ROC Curve :-

##### (a) Hard Voting :-

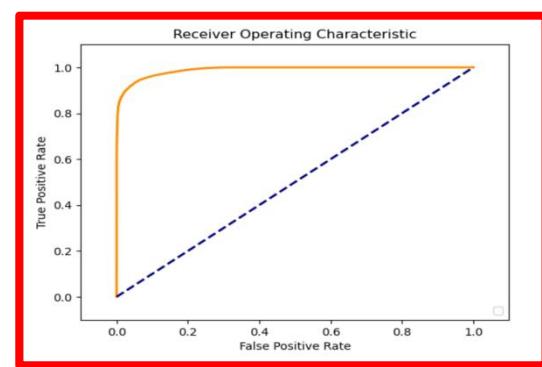


(b) Soft Voting :-

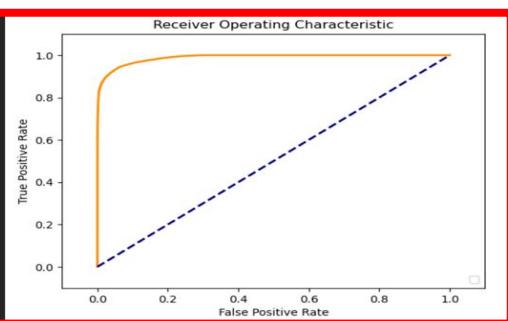
<b>16891</b> 32.78%	<b>1784</b> 3.46%
<b>1463</b> 2.84%	<b>31397</b> 60.92%

(c) Bayesian Model Average :-

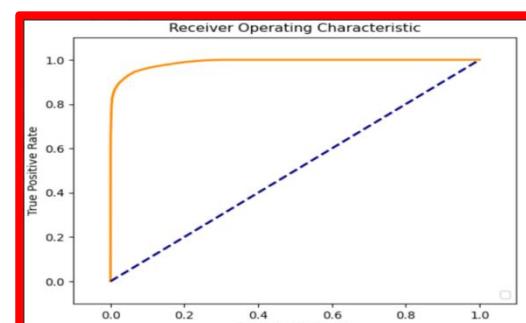
<b>17021</b> 33.03%	<b>1654</b> 3.21%
<b>1401</b> 2.72%	<b>31459</b> 61.04%

(d) Dynamic Method :-

<b>17021</b> 33.03%	<b>1654</b> 3.21%
<b>1401</b> 2.72%	<b>31459</b> 61.04%

(e) Weighted Average :-

<b>15959</b> 30.97%	<b>2716</b> 5.27%
<b>1885</b> 3.66%	<b>30975</b> 60.10%



# Chapter 12

## Final Result and Conclusion

In the pursuit of fortifying network security against the relentless tide of cyber threats, this research project has undertaken a comprehensive exploration of hybrid intrusion detection systems. The imperative to safeguard networked environments has been addressed through the evaluation of diverse machine learning models—AdaBoost, decision tree, random forest, gradient boosting, extremely randomized trees, and neural networks—leveraging datasets (CIC-IDS2017, NSL-KDD, SDN, UNSW-NB15) that mirror real-world network challenges.

Table-08 :: Hybrid Model										
Dataset	Models Used	Voting	AUC	CA	F1	Precision	Recall	MCC	Spec	LogLoss
NSL-KDD19	RF, DT, ET, AdaB, GB	Hard	0.9949	99.50	0.9952	0.9961	0.9943	0.9900	0.9958	0.0171
CICIDS-17	NN, DT, ET, AdaB, GB	Hard	0.9991	99.96	0.9998	1.0000	0.9996	0.9989	1.0000	0.4588
UNSW-NB15	RF, DT, ET, AdaB, GB	Weighted Average	0.9357	94.13	0.9187	0.9155	0.9219	0.8727	0.9522	0.1697
SDN	RF, DT, ET, AdaB, GB	Hard	0.9999	99.99	0.9999	0.9998	1.0000	0.9998	0.9998	0.0020

The meticulous evaluation of the hybrid model showcased promising advancements in intrusion detection capabilities. The ensemble model demonstrated improved performance metrics, surpassing the individual models in terms of accuracy, precision, recall, and F1 score. This achievement substantiates the efficacy of the hybrid approach in enhancing the overall detection accuracy and reducing false positives.

The selective selection procedure identified the five most effective models, forming the foundation for the hybrid ensemble. The integration of these models exhibited a synergistic effect, leveraging the unique strengths of each algorithm while mitigating their respective drawbacks. This not only led to a more robust defense against known threats but also showcased adaptability to novel intrusion techniques.

# Chapter 13

## Future Scope

### 13.1 Optimization for Computational Efficiency :-

Future research can focus on optimizing the computational efficiency of the hybrid ensemble model, especially when utilizing resource-intensive machine learning algorithms like neural networks. This optimization could facilitate real-time intrusion detection without compromising performance.

### 13.2 Scalability in Big Data Analytics :-

Evaluate the scalability of the hybrid model in the context of big data analytics. Ensure that the model can efficiently scale to handle massive volumes of network data, making it applicable and effective in large-scale and high-velocity data environments.

### 13.3 Integration with Deep Learning Architectures :-

Explore the integration of deep learning architectures within the hybrid ensemble model. Investigate the potential benefits of incorporating deep neural networks to capture intricate patterns in network traffic data, enhancing the model's ability to detect sophisticated and evolving intrusion techniques.

### 13.4 Real-Time Threat Intelligence Integration :-

Explore methods to integrate real-time threat intelligence feeds into the hybrid model. This can enhance the system's ability to adapt to the latest threats by incorporating up-to-the-minute information about known malicious entities and attack patterns.

Addressing these future research directions would not only mitigate the identified drawbacks but also position the hybrid intrusion detection system as a more potent tool in the field of network security, with applications extending into the realms of big data analytics and deep learning.

# Chapter 14

## References

- [1] HIDM : A Hybrid Intrusion Detection Model for Cloud Based Systems by Lalit Kumar Vashishtha, Akhil Pratap Singh and Kakali Chatterjee, September 2022
- [2] A Detail Analysis on Intrusion Detection Datasets by Dr Santosh Kumar Sahu, Sanjay Kumar Jena and Sauravranjan Sarangi, February 2014
- [3] Hybrid Intrusion Detection System using Machine Learning by Amar Meryem and Bouabid EL Ouahidi, May 2020
- [4] A Hybrid Intrusion Detection System based on ABC-AFS Algorithm for Misuse and Anomaly Detection by Vajiheh Hajisalem and Shahram Babaie, February 2018
- [5] Improving the Performance of Machine Learning-Based Network Intrusion Detection Systems on the UNSW-NB15 Dataset by Soulaiman Moualla, Khaldoun Khorzom and Assef Jafar, June 2021
- [6] A Hybrid Intrusion Detection System Based on Feature Selection and Weighted Stacking Classifier by RUIZHE ZHAO, YINGXUE MU, LONG ZOU AND XIUMEI WEN, June 2022
- [7] Hybrid Intrusion Detection System Using Machine Learning Techniques in Cloud Computing Environments by Ibraheem Aljamal, Ali Tekeoglu, Korkut Bekiroglu and Saumendra Sengupta, May 2019
- [8] UNSW-NB15: A Comprehensive Data set for Network Intrusion Detection Systems by Nour Moustafa and Jill Slay, March 2015
- [9] A Hybrid Machine Learning Method for Increasing the Performance of Network Intrusion Detection Systems by Achmad Akbar, Megantara and Tohari Ahmad, December 2021
- [10] Efficient Hybrid Model for Intrusion Detection Systems by Nesrine Kaaniche, Aymen Boudguiga and Gustavo Gonzalez-Granadillo, July 2022