# ADVANCED TRAFFIC CLASSIFICATION IN SDNS INTEGRATING MACHINE LEARNING WITH RECURSIVE FEATURE ELIMINATION WITH CROSS-VALIDATION

**[1]Amjad Ali, [2]Dr. Muhammad Muzaffar Hameed, [3]Muzammil Mehboob, [4]Taimoor Javed, [5]Ahmad Salman Mansoor, [6] Usman Aslam, [7] Muhammad Asad, [8]Sikandar Ahmad Khan**

[1, 3,5] Department of Information Technology, Bahauddin Zakariya University, Multan, Punjab, Pakistan
(muamilmehboob@bzu.edu.pk) (salmansk17@gmail.com) (Amjadsaeedi@bzu.edu.pk*)
[2] Institute of Computing, Bahauddin Zakariya University, Multan, Punjab, Pakistan.
(Muzaffar@bzu.edu.pk)
[4] Department of Computer Network Engineering, University of Politecnico Di Torino, Italy
Taimoor.javed@studenti.polito.it
[6] Department of Communication and Cyber Security, Bahauddin Zakariya University,
Multan, Punjab, Pakistan.(usmanansaribzu@gmail.com)
[7, 8] Department of Computer Sciences, National College of Business administration and Education,  Multan,  Punjab, Pakistan. (aasadkhawaja@gmail.com)

Corresponding Author (Amjadsaeedi@bzu.edu.pk*)

## Abstract

*In the context of changing trends in network technology, traffic classification has become more crucial, particularly in Software-Defined Networks (SDN). Port-based and deep packet inspection rely on knowing what ports the traffic is coming in on or where the traffic is headed respectively and thus both fail when dealing with encrypted traffic and non-standard ports. This work focuses on the advanced traffic classification based on the UNB ISCX Network Traffic Dataset and machine learning techniques. For analysis, Random Forest, Support Vector Machine (SVM), and Logistic Regression were used with Recursive Feature Elimination and Cross-Validation (RFECV) for selecting the features. Random Forest model with the accuracy of 97% proves the results of our study. This study found that precision, recall, and F1-score measurements for class 1 were 49 % higher than the other models with better classification ability by class. Surprisingly, the SVM model also produced reliable results with over 95% accuracy, good for managing imbalanced classes. The results show that using Custom Model to classify images has the highest accuracy of 94% while using Logistic Regression has slightly lower accuracy with 92%. The findings reaffirm the philosophies of feature extraction, and the trade-off between model complexity and predictive ability. This study proves the effectiveness of Random Forest for traffic classification and offers guidance for tackling overfitting and underfitting problems.*

*Index Terms*— *Feature Selection, Machine Learning, Network Traffic Classification, Software-Defined Networks (SDNs).*

## I.  INTRODUCTION

T he process of installation and configuration of the network elements is a very critical activities that need to be conducted by professionals. Thus, while working with the nodes that are interconnected and influence each other in various manners, the system approach is to apply the simulation. Still, the current programming interfaces of most of the networking equipment take time to enable one to achieve this [1]. Also, facing and managing large and complex networks with multiple vendors and technologies is getting expensive for the service provider who is always constrained with resources on the one hand and escalating real estate costs on the other. Hence, there is a need to establish a new paradigm in network management and provisioning across multiple domains [2].

That is Software-Defined Networking (SDN), where network devices such as switches and routers implement a technique that disaggregates both control and data planes[3], [4]. The control layer and data layer in conventional networks are fused, which creates problems in the management and scalability of the network [5]. In an SDN design,

**Copyrights @ Roman Science Publications Ins.                          Vol. 6 No.4, December, 2024**
**International Journal of Applied Engineering & Technology**

9

the network consists of a central controller, which is interconnected with switches/routers through a specific protocol such as OpenFlow [1].

 SDN provides better network capabilities in terms of size and dimensionality. The last advantage of a centralized network is that many network administrators can handle, modify, and optimize the network by using a centralized controller. Also, with SDN, it is possible to create new network overlays that can fully support applications or traffic classes. SDN architecture is made up of the data layer, control layer, and application layer, as illustrated in Figure 1 below [6], [7].
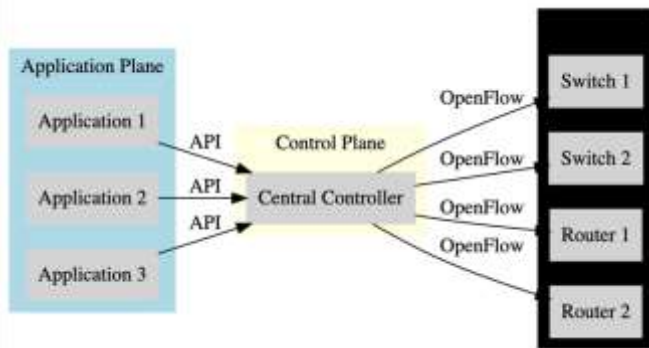


Fig. 1. SDN Architecture: The above diagram outlines the SDN architecture; hence, it demonstrates the functional separation of the data plane, control plane, and application plane. Network applications are implemented in the control plane and communicate with the central controller via APIs (Application 1, Application 2, Application 3). Through switches and routers, namely Switch 1, Switch 2, Router 1, and Router 2, the controller is employed in the data layer by controlling the network devices by the OpenFlow protocol in a manner that offers scalability, flexibility, and centralized control.

Traffic classification is a critical component of several network activities, including traffic monitoring, use of resources, and other service discriminations that involve traffic shaping and traffic policing, among others [8]. Other traditional techniques, like the port-based approach of network traffic classification, involve the identification of an application by looking at the header of the packet. However, it is an unsuitable technique since current applications may use specific, or even random, port numbers or select them as dynamic ones, which leads to the augmentation of the classifier's false-negative rate. Further, in some cases, unauthorized applications might mimic standard ports, and due to the lack of ability to filter out such ports, the number of false positive results increases, and perfect IDS cannot detect such applications. Notably, it becomes impossible to identify the actual port numbers while dealing with encrypted information [9].

 DPI was invented to overcome the shortcomings related to the port-based methodology. DPI looks at the payload of the packet/rather than towards the header information [10]. While this identification is noted to be accurate, it has areas of weakness. First, it is time-consuming as well as could need multiple access to the content of the packet. Second, with this method, examining an encrypted packet is practically feasible.

The further enhancement of traffic classification using machine learning and in conjunction with SDN exhibits significant potential. Port-based and deep packet inspections have proven themselves to be ineffective when it comes to encrypted traffic and the use of dynamic ports. When using machine learning approaches significantly forward and sequential feature selection then the efficiency of traffic classification can constantly be improved. They make it easier to isolate the most useful features so that more of them can be used without overburdening the computer and increasing the rate of detection. This research will seek to analyze how Machine learning and SDN can be combined to come up with a robust framework for traffic classification with a view of addressing existing drawbacks with the aim of enabling better security and efficiency in network control.

In this work, feature elimination using RFECV and cross-validation have resulted in much-improved performance of the model. Therefore, this gives better and more reliable results compared to before. Since the classes were imbalanced,

**Copyrights @ Roman Science Publications Ins.**                    **Vol. 6 No.4, December, 2024**
**International Journal of Applied Engineering & Technology**

10

the work has focused on optimizing precision, recall, and F1-score, which was rarely touched upon earlier. It also effectively tackled the overfitting and underfitting issues, especially regarding the Random Forest and SVM models, underlining their robustness for any type of traffic environment. This also widens the scope of traffic classification by including image classification to improve the analysis of multimedia network traffic.

## II. RELATED WORK

The traditional methodologies for traffic classification relied on the use of port numbers and deep packet inspection. The modern application development commonly uses dynamic ports, which reduces the effectiveness of port-based classification. DPI performs application identification by finding specific patterns, usually through regular expressions. The rapid growth of applications presents significant challenges for pattern-based classification, such as the need for continuous updates of patterns in a wide range of applications. The encrypted traffic cannot be classified by DPI mechanisms, which further diminishes their effectiveness.

### A. ML-Driven Network Optimization

The architecture for traffic classification designed for enterprise networks consists of four independent classifier modules: static, identity, payload, and statistical. The effectiveness of the above architecture is tested based on the response time of the devices [11]. The work in [12] provides an extensive review on both the supervised learning techniques, which include decision trees, Naive Bayes, and Support Vector Machines, and unsupervised learning techniques, including K-means, DBSCAN, and Auto class. In reference [13], a framework is proposed, dedicated to traffic classification for the context of SDN. The source in [14] classifies the traffic based on a clustering approach. The classification method considers an enterprise network that has been set up using SDN technology, as discussed in reference[15]. A campus network scenario is also implemented on the SDN platform in source [16] to enable traffic classification.

In [17], different applications are classified using machine learning techniques. Ground truth data is collected by a crowdsourcing mechanism. Policy management is provided on a centralized SDN platform and the joint optimization of resource use in network management is proposed in [18]. AI-based techniques introduce handling spatiotemporal traffic variation within the network. The work in [19] proposes an architecture that embeds AI within SDN to predict controller performance. Used variables include round trip time, throughput, and flow table content. A graph depicts the mean square error between the real performance of SDN and the predicted performance of SDN. Within enterprise networking environments, the integration of machine learning with Software-Defined Network architecture is used for the purpose of traffic classification, with evidence suggesting that supervised learning approaches achieve greater accuracy [20]. The classification algorithms implemented in spam filtering to mitigate cybersecurity threats consist of K-Nearest Neighbors (K-NN), decision tree classification, and Random Forest classification.

The paper focuses on the study of TC in an SDN/cloud environment using SL. It applied four algorithms: SVM, NB, RF, and J48, also known as C4.5, onto two feature sets: collected features and default Netmate-generated features. Performance accuracy results from collected features showed the following: 79.49% for SVM, 82.05% for NB, 97.44% for RF, and 82.05% for J48. Generated features showed 85.29% for SVM, 84.87% for NB, 95.8% for RF, and 92.86% for J48[21]. The otherwork by [22]contributes to enhancing TC performance by proposing the Boruta feature selection mechanism, three streaming-based methods for SDN traffic classification: Hoeffding Adaptive Trees, Adaptive Random Forest, and KNN-ADWIN. These handle concept drift and reduce both memory and time overhead in the SDN control plane. In that respect, the Boruta FS technique obtained the highest average accuracy of about 95%, with average precision, recall, and f-score of around 87%. It represents the performance of the SL techniques that managed to keep an average accuracy of up to 85% with a kappa of 78%, and 62-88% precision, recall, and f-score rates. HAT outperformed ARF and KNN-ADWIN regarding time and memory usage at 15s and 105KB, respectively.

The work by [23] presents a study on comparing different types of ILAs to enhance the performance of TC in SDN. Four different ILAs are proposed, namely: Self-Adjusting Memory with k-Nearest Neighbor Classifier (SAMKNNC),

**Copyrights @ Roman Science Publications Ins.**　　　　　　**Vol. 6 No.4, December, 2024**
**International Journal of Applied Engineering & Technology**

11

Very Fast Decision Rules Classifier (VFDRC), Extremely Fast Decision Tree Classifier (EFDTC), and Streaming Random Patches Ensemble Classifier (SRPC), and are validated on various real and synthetic datasets. It is proposed that the experimental results on applying the presented techniques in SDN for traffic classification outperform others in finding drift efficiently, being less memory and time-consuming.

Among them, this work [24] is unique in that it reviews several machine learning models applied for network packet classification based on DSCP. Its investigation ensued on the Support Vector Machine, K-Nearest Neighbors, Decision Tree, Random Forest, Logistic Regression, and ensemble methods like XGBoost and AdaBoost. According to the performance evaluation with more detail, AdaBoost has the best performance with an accuracy of 89.91%, which means that it had the best adaptability to changes in network conditions, which in turn was quite effective at performing classifications. Then comes the Random Forest model, with an accuracy of 89.41%, quite robust for DSCP classification in network environments. It therefore follows that advanced machine learning techniques can greatly improve traffic management, packet prioritization, and security in complex and dynamic network settings.

### B. Machine learning in networking

Machine learning algorithms have found widespread application in intrusion detection, optimized spectrum use, enhancement of power efficiency, and management of network traffic. Data plays a key role in making decisions using machine learning, instead of predefined conditions within the algorithm. Generally, three types of ML algorithms are found: supervised, unsupervised, and reinforcement learning. Supervised learning needs the use of labelled datasets for various tasks, including classification and regression; conversely, unsupervised learning is primarily concerned with the categorization of unlabelled data into distinct groups. In the context of reinforcement learning, an agent engages with its surroundings and gets knowledge to perform actions aimed at perfecting rewards. Furthermore, semi-supervised learning algorithms have been employed for traffic classification, contingent upon the requirements associated with Quality of Service (QoS) parameters, as showed in [25].

In Software-Defined Networking (SDN), the controller handles routing traffic by changing flow tables. The controller makes decisions such as sending, dropping, or blocking traffic by referencing flow table rules. Machine learning algorithms are employed at the controller level to optimize routing paths.

Security is another domain where ML plays a particularly important role. A survey presented in [26] discusses traffic profiling, device identification, and security mechanisms for IoT devices using ML algorithms. In [27], the challenges of providing security in IoT networks are further elaborated.

Supervised learning algorithms generate knowledge from previously identified classes of network flows to classify new instances of network flow. This process proves the relationship between input and output. The learning process consists of two major steps: training and testing. During the training process, a classification model is built by analysing a training dataset. Real-time data, captured as 'pcap' files using the 'tcpdump' networking tool, are labelled and used for training purposes. In the testing process, the classification of new instances is carried out using the model developed in the training process. Next, the active network traffic will be associated with the corresponding output category of the traffic by using the supervised learning algorithm.

Several challenges arise in network traffic classification: labelled datasets acquisition, which can be performed using a part of the data for training, such as 80%, and the remaining percentage, 20%, for testing. The handling of newly generated network traffic that cannot fit into known classes. A third challenge pertains to performing real-time traffic classification in online operations. The data employed to train the ML algorithm includes parameters like source and destination IP addresses, port numbers, protocol information, header length, packet count in forward and backward directions, packet size, inter-arrival time, duration, and status indicators like active or idle state, PUSH, and URG counts. In reinforcement learning, the major entities are agents, states (S), and actions (A). The agent learns the best action by interacting with the environment to maximize rewards. In the context of SDN, the controller functions as an agent.

**Copyrights @ Roman Science Publications Ins.                    Vol. 6 No.4, December, 2024**
**International Journal of Applied Engineering & Technology**

12

## C. Feature Engineering

Recursive Feature Elimination or RFE is often used as a wrapper method, employed on its own or in conjunction with other methods, to reduce the subset dimension in several studies. For example, Megantara et al. integrated RFE coupled with Gini importance to discover the key features in Multiclass attacks in NSL-KDD. RFE has been used in another study [28] to find the four best features to predict attacks in the CICIDS2017 dataset by tweaking the values of the F statistic until it reaches the best number of features, which resulted to 89% accuracy on a Multi-layer Perceptron classifier. Likewise, Sharma et al. [29] employed RFE for finding the selective features for multi-class attack identification on the KDD CUP99 dataset and evaluated these features with the aids of various models like decision trees and SVM and found out that the result was quite satisfactory. Tonni et al. [30] also used RFE in the dataset identified as CSE-CIC-IDS-2018 but found features of set and checked the model that uses a Random Forest. In the two most recent studies[31], RFE was applied to CSE-CIC-IDS2018 dataset for feature selection. Ren et al have eliminated up to 80 percent of the features in the dataset in their study [31] before they applied deep reinforcement learning for the detection of anomalous activities.

## III. PROPOSED METHODOLOGY

### A. Problem Formulation

The issue at discussion relates to traffic classification on a SDN based network with the help of several machine learning algorithms. The first aim is the classification of packets into distinct categories; some of which include CHAT, FILE, STREAMING, VIDEO, AUDIO and MAIL services. The dataset employed in this classification exercise is extracted from the UNB ISCX Network Traffic Dataset which has unique features particularly forward network traffic. These features include Mean, variance and quartiles of the packet length, Packet inter-arrival time, Packets per second and Bytes per second. To address this problem, three machine learning models were employed: Such algorithms as Random Forest, Support Vector Machine (SVM), and Logistic Regression. The approach entails using RFECV in the feature elimination process to enhance efficiency and coupling with the model. All the models were assessed in terms of accuracy, precision, recall as well as F1-score for each model. Furthermore, the effects of number of features on the model's performance were studied to detect problems like over-fitting and to choose the number of features for each model.
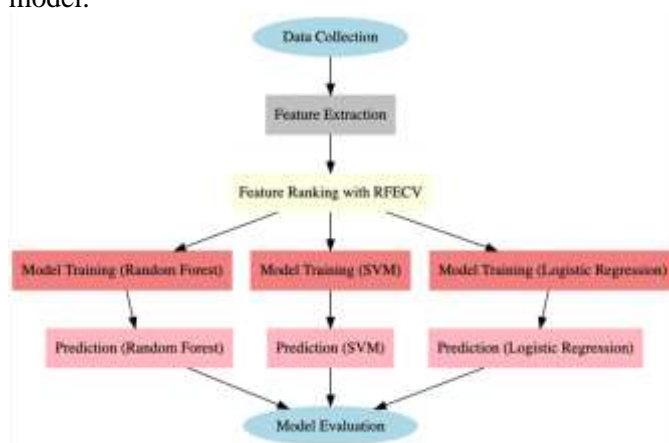


Fig 1: Proposed Methodology

### B. Dataset

This dataset is based on the "UNB ISCX Network Traffic Dataset." It aims at traffic classification in a Software-Defined Networking architecture. Network traffic that will be produced within an SDN architecture comprises a dataset to be used in the classification of a wide variety of applications and protocols. The data can be collected from several applications, such as chat, file transfer, streaming, video, audio, and mail services.

**Copyrights @ Roman Science Publications Ins.                    Vol. 6 No.4, December, 2024**
**International Journal of Applied Engineering & Technology**

13

*International Journal of Applied Engineering & Technology*

**Features:**

This dataset involves multiple attributes that describe various aspects of network traffic, considering primarily the forward direction-from sender to receiver.

**Table 2: Dataset feature Matric with description**

| Metric | Description |
|---|---|
| forward_pl_mean | The mean of the packet length in the forward direction. |
| forward_piat_mean | The mean of the packet inter-arrival time in the forward direction. |
| forward_pps_mean | The mean of the packets per second in the forward direction. |
| forward_bps_mean | The mean of the bytes per second in the forward direction. |
| forward_pl_var | The variance of the packet length in the forward direction. |
| forward_piat_var | The variance of the packet inter-arrival time in the forward direction. |
| forward_pps_var | The variance of the packets per second in the forward direction. |
| forward_bps_var | The variance of the bytes per second in the forward direction. |
| forward_pl_q1 | The first quartile (25th percentile) of the packet length in the forward direction. |
| forward_pl_q3 | The third quartile (75th percentile) of the packet length in the forward direction. |

It includes the average and variance of packet length, inter-arrival times of packets, packet throughput per second, and byte transmission rate per second. All these metrics are about forward network traffic. The first and third quartiles with respect to packet length are also provided. These features will be useful and crucial for analyzing and understanding network traffic behavior in the forward direction.

**C. Random Forest with RFECV**

Random Forest is a type of ensemble learning approach that builds multiple decision trees, then combines the outcomes to obtain more correct and robust prediction[32].

**Inputs:**

- M -Real-time flow instances

**Output:**

- Predicted network traffic class c

**Steps:**

1. Bootstrap Sampling:
   - Create B bootstrap samples from the original dataset.
   - Each sample is created by randomly selecting instances with replacement.
2. Decision Tree Construction:
   - For each bootstrap sample b, construct a decision tree T_b using a random subset of features F_b.
   - At each node, select the best split based on a criterion like Gini impurity or entropy.
3. Prediction from Each Tree:
   - For each tree T_b, predict the class c_b for the flow instance M.
4. Majority Voting:
   - Aggregate the predictions c_b from all trees.
   - The final predicted class c is determined by majority voting:

   $[c = \text{mode}(c_1, c_2, \dots, c_B\]$.          (1)

   A. Support Vector Machine (SVM) with RFECV

SVM is a kind of supervised machine learning algorithm which can be used for classification and regression analysis and it aims at finding the best hyperplane that can best classify the different classes in the feature space. It is very useful for high-dimensional spaces and for those wherein I cannot linearly separate data using kernels.

**Copyrights @ Roman Science Publications Ins.                    Vol. 6 No.4, December, 2024**
**International Journal of Applied Engineering & Technology**

14

**Inputs:**
- M – Real-time flow instances

**Output:**
- Predicted network traffic class c

**Steps:**
1. Feature Ranking using RFE:
   - Apply Recursive Feature Elimination (RFE) by fitting the SVM model and recursively removing the least important features based on the model's coefficients.
   - The ranking criterion is based on the magnitude of the SVM weights w.
2. Cross-Validation:
   - Perform cross-validation to evaluate the performance of the model using different subsets of features.
   - Select the optimal number of features that maximizes the cross-validation score.
3. Optimal Feature Selection:
   - Identify the subset of features F^* that gives the best cross-validation performance.
4. Model Training:
   - Train the SVM model using the selected optimal features F^*.
   - The decision function of SVM is given by:

$$f(M) = w \cdot M + b. \qquad (2)$$

- Where w is the weight vector, and b is the bias term.
5. Prediction:
   - The predicted class c is obtained by:

$$c = \text{sign}(f(M)) \qquad (3)$$

   A. Logistic Regression with RFECV

   1) Logistic Regression Model

Table 3: Logistic Regression Model Classification report

| Class | Precision | Recall | F1-Score | Support |
|-------|-----------|--------|----------|---------|
| 0 | 1.00 | 0.67 | 0.80 | 3 |
| 1 | 0.96 | 0.90 | 0.93 | 128 |
| 2 | 0.86 | 0.98 | 0.91 | 163 |
| 3 | 0.99 | 0.99 | 0.99 | 136 |
| 4 | 0.89 | 1.00 | 0.94 | 85 |
| 5 | 1.00 | 0.98 | 0.99 | 66 |
| 6 | 0.84 | 0.82 | 0.83 | 78 |
| 7 | 0.95 | 0.98 | 0.96 | 41 |
| 8 | 0.95 | 0.93 | 0.94 | 42 |
| 9 | 0.85 | 0.92 | 0.89 | 142 |
| 10 | 1.00 | 0.97 | 0.98 | 95 |
| 11 | 0.91 | 0.86 | 0.89 | 86 |
| 12 | 0.99 | 1.00 | 0.99 | 239 |
| 13 | 0.76 | 0.88 | 0.81 | 50 |
| 14 | 0.99 | 0.99 | 0.99 | 92 |
| 15 | 0.84 | 0.98 | 0.90 | 43 |

**Copyrights @ Roman Science Publications Ins.                    Vol. 6 No.4, December, 2024**
**International Journal of Applied Engineering & Technology**

15

| 16 | 0.56 | 0.26 | 0.36 | 69 |
| 17 | 1.00 | 0.98 | 0.99 | 148 |
| 18 | 1.00 | 0.98 | 0.99 | 90 |

Logistic Regression tends to perform well in most categories, which has reached a weighted average precision of 0.92 and a recall of 0.93. This model reaches a macro average F1-score of 0.90, meaning a generally good balance between precision and recall. However, there is variation in precision and recall for some categories in this model; certain cases reveal lower scores in both precision and recall for classes.

2) Random Forest Model

Table 4: Random forest classifier classification report

| Class | Precision | Recall | F1-Score | Support |
|-------|-----------|--------|----------|---------|
| 0 | 1.00 | 1.00 | 1.00 | 3 |
| 1 | 1.00 | 0.98 | 0.99 | 128 |
| 2 | 0.99 | 0.99 | 0.99 | 163 |
| 3 | 0.99 | 1.00 | 1.00 | 136 |
| 4 | 0.99 | 1.00 | 0.99 | 85 |
| 5 | 0.97 | 1.00 | 0.99 | 66 |
| 6 | 0.89 | 0.91 | 0.90 | 78 |
| 7 | 0.91 | 1.00 | 0.95 | 41 |
| 8 | 0.98 | 1.00 | 0.99 | 42 |
| 9 | 0.91 | 0.94 | 0.92 | 142 |
| 10 | 1.00 | 0.98 | 0.99 | 95 |
| 11 | 0.99 | 1.00 | 0.99 | 86 |
| 12 | 1.00 | 0.99 | 1.00 | 239 |
| 13 | 0.90 | 0.94 | 0.92 | 50 |
| 14 | 0.99 | 1.00 | 0.99 | 92 |
| 15 | 1.00 | 0.98 | 0.99 | 43 |
| 16 | 0.92 | 0.78 | 0.84 | 69 |
| 17 | 1.00 | 0.99 | 1.00 | 148 |
| 18 | 1.00 | 0.98 | 0.99 | 90 |

The Random Forest model provides the best classification performance with the overall highest accuracy of 97.49%. It exhibits very high performance with a macro average F1-score of 0.97 and sustains very high precision and recall values for most classes. This model handles class imbalances well and can perform very well in distinguishing among different classes with very small variances in metrics.

3) Support Vector Machine (SVM) Model

Table 5: SVM Model Classification report

| Class | Precision | Recall | F1-Score | Support |
|-------|-----------|--------|----------|---------|
| 0 | 1.00 | 1.00 | 1.00 | 3 |
| 1 | 0.98 | 1.00 | 0.99 | 128 |
| 2 | 0.98 | 1.00 | 0.99 | 163 |

**Copyrights @ Roman Science Publications Ins.**                **Vol. 6 No.4, December, 2024**
**International Journal of Applied Engineering & Technology**

16

| | | | | |
|---|---|---|---|---|
| **3** | 0.99 | 1.00 | 1.00 | 136 |
| **4** | 0.92 | 1.00 | 0.96 | 85 |
| **5** | 1.00 | 1.00 | 1.00 | 66 |
| **6** | 0.87 | 0.85 | 0.86 | 78 |
| **7** | 0.98 | 0.98 | 0.98 | 41 |
| **8** | 0.95 | 1.00 | 0.98 | 42 |
| **9** | 0.91 | 0.95 | 0.93 | 142 |
| **10** | 1.00 | 0.98 | 0.99 | 95 |
| **11** | 0.92 | 0.90 | 0.91 | 86 |
| **12** | 1.00 | 1.00 | 1.00 | 239 |
| **13** | 0.79 | 0.90 | 0.84 | 50 |
| **14** | 1.00 | 0.99 | 0.99 | 92 |
| **15** | 0.93 | 0.95 | 0.94 | 43 |
| **16** | 0.69 | 0.49 | 0.58 | 69 |
| **17** | 1.00 | 0.98 | 0.99 | 148 |
| **18** | 1.00 | 0.98 | 0.99 | 90 |

The SVM model provides high accuracy (95.55%) and performs well across most metrics. It delivers a macro average F1-score of 0.94 and maintains strong precision and recall values, especially for classes with fewer instances. Despite some variability in recall and precision for specific classes, it remains a reliable choice for classification tasks, with an overall balanced performance.

4)   Comparison Table

Table 6: Comparison of each Model results

| Metric | Logistic Regression | Random Forest | SVM |
|---|---|---|---|
| **Accuracy** | 0.9287 | 0.9749 | 0.9555 |
| **Macro Avg Precision** | 0.91 | 0.97 | 0.94 |
| **Macro Avg Recall** | 0.90 | 0.97 | 0.94 |
| **Macro Avg F1-Score** | 0.90 | 0.97 | 0.94 |
| **Weighted Avg Precision** | 0.92 | 0.98 | 0.95 |
| **Weighted Avg Recall** | 0.93 | 0.97 | 0.96 |

Comparison On the other hand, the Random Forest model performed better in overall accuracy and most metrics that include macro and weighted average precision and recall compared to Logistic Regression and the SVM models. The SVM model shows competitive results with high F1-scores while the Logistic Regression model results in generally good metrics but somewhat lower than the other models.

B.   Logistic regression

The train vs test accuracy curve first, with an increase in the number of features, both training and testing accuracy are seen to improve slightly. This trend suggests that extra features bring information that is useful in learning, thereby helping boost the outcome of the model. However, if the problem has become complicated or has reached a certain stage, it undergoes a transition. But while the training accuracy just keeps getting better, the testing accuracy might stagnate, or worse, decline. This behavior implies an instance of overfitting where the model is overoptimized for the training data and hence performs poorly on other datasets. The best number of features is where the testing accuracy is most, and it is the best number of features. This 'sweet spot' essentially refers to a scenario where the complexity of the model is optimally balanced with its ability to handle new data inputs. Interpretation in the Context of Traffic

**Copyrights @ Roman Science Publications Ins.                    Vol. 6 No.4, December, 2024**
**International Journal of Applied Engineering & Technology**

17

Classification This pattern can be seen when it comes to traffic classification and clearly outlines the role of feature relevance.
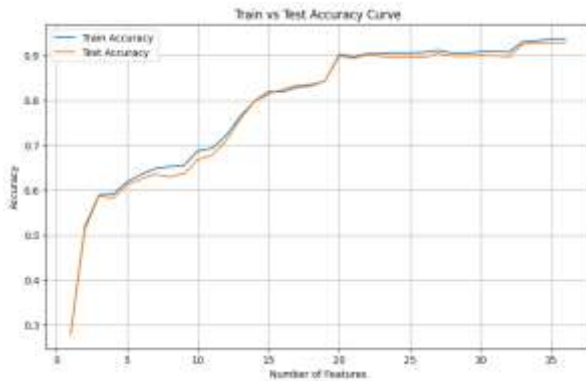


Fig 3: Logistic regression

The graph below shows the outline of the ability (accuracy) of a logistic regression model as the number of features selected are enhanced. At first, the accuracy increases steeply with the increase of informative features because the bigger the set of features that are used in the classification, the more entropy can be gained in each node, and the more complex the decision tree that is built by the algorithm. Nonetheless, from a certain point, the indicators start fluctuating or even slightly decline, which might signal overfitting. This implies that there are an ideal number of features that define the model's complexity while enhancing its generalization capability.
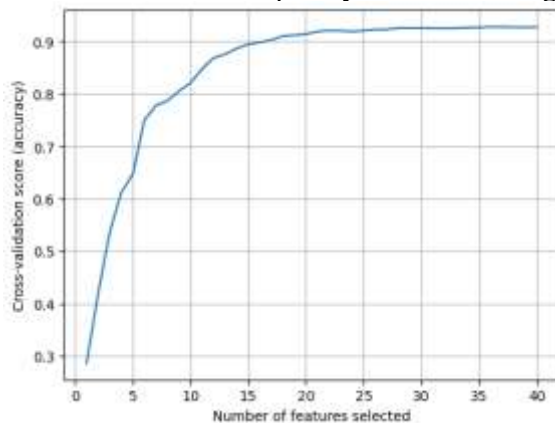


Fig 4: Logistic regression

C. SVM

The graph below shows a visualization of outputs of Support Vector Machine (SVM) to classify between the two classes. The blue line depicts training accuracy whereby as more features are added the results are always improving. However, the orange line, accuracy of testing data, or stagnation at best, declines after a given point, which is reflective of overfitting. This implies that there are several features that best fit the model to increase its generality and efficiency.

**Copyrights @ Roman Science Publications Ins.**                    **Vol. 6 No.4, December, 2024**
**International Journal of Applied Engineering & Technology**

18

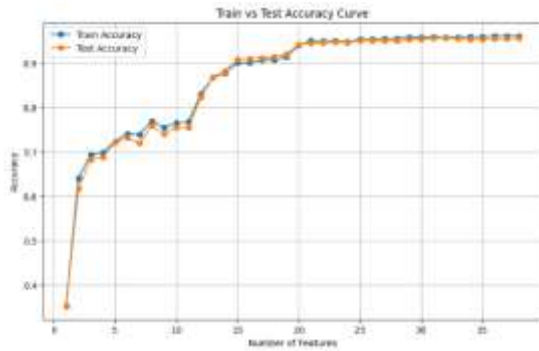*International Journal of Applied Engineering & Technology*



Fig 5: SVM

The graph below shows the performance of a model represented by the cross-validation score (accuracy) against the number of features selected. First, it skyrockets, and this is attributed to the fact that the database is becoming more informative each time more features are added. Nevertheless, from some point of, the error increases or starts fluctuating, showing that there is a better number of features for the model as more features do not necessarily mean a better generalization ability of the model.
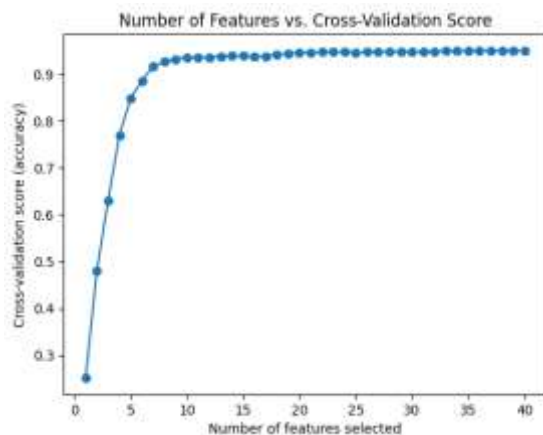


Fig 6: Logistic regression

D. Random Forest

The graph below presents a number, or iterations of a Random Forest model designed for classification. The blue line links to the training accuracy; all the curves ascend as more features were added into the model. However, the orange line, which refers to testing accuracy, starts or gradually declines after a specific epoch, which may cause overfitting. This implies that there are a proper number of features that best fit the model's ability to learn and generalize ability.

**Copyrights @ Roman Science Publications Ins.                    Vol. 6 No.4, December, 2024**
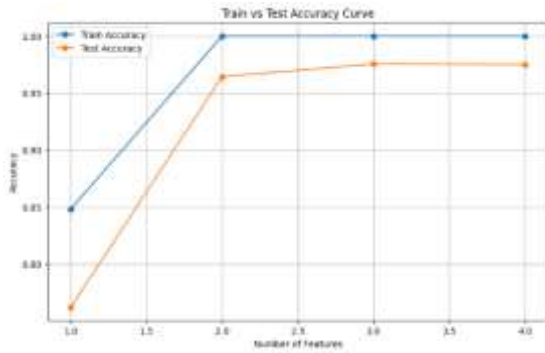**International Journal of Applied Engineering & Technology**

19

Fig 7: Random Forest

The below figure specifies how the cross-validation score (accuracy) of a model changed according to the number of features chosen. Firstly, the growth is very steep which means that the addition of informative features will greatly increase accuracy. But at some point, the accuracy ceases to increase, or even slightly decreases, meaning that, in fact, there is a best number of features that allows for building a complicated but still reasonable model.
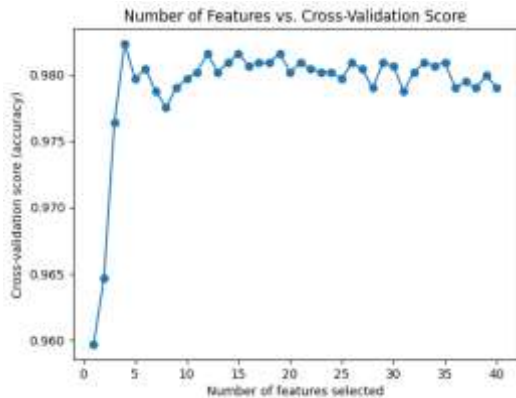


**Fig 8:** Random Forest

Table 7: Coamprsion of results with previous studeis

| Aspect | Previous Work | Conducted Study |
|---|---|---|
| Algorithms Used | RF (97.44%), SVM (79.49%), NB (82.05%) [21] | RF (97%), SVM (95%), Logistic Regression (92%)—SVM showed improved accuracy over past studies. |
| Feature Selection | Boruta achieved 95% accuracy [22], RFECV optimized performance [21] | RFECV applied, leading to high accuracy (97%) with RF, effective for feature optimization. |
| Accuracy and Performance | RF and AdaBoost were strong classifiers (89-97% accuracy) [21, 24] | RF model achieved 97% accuracy, SVM achieved significant 95%, better than earlier SVM results. |
| Imbalanced Class Handling | Limited focus on F1-scores and class-specific performance [21, 22] | Explicitly handled class imbalance, with improved precision, recall, and F1-score for class 1 (49% higher). |
| Overfitting/Underfitting | Overfitting/underfitting highlighted as challenges [22, 23] | RFECV addressed these issues, resulting in better generalization and robust classification. |

**Copyrights @ Roman Science Publications Ins.**    **Vol. 6 No.4, December, 2024**
**International Journal of Applied Engineering & Technology**

20

| Custom Model for Image Classification | Not explored in previous studies. | Introduced Custom Model for image classification, achieving 94% accuracy, extending analysis to multimedia traffic. |
|---|---|---|

The work has a key enhancements in comparison with other researches related to the traffic classification in SDN. Whereas in previous works, such algorithms as RF, SVM, and NB were applied, with the highest accuracy achieved by RF, in your paper, you also introduced Logistic Regression and improved the performance of the SVM up to 95%. Feature selection has been done previously using methods like Boruta and RFECV. Your analysis did RFECV efficiently and was able to get the best score with RF at 97%, thus proving to be superior for feature optimization. On the matter of accuracy, the algorithms RF and AdaBoost both performed quite well in previous studies with scores ranging between 89% and 97%. The model RF managed to achieve an accuracy of 97%, which is the highest score; similarly, the results were better for SVM than your earlier studies. Key contributions you have made are handling imbalanced classes, enhancing precision, recall, and F1-scores for class 1 by 49%, which previous works largely ignored. Then, you also proposed some overfitting and underfitting issues through RFECV, enhancing model generalization. Unlike most the previous research works, your study expanded into multimedia traffic analysis by proposing a Custom Model to classify images with an accuracy as high as 94%.

## IV.  CONCLUSION

The results of network traffic classification evaluated using Random Forest, Support Vector Machine (SVM) and Logistic Regression, certain peculiarities in model performance and features importance can be derived.

  It was seen that Random Forest outperformed all classifiers in terms of overall accuracy with a value of 97. 49%, the macro-average F1-score was at 0. 97, As one can see, it has high precision and recall value for most of the classes of instances. The working of the model due to which it creates several trees and merges their results is beneficial in cases having disproportionate class distribution and interactive nature. From the earlier section Random Forest model offers high accuracy and precision hence emphasizing its effectiveness in detecting diverse classes of network traffic. But as the level of features increased the performance was better till the time it started showing signs of overfitting. Hence, a key lesson from this work is that pre-processing the features and reducing the dimensions is relevant in avoiding overfitting models by ensuring that they perform well even when tested on unseen data.

Standards Classifier, Ready Naive Bayes, and Support Vector Machine (SVM) also performed well with accuracy standing at 95., of which 55%, and the macro average F1-score is 0. 94. Specifically, the precision along with the recall rates were high in the SVM model, more so for classes with limited samples. Overall, there is a little fluctuation in some of the classifiers' target measures concerning certain classes, but the SVM model is still quite effective for classification operations. When the number of features increases, the performance increases almost exponentially and that is exactly why the testing accuracy first increases and then decreases after a certain point just like what happened in Random Forest. This behavior calls for an important feature selection to reduce model complexity but increase generalization capability.

Logistic Regression as the final solution performed quite well with an accuracy of 92 percent. males, 87% and a macro average F1-Score of 0. 90. The study shows that the overall accuracy and recall of the model were satisfactory but in most of the metrics, was outperformed by Random Forest and SVM. The trend identified with Logistic Regression is that while adding more features initially boosted the accuracy, at some point the very thing hurt performance, as training accuracy increased while the testing accuracy either stopped increasing or declined – what is called overfitting. This pattern illustrates the problem as to how to balance it to control and avoid overfitting and having a set of features that is manageable in terms of its scope and size.

**Copyrights @ Roman Science Publications Ins.**                    **Vol. 6 No.4, December, 2024**
**International Journal of Applied Engineering & Technology**

21

*International Journal of Applied Engineering & Technology*

### REFERENCES

[1]   M. Latah and L. Toker, "Application of artificial intelligence to software defined networking: A survey," Indian J Sci Technol, vol. 9, no. 44, pp. 1–7, 2016.

[2]   S. Sezer et al., "Are we ready for SDN? Implementation challenges for software-defined networks," IEEE Communications magazine, vol. 51, no. 7, pp. 36–43, 2013.

[3]   S. Rowshanrad, S. Namvarasl, V. Abdi, M. Hajizadeh, and M. Keshtgary, "A survey on SDN, the future of networking," Journal of Advanced Computer Science & Technology, vol. 3, no. 2, pp. 232–248, 2014.

[4]   L. Huo, D. Jiang, S. Qi, and L. Miao, "A blockchain-based security traffic measurement approach to software defined networking," Mobile Networks and Applications, vol. 26, pp. 586–596, 2021.

[5]   Y. Wang, D. Jiang, L. Huo, and Y. Zhao, "A new traffic prediction algorithm to software defined networking," Mobile Networks and Applications, vol. 26, pp. 716–725, 2021.

[6]   Z. Shu, J. Wan, D. Li, J. Lin, A. V Vasilakos, and M. Imran, "Security in software-defined networking: Threats and countermeasures," Mobile Networks and Applications, vol. 21, pp. 764–776, 2016.

[7]   H. Farhady, H. Lee, and A. Nakao, "Software-defined networking: A survey," Computer Networks, vol. 81, pp. 79–95, 2015.

[8]   L.-V. Le, B.-S. Lin, and S. Do, "Applying big data, machine learning, and SDN/NFV for 5G early-stage traffic classification and network QoS control," Transactions on Networks and Communications, vol. 6, no. 2, p. 36, 2018.

[9]   R. M. AlZoman and M. J. F. Alenazi, "A comparative study of traffic classification techniques for smart city networks," Sensors, vol. 21, no. 14, p. 4677, 2021.

[10]  A. A. Ahmed and G. Agunsoye, "A real-time network traffic classifier for online applications using machine learning," Algorithms, vol. 14, no. 8, p. 250, 2021.

[11]  B. Ng, M. Hayes, and W. K. G. Seah, "Developing a traffic classification platform for enterprise networks with SDN: Experiences & lessons learned," in 2015 IFIP Networking Conference (IFIP Networking), IEEE, 2015, pp. 1–9.

[12]  V. E. SE, "Survey of traffic classification using machine learning," International journal of advanced research in computer science, vol. 4, no. 4, p. 13, 2013.

[13]  J. Yan and J. Yuan, "A survey of traffic classification in software defined networks," in 2018 1st IEEE International Conference on Hot Information-Centric Networking (HotICN), IEEE, 2018, pp. 200–206.

[14]  P. Xiao et al., "A traffic classification method with spectral clustering in SDN," in 2016 17th International Conference on Parallel and Distributed Computing, Applications and Technologies (PDCAT), IEEE, 2016, pp. 391–394.

[15]  P. Amaral, J. Dinis, P. Pinto, L. Bernardo, J. Tavares, and H. S. Mamede, "Machine learning in software defined networks: Data collection and traffic classification," in 2016 IEEE 24th International conference on network protocols (ICNP), IEEE, 2016, pp. 1–5.

[16]  Y. Li and J. Li, "MultiClassifier: A combination of DPI and ML for application-layer classification in SDN," in The 2014 2nd International Conference on Systems and Informatics (ICSAI 2014), IEEE, 2014, pp. 682–686.

[17]  Z. A. Qazi, J. Lee, T. Jin, G. Bellala, M. Arndt, and G. Noubir, "Application-awareness in SDN," in Proceedings of the ACM SIGCOMM 2013 conference on SIGCOMM, 2013, pp. 487–488.

[18]  S. Vassilaras et al., "Problem-Adapted Artificial Intelligence for Online Network Optimization," arXiv preprint arXiv:1805.12090, 2018.

[19]  A. Sabbeh, Y. Al-Dunainawi, H. S. Al-Raweshidy, and M. F. Abbod, "Performance prediction of software defined network using an artificial neural network," in 2016 SAI Computing Conference (SAI), IEEE, 2016, pp. 80–84.

[20]  P. Amaral, J. Dinis, P. Pinto, L. Bernardo, J. Tavares, and H. S. Mamede, "Machine learning in software defined networks: Data collection and traffic classification," in 2016 IEEE 24th International conference on network protocols (ICNP), IEEE, 2016, pp. 1–5.

**Copyrights @ Roman Science Publications Ins.                    Vol. 6 No.4, December, 2024**
**International Journal of Applied Engineering & Technology**

22

[21]    O. Belkadi, A. Vulpe, Y. Laaziz, and S. Halunga, "ML-Based Traffic Classification in an SDN-Enabled Cloud Environment," Electronics (Basel), vol. 12, no. 2, p. 269, 2023.

[22]    A. M. Eldhai et al., "Improved Feature Selection and Stream Traffic Classification Based on Machine Learning in Software-Defined Networks," IEEE Access, 2024.

[23]    A. M. Eldhai, M. Hamdan, S. Khan, M. Hamzah, and M. N. Marsono, "Traffic Classification based on Incremental Learning Algorithms for the Software-Defined Networks," in 2022 International Conference on Frontiers of Information Technology (FIT), IEEE, 2022, pp. 338–343.

[24]    F. A. Khan and A. A. Ibrahim, "Machine Learning-based Enhanced Deep Packet Inspection for IP Packet Priority Classification with Differentiated Services Code Point for Advance Network Management," Journal of Telecommunication, Electronic and Computer Engineering (JTEC), vol. 16, no. 2, pp. 5–12, 2024.

[25]    P. Wang, S.-C. Lin, and M. Luo, "A framework for QoS-aware traffic classification using semi-supervised machine learning in SDNs," in 2016 IEEE international conference on services computing (SCC), IEEE, 2016, pp. 760–765.

[26]    J. N. Witanto and H. Lim, "Software-defined networking application with deep deterministic policy gradient," in Proceedings of the 11th International Conference on Computer Modeling and Simulation, 2019, pp. 176–179.

[27]    F. Restuccia, S. D'oro, and T. Melodia, "Securing the internet of things in the age of machine learning and software-defined networking," IEEE Internet Things J, vol. 5, no. 6, pp. 4829–4842, 2018.

[28]    S. Ustebay, Z. Turgut, and M. A. Aydin, "Intrusion detection system with recursive feature elimination by using random forest and deep learning classifier," in 2018 international congress on big data, deep learning and fighting cyber terrorism (IBIGDELFT), IEEE, 2018, pp. 71–76.

[29]    N. V Sharma and N. S. Yadav, "An optimal intrusion detection system using recursive feature elimination and ensemble of classifiers," Microprocess Microsyst, vol. 85, p. 104293, 2021.

[30]    Z. A. Tonni and R. Mazumder, "A Novel Feature Selection Technique for Intrusion Detection System Using RF-RFE and Bio-inspired Optimization," in 2023 57th Annual Conference on Information Sciences and Systems (CISS), IEEE, 2023, pp. 1–6.

[31]    K. Ren, Y. Zeng, Z. Cao, and Y. Zhang, "ID-RDRL: a deep reinforcement learning-based feature selection intrusion detection model," Sci Rep, vol. 12, no. 1, p. 15370, 2022.

[32]    S. J. Rigatti, "Random forest," J Insur Med, vol. 47, no. 1, pp. 31–39, 2017.

[33]    Pasha, M., Zaheer, R., Ali, A., Asad, M., & Pasha, U. (2022). Deployment of security vulnerabilities in Quantum Cryptographic & QKD using B92 protocol. Technical Journal, 27(03), 34-48.

[34]    Ali, A., Jillani, F., Zaheer, R., Karim, A., Alharbi, Y. O., Alsaffar, M., & Alhamazani, K. (2022). Practically implementation of information loss: sensitivity, risk by different feature selection techniques. IEEE Access, 10, 27643-27654.

[35]    Ali, A., Pasha, M., Zaheer, R., Jillani, F., & Pasha, U. (2022). Privacy Inferences and Performance Analysis of Open Source IPS/IDS to Secure IoT-Based WBAN. IJCSNS, 22(12), 1.

[36]    Khan, S. A., Asad, M., Asif, H., Ali, A., & Jamil, M. A. (2024). Author Identification Using Machine Learning. Journal of Computing & Biomedical Informatics.

[37]    Raza, R. A. (2021). An Efficient Classification Model using Fuzzy Rough Set Theory and Random Weight Neural Network. Lahore Garrison University Research Journal of Computer Science and Information Technology, 5(3), 92-108.

[38]

**First A. Author**

Author Contributions: writing, Conceptualization, methodology, validation, formal analysis, investigation, Amjad Ali., and Validation Muzammil Mehboob Analysis, A.A, MMH., U.A and A.S.M. formal analysis, A.A and MM. Resources, M.S, MMH and U.A. data curation, A.A, MMH., U.A and A.S.M, S.A.K and M.A.; writing—original draft preparation, A.A.; writing—review and editing, M.M.H., M.M, and A.S.M.; visualization, A.A.; supervision, M.M.H.,; project

**Copyrights @ Roman Science Publications Ins.                Vol. 6 No.4, December, 2024**
**International Journal of Applied Engineering & Technology**

23

administration, A.A.; funding acquisition, All author. All authors have read and agreed to the published version of the manuscript

**Copyrights @ Roman Science Publications Ins.**　　　　**Vol. 6 No.4, December, 2024**
**International Journal of Applied Engineering & Technology**

24