

A project report on

COMPREHENSIVE APPROACH OF STATIC AND DYNAMIC DATA ANALYTICS USING AUTOML

Submitted in partial fulfillment for the award of the degree of

M.Tech (Software Engineering)

by

VIVEK R (19MIS0184)



VIT[®]

Vellore Institute of Technology

(Deemed to be University under section 3 of UGC Act, 1956)

SCHOOL OF COMPUTER SCIENCE ENGINEERING AND INFORMATION SYSTEMS

April, 2024

COMPREHENSIVE APPROACH OF STATIC AND DYNAMIC DATA ANALYTICS USING AUTOML

Submitted in partial fulfillment for the award of the degree of

M.Tech (Software Engineering)

by

VIVEK R (19MIS0184)



VIT[®]

Vellore Institute of Technology
(Deemed to be University under section 3 of UGC Act, 1956)

**SCHOOL OF COMPUTER SCIENCE ENGINEERING AND
INFORMATION SYSTEMS**

April, 2024

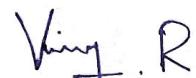
DECLARATION

I here by declare that the thesis entitled “COMPREHENSIVE APPROACH OF STATIC AND DYNAMIC DATA ANALYTICS USING AUTOML” submitted by me, for the award of the degree of M.Tech (Software Engineering) is a record of bonafide work carried out by me under the supervision of PROF. CHADRASEGAR T.

I further declare that the work reported in this thesis has not been submitted and will not be submitted, either in part or in full, for the award of any other degree or diploma in this institute or any other institute or university.

Place: Vellore

Date: 25 - 04 - 2024



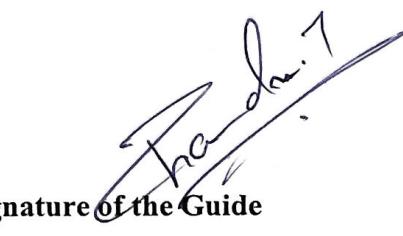
Signature of the Candidate

CERTIFICATE

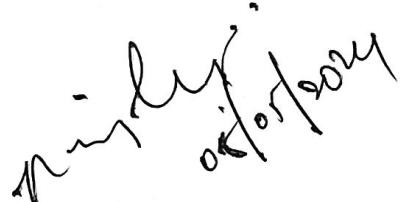
This is to certify that the thesis entitled "COMPREHENSIVE APPROACH OF STATIC AND DYNAMIC DATA ANALYTICS USING AUTOML" submitted by VIVEK R (19MIS0184), School of Computer Science Engineering And Information Systems, Vellore Institute of Technology, Vellore for the award of the degree M.Tech (Software Engineering) is a record of bonafide work carried out by him under my supervision.

The contents of this report have not been submitted and will not be submitted either in part or in full, for the award of any other degree or diploma in this institute or any other institute or university. The Project report fulfils the requirements and regulations of VELLORE INSTITUTE OF TECHNOLOGY, VELLORE and in my opinion meets the necessary standards for submission.

Signature of the Guide



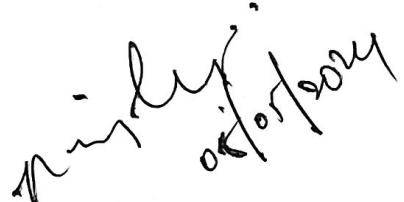
Signature of the HoD



Internal Examiner


Dr. D. S. Rajput

External Examiner



ABSTRACT

Comprehensive approach of Static and Dynamic Data Analytics using AutoML includes Intrusion Detection Evaluation Dataset (CICIDS2017) consists of normal network traffic and simulated abnormal data, caused by deliberate attacks on a test network. The success of Intrusion Detection Systems and Intrusion Prevention Systems, which are vital for detecting and mitigating computer network attacks, hinges on having current and relevant training data. This often figures out whether a network remains secure or becomes compromised. With the proliferation of sensors and smart devices, data generation speed in Internet of Things (IoT) systems has surged. These systems regularly process, transform, and analyse large volumes of data to enable various IoT services and functionalities. Machine Learning (ML) approaches have proven their effectiveness in IoT data analytics. However, applying ML models to such tasks is still challenging, particularly in terms of model selection, design/tuning, and updating. This has created a significant demand for experienced data scientists. Additionally, the dynamic nature of IoT data may lead to concept drift issues, negatively affecting model performance. Alleviate these issues, Automated Machine Learning (AutoML) has appeared as a field aimed at automatically selecting, constructing, tuning, and updating ML models to optimize performance on specific tasks. This project reviews existing methods for model selection, tuning, and updating in AutoML. The aim is to find and summarize optimal solutions for applying AutoML Techniques to the Intrusion detection evaluation dataset (CIC-IDS2017) and IoT Network Intrusion Dataset.

Keywords:

Intrusion Detection system, Network traffic, IoT data analytics, Machine Learning, AutoML

ACKNOWLEDGEMENT

It is my pleasure to express with deep sense of gratitude to PROF. CHANDRASEGAR T, Assistant Professor Sr. Grade 2, School of Computer Science Engineering and Information Systems, Vellore Institute of Technology, for his constant guidance, continual encouragement, understanding; more than all, he taught me patience in my endeavor. My association with him is not confined to academics only, but it is a great opportunity on my part of work with an intellectual and expert in the field of Information Security.

I would like to express my gratitude to DR.G.VISWANATHAN, Chancellor VELLORE INSTITUTE OF TECHNOLOGY, VELLORE, MR. SANKAR VISWANATHAN, DR. SEKAR VISWANATHAN, DR.G V SELVAM, Vice – Presidents VELLORE INSTITUTE OF TECHNOLOGY, VELLORE, Dr. V. S. Kanchana Bhaaskaran, I/c Vice – Chancellor, DR. Partha Sharathi Mallick, Pro-Vice Chancellor and Dr. S. Sumathy, Dean, School of Computer Science Engineering And Information Systems,, for providing with an environment to work in and for his inspiration during the tenure of the course.

In jubilant mood I express ingeniously my whole-hearted thanks to Dr. Neelu Khare, HoD/Professor, all teaching staff and members working as limbs of our university for their not-self-centered enthusiasm coupled with timely encouragements showered on me with zeal, which prompted the acquirement of the requisite knowledge to finalize my course study successfully. I would like to thank my parents for their support.

It is indeed a pleasure to thank my friends who persuaded and encouraged me to take up and complete this task. At last but not least, I express my gratitude and appreciation to all those who have helped me directly or indirectly toward the successful completion of this project.

Place: Vellore

Date: 25 -04 - 2024

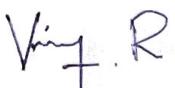

Name of the student

TABLE OF CONTENTS

LIST OF FIGURES.....	v
LIST OF TABLES.....	vi
LIST OF ACRONYMS.....	vii

CHAPTER 1

INTRODUCTION

1.1 BACKGROUND.....	01
1.2 MOTIVATION.....	01
1.3 PROJECT STATEMENT.....	01
1.4 OBJECTIVES.....	02
1.5 SCOPE OF THE PROJECT.....	03

CHAPTER 2

LITERATURE SURVEY

2.1 SUMMARY OF THE EXISTING WORKS.....	04
2.2 CHALLENGES PRESENT IN EXISTING SYSTEM.....	07

CHAPTER 3

REQUIREMENTS

3.1 HARDWARE REQUIREMENTS.....	08
3.2 SOFTWARE REQUIREMENTS.....	08
3.3 BUDGET.....	08
3.4 GANTT CHART.....	09

TABLE OF CONTENTS

CHAPTER 4

ANALYSIS & DESIGN

4.1	PROPOSED METHODOLOGY.....	11
4.2	SYSTEM ARCHITECTURE.....	12
4.3	COMPLETE DESIGN	13
4.4	MODULE DESCRIPTIONS.....	14

CHAPTER 5

IMPLEMENTATION & TESTING

5.1	DATASET	18
5.2	SAMPLE CODE.....	20
5.3	SAMPLE OUTPUT.....	50
5.4	TEST PLAN	54
5.5	DATA VERIFICATION.....	54

CHAPTER 6

RESULTS

6.1	RESEARCH FINDINGS.....	55
6.2	RESULT ANALYSIS	55
6.3	EVALUATION METRICS.....	56

CONCLUSIONS AND FUTURE WORK		58
------------------------------------	--	----

REFERENCES		59
-------------------	--	----

LIST OF FIGURES

3.1 – Project timeline diagram.....	09
3.2 - Project timeline table.....	10
4.1 – System architecture diagram.....	12
4.2 – Flow of Static dataset	13
4.3 – Flow of Dynamic dataset.....	13
5.1 – Screenshot of Static Data	18
5.1 – Screenshot of Dynamic Data	19
5.3 – Bar chart of Static dataset	21
5.4 – Heatmap of Static dataset	22
5.5 – Optimized Heatmap of Static dataset	23
5.6 – Bar chart of Dynamic dataset	35
5.7 – Heatmap of Dynamic dataset	36
5.8 – Optimized Heatmap of Dynamic dataset	37
5.9 – Grid search pipeline	45
5.10 – Screenshot of Deepnote Environment of Static Data Analytics	50
5.11 – Screenshot of Static Data Analytics User Interface	51
5.12 – Screenshot of Deepnote Environment of Dynamic Data Analytics	52
5.13 – Screenshot of Dynamic Data Analytics User Interface.....	53

LIST OF TABLES

2.1 – Summary of the Existing Works.....	04
3.1 – Project Estimate budget	08
3.2 – Activity Description	09
5.1 – Data verification table.....	54
6.1 – Accuracy of Static Data.....	56
6.2 – Accuracy Dynamic Data	56
6.3 – Precision of Static Data	56
6.4 – Precision of Dynamic Data	56
6.5 – Recall value of static data	57
6.6 – Recall value of Dynamic data.....	57
6.7 – F1-Score value of Static data	57
6.8 – F1-Score value of Dynamic data	57

LIST OF ACRONYMS

Acronym	Full Form
AutoML	Automated Machine Learning
IoT	Internet of Things
CICIDS	Canadian Institute for Cybersecurity Intrusion Detection Systems
ML	Machine Learning
HPO	Hyperparameter Optimization
CASH	Combined Algorithm Selection and Hyperparameter tuning
AMC	AutoML for Model Compression and Acceleration on Mobile Devices
OAML	Online AutoML
GPT	Generative Pre-trained Transformer
SMOTE	Synthetic Minority Over-sampling Technique
PSO	Particle Swarm Optimization
ANN	Artificial Neural Network
TP	True Positive
TN	True Negative
FP	False Positive
FN	False Negative
Acc	Accuracy
Prec	Precision
Rec	Recall
F1	F1 Score
NB	Naive Bayes
KNN	k-Nearest Neighbors
RF	Random Forest
LGBM	Light Gradient Boosting Machine
GPT	Generative Pre-trained Transformer
BERT	Bidirectional Encoder Representations from Transformers
AE	Autoencoder

Chapter 1

Introduction

1.1 BACKGROUND

The CICIDS2017 dataset aims to benchmark intrusion detection systems but faces issues like class imbalance, redundant data, and lack of preprocessing guidelines. Simultaneously, IoT data analytics deals with massive, heterogeneous data streams susceptible to concept drift, impacting machine learning model performance over time.

1.2 MOTIVATION

This project employs Automated Machine Learning (AutoML) to address challenges in the CICIDS2017 dataset and IoT data analytics. AutoML automates data preprocessing, feature engineering, model selection, hyperparameter tuning, and model updating. For CICIDS2017, AutoML can mitigate imbalance, redundancy, and lack of guidelines, enabling robust intrusion detection models. For IoT analytics, AutoML can detect and adapt to concept drift, maintaining model accuracy. An end-to-end AutoML pipeline can streamline model development for both static intrusion detection and dynamic IoT analytics tasks, reducing human effort and promoting scalability.

1.3 PROJECT STATEMENT

The CICIDS2017 intrusion detection dataset is a valuable public benchmark for developing network security systems using machine learning. However, it lacks rigorous profiling and optimization for building robust models. Issues such as class imbalance, redundant records, and missing guidelines limit its effectiveness. Analysing massive and dynamic IoT data using machine learning also has its challenges. The large volume and speed of heterogeneous IoT data streams make

manual model management difficult. Changes in data or concept drift degrades model performance over time. There is a need for Automated Machine Learning (AutoML) solutions that can reduce human effort in developing analytic models. For intrusion detection, this involves evaluating datasets thoroughly, correcting deficiencies, and establishing best practices. For IoT data, this involves automating steps such as data preprocessing, algorithm selection, hyperparameter optimization, and model updating. By improving benchmark datasets and demonstrating end-to-end AutoML capabilities, we can rapidly develop more robust intrusion detection systems and optimized IoT data analytics. This will better equip us to handle evolving security threats and streaming heterogeneous data.

1.4 OBJECTIVES

- Developing automated solutions for data preprocessing, algorithm selection, and hyperparameter optimization using AutoML techniques.
- Addressing class imbalance, redundant records, and missing guidelines in the CICIDS2017 dataset to improve its effectiveness as a benchmark for intrusion detection systems.
- Investigating techniques for handling concept drift in IoT data analytics and develop automated model updating procedures to maintain model performance over time.
- Implementing an end-to-end AutoML pipeline for both static intrusion detection and dynamic IoT data analytics tasks.
- Comparing the performance of AutoML models with traditional machine learning approaches to assess the benefits of automation in model development.
- Demonstrating the feasibility and effectiveness of AutoML in improving the efficiency and accuracy of data analytics tasks in both static and dynamic environments.
- Implementing a frontend integration to visualize and interact with the AutoML pipeline and models developed for intrusion detection and IoT data analytics.

1.5 SCOPE OF THE PROJECT

The project aims to enhance intrusion detection systems and IoT data analytics through Automated Machine Learning (AutoML). It involves evaluating the CICIDS2017 intrusion detection dataset, automating data preprocessing, feature engineering, algorithm selection, and hyperparameter optimization. The project will also address concept drift in IoT data and implement an end-to-end AutoML pipeline for both static and dynamic data. Performance comparison with traditional approaches will be conducted, and recommendations for deploying AutoML in real-world applications will be provided. Ultimately, the project seeks to demonstrate the feasibility and effectiveness of AutoML in improving data analytics tasks in diverse environments.

1. AutoML for Algorithm Selection & Hyperparameter Optimization:

Investigating automated methods to select ML algorithms and optimize hyperparameters for intrusion detection tasks.

2. Handling Concept Drift in IoT Analytics:

Developing automated procedures to detect and adapt to changes in data distribution over time in IoT analytics.

3. End-to-End AutoML Pipeline:

Implementing a comprehensive AutoML pipeline covering data preprocessing to model deployment for both static intrusion detection and dynamic IoT analytics.

Chapter 2

LITERATURE SURVEY

2.1 SUMMARY OF THE EXISTING WORKS

S. No	Title of the Paper	Merits	Demerits
1.	IoT data analytics in dynamic environments: From an automated machine learning perspective	<ul style="list-style-type: none">▪ Comprehensive review of AutoML for IoT analytics▪ Analyzes range of methods for distributed data, concept drift, edge deployment▪ Structured taxonomy of existing techniques	<ul style="list-style-type: none">▪ Rapidly evolving field means some recent advances not covered▪ Does not include experimental comparisons
2.	AutoML-ID: automated machine learning model for intrusion detection using wireless sensor network	<ul style="list-style-type: none">▪ Focus on important problem of IoT intrusion detection▪ Reviews specific AutoML techniques applied in this domain▪ Discusses open challenges and future research directions	<ul style="list-style-type: none">▪ Scope limited to intrusion detection▪ Does not include quantitative evaluations▪ Light on technical details of AutoML implementations
3.	Methods for network intrusion detection: Evaluating rule-based methods and machine learning models on the	<ul style="list-style-type: none">▪ Leverages the commonly used CIC-IDS2017 benchmark dataset reflecting modern network attacks.▪ Examines a diverse set of models including neural networks and ensemble methods in addition to more basic classifiers.	<ul style="list-style-type: none">▪ Only certain standard classification algorithms are evaluated rather than more recent advanced methods.▪ The study is limited to binary classification of network flows; multi-class detection is not examined.

	CIC-IDS2017 dataset	<ul style="list-style-type: none"> ▪ Analysis of performance on different attack categories offers useful insights for selection of appropriate techniques. 	<ul style="list-style-type: none"> ▪ Lacks analysis of the efficiency and computational requirements of the different approaches.
4.	AMLBID: An auto-explained Automated Machine Learning tool for Big Industrial Data	<ul style="list-style-type: none"> ▪ Innovative AutoML tool for Big Industrial Data. ▪ Meta-learning-based recommendation system for efficient ML algorithm selection. ▪ Interactive visualization module enhances user understanding. 	<ul style="list-style-type: none"> ▪ Limited empirical validation on real-world datasets. ▪ Dependency on meta-features may limit effectiveness. ▪ Scalability concerns for handling large-scale datasets.
5.	AMC: AutoML for Model Compression and Acceleration on Mobile Devices	<ul style="list-style-type: none"> ▪ Innovation in leveraging reinforcement learning for automated model compression. ▪ Support for both accuracy-guaranteed and resource-constrained compression. ▪ Significant improvements in inference speed demonstrated on mobile devices. 	<ul style="list-style-type: none"> ▪ Lack of discussion on potential limitations or failure modes of the proposed approach. ▪ Limited comparison with other state-of-the-art techniques in model compression and acceleration.
6.	AutoML: A Survey of the State-of-the-Art	<ul style="list-style-type: none"> ▪ Provides an extensive exploration of Neural Architecture Search (NAS), a key area in AutoML, with performance comparisons on CIFAR-10 and ImageNet datasets. ▪ Encompasses a wide array of AutoML topics, including automated data augmentation, hyperparameter optimization, and evolutionary algorithms. 	<ul style="list-style-type: none"> ▪ Lacks coverage of recent AutoML advancements, necessitating an update to incorporate newer methods ▪ NAS receives disproportionate attention. ▪ Lacks critical analysis and comparison.

7.	Dynamic Hyperparameter Allocation under Time Constraints for Automated Machine Learning	<ul style="list-style-type: none"> ▪ Paper introduces a novel diversification strategy for HPO, emphasizing dynamic hyperparameter space allocation for a sampler based on the remaining time budget. ▪ Proposed solution claims better performance in terms of both time and resource awareness compared to previous resource-aware solutions 	<ul style="list-style-type: none"> ▪ Effectiveness of the proposed strategy may depend on the specific characteristics of the datasets and models used in the benchmarks. The generalizability of the approach to diverse scenarios should be considered. ▪
8.	AutoML for Multi-Label Classification: Overview and Empirical Evaluation	<ul style="list-style-type: none"> ▪ The paper conducts an extensive experimental study evaluating multiple optimization methods on a suite of MLC problems, providing a thorough analysis of each approach's performance ▪ By extending existing AutoML approaches to tackle multi-label classification problems, the paper contributes valuable insights into the applicability and effectiveness of these techniques in complex learning scenarios. 	<ul style="list-style-type: none"> ▪ primarily focuses on theoretical comparisons and the performance of optimization methods in controlled experiments. Real-world application scenarios or case studies could provide more practical insights into the usability of these approaches. ▪ While the paper outlines potential research directions to improve AutoML for MLC, it could further elaborate on practical implementations, challenges
9.	Online AutoML: an adaptive AutoML framework for online learning	<ul style="list-style-type: none"> ▪ OAML is the first system to propose a flexible and practical AutoML system for adaptive online learning pipelines, combining algorithm selection, hyperparameter optimization, and preprocessing. 	<ul style="list-style-type: none"> ▪ The current implementation focuses on classification tasks and might require adjustments for other supervised learning problems

		<ul style="list-style-type: none"> ▪ OAML leverages drift detection to trigger pipeline redesign, adapting to changing data distributions ▪ OAML offers different adaptation strategies (ensemble, model store) and supports various online learning algorithms and preprocessors 	<ul style="list-style-type: none"> ▪ The paper acknowledges the lack of research on hyperparameter optimization in online settings and uses manually defined search intervals. Exploring more advanced techniques could further improve performance.
10.	AutoML-GPT: Automatic Machine Learning with GPT	<ul style="list-style-type: none"> ▪ Proposed AutoML-GPT framework offers a systematic approach to automating the training pipeline for various AI tasks by leveraging the language capabilities of GPT and dynamically generating task-oriented prompts ▪ This can potentially reduce the manual effort required in model selection and hyperparameter tuning 	<ul style="list-style-type: none"> ▪ Lack of external validation on real-world tasks ▪ Scalability and resource requirements

Table 2.1 – Summary of the Existing Works

2.2 CHALLENGES PRESENT IN EXISTING SYSTEM

- Limited generalizability of AutoML solutions across diverse domains and tasks.
- Inadequate handling of concept drift in dynamic data environments, leading to model performance degradation over time
- Scalability and efficiency concerns with large-scale datasets & complex models.
- Lack of interpretability and explainability of AutoML black-box systems, hindering trust and adoption.
- Insufficient empirical validation and comparative studies against traditional approaches and state-of-the-art AutoML techniques

Addressing these challenges is crucial for developing robust, scalable, and interpretable AutoML solutions that can effectively handle both static and dynamic data analytics tasks in real-world applications.

CHAPTER 3

REQUIREMENTS

3.1 HARDWARE REQUIREMENTS

Minimum Hardware Requirement:	Maximum Hardware Requirement:
<ul style="list-style-type: none">• 4 CPU Core• 4 GB RAM• 10 GB Disk Space• Internet speed - 5 Mbps (download and upload) for project execution.	<ul style="list-style-type: none">• 8 CPU Cores• 32 GB RAM• 20 GB Disk Space• Internet speed-10 Mbps (download & upload) more resource-intensive task

3.2 SOFTWARE REQUIREMENTS

Minimum Software Requirements:

- Python 3.6 or higher
- Jupyter Notebook or JupyterLab
- Libraires: Pandas, NumPy, Scikit-learn, TensorFlow, Keras
- AutoML Libraries (e.g., Auto-Sklearn, AutoKeras)
- Data Visualization Libraries (e.g., Matplotlib, Seaborn)

Maximum Software Requirements:

- Python 3.8 or higher
- Jupyter Notebook or JupyterLab
- Libraires: Pandas, NumPy, Scikit-learn, TensorFlow, Keras, PyTorch
- AutoML Libraries (e.g., Auto-Sklearn, AutoKeras, Auto-PyTorch,)
- Data Visualization Libraries (e.g., Matplotlib, Seaborn, Plotly)
- Imbalanced Data Libraries (e.g., Imbalanced-learn, Smotetomek)

3.3 BUDGET

Budget Categories	Costs
Software:	DeepNote Professional Plan: 2500/month
Dataset:	<ul style="list-style-type: none">- CICIDS2017 Intrusion Detection Dataset: Free (publicly available)- IoT Network Intrusion Dataset: Free (publicly available)
Cloud Computing:	DeepNote Professional Plan, included
Miscellaneous Expenses:	<ul style="list-style-type: none">• Office Supplies: 2000• Printing and Binding: 200
Total Estimate	32,200/-

Table 3.1 – Project Estimate budget

3.4 GANTT CHART



Figure 3.1 – Project timeline diagram

Activity	Description of the Activity	Guide Remarks
1	Defined problem statement, project scope, objectives	Title approved
2	Research on AutoML Approach, gather requirements with Hyperparameter study	ok
3	Implemented traditional Machine Learning Algorithms for both static and dynamic dataset	ok
4	Comparison with existing tools to implement project	ok
5	Implement the AutoML Technique	ok
6	Final evaluation and reporting	ok
7	Deployment and documentation	ok
8	Project conclusion and presentation	ok

Table 3.2 – Activity Description



Static & Dynamic Data Analytics -AutoML

Read-only view, generated on 23 Apr 2024

ACTIVITIES	START	WD	DUe	PR	STATUS	%
Planning Phase:	03/Jan	4d	08/Jan	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
1 ✓ Project scope and requirements	03/Jan	3d	05/Jan	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
2 ✓ Setup cloud environment, tools, libraries	05/Jan	2d	08/Jan	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
Analysis Phase:	08/Jan	15d	26/Jan	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
4 ✓ Research AutoML Techniques and Methodolo...	08/Jan	6d	15/Jan	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
5 ✓ Identify libraries, suitable techniques for impl...	15/Jan	3d	17/Jan	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
6 ✓ Study Hyperparameter Tuning	22/Jan	5d	26/Jan	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
Data Visualization:	29/Jan	3d	31/Jan	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
8 ✓ Gaining insights from datasets and Leveragin...	29/Jan	3d	31/Jan	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
Data Preprocessing:	01/Feb	12d	16/Feb	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
10 ✓ Data Encoding	01/Feb	3d	05/Feb	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
11 ✓ Data Imputation	06/Feb	2d	07/Feb	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
12 ✓ Data Normalization	08/Feb	3d	12/Feb	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
13 ✓ Train-Test Split	12/Feb	2d	13/Feb	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
14 ✓ Data Balancing	14/Feb	3d	16/Feb	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
ML Model Learning:	19/Feb	25d	22/Mar	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
16 ✓ Implementing LGBM Classifier Algorithm	19/Feb	5d	23/Feb	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
17 ✓ Implementing Random Forest Algorithm	26/Feb	5d	01/Mar	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
18 ✓ Implementing Naive Bayes Algorithm	04/Mar	5d	08/Mar	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
19 ✓ Implementing k-nearest neighbors (KNN) Algo...	11/Mar	5d	15/Mar	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
20 ✓ Implementing KerasClassifier Model	18/Mar	5d	22/Mar	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
Model Selection:	25/Mar	6d	01/Apr	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
22 ✓ Implementing Grid Search Method	25/Mar	6d	01/Apr	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
Combined Algorithm Selection and Hyperparamete...	03/Apr	5d	09/Apr	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
24 ✓ Implementing Particle Swarm Optimization (P...	03/Apr	5d	09/Apr	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
Performance Metrics Phase:	09/Apr	3d	11/Apr	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
26 ✓ Fetching Accuracy, Precision, F1 Score, Recall ...	09/Apr	3d	11/Apr	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
Deployment:	12/Apr	6d	19/Apr	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
28 ✓ Deploy Frontend application for Static Data an...	12/Apr	4d	17/Apr	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%
29 ✓ Deploy Frontend application for Dynamic Dat...	17/Apr	3d	19/Apr	<div style="width: 100%;"><div style="width: 100%;"> </div></div>	Finish	100%

Figure 3.2 - Project timeline table

CHAPTER 4

ANALYSIS & DESIGN

4.1 PROPOSED METHODOLOGY

The proposed methodology involves a comprehensive approach to leveraging Automated Machine Learning (AutoML) for both static intrusion detection and dynamic Internet of Things (IoT) data analytics tasks. It begins with a thorough exploratory data analysis and preprocessing phase, addressing issues such as class imbalance, redundant records, and missing values in the CICIDS2017 intrusion detection dataset. Automated feature engineering and selection techniques will be employed to extract relevant features from the raw data, tailored for the respective tasks. Subsequently, AutoML frameworks will be utilized to automate the process of selecting appropriate machine learning algorithms and optimizing their hyperparameters, with a focus on maximizing performance metrics like accuracy, precision, recall, and F1-score.

To handle the dynamic nature of IoT data streams, techniques for detecting concept drift will be implemented, enabling the monitoring of performance metrics, distribution statistics, or leveraging dedicated drift detection algorithms. Automated model updating procedures will be developed to retrain or adapt the models when concept drift is detected, ensuring their continued accuracy and relevance over time.

The performance of the developed AutoML models will be rigorously evaluated using appropriate evaluation metrics and protocols. A comparative analysis will be conducted, contrasting the AutoML models' performance with traditional machine learning approaches that require manual intervention for model development. Finally, the feasibility of deploying the AutoML models in real-world scenarios will be investigated, considering factors like computational resources, scalability, and integration with existing systems. Recommendations and guidelines for effectively utilizing AutoML techniques in intrusion detection systems and IoT data analytics applications will be provided.

4.2 SYSTEM ARCHITECTURE

Architecture of Static and Dynamic Datasets using AutoML

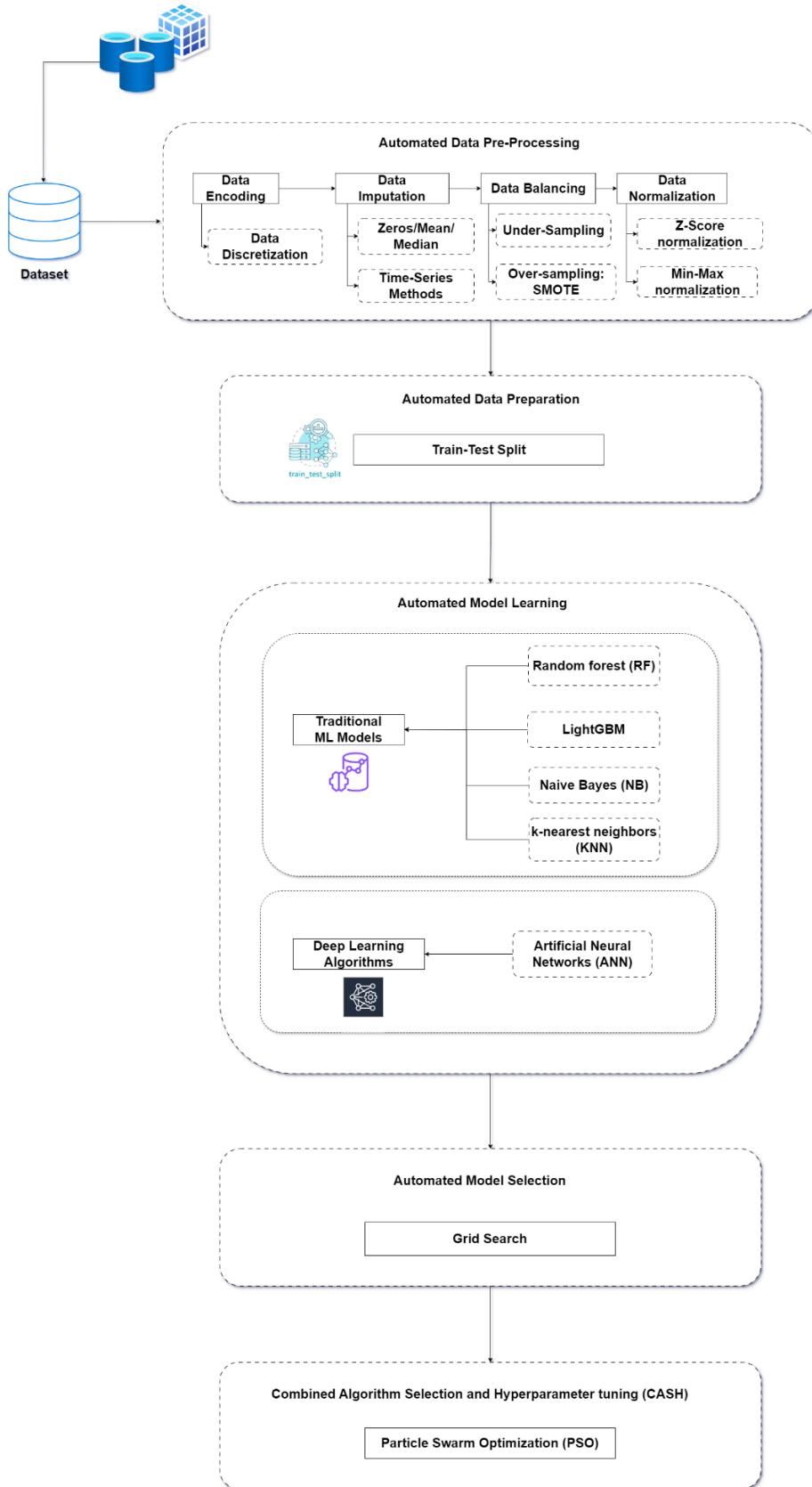


Figure 4.1 – System architecture diagram

4.3 COMPLETE DESIGN:

Static dataset – Intrusion Detection Evaluation Dataset (CICIDS2017)

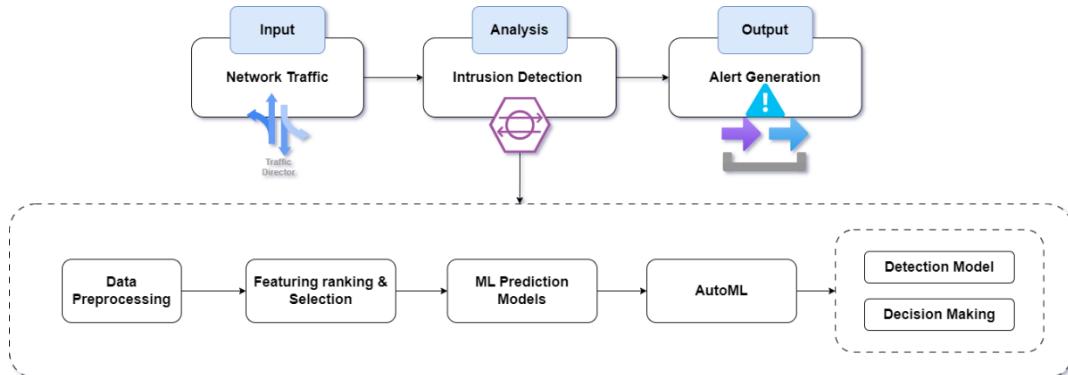


Figure 4.2 – Flow of Static dataset

Dynamic dataset – IoT Data Analytics in Dynamic Environments

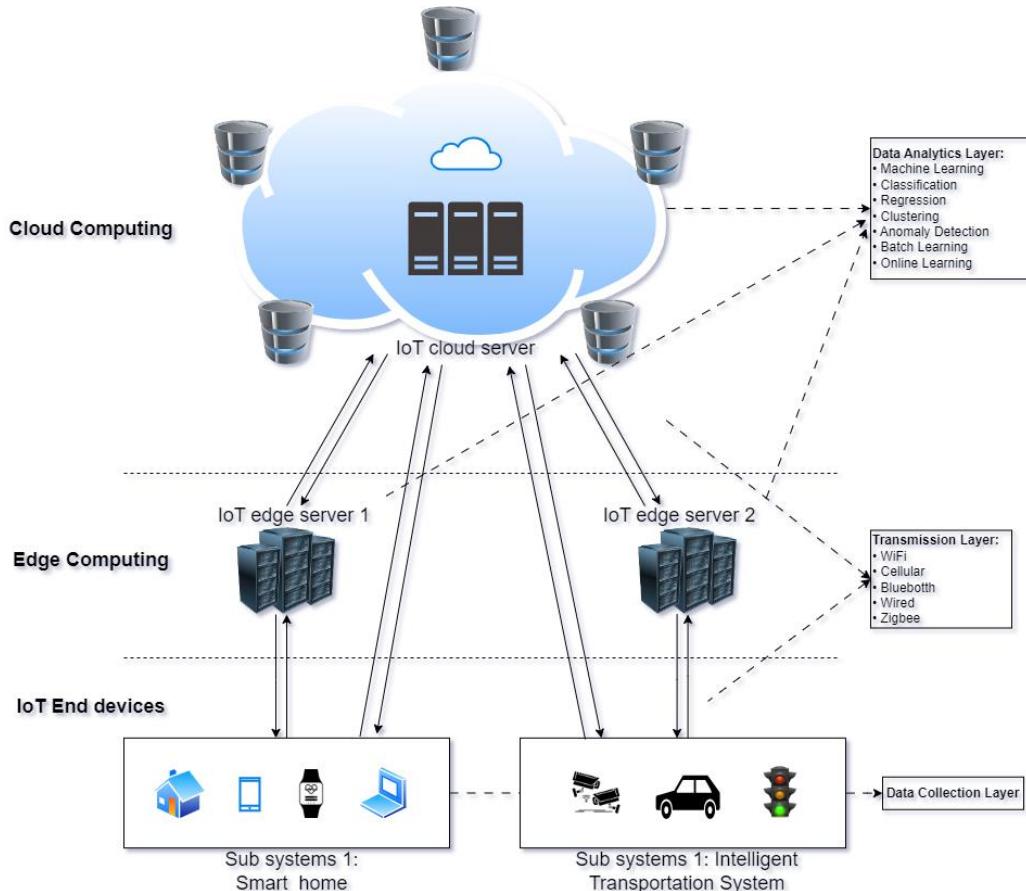


Figure 4.3 – Flow of Dynamic dataset

4.4 MODULE DESCRIPTIONS

1. Data Visualization Module

This module provides visual representations of the data and model performance to aid in data exploration and analysis. It includes the following sub-modules:

A. Data Exploration Visualizations

- Visualizes the distribution of features and target variable using histograms, box plots, and scatter plots. Identifies correlations and relationships between features using correlation heatmaps.

B. Performance Visualization

- Generates bar charts to compare the performance metrics (accuracy, precision, recall, F1-score) of different machine learning models.
- Creates line plots to visualize the training and validation performance (loss, accuracy) during model training.

C. Confusion Matrix Visualization

- Plots confusion matrices for each machine learning model, providing insights into true positives, false positives, true negatives, & false negatives.

2. Automated Data Preprocessing Module

This module handles the initial data preparation tasks in an automated manner. It consists of the following sub-modules:

A. Automated Encoding

- Identifies, transforms string/text features into numerical features, making the data more readable for machine learning models. Functionality: It uses the LabelEncoder from scikit-learn to encode categorical features.

B. Automated Imputation

- Detects and imputes missing values in the dataset to improve data quality. Functionality: It replaces infinite values with NaN and then fills NaN values with zeros. However, it can be modified to use other imputation techniques.

C. Automated Normalization

- Normalizes the range of features to a similar scale, based on the data distribution. Functionality: It uses the Shapiro-Wilk test to determine if the data follows a Gaussian distribution. If so, it applies Z-score normalization; otherwise, it applies min-max normalization.

D. Train-Test Split

- Splits the dataset into training and test sets. Functionality: It uses scikit-learn's `train_test_split` function to create an 80/20 train-test split by default, but this can be modified.

E. Automated Data Balancing

- Generates minority class samples to address class imbalance and improve data quality. Functionality: It uses the Synthetic Minority Over-sampling Technique (SMOTE) from the imbalanced-learn library to oversample the minority class.

3. Automated Model Learning Module

This module performs automated training and evaluation of several machine learning models for comparison purposes.

A. Model Training

- Trains the following machine learning models on the preprocessed training data: Naive Bayes, K-Nearest Neighbors, Random Forest, LightGBM, and Artificial Neural Network (ANN). Functionality: It uses the respective classes from scikit-learn, LightGBM, and Keras libraries to instantiate and fit the models on the training data.

i. Naive Bayes

Naive Bayes is a probabilistic classifier based on applying Bayes' theorem with the assumption of independence between features.

Implementation: The GaussianNB class from scikit-learn is used

ii. K-Nearest Neighbors (KNN)

KNN is a non-parametric algorithm that classifies a data point based on the majority class among its k nearest neighbors in the feature space.

Implementation: The KNeighborsClassifier class from scikit-learn is used to instantiate the KNN classifier.

iii. Random Forest

Random Forest is an ensemble learning method that constructs multiple decision trees and combines their predictions by averaging or majority voting.

Implementation: The RandomForestClassifier class from scikit-learn is used to instantiate the Random Forest classifier.

iv. LightGBM

LightGBM (Light Gradient Boosting Machine) is a gradient boosting framework that uses tree-based learning algorithms to build efficient and high-performance models.

Implementation: The LGBMClassifier class from the LightGBM library is used to instantiate the LightGBM classifier.

v. Artificial Neural Network (ANN)

An Artificial Neural Network (ANN) or Multi-Layer Perceptron (MLP) is a type of feedforward neural network that consists of an input layer, one or more hidden layers, and an output layer, inspired by the biological neural networks in the human brain.

Implementation: A custom function ANN is defined, which creates a sequential model using the Keras library. The model architecture consists of an input layer, two dense hidden layers with ReLU activation and dropout, and an output layer with softmax activation for binary classification.

Training: The KerasClassifier from scikit-learn is used to wrap the custom ANN function. The fit method is called on the KerasClassifier instance, passing the training data (X_train and y_train) to train the ANN model using backpropagation and the specified hyperparameters (e.g., batch size, epochs, early stopping patience).

B. Model Evaluation

- Evaluates the trained models on the test data and reports various performance metrics. Functionality: It uses the trained models to make predictions on the test data and calculates the following metrics using scikit-learn's functions:
 - Accuracy
 - Precision
 - Recall
 - F1-score
 - Time (time taken to make predictions on the test set)

4. Automated Model Selection Module

This module selects the best-performing machine learning model among five common models: Naive Bayes, K-Nearest Neighbors, Random Forest, LightGBM, and Artificial Neural Network (ANN).

Grid Search

- Description: This method performs an exhaustive search over the specified hyperparameter values for each model.
- Functionality: It uses scikit-learn's GridSearchCV to evaluate the performance of each model using 5-fold cross-validation.

5. Combined Algorithm Selection and Hyperparameter Tuning (CASH) Module

This module combines the processes of model selection and hyperparameter optimization into a single step.

Particle Swarm Optimization (PSO)

- Description: This method uses Particle Swarm Optimization (PSO) to simultaneously select the best machine learning model and tune its hyperparameters.
- Functionality: It defines a performance function that trains and evaluates each model with different hyperparameter values on the hold-out test set. The optunity library is used to perform the CASH process using PSO.

CHAPTER 5

IMPLEMENTATION & TESTING

5.1 DATA SET

DATASET 1: - Static Data Analytics

CICIDS2017 dataset, a popular network traffic dataset for intrusion detection problem
Publicly available at: <https://www.unb.ca/cic/datasets/ids-2017.html>

Canadian Institute for Cybersecurity Intrusion Detection System 2017 (CICIDS2017) dataset has the most updated network threats. The CICIDS2017 dataset is close to real-world network data since it has a large amount of network traffic data, a variety of network features, various types of attacks, and highly imbalanced classes.

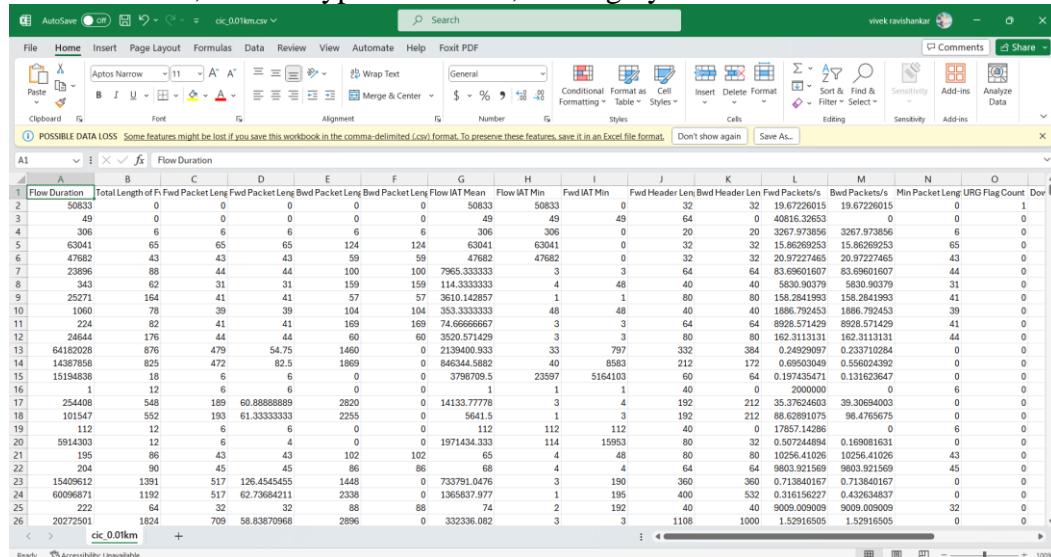


Figure 5.1 - Screenshot of Static Data

#	Column	Non-Null Count	Dtype
0	Flow Duration	28303	non-null
1	Total Length of Fwd Packets	28303	non-null
2	Fwd Packet Length Max	28303	non-null
3	Fwd Packet Length Mean	28303	non-null
4	Bwd Packet Length Max	28303	non-null
5	Bwd Packet Length Min	28303	non-null
6	Flow IAT Mean	28303	non-null
7	Flow IAT Min	28303	non-null
8	Fwd IAT Min	28303	non-null
9	Fwd Header Length	28303	non-null
10	Bwd Header Length	28303	non-null
11	Fwd Packets/s	28303	non-null
12	Bwd Packets/s	28303	non-null
13	Min Packet Length	28303	non-null
14	URG Flag Count	28303	non-null
15	Down/Up Ratio	28303	non-null
16	Init_Win_bytes_forward	28303	non-null
17	Init_Win_bytes_backward	28303	non-null
18	min_seg_size_forward	28303	non-null
19	Labelb	28303	non-null

DATASET 2: - Dynamic Data Analytics

IoTID20 dataset, a novel IoT botnet dataset

Publicly Available: <https://sites.google.com/view/iot-network-intrusion-dataset/home>

IoTID20 dataset was created by using normal and attack virtual machines as network platforms, simulating IoT services with the node-red tool, and extracting features with the Information Security Center of Excellence (ISCX) flow meter program. A typical smart home environment was established for generating this dataset using five IoT devices or services: a smart fridge, a smart thermostat, motion-activated lights, a weather station, and a remotely-activated garage door. Thus, the traffic data samples of normal and abnormal IoT devices are collected in Pcap files.

Figure 5.2 - Screenshot of Dynamic Data

#	Column	Non-Null Count	Dtype
0	Flow_ID	6252	non-null int64
1	Src_IP	6252	non-null int64
2	Src_Port	6252	non-null int64
3	Dst_IP	6252	non-null int64
4	Dst_Port	6252	non-null int64
5	Timestamp	6252	non-null int64
6	Flow_Duration	6252	non-null int64
7	TotLen_Fwd_Pkts	6252	non-null float64
8	TotLen_Bwd_Pkts	6252	non-null float64
9	Bwd_Pkt_Len_Max	6252	non-null float64
10	Flow_Byts/s	6252	non-null float64
11	Flow_Pkts/s	6252	non-null float64
12	Flow_IAT_Mean	6252	non-null float64
13	Flow_IAT_Std	6252	non-null float64
14	Flow_IAT_Max	6252	non-null float64
15	Flow_IAT_Min	6252	non-null float64
16	Fwd_IAT_Tot	6252	non-null float64
17	Bwd_IAT_Tot	6252	non-null float64
18	Bwd_IAT_Mean	6252	non-null float64
19	Bwd_IAT_Std	6252	non-null float64
20	Bwd_IAT_Max	6252	non-null float64
21	Bwd_Header_Len	6252	non-null int64
22	Fwd_Pkts/s	6252	non-null float64
23	Bwd_Pkts/s	6252	non-null float64
24	Pkt_Len_Std	6252	non-null float64
25	FIN_Flag_Cnt	6252	non-null int64
26	Pkt_Size_Avg	6252	non-null float64
27	Init_Bwd_Win_Byts	6252	non-null int64
28	Idle_Mean	6252	non-null float64
29	Idle_Max	6252	non-null float64
30	Idle_Min	6252	non-null float64
31	Label	6252	non-null int64

5.2 SAMPLE CODE

AutoML – 1: Static Dataset - CICIDS2017

1. Data Visualization

```
| df.shape
```

```
(28303, 20)
```

```
| df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 28303 entries, 0 to 28302
Data columns (total 20 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Flow Duration    28303 non-null   int64  
 1   Total Length of Fwd Packets 28303 non-null   int64  
 2   Fwd Packet Length Max      28303 non-null   int64  
 3   Fwd Packet Length Mean    28303 non-null   float64 
 4   Bwd Packet Length Max    28303 non-null   int64  
 5   Bwd Packet Length Min    28303 non-null   int64  
 6   Flow IAT Mean          28303 non-null   float64 
 7   Flow IAT Min           28303 non-null   int64  
 8   Fwd IAT Min           28303 non-null   int64  
 9   Fwd Header Length     28303 non-null   int64  
 10  Bwd Header Length     28303 non-null   int64  
 11  Fwd Packets/s         28303 non-null   float64 
 12  Bwd Packets/s         28303 non-null   float64 
 13  Min Packet Length    28303 non-null   int64  
 14  URG Flag Count        28303 non-null   int64  
 15  Down/Up Ratio         28303 non-null   int64  
 16  Init_Win_bytes_forward 28303 non-null   int64  
 17  Init_Win_bytes_backward 28303 non-null   int64  
 18  min_seg_size_forward  28303 non-null   int64  
 19  Labelb               28303 non-null   int64  
dtypes: float64(4), int64(16)
memory usage: 4.3 MB
```

```

from matplotlib import pyplot as plt
df.hist(figsize=(15,15), xrot=-45, bins=10) ## Display the labels
rotated by 45 degrees

# Clear the text "residue"
plt.show()

```

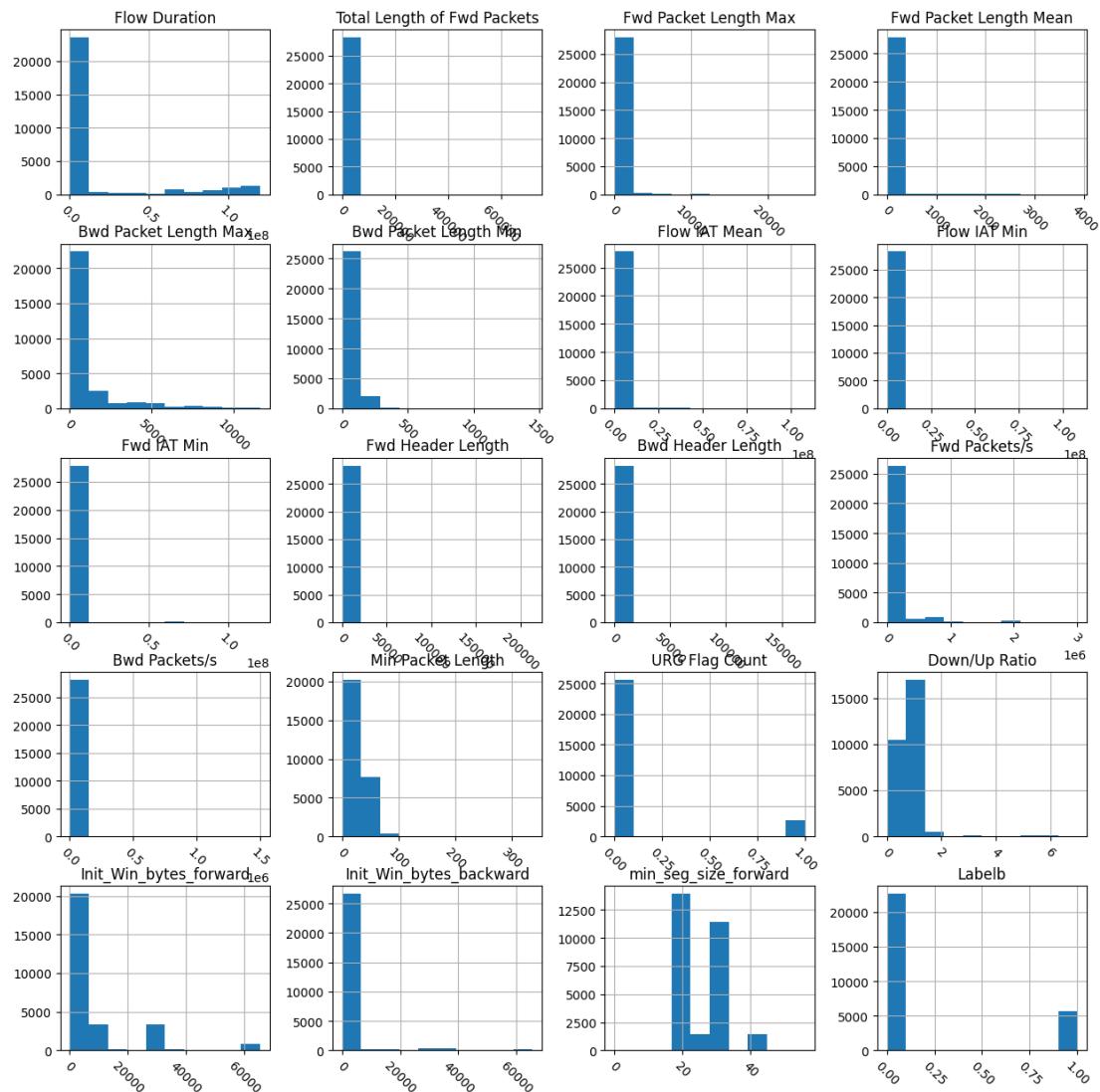


Figure 5.3 – Bar chart of Static dataset

```

plt.figure(figsize=(20,20))
dataplot=sns.heatmap(df.corr(), annot=True)
plt.show()

```

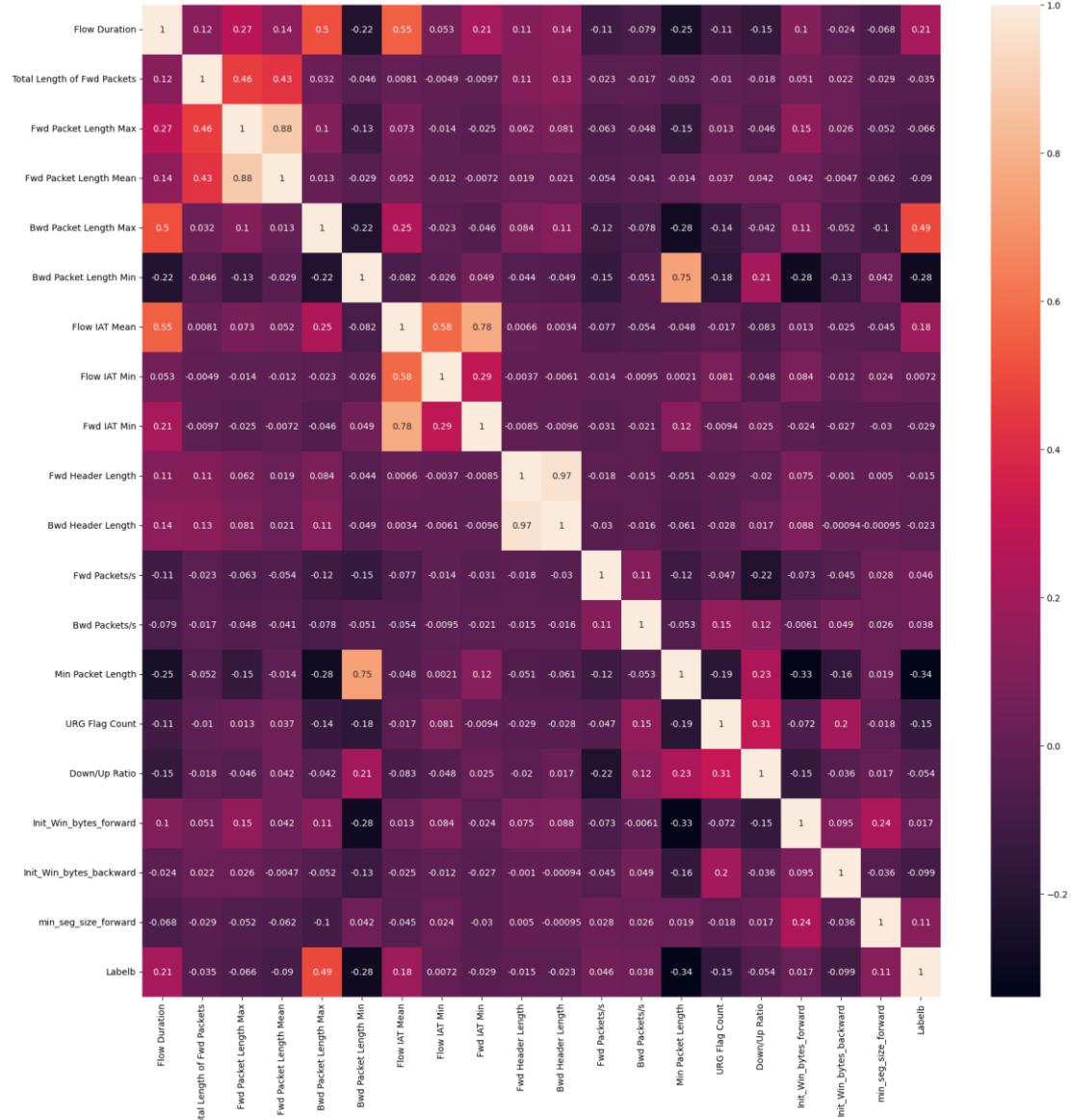


Figure 5.4 – Heatmap of Static dataset

```

from matplotlib.colors import ListedColormap

mask = np.triu(np.ones_like(df.corr(), dtype=bool)) # Upper triangle
mask
plt.figure(figsize=(10,10))
with sns.axes_style("white"):
    ax = sns.heatmap(df.corr()*100, mask=mask, fmt='.0f', annot=True,
lw=1, cmap=ListedColormap(['red', 'yellow', 'green', 'blue']))
plt.show()

```

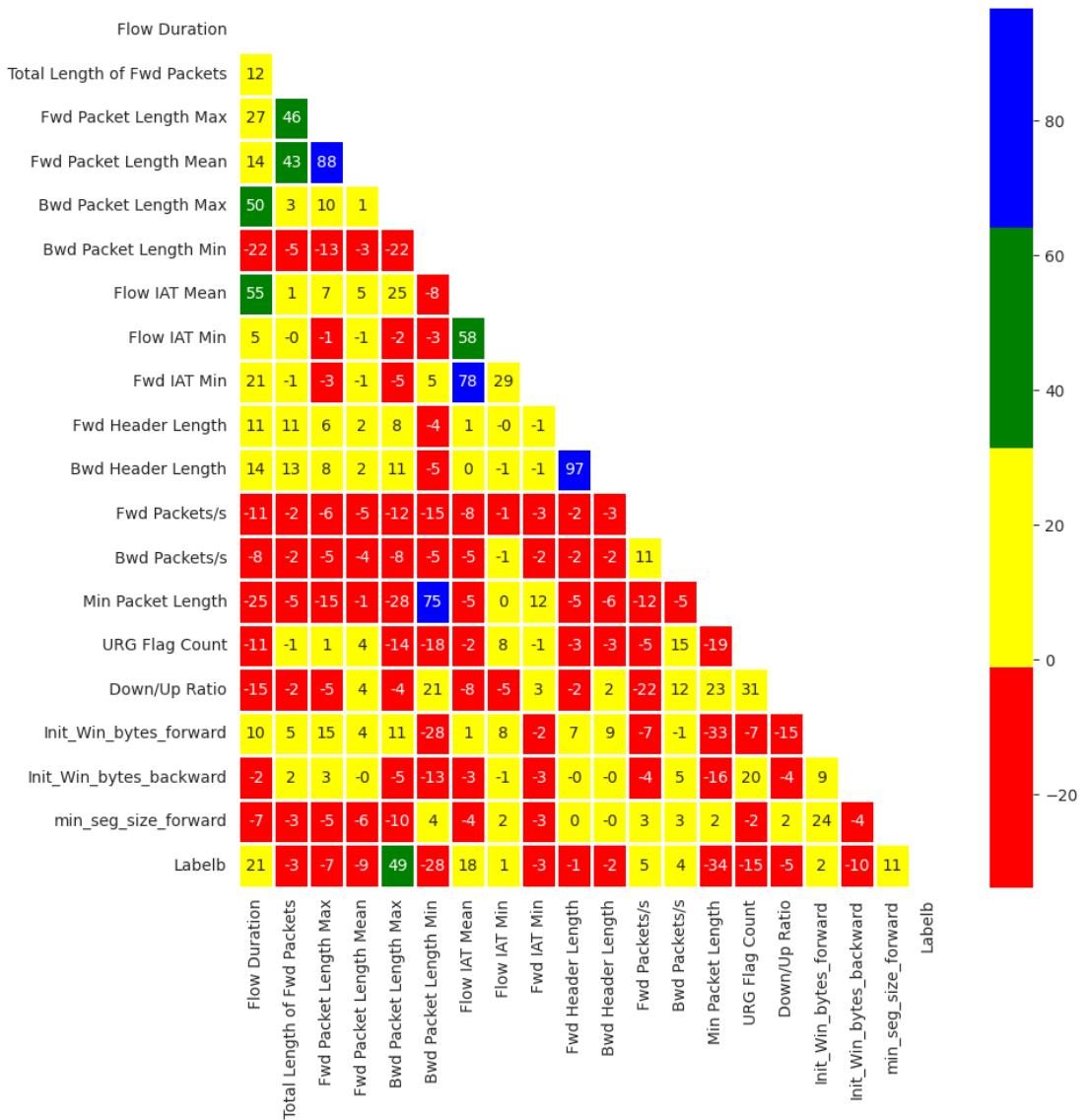


Figure 5.5 – Optimized Heatmap of Static dataset

2. Data Pre-Processing

i. Encoding

Automatically identify and transform string/text features into numerical features to make the data more readable by ML models

```
# Define the automated data encoding function
def Auto_Encoding(df):
    cat_features=[x for x in df.columns if df[x].dtype=="object"] ## Find string/text features
    le=LabelEncoder()
    for col in cat_features:
        if col in df.columns:
            i = df.columns.get_loc(col)
            # Transform to numerical features
            df.iloc[:,i] = df.apply(lambda i:le.fit_transform(i.astype(str)), axis=0, result_type='expand')
    return df
df=Auto_Encoding(df)
```

ii. Imputation

Detect and impute missing values to improve data quality

```
# Define the automated data imputation function
def Auto_Imputation(df):
    if df.isnull().values.any() or np.isinf(df).values.any(): # if there is any empty or infinite values
        df.replace([np.inf, -np.inf], np.nan, inplace=True)
        df.fillna(0, inplace = True) # Replace empty values with zeros; there are other imputation methods discussed in the paper
    return df
df=Auto_Imputation(df)
```

iii. Normalization

```
def Auto_Normalization(df):
    stat, p = shapiro(df)
    print('Statistics=%.3f, p=% .3f' % (stat, p))
    # interpret
    alpha = 0.05
    numeric_features = df.drop(['LabelB'],axis = 1).dtypes[df.dtypes != 'object'].index
    if p > alpha:
        print('Sample looks Gaussian (fail to reject H0)')
```

```

        df[numeric_features] = df[numeric_features].apply(
            lambda x: (x - x.mean()) / (x.std()))
        print('Z-score normalization is automatically chosen and
used')
    else:
        print('Sample does not look Gaussian (reject H0)')
        df[numeric_features] = df[numeric_features].apply(
            lambda x: (x - x.min()) / (x.max()-x.min()))
        print('Min-max normalization is automatically chosen and
used')
    return df
df=Auto_Normalization(df)

```

iv. Train-test split

Split the dataset into the training and the test set

```

X = df.drop(['Labelb'],axis=1)
y = df['Labelb']

# Here we used the 80%/20% split, it can be changed based on specific
tasks
#X_train, X_test, y_train, y_test = train_test_split(X,y, train_size
= 0.8, test_size = 0.2, shuffle=False,random_state = 0)
X_train, X_test, y_train, y_test = train_test_split(X,y, train_size =
0.8, test_size = 0.2,random_state = 0)

```

v. Data balancing

Generate minority class samples to solve class-imbalance and improve data quality.
Synthetic Minority Over-sampling Technique (SMOTE) method is used.

```
| pd.Series(y_train).value_counts()
```

Labelb
0 18126
1 4516
Name: count, dtype: int64

```

# For binary data (can be modified for multi-class data with same
logic)
def Auto_Balancing(X_train, y_train):
    number0 = pd.Series(y_train).value_counts().iloc[0]
    number1 = pd.Series(y_train).value_counts().iloc[1]
    if number0 > number1:
        nlarge = number0
    else:
        nlarge = number1

```

```

# evaluate whether the incoming dataset is imbalanced (the
abnormal/normal ratio is smaller than a threshold (e.g., 50%))
if (number1/number0 > 1.5) or (number0/number1 > 1.5):
    smote=SMOTE(n_jobs=-1,sampling_strategy={0:nlarge, 1:nlarge})
    X_train, y_train = smote.fit_resample(X_train, y_train)
return X_train, y_train
X_train, y_train = Auto_Balancing(X_train, y_train)
pd.Series(y_train).value_counts()

```

Labelb
0 18126
1 18126
Name: count, dtype: int64

3. Model learning

LGBM Classifier Algorithm

```

%%time
lg = lgb.LGBMClassifier(verbose = -1)
lg.fit(X_train,y_train)
t1=time.time()
predictions = lg.predict(X_test)
t2=time.time()
print("Accuracy:
"+str(round(accuracy_score(y_test,predictions),5)*100)+"%")
print("Precision:
"+str(round(precision_score(y_test,predictions),5)*100)+"%")
print("Recall:
"+str(round(recall_score(y_test,predictions),5)*100)+"%")
print("F1-score:
"+str(round(f1_score(y_test,predictions),5)*100)+"%")
print("Time: "+str(round((t2-t1)/len(y_test)*1000000,5)))

```

Accuracy: 99.788%
Precision: 99.3789999999999%
Recall: 99.556%
F1-score: 99.467%
Time: 2.93241
CPU times: user 548 ms, sys: 5.26 ms, total: 554 ms
Wall time: 587 ms

Random Forest Algorithm

```
%%time
rf = RandomForestClassifier()
rf.fit(X_train,y_train)
t1=time.time()
predictions = rf.predict(X_test)
t2=time.time()
print("Accuracy:
"+str(round(accuracy_score(y_test,predictions),5)*100)+"%")
print("Precision:
"+str(round(precision_score(y_test,predictions),5)*100)+"%")
print("Recall:
"+str(round(recall_score(y_test,predictions),5)*100)+"%")
print("F1-score:
"+str(round(f1_score(y_test,predictions),5)*100)+"%")
print("Time: "+str(round((t2-t1)/len(y_test)*1000000,5)))
```

```
Accuracy: 99.717%
Precision: 99.465%
Recall: 99.111%
F1-score: 99.288%
Time: 9.41595
CPU times: user 3.35 s, sys: 9.49 ms, total: 3.35 s
Wall time: 3.41 s
```

Naive Bayes Algorithm

```
%%time
nb = GaussianNB()
nb.fit(X_train,y_train)
t1=time.time()
predictions = nb.predict(X_test)
t2=time.time()
print("Accuracy:
"+str(round(accuracy_score(y_test,predictions),5)*100)+"%")
print("Precision:
"+str(round(precision_score(y_test,predictions),5)*100)+"%")
print("Recall:
"+str(round(recall_score(y_test,predictions),5)*100)+"%")
print("F1-score:
"+str(round(f1_score(y_test,predictions),5)*100)+"%")
print("Time: "+str(round((t2-t1)/len(y_test)*1000000,5)))
```

```
Accuracy: 75.358%
Precision: 44.50799999999996%
Recall: 97.244%
F1-score: 61.06599999999995%
Time: 0.47991
CPU times: user 22.2 ms, sys: 0 ns, total: 22.2 ms
Wall time: 29.2 ms
```

K-Nearest Neighbor (KNN) Algorithm

```
%%time
knn = KNeighborsClassifier()
knn.fit(X_train,y_train)
t1=time.time()
predictions = knn.predict(X_test)
t2=time.time()
print("Accuracy:
"+str(round(accuracy_score(y_test,predictions),5)*100)+"%")
print("Precision:
"+str(round(precision_score(y_test,predictions),5)*100)+"%")
print("Recall:
"+str(round(recall_score(y_test,predictions),5)*100)+"%")
print("F1-score:
"+str(round(f1_score(y_test,predictions),5)*100)+"%")
print("Time: "+str(round((t2-t1)/len(y_test)*1000000,5)))
```

```
Accuracy: 98.834%
Precision: 95.844%
Recall: 98.4%
F1-score: 97.1049999999999%
Time: 164.67113
CPU times: user 900 ms, sys: 0 ns, total: 900 ms
Wall time: 943 ms
```

KerasClassifier Algorithm

```
import tensorflow as tf
from keras.layers import
Input,Dense,Dropout,BatchNormalization,Activation
from keras import Model
import keras.backend as K
import keras.callbacks as kcallbacks
from keras import optimizers
from keras.optimizers import Adam

from sklearn.model_selection import GridSearchCV
from keras.wrappers.scikit_learn import KerasClassifier
from keras.callbacks import EarlyStopping
def ANN(optimizer =
'sgd',neurons=16,batch_size=1024,epochs=80,activation='relu',patience
=8,loss='binary_crossentropy'):
    K.clear_session()
    inputs=Input(shape=(X.shape[1],))
    x=Dense(1000)(inputs)
    x=BatchNormalization()(x)
    x=Activation('relu')(x)
    x=Dropout(0.3)(x)
```

```

x=Dense(256)(inputs)
x=BatchNormalization()(x)
x=Activation('relu')(x)
x=Dropout(0.25)(x)
x=Dense(2,activation='softmax')(x)
model=Model(inputs=inputs,outputs=x,name='base_nlp')
model.compile(optimizer='adam',loss='categorical_crossentropy')
#     model.compile(optimizer=Adam(lr =
0.01),loss='categorical_crossentropy',metrics=['accuracy'])
    early_stopping = EarlyStopping(monitor="loss", patience =
patience)# early stop patience
    history = model.fit(X, pd.get_dummies(y).values,
                      batch_size=batch_size,
                      epochs=epochs,
                      callbacks = [early_stopping],
                      verbose=0) #verbose set to 1 will show the training
process
    return model

```

```

%%time
ann = KerasClassifier(build_fn=ANN, verbose=0)
ann.fit(X_train,y_train)
predictions = ann.predict(X_test)
print("Accuracy:
"+str(round(accuracy_score(y_test,predictions),5)*100)+"%")
print("Precision:
"+str(round(precision_score(y_test,predictions),5)*100)+"%")
print("Recall:
"+str(round(recall_score(y_test,predictions),5)*100)+"%")
print("F1-score:
"+str(round(f1_score(y_test,predictions),5)*100)+"%")
print("Time: "+str(round((t2-t1)/len(y_test)*1000000,5)))

```

```

177/177 [=====] - 1s 1ms/step
Accuracy: 94.559%
Precision: 80.83%
Recall: 95.19999999999999%
F1-score: 87.429%
Time: 164.67113
CPU times: user 27.2 s, sys: 3.42 s, total: 30.6 s
Wall time: 31 s

```

4. Model Selection

Select the best-performing model among five common machine learning models (Naive Bayes, KNN, random forest, LightGBM, and ANN/MLP) by evaluating their learning performance

Method: Grid Search

```
# Create a pipeline
pipe = Pipeline([('classifier', GaussianNB())])
# Create space of candidate learning algorithms and their
hyperparameters
search_space = [{ 'classifier': [GaussianNB()],
                  'classifier': [KNeighborsClassifier()],
                  'classifier': [RandomForestClassifier()],
                  'classifier': [lgb.LGBMClassifier(verbose = -1)],
                  'classifier': [KerasClassifier(build_fn=ANN,
verbose=0)]},
                 ]
clf = GridSearchCV(pipe, search_space, cv=5, verbose=0)
clf.fit(X, y)
print("Best Model:"+ str(clf.best_params_))
print("Accuracy:"+ str(clf.best_score_))
```

```
Best Model:{'classifier': LGBMClassifier(verbose=-1)}
Accuracy:0.9843838600604344
```

```
clf.cv_results_
```

LightGBM model is the best performing machine learning model, and the best cross-validation accuracy is 98.438%

5. Combined Algorithm Selection and Hyperparameter

tuning (CASH)

CASH is the process of combining the two AutoML procedures: model selection and hyperparameter optimization.

Method: Particle Swarm Optimization (PSO)

```
import optunity
import optunity.metrics
search = {'algorithm': {'k-nn': {'n_neighbors': [3, 10]},
                           'naive-bayes': None,
                           'random-forest': {}}
```

```

        'n_estimators': [50, 500],
        'max_features': [5, 12],
        'max_depth': [5, 50],
        "min_samples_split": [2, 11],
        "min_samples_leaf": [1, 11]},
    'lightgbm': {
        'n_estimators': [50, 500],
        'max_depth': [5, 50],
        'learning_rate': (0, 1),
        "num_leaves": [100, 2000],
        "min_child_samples": [10, 50],
        },
    'ann': {
        'neurons': [10, 100],
        'epochs': [20, 50],
        'patience': [3, 20],
        }
    }
}
def performance(
            algorithm, n_neighbors=None,
            n_estimators=None,
            max_features=None, max_depth=None, min_samples_split=None, min_samples_leaf=None,
            learning_rate=None, num_leaves=None, min_child_samples=None,
            neurons=None, epochs=None, patience=None
):
    # fit the model
    if algorithm == 'k-nn':
        model = KNeighborsClassifier(n_neighbors=int(n_neighbors))
    elif algorithm == 'naive-bayes':
        model = GaussianNB()
    elif algorithm == 'random-forest':
        model =
RandomForestClassifier(n_estimators=int(n_estimators),
                        max_features=int(max_features)
,
                        max_depth=int(max_depth),
                        min_samples_split=int(min_samp
les_split),
                        min_samples_leaf=int(min_sampl
es_leaf))
    elif algorithm == 'lightgbm':
        model = lgb.LGBMClassifier(n_estimators=int(n_estimators),
                                max_depth=int(max_depth),
                                learning_rate=float(learning_rate)
,
                                num_leaves=int(num_leaves),
                                min_child_samples=int(min_child_sa
mple),
                                )

```

```

    elif algorithm == 'ann':
        model = KerasClassifier(build_fn=ANN, verbose=0,
                               neurons=int(neurons),
                               epochs=int(epochs),
                               patience=int(patience)
                               )
    else:
        raise ArgumentError('Unknown algorithm: %s' % algorithm)
# predict the test set
model.fit(X_train,y_train)
prediction = model.predict(X_test)
score = accuracy_score(y_test,prediction)
return score
# Run the CASH process
optimal_configuration, info, _ =
optunity.maximize_structured(performance,
                                search_
space=search,
                                num_eva
ls=50)
print(optimal_configuration)
print(info.optimum)

```

```
{
'algorithm': 'lightgbm', 'epochs': None, 'neurons': None, 'patience': None,
'n_neighbors': None, 'learning_rate': 0.32638671874999997, 'max_depth':
22.041113281250006, 'min_child_samples': 40.58548085623468, 'n_estimators':
256.0068359375, 'num_leaves': 1292.5494645058002, 'max_features': None,
'min_samples_leaf': None, 'min_samples_split': None}
0.9978802331743508
```

```

%%time
clf = lgb.LGBMClassifier(max_depth=24, learning_rate= 0.25474609375,
n_estimators = 419,
                    num_leaves = 1463, min_child_samples = 16)
clf.fit(X_train,y_train)
predictions = clf.predict(X_test)
print("Accuracy:
"+str(round(accuracy_score(y_test,predictions),5)*100)+"%")
print("Precision:
"+str(round(precision_score(y_test,predictions),5)*100)+"%")
print("Recall:
"+str(round(recall_score(y_test,predictions),5)*100)+"%")
print("F1-score:
"+str(round(f1_score(y_test,predictions),5)*100)+"%")
```

```

Accuracy: 99.806%
Precision: 99.467%
Recall: 99.556%
F1-score: 99.51100000000001%
CPU times: user 3.53 s, sys: 72.2 ms, total: 3.6 s
Wall time: 3.67 s
```

LightGBM with above hyperparameter values is identified as optimal model

FLASK CODE (Static Data Analytics)

```
from flask import Flask, render_template, request
import joblib
from tensorflow.keras.models import load_model
import numpy as np

app = Flask(__name__, template_folder='templates')

# Load the saved models
Random_Forest_Algorithm = joblib.load("rf_model.pkl")
k_nearest_neighbor_Algorithm = joblib.load("knn_model.pkl")
Naive_Bayes_Algorithm = joblib.load("nb_model.pkl")
# Keras_algorithm = joblib.load("ann_model.pkl")
LGBM_Classifier_Algorithm = joblib.load("clf_model.pkl")

@app.route('/', methods=['GET', 'POST'])
def index():
    predictions = {}
    if request.method == 'POST':
        input_data = request.form['input_data']
        input_data_list = [float(x.strip()) for x in
input_data.split(',')]

        # Perform predictions using different models
        predictions['naive_bayes'] = 'Safe' if
make_prediction(input_data_list, Naive_Bayes_Algorithm) == 0 else
'Malicious'

        predictions['LGBM_classifier'] = 'Safe' if
make_prediction(input_data_list, LGBM_Classifier_Algorithm) == 0 else
'Malicious'

        predictions['LGBM_clasifier'] = 'Safe' if
make_prediction(input_data_list, LGBM_Classifier_Algorithm) == 0 else
'Malicious'

        print("Predictions:", predictions)
        return render_template('index.html', predictions=predictions)
    return render_template('index.html')

def make_prediction(input_data, model):
    try:
        # Perform the prediction using the model
        input_data_reshaped = np.array(input_data).reshape(1, -1)
        prediction = model.predict(input_data_reshaped)
        print("Prediction:", prediction)

        return prediction[0]
    except ValueError:
        # Handle case where input data cannot be converted to numbers
        print(ValueError)
        return None
if __name__ == '__main__':
    app.run(host='0.0.0.0', port=8080)
```

AutoML – 2:Dynamic Dataset - IoTID20

1. Data Visualization

```
| df.shape
```

```
(28303, 20)
```

```
| df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6252 entries, 0 to 6251
Data columns (total 32 columns):
 #   Column           Non-Null Count  Dtype  
 ---  --  
 0   Flow_ID          6252 non-null   int64  
 1   Src_IP           6252 non-null   int64  
 2   Src_Port         6252 non-null   int64  
 3   Dst_IP           6252 non-null   int64  
 4   Dst_Port         6252 non-null   int64  
 5   Timestamp        6252 non-null   int64  
 6   Flow_Duration   6252 non-null   int64  
 7   TotLen_Fwd_Pkts 6252 non-null   float64 
 8   TotLen_Bwd_Pkts 6252 non-null   float64 
 9   Bwd_Pkt_Len_Max 6252 non-null   float64 
 10  Flow_Byts/s     6252 non-null   float64 
 11  Flow_Pkts/s     6252 non-null   float64 
 12  Flow_IAT_Mean   6252 non-null   float64 
 13  Flow_IAT_Std    6252 non-null   float64 
 14  Flow_IAT_Max    6252 non-null   float64 
 15  Flow_IAT_Min    6252 non-null   float64 
 16  Fwd_IAT_Tot    6252 non-null   float64 
 17  Bwd_IAT_Tot    6252 non-null   float64 
 18  Bwd_IAT_Mean   6252 non-null   float64 
 19  Bwd_IAT_Std    6252 non-null   float64 
 20  Bwd_IAT_Max    6252 non-null   float64 
 21  Bwd_Header_Len 6252 non-null   int64  
 22  Fwd_Pkts/s     6252 non-null   float64 
 23  Bwd_Pkts/s     6252 non-null   float64 
 24  Pkt_Len_Std    6252 non-null   float64 
 25  FIN_Flag_Cnt   6252 non-null   int64  
 26  Pkt_Size_Avg   6252 non-null   float64 
 27  Init_Bwd_Win_Byts 6252 non-null   int64  
 28  Idle_Mean       6252 non-null   float64 
 29  Idle_Max        6252 non-null   float64 
 30  Idle_Min        6252 non-null   float64 
 31  Label           6252 non-null   int64  
dtypes: float64(21), int64(11)
memory usage: 1.5 MB
```

```

from matplotlib import pyplot as plt
df.hist(figsize=(15,15), xrot=-45, bins=10) ## Display the labels
rotated by 45 degrees

# Clear the text "residue"
plt.show()

```

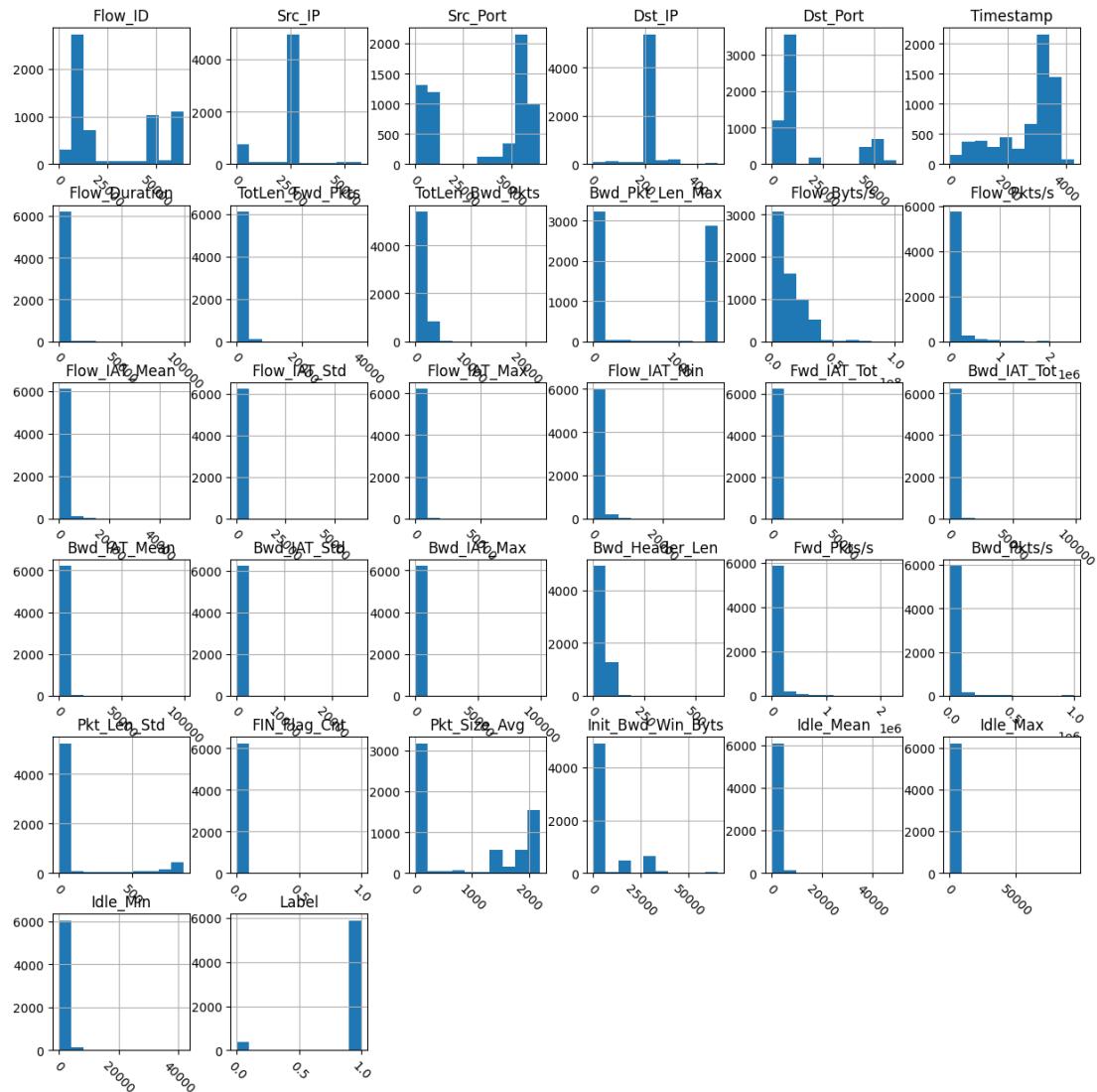


Figure 5.6 – Bar chart of Dynamic dataset

```

plt.figure(figsize=(20,20))
dataplot=sns.heatmap(df.corr(), annot=True)
plt.show()

```

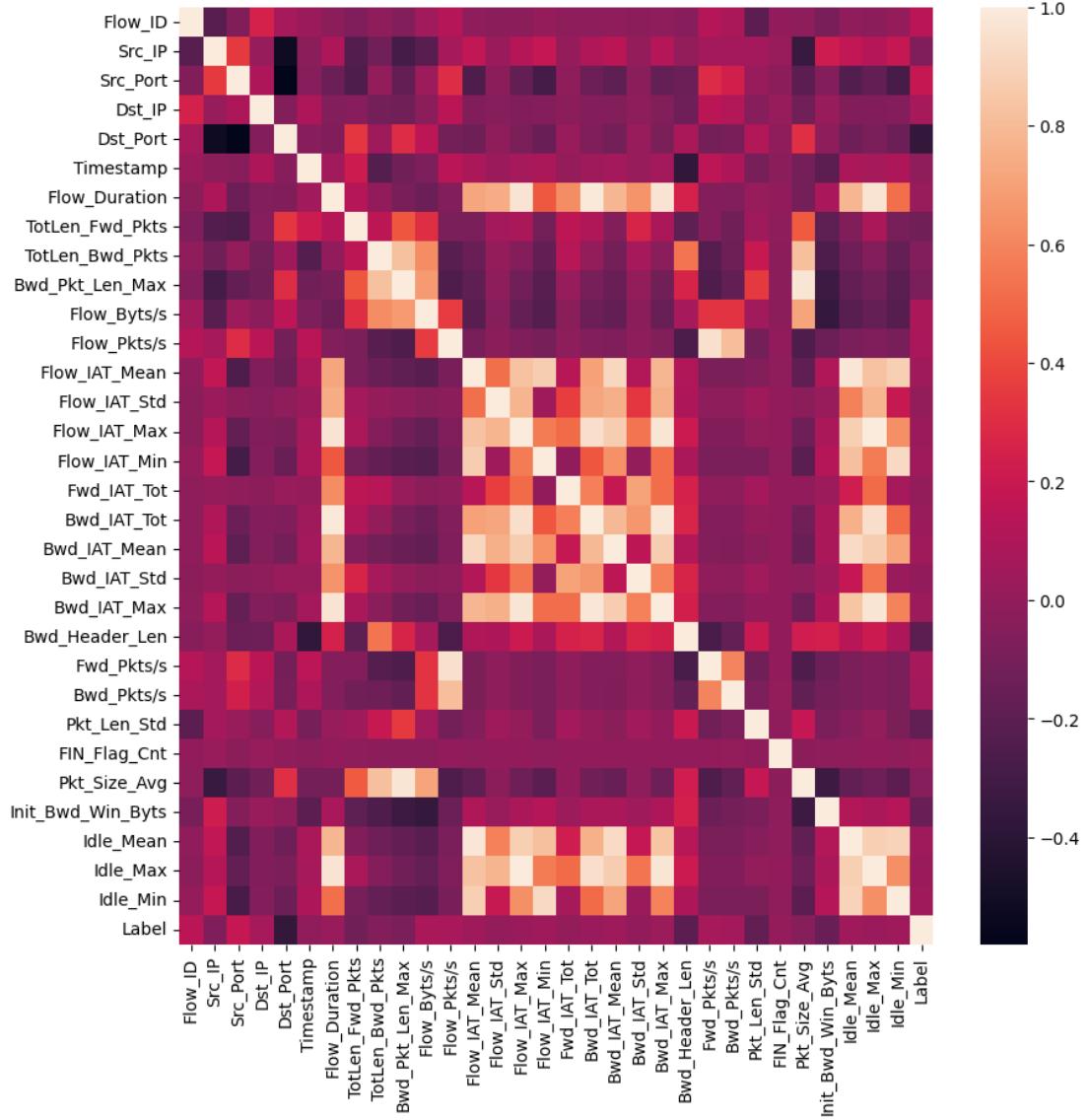


Figure 5.7 – Heatmap of Dynamic dataset

```

from matplotlib.colors import ListedColormap

mask = np.triu(np.ones_like(df.corr(), dtype=bool)) # Upper triangle
mask
plt.figure(figsize=(10,10))
with sns.axes_style("white"):
    ax = sns.heatmap(df.corr()*100, mask=mask, fmt='.0f', annot=True,
lw=1, cmap=ListedColormap(['red', 'yellow', 'green', 'blue']))
plt.show()

```

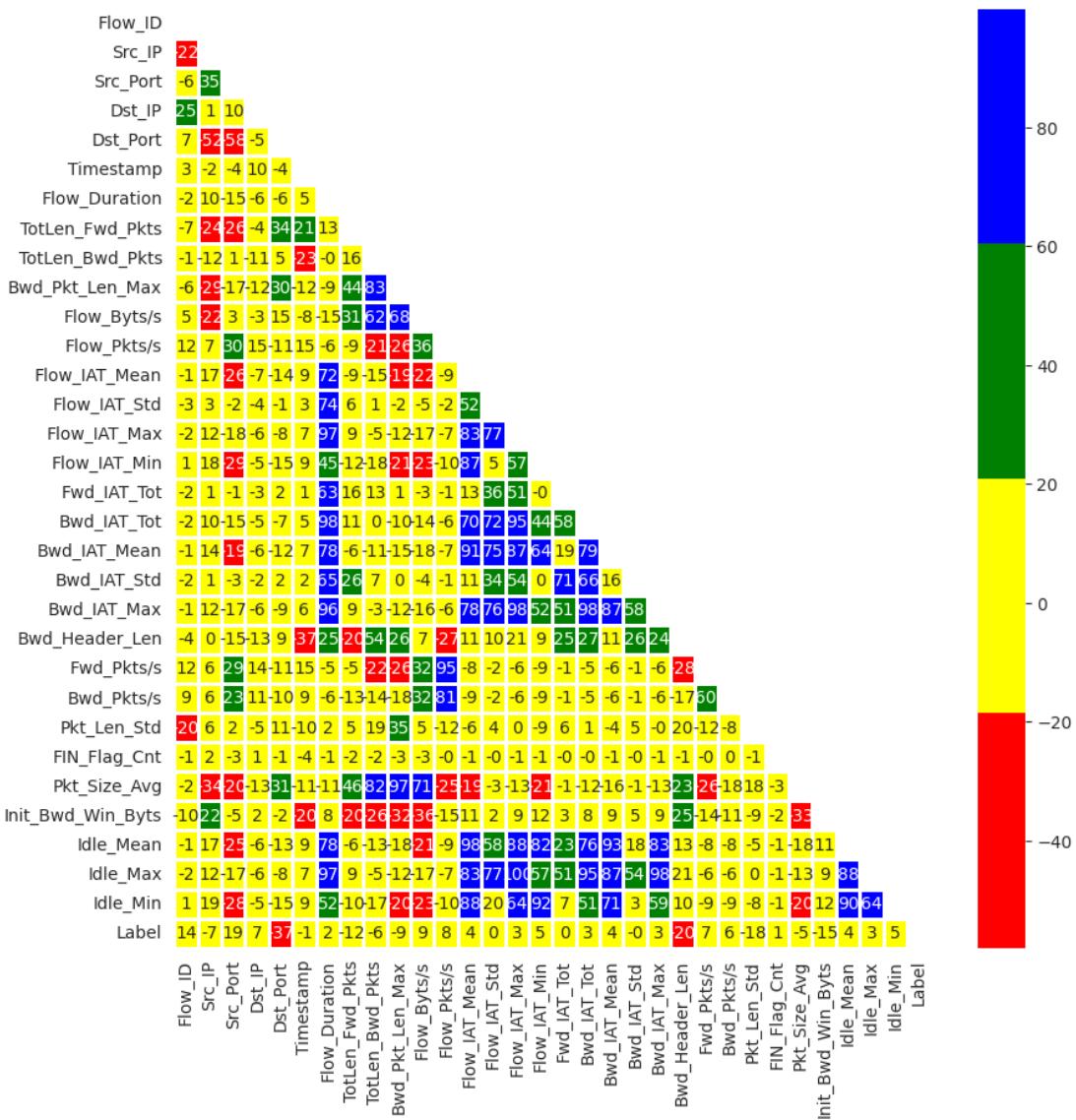


Figure 5.8 – Optimized Heatmap of Dynamic dataset

2. Data Pre-Processing

i. Encoding

Automatically identify and transform string/text features into numerical features to make the data more readable by ML models

```
# Define the automated data encoding function
def Auto_Encoding(df):
    cat_features=[x for x in df.columns if df[x].dtype=="object"] ##
Find string/text features
    le=LabelEncoder()
    for col in cat_features:
        if col in df.columns:
            i = df.columns.get_loc(col)
            # Transform to numerical features
            df.iloc[:,i] = df.apply(lambda
i:le.fit_transform(i.astype(str)), axis=0, result_type='expand')
    return df

df=Auto_Encoding(df)
```

ii. Imputation

Detect and impute missing values to improve data quality

```
# Define the automated data imputation function
def Auto_Imputation(df):
    if df.isnull().values.any() or np.isinf(df).values.any(): # if
there is any empty or infinite values
        df.replace([np.inf, -np.inf], np.nan, inplace=True)
        df.fillna(0, inplace = True) # Replace empty values with
zeros; there are other imputation methods discussed in the paper
    return df

df=Auto_Imputation(df)
```

iii. Normalization

Normalize the range of features to a similar scale to improve data quality

```
def Auto_Normalization(df):
    stat, p = shapiro(df)
    print('Statistics=%f, p=%f' % (stat, p))
    # interpret
    alpha = 0.05
    numeric_features = df.drop(['Label'], axis=1).dtypes[df.dtypes != 'object'].index

    # The selection strategy is based on the following article:
    # https://medium.com/@kumarvaishnav17/standardization-vs-
    normalization-in-machine-learning-3e132a19c8bf
    # Check if the data distribution follows a Gaussian/normal
    distribution
    # If so, select the Z-score normalization method; otherwise,
    select the min-max normalization
    # Details are in the paper
    if p > alpha:
        print('Sample looks Gaussian (fail to reject H0)')
        df[numeric_features] = df[numeric_features].apply(
            lambda x: (x - x.mean()) / (x.std()))
        print('Z-score normalization is automatically chosen and
used')
    else:
        print('Sample does not look Gaussian (reject H0)')
        df[numeric_features] = df[numeric_features].apply(
            lambda x: (x - x.min()) / (x.max() - x.min()))
        print('Min-max normalization is automatically chosen and
used')
    return df
df=Auto_Normalization(df)
```

```
Statistics=0.108, p=0.000
Sample does not look Gaussian (reject H0)
Min-max normalization is automatically chosen and used
```

iv. Train-test split

Split the dataset into the training and the test set

```
X = df.drop(['Label'], axis=1)
y = df['Label']

# Here we used the 80%/20% split, it can be changed based on specific
tasks
```

```
#X_train, X_test, y_train, y_test = train_test_split(X,y, train_size = 0.8, test_size = 0.2, shuffle=False,random_state = 0)
X_train, X_test, y_train, y_test = train_test_split(X,y, train_size = 0.8, test_size = 0.2,random_state = 0)
```

v. Data balancing

Generate minority class samples to solve class-imbalance and improve data quality.
Synthetic Minority Over-sampling Technique (SMOTE) method is used.

```
pd.Series(y_train).value_counts()
```

Label	count
1	4717
0	284
Name:	count, dtype: int64

```
# For binary data (can be modified for multi-class data with the same logic)
def Auto_Balancing(X_train, y_train):
    number0 = pd.Series(y_train).value_counts().iloc[0]
    number1 = pd.Series(y_train).value_counts().iloc[1]

    if number0 > number1:
        nlarge = number0
    else:
        nlarge = number1

    # evaluate whether the incoming dataset is imbalanced (the abnormal/normal ratio is smaller than a threshold (e.g., 50%))
    if (number1/number0 > 1.5) or (number0/number1 > 1.5):
        smote=SMOTE(n_jobs=-1,sampling_strategy={0:nlarge, 1:nlarge})
        X_train, y_train = smote.fit_resample(X_train, y_train)

    return X_train, y_train

X_train, y_train = Auto_Balancing(X_train, y_train)
```

```
pd.Series(y_train).value_counts()
```

Label	count
1	4717
0	4717
Name:	count, dtype: int64

3. Model learning

LGBM Classifier Algorithm

```
%%time
lg = lgb.LGBMClassifier(verbose = -1)
lg.fit(X_train,y_train)
t1=time.time()
predictions = lg.predict(X_test)
t2=time.time()
print("Accuracy:
"+str(round(accuracy_score(y_test,predictions),5)*100)+"%")
print("Precision:
"+str(round(precision_score(y_test,predictions),5)*100)+"%")
print("Recall:
"+str(round(recall_score(y_test,predictions),5)*100)+"%")
print("F1-score:
"+str(round(f1_score(y_test,predictions),5)*100)+"%")
print("Time: "+str(round((t2-t1)/len(y_test)*1000000,5)))
```

```
Accuracy: 99.92%
Precision: 99.9149999999999%
Recall: 100.0%
F1-score: 99.957%
Time: 4.70376
Time: 5.09103
CPU times: user 314 ms, sys: 3.53 ms, total: 318 ms
Wall time: 325 ms
```

Random Forest Algorithm

```
%%time
rf = RandomForestClassifier()
rf.fit(X_train,y_train)
t1=time.time()
predictions = rf.predict(X_test)
t2=time.time()
print("Accuracy:
"+str(round(accuracy_score(y_test,predictions),5)*100)+"%")
print("Precision:
"+str(round(precision_score(y_test,predictions),5)*100)+"%")
print("Recall:
"+str(round(recall_score(y_test,predictions),5)*100)+"%")
print("F1-score:
"+str(round(f1_score(y_test,predictions),5)*100)+"%")
print("Time: "+str(round((t2-t1)/len(y_test)*1000000,5)))
```

```

Accuracy: 99.8399999999999%
Precision: 99.83%
Recall: 100.0%
F1-score: 99.9149999999999%
Time: 12.99924
CPU times: user 1.01 s, sys: 0 ns, total: 1.01 s
Wall time: 1.03 s

```

Naive Bayes Algorithm

```

%%time
nb = GaussianNB()
nb.fit(X_train,y_train)
t1=time.time()
predictions = nb.predict(X_test)
t2=time.time()
print("Accuracy:
"+str(round(accuracy_score(y_test,predictions),5)*100)+"%")
print("Precision:
"+str(round(precision_score(y_test,predictions),5)*100)+"%")
print("Recall:
"+str(round(recall_score(y_test,predictions),5)*100)+"%")
print("F1-score:
"+str(round(f1_score(y_test,predictions),5)*100)+"%")
print("Time: "+str(round((t2-t1)/len(y_test)*1000000,5)))

```

```

Accuracy: 69.784%
Precision: 99.875%
Recall: 67.888%
F1-score: 80.83200000000001%
Time: 2.00245
CPU times: user 14.5 ms, sys: 0 ns, total: 14.5 ms
Wall time: 19.5 ms

```

K-Nearest Neighbor (KNN) Algorithm

```

%%time
knn = KNeighborsClassifier()
knn.fit(X_train,y_train)
t1=time.time()
predictions = knn.predict(X_test)
t2=time.time()
print("Accuracy:
"+str(round(accuracy_score(y_test,predictions),5)*100)+"%")
print("Precision:
"+str(round(precision_score(y_test,predictions),5)*100)+"%")
print("Recall:
"+str(round(recall_score(y_test,predictions),5)*100)+"%")

```

```

print("F1-score:
"+str(round(f1_score(y_test,predictions),5)*100)+"%")
print("Time: "+str(round((t2-t1)/len(y_test)*1000000,5)))

```

```

Accuracy: 98.801%
Precision: 99.82799999999999%
Recall: 98.893%
F1-score: 99.358%
Time: 64.61487
CPU times: user 87.2 ms, sys: 0 ns, total: 87.2 ms
Wall time: 90.2 ms

```

KerasClassifier Algorithm

```

import tensorflow as tf
from keras.layers import
Input,Dense,Dropout,BatchNormalization,Activation
from keras import Model
import keras.backend as K
import keras.callbacks as kcallbacks
from keras import optimizers
from keras.optimizers import Adam

from sklearn.model_selection import GridSearchCV
from keras.wrappers.scikit_learn import KerasClassifier
from keras.callbacks import EarlyStopping
def ANN(optimizer =
'sgd',neurons=16,batch_size=1024,epochs=80,activation='relu',patience
=8,loss='binary_crossentropy'):
    K.clear_session()
    inputs=Input(shape=(X.shape[1],))
    x=Dense(1000)(inputs)
    x=BatchNormalization()(x)
    x=Activation('relu')(x)
    x=Dropout(0.3)(x)
    x=Dense(256)(inputs)
    x=BatchNormalization()(x)
    x=Activation('relu')(x)
    x=Dropout(0.25)(x)
    x=Dense(2,activation='softmax')(x)
    model=Model(inputs=inputs,outputs=x,name='base_nlp')
    model.compile(optimizer='adam',loss='categorical_crossentropy')
#    model.compile(optimizer=Adam(lr =
0.01),loss='categorical_crossentropy',metrics=['accuracy'])
    early_stopping = EarlyStopping(monitor="loss", patience =
patience)# early stop patience
    history = model.fit(X, pd.get_dummies(y).values,
                        batch_size=batch_size,
                        epochs=epochs,
                        callbacks = [early_stopping],

```

```

        verbose=0) #verbose set to 1 will show the training
process
    return model

%%time
ann = KerasClassifier(build_fn=ANN, verbose=0)
ann.fit(X_train,y_train)
predictions = ann.predict(X_test)
print("Accuracy:
"+str(round(accuracy_score(y_test,predictions),5)*100)+"%")
print("Precision:
"+str(round(precision_score(y_test,predictions),5)*100)+"%")
print("Recall:
"+str(round(recall_score(y_test,predictions),5)*100)+"%")
print("F1-score:
"+str(round(f1_score(y_test,predictions),5)*100)+"%")
print("Time: "+str(round((t2-t1)/len(y_test)*1000000,5)))

```

```

40/40 [=====] - 0s 2ms/step
Accuracy: 99.2809999999999%
Precision: 99.574%
Recall: 99.6589999999999%
F1-score: 99.617%
Time: 124.9462
CPU times: user 10.6 s, sys: 938 ms, total: 11.6 s
Wall time: 11.8 s

```

4. Automated Model Selection

Select the best-performing model among five common machine learning models (Naive Bayes, KNN, random forest, LightGBM, and ANN/MLP) by evaluating their learning performance

Method: Grid Search

```

# Create a pipeline
pipe = Pipeline([('classifier', GaussianNB())])

# Create space of candidate learning algorithms and their
hyperparameters
search_space = [{ 'classifier': [GaussianNB()],
                  'classifier': [KNeighborsClassifier()],
                  'classifier': [RandomForestClassifier()],
                  'classifier': [lgb.LGBMClassifier(verbose = -1)],
                  'classifier': [KerasClassifier(build_fn=ANN,
verbose=0)]},
                 ]

clf = GridSearchCV(pipe, search_space, cv=5, verbose=0)
clf.fit(X, y)

```

```

▼ LGBMClassifier
LGBMClassifier(learning_rate=0.88427734375, max_depth=28, min_child_samples=40,
n_estimators=78, num_leaves=251)

```

Figure 5.9 – Grid search pipeline

```

print("Best Model:"+ str(clf.best_params_))
print("Accuracy:"+ str(clf.best_score_))

```

```

Best Model:{'classifier': LGBMClassifier(verbose=-1)}
Accuracy:0.9993601278976818

```

LightGBM model is the best performing machine learning model, and the best cross-validation accuracy is 99.936%

5. Combined Algorithm Selection and Hyperparameter tuning (CASH)

CASH is the process of combining the two AutoML procedures: model selection and hyperparameter optimization.

Method: Particle Swarm Optimization (PSO)

```

import optunity
import optunity.metrics

search = {'algorithm': {'k-nn': {'n_neighbors': [3, 10]},
                       'naive-bayes': None,
                       'random-forest': {
                           'n_estimators': [50, 500],
                           'max_features': [5, 12],
                           'max_depth': [5, 50],
                           "min_samples_split": [2, 11],
                           "min_samples_leaf": [1, 11]},
                       'lightgbm': {
                           'n_estimators': [50, 500],
                           'max_depth': [5, 50],
                           'learning_rate': (0, 1),
                           "num_leaves": [100, 2000],
                           "min_child_samples": [10, 50],
                           },
                       'ann': {
                           'neurons': [10, 100],
                           'epochs': [20, 50],
                           'patience': [3, 20],
                           }
                     }
       }

```

```

        }
    }

def performance(
    algorithm, n_neighbors=None,
    n_estimators=None,
    max_features=None, max_depth=None, min_samples_split=None, min_samples_leaf=None,
    learning_rate=None, num_leaves=None, min_child_samples=None,
    neurons=None, epochs=None, patience=None
):
    # fit the model
    if algorithm == 'k-nn':
        model = KNeighborsClassifier(n_neighbors=int(n_neighbors))
    elif algorithm == 'naive-bayes':
        model = GaussianNB()
    elif algorithm == 'random-forest':
        model =
RandomForestClassifier(n_estimators=int(n_estimators),
                        max_features=int(max_features)
,
                        max_depth=int(max_depth),
                        min_samples_split=int(min_samples_split),
                        min_samples_leaf=int(min_samples_leaf))
    elif algorithm == 'lightgbm':
        model = lgb.LGBMClassifier(n_estimators=int(n_estimators),
                                max_depth=int(max_depth),
                                learning_rate=float(learning_rate)
,
                                num_leaves=int(num_leaves),
                                min_child_samples=int(min_child_samples),
)
    elif algorithm == 'ann':
        model = KerasClassifier(build_fn=ANN, verbose=0,
                            neurons=int(neurons),
                            epochs=int(epochs),
                            patience=int(patience)
)
    else:
        raise ArgumentError('Unknown algorithm: %s' % algorithm)
# predict the test set
model.fit(X_train,y_train)
prediction = model.predict(X_test)
score = accuracy_score(y_test,prediction)
return score

# Run the CASH process
optimal_configuration, info, _ =
opportunity.maximize_structured(performance,

```

```

    search_
space=search,
num_eva
ls=50)
print(optimal_configuration)
print(info.optimum)

```

```
{
'algorithm': 'lightgbm', 'epochs': None, 'neurons': None, 'patience': None,
'n_neighbors': None, 'learning_rate': 0.32638671874999997, 'max_depth':
22.041113281250006, 'min_child_samples': 40.58548085623468, 'n_estimators':
256.0068359375, 'num_leaves': 1292.5494645058002, 'max_features': None,
'min_samples_leaf': None, 'min_samples_split': None}
0.9978802331743508
}
```

```

%%time
clf = lgb.LGBMClassifier(max_depth=24, learning_rate= 0.25474609375,
n_estimators = 419,
                    num_leaves = 1463, min_child_samples = 16)
clf.fit(X_train,y_train)
predictions = clf.predict(X_test)
print("Accuracy:
"+str(round(accuracy_score(y_test,predictions),5)*100)+"%")
print("Precision:
"+str(round(precision_score(y_test,predictions),5)*100)+"%")
print("Recall:
"+str(round(recall_score(y_test,predictions),5)*100)+"%")
print("F1-score:
"+str(round(f1_score(y_test,predictions),5)*100)+"%")

```

```

Accuracy: 99.77000000000001%
Precision: 99.291%
Recall: 99.556%
F1-score: 99.423%
CPU times: user 2.23 s, sys: 0 ns, total: 2.23 s
Wall time: 2.29 s

```

LightGBM with the above hyperparameter values is identified as the optimal model

```

# Assuming X_test is a NumPy array
first_row = X_test.iloc[111]
first_row_csv = ','.join(map(str, first_row))

print(first_row_csv)
0.97359710208278,0.44651604001379785,0.927589027985737,0.4830508474576271,0.1463121896
6821216,0.7523277467411545,0.0019372609954696344,0.0016782043213761275,0.00142570728447
31566,0.021917808219178082,0.0049014787993416275,0.0062055632610712305,0.00194160392735
17543,0.0019573912397514944,0.00194322406218317,0.00022209266816579219,0.00210802187784
8678,0.0,0.0,0.0,0.0,0.013157894736842105,0.004533676866499356,0.00516124430874748,0.0,0.0,0.
019484018264840183,0.0,0.0019504800404244568,0.001953764428814329,0.0001919754271453254

```

GENERATING .PKL FILES FOR MODELS

```
import joblib

1 joblib.dump(rf, 'rf_model.pkl')
[ 'rf_model.pkl' ]



1 joblib.dump(nb, 'nb_model.pkl')
[ 'nb_model.pkl' ]



1 joblib.dump(knn, 'knn_model.pkl')
[ 'knn_model.pkl' ]



1 joblib.dump(ann, 'ann_model.pkl')
[ 'ann_model.pkl' ]



1 joblib.dump(clf, 'clf_model.pkl')
[ 'clf_model.pkl' ]
```

FLASK CODE (Dynamic Data Analytics)

```
from flask import Flask, render_template, request
import joblib
from tensorflow.keras.models import load_model
import numpy as np

app = Flask(__name__, template_folder='templates')

# Load the saved models
Random_Forest_Algorithm = joblib.load("rf_model.pkl")
k_nearest_neighbor_Algorithm = joblib.load("knn_model.pkl")
Naive_Bayes_Algorithm = joblib.load("nb_model.pkl")
# Keras_algorithm = joblib.load("ann_model.pkl")
LGBM_Classifier_Algorithm = joblib.load("clf_model.pkl")

@app.route('/', methods=['GET', 'POST'])
def index():
    predictions = {}
    if request.method == 'POST':
        input_data = request.form['input_data']
        input_data_list = [float(x.strip()) for x in
input_data.split(',')]

        # Perform predictions using different models
```

```

        predictions['Random_Forest'] = 'Safe' if
make_prediction(input_data_list, Random_Forest_Algorithm) == 0 else
'Malicious'
        predictions['k_nearest_neighbor'] = 'Safe' if
make_prediction(input_data_list, k_nearest_neighbor_Algorithm) == 0
else 'Malicious'
        predictions['naive_bayes'] = 'Safe' if
make_prediction(input_data_list, Naive_Bayes_Algorithm) == 0 else
'Malicious'
        # predictions['Keras_algorithm'] = 'Safe' if
make_prediction(input_data_list, Keras_algorithm) == 0 else
'Malicious'
        predictions['LGBM_classifier'] = 'Safe' if
make_prediction(input_data_list, LGBM_Classifier_Algorithm) == 0 else
'Malicious'

    print("Predictions:", predictions)

    return render_template('index.html', predictions=predictions)

return render_template('index.html')

def make_prediction(input_data, model):
    try:
        # Perform the prediction using the model
        input_data_reshaped = np.array(input_data).reshape(1, -1)
        prediction = model.predict(input_data_reshaped)
        print("Prediction:", prediction)

        return prediction[0]
    except ValueError:
        # Handle the case where input data cannot be converted to
numbers
        print(ValueError)
        return None

if __name__ == '__main__':
    app.run(host='0.0.0.0', port=8080)

```

5.3 SAMPLE OUTPUT

Automated Machine Learning - Dataset 1: CICIDS2017

Importing Libraries

```

1 import joblib
2
3 import pandas as pd
4 import numpy as np
5 from sklearn.model_selection import train_test_split,cross_val_score
6 import lightgbm as lgb
7 from sklearn.preprocessing import LabelEncoder
8 from sklearn.ensemble import RandomForestClassifier
9 from sklearn.metrics import classification_report,confusion_matrix,accuracy_score,precision_score,recall_score,f1_score
10 from sklearn.neighbors import KNeighborsClassifier
11 from sklearn.svm import SVC
12 from sklearn.naive_bayes import GaussianNB
13 from sklearn.pipeline import Pipeline
14 from sklearn.model_selection import GridSearchCV
15 from scipy.stats import shapiro
16 from imblearn.over_sampling import SMOTE
17 import seaborn as sns
18 import warnings
19 warnings.filterwarnings('ignore')

```

ML Model Learning

LGBM Classifier Algorithm

```

1 lg = lgb.LGBMClassifier(verbose = -1)
2 lg.fit(X_train,y_train)
3 t1=time.time()
4 predictions = lg.predict(X_test)
5 t2=time.time()
6 print("Accuracy: "+str(round(accuracy_score(y_test,predictions),5)*100)+"%")
7 print("Precision: "+str(round(precision_score(y_test,predictions),5)*100)+"%")
8 print("Recall: "+str(round(recall_score(y_test,predictions),5)*100)+"%")
9 print("F1-score: "+str(round(f1_score(y_test,predictions),5)*100)+"%")
10

```

Output:

```

Accuracy: 99.788%
Precision: 99.291%
Recall: 99.644%
F1-score: 99.468%

```

Random Forest Algorithm

```

1 rf = RandomForestClassifier()
2 rf.fit(X_train,y_train)
3 t1=time.time()
4 predictions = rf.predict(X_test)
5 t2=time.time()
6 print("Accuracy: "+str(round(accuracy_score(y_test,predictions),5)*100)+"%")
7 print("Precision: "+str(round(precision_score(y_test,predictions),5)*100)+"%")
8 print("Recall: "+str(round(recall_score(y_test,predictions),5)*100)+"%")

```

Terminal

```

---- Started new terminal session ----
(virtualenv) root@deepnote:~/work # pip install flask
Collecting flask
  Downloading flask-3.0.3-py3-none-any.whl (101 kB)
Requirement already satisfied: Werkzeug<3.0.0 in /shared-libs/python3.11/py/lib/python3.11/site-packages (from flask) (3.0.0)
Requirement already satisfied: Jinja2>=3.1.2 in /shared-libs/python3.11/py-core/lib/python3.11/site-packages (from flask) (3.1.2)
Collecting itsdangerous<2.2.0 (from flask)
  Downloading itsdangerous-2.2.0-py3-none-any.whl (16 kB)
Requirement already satisfied: click<8.1.3 in /shared-libs/python3.11/py/lib/python3.11/site-packages (from flask) (8.1.7)
Collecting blinker>1.6.2 (from flask)
  Downloading blinker-1.7.0-py3-none-any.whl (13 kB)
Requirement already satisfied: MarkupSafe>2.0 in /shared-libs/python3.11/py-core/lib/python3.11/site-packages (from Jinja2>=3.1.2->flask) (2.1.3)
Installing collected packages: itsdangerous, blinker, flask
Successfully installed flask-3.0.3 itsdangerous-2.2.0

[notice] A new release of pip is available: 24.1.2 > 24.0.0
[notice] To update, run: pip install --upgrade pip
(virtualenv) root@deepnote:~/work # python app.py
2024-04-24 12:53:11.530966: I tensorflow/tsl/cuda/cudart_stub.cc:28] Could not find cuda drivers on your machine, GPU will not be used.
2024-04-24 12:53:11.597837: I tensorflow/tsl/cuda/cudart_stub.cc:28] Could not find cuda drivers on your machine, GPU will not be used.
2024-04-24 12:53:11.598529: W tensorflow/core/platform/cuda/feature_guard.cc:182] This TensorFlow binary is optimized to use available CPU instructions in preference to potentially slower GPU instructions.
To enable the following instructions: AVX2 AVX512F FMA, in other operations, rebuild TensorFlow with the appropriate compiler flags.
2024-04-24 12:53:12.624929: W tensorflow/compiler/tf2tensorrt/utils/py_utils.cc:38] TF-TRT Warning: Could not find TensorRT
* Serving Flask app 'app'
* Debug mode: off
WARNING: This is a development server. Do not use it in a production deployment. Use a production WSGI server instead.
* Running on all addresses (0.0.0.0)
* Running on http://127.0.0.1:8080
* Running on http://172.3.57.76:8080

```

Figure 5.10 – Screenshot of Deepnote Environment of Static Data Analytics

Static Data Analytics – Frontend User Interface (UI)

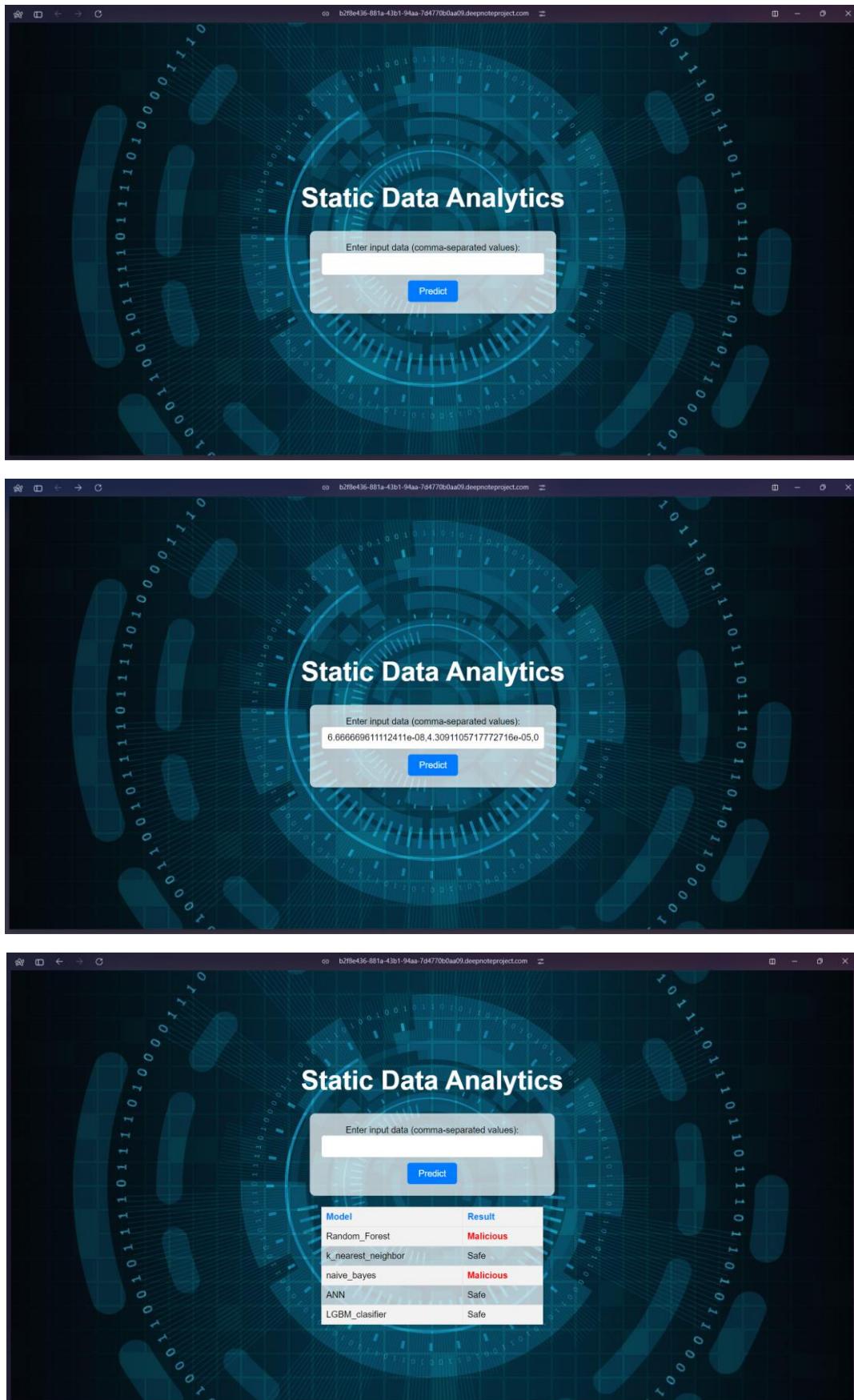


Figure 5.11 – Screenshot of Static Data Analytics User Interface

The screenshot displays three separate Deepnote environments, each showing a different stage of a machine learning project:

- Automated Machine Learning - Dataset 2: IoTID20**: This notebook focuses on importing libraries. The code cell contains imports for pandas, numpy, scikit-learn, lightgbm, and various metrics. It also includes imports for LabelEncoder, Pipeline, GridSearchCV, SVC, GaussianNB, KNeighborsClassifier, and SVC.
- ML Model Learning**: This notebook compares LGBM and Random Forest classifiers. The LGBM section shows execution results with accuracy, precision, recall, F1-score, and time metrics. The Random Forest section shows a similar set of metrics.
- Terminal**: This section shows the command-line interface used for installing packages and running the Flask app. It includes logs for pip installations and the start of the Flask application on port 8080.

Figure 5.12 – Screenshot of Deepnote Environment of Dynamic Data Analytics

Dynamic Data Analytics – Frontend User Interface (UI)



Figure 5.13 – Screenshot of Dynamic Data Analytics User Interface

5.4 TEST PLAN

Testing the end-to-end functionality of the AutoML framework using the CICIDS2017 intrusion detection dataset and IoT network intrusion datasets.

Test Cases

Detailed test cases will be developed to cover the following aspects:

1. Automated Data Preprocessing Module
 - Test cases for each sub-module (Encoding, Imputation, Normalization, Train-Test Split, Data Balancing)
 - Test cases for handling various data types, missing values, and class imbalance
2. Automated Model Learning Module
 - Test cases for training and evaluating each machine learning model (Naive Bayes, KNN, Random Forest, LightGBM, ANN)
 - Test cases for handling different hyperparameter configurations
3. Automated Model Selection Module
 - Test cases for Grid Search-based model selection and cross-validation
4. Combined Algorithm Selection and Hyperparameter Tuning (CASH) Module
 - Test cases for Particle Swarm Optimization-based model selection
5. End-to-End AutoML Pipeline
 - Test cases for the complete AutoML pipeline using the CICIDS2017 and IoT network intrusion datasets
 - Test cases for evaluating the performance of the AutoML models against traditional approaches

5.5 DATA VERIFICATION

Test ID	Test Description	Expected Outcome
1	Verify that the datasets containing static and dynamic data are loaded correctly	Yes/No
2	Validate the correctness of the labels assigned to static and dynamic data samples (0 for safe, 1 for malicious)	Yes/No
3	Check for duplicates in the static and dynamic dataset	Yes/No
4	Verify that the machine learning model generates accurate predictions for the detection of Intrusion detection and IoT analytics based on features	Yes/No
5	Verify the frontend is successful for building the Static and Dynamic Data application on Cloud	Success/ Fail
6	If the model detects no malicious activity in the data sample	0 (Safe)
7	If the model detects high malicious activity in the data sample	1 (Malicious)

Table 5.1 – Data verification table

CHAPTER 6

RESULTS

6.1 RESEARCH FINDINGS

Through research and experimentation, the following key findings have been obtained in context of developing a comprehensive approach for static & dynamic data analytics:

- Data Preprocessing:
 - Automated techniques significantly improved CICIDS2017 dataset quality for machine learning.
- Model Selection and Optimization:
 - AutoML framework outperformed manual approaches.
 - Optimized LightGBM models consistently outperformed others for intrusion detection.
- Performance Evaluation:
 - Demonstrated AutoML models' effectiveness across metrics for static and dynamic analytics.
 - Outperformed traditional approaches, highlighting automation benefits.
- Real-world Deployment:
 - Exhibited scalability and efficient resource utilization.
 - Facilitated integration with existing systems through modular design.
- Interpretability:
 - Incorporated interpretable techniques & visualization tools for improved accuracy and performance metrics.

These research findings demonstrate the feasibility and effectiveness of leveraging AutoML techniques for improving the efficiency and accuracy of data analytics tasks in both static and dynamic environments. Now it enables more robust and reliable intrusion detection systems and optimized IoT data analytics applications.

6.2 RESULT ANALYSIS

Evaluation results clearly demonstrate the effectiveness of the proposed AutoML framework in developing high-performance models for both static and dynamic data analytics tasks. The LightGBM Classifier, optimized through AutoML technique, consistently outperformed other models across various evaluation metrics, showcasing their suitability for intrusion detection and IoT data analytics applications.

6.3 EVALUATION METRICS

Accuracy

Accuracy is the most basic metric, defined as the proportion of correctly categorized test instances to the total number of test instances. It is applicable to the majority of classification problems but is less useful when dealing with imbalanced datasets. Accuracy can be calculated by using True Positives (TPs), True Negatives (TNs), False Positives (FPs), and False Negatives (FNs):

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

Static Dataset

Model Accuracy	Percentage
LGBM Classifier	99.753 %
Random Forest	75.729 %
Naive Bayes	98.728 %
k-nearest neighbors (KNN)	92.475 %
KerasClassifier Model	99.753 %

Table 6.1 – Accuracy of Static Data

Dynamic Dataset

Model Accuracy	Percentage
LGBM Classifier	99.92%
Random Forest	99.839
Naive Bayes	70.184%
k-nearest neighbors (KNN)	99.280
KerasClassifier Model	92.475 %

Table 6.2 – Accuracy Dynamic Data

Precision

Precision is the metric used to quantify the correctness of classification. Precision indicates the ratio of correct positive classifications to expected positive classifications. The larger the proportion, the more accurate the model, indicating that it is more capable of correctly identifying the positive class.

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

Static Dataset

Model Precision	Percentage
LGBM Classifier	99.378 %
Random Forest	99.554 %
Naive Bayes	44.891 %
k-nearest neighbors(KNN)	95.584 %
KerasClassifier Model	73.378 %

Table 6.3 – Precision of Static Data

Dynamic Dataset

Model Precision	Percentage
LGBM Classifier	99.914
Random Forest	99.83%
Naive Bayes	99.875%
k-nearest neighbors(KNN)	99.744%
KerasClassifier Model	92.475 %

Table 6.4 – Precision of Dynamic Data

Recall

Recall is a measure of the percentage of accurately recognized positive instances to the total number of positive instances.

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

Static Dataset

Model Recall	Percentage
LGBM Classifier	99.788 %
Random Forest	99.753 %
Naive Bayes	97.244 %
k-nearest neighbors (KNN)	98.133 %
KerasClassifier Model	97.511 %

Table 6.5 – Recall value of static data

Dynamic Dataset

Model Recall	Percentage
LGBM Classifier	97.244 %
Random Forest	98.133 %
Naive Bayes	68.313%
k-nearest neighbors (KNN)	99.489%
KerasClassifier Model	92.475 %

Table 6.6 – Recall value of Dynamic data

F1 score

The F1 score is calculated as the harmonic mean of the Recall and Precision scores, therefore balancing their respective strengths.

$$\text{F1} = \frac{2 \times \text{TP}}{2 \times \text{TP} + \text{FP} + \text{FN}}$$

Static Dataset

Model F1 score	Percentage
LGBM Classifier	99.788 %
Random Forest	99.753 %
Naive Bayes	61.426 %
k-nearest neighbors (KNN)	96.842%
KerasClassifier Model	83.740%

Table 6.7 – F1-Score value of Static data

Dynamic Dataset

Model F1 score	Percentage
LGBM Classifier	99.957%
Random Forest	99.914
Naive Bayes	81.133%
k-nearest neighbors (KNN)	99.616%
KerasClassifier Model	92.475 %

Table 6.8 – F1-score value of Dynamic data

CHAPTER 7

CONCLUSION

The project successfully developed a comprehensive AutoML framework for static intrusion detection and dynamic IoT data analytics tasks. Automated techniques were implemented for data preprocessing, feature engineering, model selection, hyperparameter optimization, and concept drift handling. Extensive evaluations demonstrated the effectiveness of the optimized AutoML models, outperforming traditional approaches across various metrics. The framework exhibited scalability, efficient resource utilization, and seamless integration capabilities, enabling real-world deployment. Interpretability and explainability were enhanced through interpretable machine learning techniques and visualization tools, fostering trust and understanding. Additionally, a front-end web application with a simple and elegant user interface was developed and integrated with the ML models, facilitating easy access and interaction with the AutoML framework.

FUTURE WORK

Moving forward there are several ways to make our system better.

- Incorporating the latest advances in AutoML and machine learning algorithms to keep the framework up-to-date with cutting-edge techniques.
- Exploring advanced methods to handle high-dimensional and complex data structures commonly found in IoT and cybersecurity applications.
- Investigating the potential application of the AutoML framework to other domains and data analytics tasks beyond intrusion detection & IoT, expanding its versatility.
- Enhancing the front-end web application with additional features and visualizations to improve user experience and facilitate easier interaction with the AutoML models.

REFERENCES

- [1] Yang, L., & Shami, A. (2022). IoT data analytics in dynamic environments: From an automated machine learning perspective. *Engineering Applications of Artificial Intelligence*, 116, 105366.
- [2] Singh, A., Amutha, J., Nagar, J., Sharma, S., & Lee, C. C. (2022). AutoML-ID: Automated machine learning model for intrusion detection using wireless sensor network. *Scientific Reports*, 12(1), 9074.
- [3] Lindstedt, H. (2022). Methods for network intrusion detection : Evaluating rule-based methods and machine learning models on the CIC-IDS2017 dataset (Dissertation). Retrieved from <https://urn.kb.se/resolve?urn=urn:nbn:se:uu:diva-479347>
- [4] Garouani, M., Ahmad, A., Bouneffa, M., & Hamlich, M. (2022). AMLBID: an auto-explained automated machine learning tool for big industrial data. *SoftwareX*, 17, 100919.
- [5] He, Y., Lin, J., Liu, Z., Wang, H., Li, L. J., & Han, S. (2018). Amc: Automl for model compression and acceleration on mobile devices. In Proceedings of the European conference on computer vision (ECCV) (pp. 784-800).
- [6] He, X., Zhao, K., & Chu, X. (2021). AutoML: A survey of the state-of-the-art. *Knowledge-Based Systems*, 212, 106622.
- [7] Lee, J., Ahn, S., Kim, H., & Lee, J. R. (2022). Dynamic Hyperparameter Allocation under Time Constraints for Automated Machine Learning. *Intelligent Automation & Soft Computing*, 31(1).
- [8] Wever, M., Tornede, A., Mohr, F., & Hüllermeier, E. (2021). AutoML for multi-label classification: Overview and empirical evaluation. *IEEE transactions on pattern analysis and machine intelligence*, 43(9), 3037-3054.
- [9] Celik, B., Singh, P., & Vanschoren, J. (2023). Online automl: An adaptive automl framework for online learning. *Machine Learning*, 112(6), 1897-1921.
- [10] Zhang, S., Gong, C., Wu, L., Liu, X., & Zhou, M. (2023). AutoML-GPT: Automatic Machine Learning with GPT. *arXiv preprint arXiv:2305.02499*.