

NETWORK INTRUSION DETECTION USING BIG DATA ANALYTICS: A PYSPARK AND HIVE APPROACH FOR UNSW-NB15

Daniel U. Okon

Article Info

Keywords: Cybersecurity; Network Intrusion Detection; Big Data Analytics; Apache Spark; Apache Hive; UNSW-NB15; Random Forest; Machine Learning

DOI

10.5281/zenodo.17135943

Abstract

Network intrusion detection remains a critical challenge as cyber threats continue to evolve in complexity and scale. This study investigates the application of big data analytics for intrusion detection using the UNSW-NB15 dataset. Apache Hive was used for large-scale querying and feature analysis, while PySpark was used for advanced analytics, including descriptive statistics, correlation, hypothesis testing, and dimensionality reduction. A RF classifier was developed and evaluated for both binary and multi-class intrusion detection tasks. The experimental results demonstrate a 99.99% accuracy in binary classification and 98.62% in multi-class classification, highlighting the effectiveness of combining Hive and PySpark for scalable intrusion detection. These findings underscore the importance of big data frameworks in strengthening cybersecurity defence systems.

1. Introduction

The proliferation of digital technology and internet connectivity has resulted in a huge increase in cyber risks, with network intrusions being the most serious and damaging threat faced by modern enterprises (Chen et al., 2019). Figure 1 depicts the amount of damage caused by cybercrime reported to the IC3 between 2001 and 2020 (Federal Bureau of Investigation, 2021). The annual loss of complaints referred to the IC3 in the most recent period was US\$4.2 billion. These alarming statistics highlight the crucial need to install effective intrusion detection systems (IDS) to protect important data and resources from various cyber-attacks. As described by Aminanto et al. (2017), by identifying malicious activity, providing early warnings, and permitting fast countermeasures, intrusion detection systems (IDS) serve a critical role in detecting and minimising the impacts of cyber-attacks. Due to the obvious reason that cyber threats are becoming more sophisticated, ongoing developments in IDS are required to keep up with the dynamic threat landscape, assuring the security of enterprises and individuals against the ever-increasing risk of cybercrime (Wang et al., 2018). Figure 2 displays intrusion detection, in which the firewall operates as an intrusion detector located among several networks,

filtering data that may constitute a threat (Alrawashdeh & Alhamid, 2020). Although it is feasible to improve network security, by employing the use of IDS, the ability of traditional intrusion detection systems (IDS) to handle and analyse large amounts of network traffic data in real time is limited. As a result, there is a growing demand for more resilient and adaptive network security solutions, which can be accomplished by using big data analytics methods (Mell, 2015).

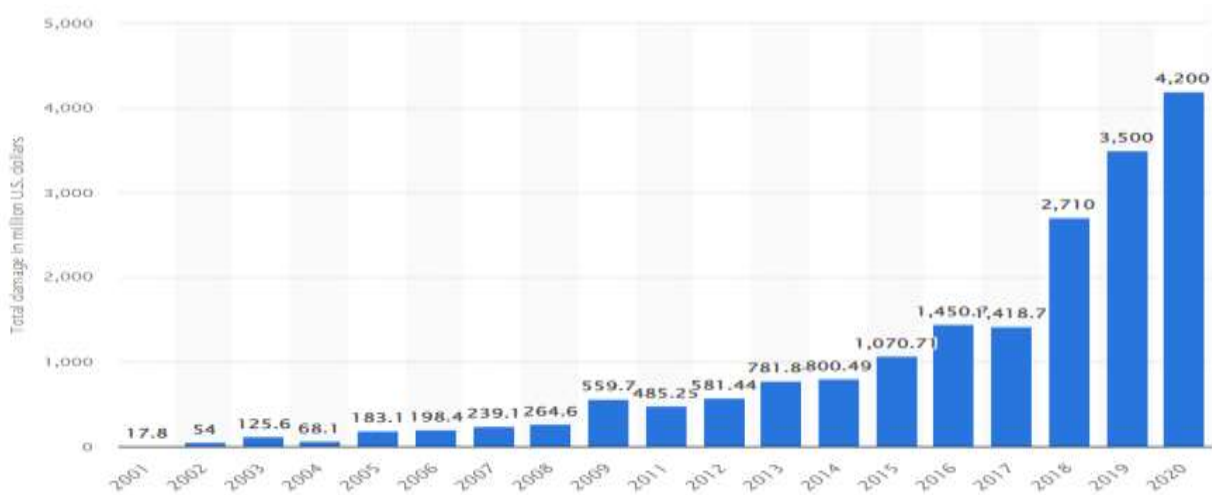


Figure 1 the Internet Crime Complaint Centre (IC3) from 2001 to 2020 Statistics

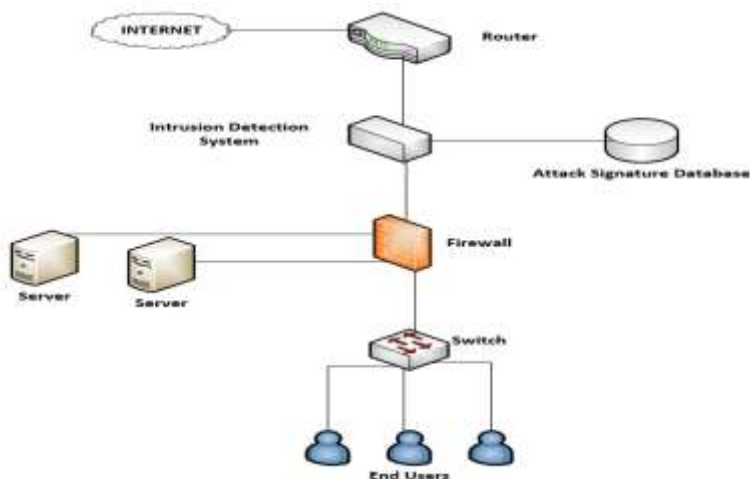


Figure 2 Network security system

Big data analytics have emerged as an indispensable element in diverse domains, encompassing network security. The introduction of distributed computing platforms, such as Apache Hadoop and Apache Spark, has altered the approach of researchers and professionals to the administration of large datasets, allowing for unequalled speed and efficacy in data processing and analysis. Furthermore, Apache Hive provides a solid foundation for running SQL-like queries on large datasets, easing the process of deriving important insights from big data. The objective of this study is to perform a thorough examination of network intrusion detection using big data analytics. The study will concentrate on the execution of Apache Spark and Apache Hive on the UNSW-NB15 dataset. This study aims to utilise various analytical techniques to extract meaningful patterns and associations from the dataset.

The obtained findings will be applied in constructing a classifier model for binary and multi-classification of the identified attack types.

1.1 Objectives

This study aims to explore the potential of big data analytics to improve the identification and categorisation of network intrusions. Our objective is to achieve the following goals:

- We employed Apache Hive for big data queries and analysis on the UNSW-NB15 dataset to generate insights for end users. This requires a thorough comprehension of the dataset, the development of Hive queries, the use of relevant visualisation tools, and the presentation of findings both numerically and graphically, as well as brief interpretations.
- Analyse and interpret UNSW data using at least four analytical methods, such as descriptive statistics, correlation, hypothesis testing, and density estimation, to acquire insights into the dataset's underlying patterns and relationships.
- Design and develop a binary and multi-classifier RF model for attack type classification, leveraging the insights gained from the analytical methods.
- Evaluation of both binary and multi-class classifier models using different evaluation metrics, such as accuracy, precision, recall, and F1-score.

1.2 Apache Spark

Apache Spark is a powerful distributed computing solution that analyses massive amounts of data efficiently by using in-memory caching and optimised query execution. Spark is an adaptable and scalable solution for big data processing, with a rich library suite that includes MLlib for machine learning, GraphX for graph processing, and Spark Streaming for real-time data processing. It is accessible to developers and academics worldwide because it is open-source, enabling continual progress and innovation in the field of distributed computing. It is an excellent platform for big data analytics in network intrusion detection because of its inherent fault tolerance, scalability, and extensibility (Chowdhury et al., 2011). Apache Spark was used to execute the analytical techniques and construct the classifier model. The MLlib library offered by Spark presents an extensive range of ML algorithms that enable diverse analytical methods for the UNSW-NB15 dataset and the creation of the RF classifier model.



Figure 3 Advantages of Spark

1.3 Apache Hive

Apache Hive is a data warehousing solution that has been developed on the Hadoop Distributed File System (HDFS). Facilitates the storage and administration of voluminous datasets in a distributed setting. Hive provides a query language, referred to as HiveQL, that bears resemblance to SQL (Apache Hive TM, 2011). This language enables users to perform intricate queries on vast datasets without the need for advanced programming proficiency. The combination of Apache Hive and Apache Spark improves the efficacy and expedites the execution of big data analytical assignments, making it an indispensable instrument in our scholarly inquiry. This study employs Apache Hive as a tool for large-scale data querying and analysis on the UNSW-NB15 dataset. Multiple Hive queries are executed to extract pertinent features, preprocess the data, and conduct exploratory data analysis. The insights derived from these queries serve as the foundation for the subsequent implementation of the analytical methods and the development of the classifier model.



Figure 4 Apache Hadoop Ecosystem

2. UNSW Dataset

Researchers from the University of New South Wales and the Australian Centre for Cyber Security created the UNSW-NB15 dataset to provide a thorough network traffic dataset that is tailored to assess the efficacy of network intrusion detection systems (Moustafa, n.d.). The dataset comprises more than 2.5 million records, each of which corresponds to a network connection. The records have 49 features, including the source and destination IP addresses, protocols, and payload characteristics, that encompass different aspects of the connection. The UNSW-NB15 dataset includes various attack types, including denial of service (DoS), remote code execution (RCE), and brute force attacks. The dataset includes a label that indicates the type of attack and a binary label that distinguishes between ordinary and malicious connections. The dataset's intricate and varied character makes it a strong option for evaluating big data analytics approaches in the network intrusion detection domain. Figure 4 shows the normal and attack distributions in the dataset. From the chart, we can see that the most common attack is the generic attack, followed by the exploit, fuzzer, and DDO attacks.

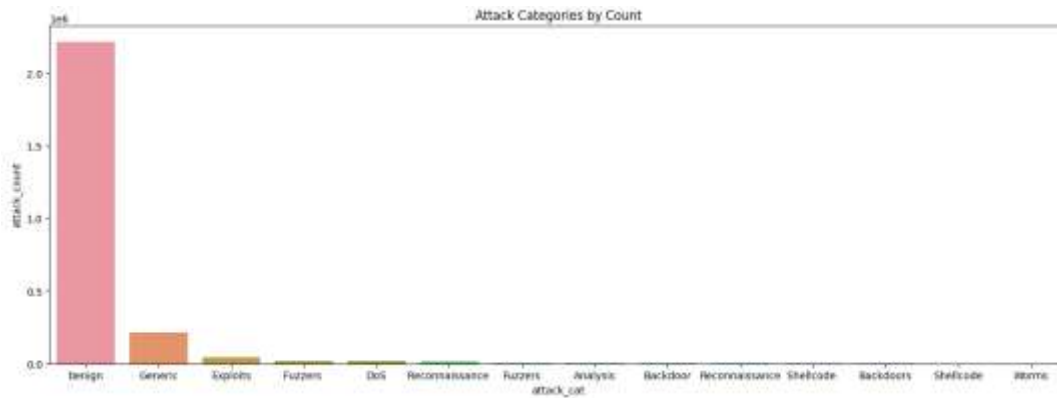


Figure 5 Distribution of normal and attack labels

2.1 Attack Types

The UNSW-NB15 dataset contains a diverse range of attack types classified into nine main classes. The aforementioned categories represent many expressions of cyber dangers, providing a thorough portrayal of the challenges experienced by computer networks' intrusion detection systems. The following sections provide a concise summary of each assault type found in the dataset:

- **Fuzzers:** Fuzzers are attacks that produce and transmit random, flawed, or invalid inputs to a target system to cause unexpected behaviour, crashes, or security vulnerabilities.
- **Analysis:** Analysis attacks focus on identifying and exploiting weaknesses in target systems. In these assaults, scanning the network for open ports, services, or applications, as well as testing systems for flaws, is common. Examples of analysis attacks include port scanning, OS fingerprinting, and vulnerability scanning.
- **DoS:** Denial of Service attacks seek to impair the availability of a network, service, or application by overwhelming it with excessive traffic or resource usage. DoS attacks can be carried out via various techniques, including SYN flood, ICMP flood, and HTTP flood. The primary purpose of these assaults is to make the target system inaccessible or ineffective, causing inconvenience to its users.
- **Exploit:** Exploit attacks employ known weaknesses in software, hardware, or protocols to gain unauthorised access, control, or privileges within a target system. These attacks frequently rely on the availability of public or private exploit code and can result in arbitrary code execution, data exfiltration, or system penetration.
- **Generic:** Generic attacks cover various approaches that do not fit into any particular category. These attacks may include network traffic manipulation, unauthorised data access, or the exploitation of legitimate system capabilities. Generic attacks can provide particular challenges to network intrusion detection systems due to their diversity.
- **Reconnaissance:** Reconnaissance attacks gather information about a target network, system, or organisation in order to discover possible targets for future assaults. Passive techniques, such as eavesdropping on network traffic, and active approaches, such as DNS enumeration or social engineering, can be used in these assaults.

Comprehending the various attack types and their unique characteristics is crucial in the creation of effective NIDS. This enables the development of models that can accurately differentiate between benign and malicious network traffic patterns and classify attacks according to their respective categories. This study explores the diverse attack categories present in the UNSW-NB15 dataset to establish a RF framework for both binary and

multi-classification objectives. This study aims to enhance the ability to detect and classify network intrusions for rapid identification.

3. Big Data Query and Analysis using Apache Hive

Big data analytics are critical in identifying and mitigating cyber threats, particularly in the domain of network intrusion detection. This study presents the results of our data analysis on the UNSW-NB15 dataset using Apache Hive queries. The goal of this study is to convert unprocessed data into useful information for end recipients and to draw meaningful conclusions from the data. We use Hive queries to perform various analyses and provide the results in both numerical and graphical representations.

3.1 Query 1: Top 5 attack categories by count

The primary objective of our preliminary investigation is to determine the five most common attack categories based on their frequency. To accomplish this task, we execute the following query on Apache Hive:

Query 1: Top 5 attack categories by count

```
query1 = """
SELECT ATTACK_CAT, COUNT(*) as ATTACK_COUNT
FROM unswnb15
WHERE attack_cat != 'benign'
GROUP BY attack_cat
ORDER BY attack_count DESC
LIMIT 5
"""
result1 = hive_context.sql(query1).toPandas()
```

Figure 6 Query: Top 5 Code Attacks

This query selects the attack category and enumerates the occurrences of each attack type. It also identifies the attack category and records the number of instances of each attack type. We deliberately ignore “benign” traffic in our query because it does not represent harmful activities. The query then groups the results by attack category and displays them in descending order according to the attack count. The output only includes the top five results.

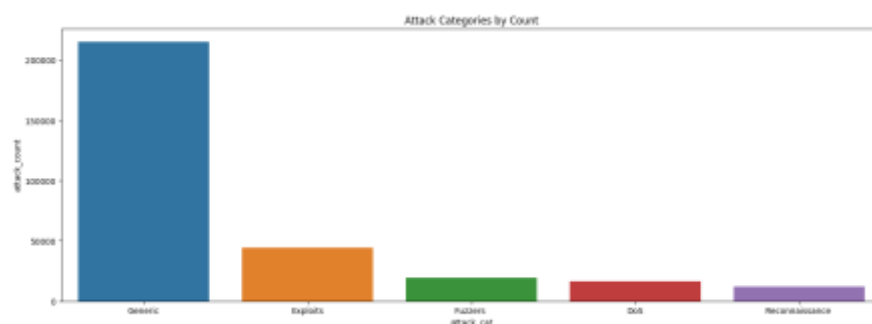


Figure 7 Top 5 Attack Labels

According to the results of this query, the “generic” attack category is the most common, accounting for 215,481 attacks. Following this are “Exploits” (44,525), “Fuzzers” (19,195), “DoS” (16,353), and “Reconnaissance” (12,228). Figure 8 shows a bar chart created with Python’s Seaborn package to visually depict the findings. In the UNSW-NB15 dataset, the “Generic” attack category is by far the most common type of attack.

3.2 Query 2: Top 10 source group with IP involved in the attacks

Our analysis aims to identify and pinpoint the top 10 source group with IP implicated in attacks by analysing the frequency of the observed source IP addresses. The following query was implemented to accomplish this task:

Query 2: Top 10 source IPs involved in attacks

```

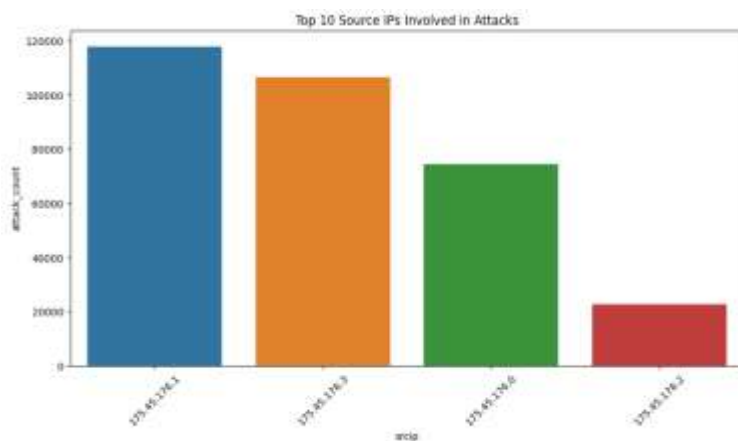
[] query2 = """
SELECT srcip, COUNT(*) as attack_count
FROM unsw_hive
WHERE Label = 1
GROUP BY srcip
ORDER BY attack_count DESC
LIMIT 10
"""

result2 = hive_context.sql(query2).toPandas()

```

Figure 8 Query: Top 10 source IPs involved in attacks

This query retrieves the source IP addresses and counts each occurrence. By filtering out non-attack traffic with the label field, we ensure that only malicious activities are considered. The results are then grouped by source IP, sorted by attack count in descending order, and limited to the top IP entries. The results of this query show that the source IP '175.45.176.1' is the most common, with 117,908 assaults. '175.45.176.3' (106,418), '175.45.176.0' (74,279), and '175.45.176.2' (22,678) follow.

*Figure 9 Top 10 source IPs involved in attacks*

To visually present these findings, a bar chart created with the Seaborn library in Python is employed, as displayed in Figure 9. The graph highlights that the top three source group with IP contribute to a significantly higher number of attacks than the remaining source group with IP in the dataset. This information is crucial in understanding the behaviour of the most active attackers and devising strategies to mitigate their occurrences.

3.3 Query 3: Attack distribution according to protocol

This query examines the attack distribution by protocol. This query focuses on the protocol column and counts the number of attacks for each protocol by excluding malicious traffic with a Label of 1 (i.e., malicious traffic with a Label of 1). The results are then categorised by protocol and based on the number of attacks, ordered in descending order. Finally, the output is restricted to the top 25 protocols with the most attacks.

Query 3: Attacks distribution by protocol

```

query3 = """
SELECT proto, COUNT(*) as attack_count
FROM unsw_hive
WHERE Label = 1
GROUP BY proto
ORDER BY attack_count DESC Limit 25
"""

result3 = hive_context.sql(query3).toPandas()

```

Figure 10 Query: Attack distribution by protocol

UDP is the most targeted protocol, with 223,750 reported attacks, followed by TCP (58,184) and UNAS (16,202). These findings are visually represented in Figure 11, which displays a bar plot of the attack distribution by protocol. The results of this analysis provide valuable insights into the behaviour of attackers targeting specific network protocols. These data can be used to inform the design and setup of intrusion detection systems, which can be tailored to detect and mitigate assaults on various protocols. Furthermore, this analysis lays the groundwork for future research into the vulnerabilities and hazards associated with various network protocols, allowing for the creation of more effective security measures.

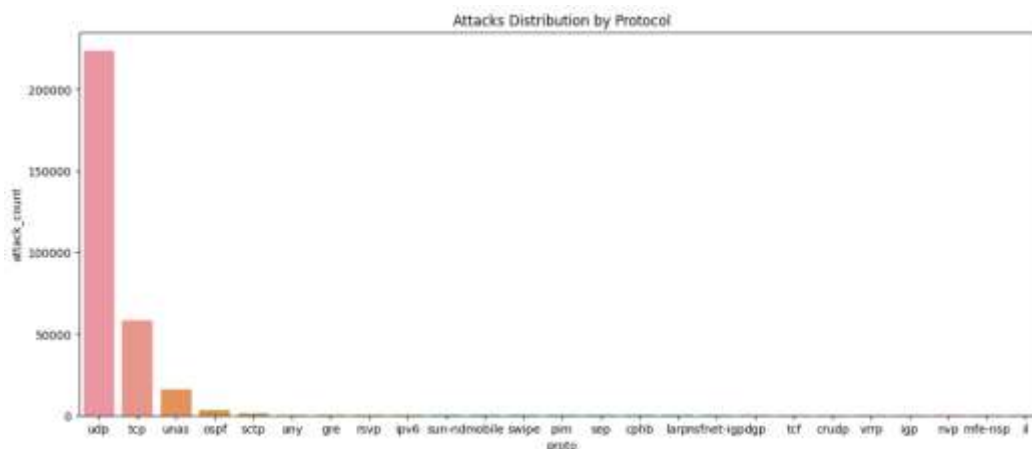


Figure 11 Attack distribution by protocol

3.4 Query 4: Average duration and total byte count per category of attack

This analysis aims to ascertain the average duration and total byte count of attacks in the UNSW-NB15 dataset for each attack category, as shown in Figure 12. This query selects the attack type and calculates the average duration and total byte count for each category while eliminating benign traffic using a WHERE clause. The findings are then categorised by attack category, sorted in descending order based on total byte count, and displayed in tabular form.

Query 4: Average duration and total byte count per attack category

```

[ ] query4 = """
SELECT attack_cat, AVG(dur) as avg_duration, SUM(sbytes + dbytes) as total_bytes
FROM unsw_hive
WHERE attack_cat != 'benign'
GROUP BY attack_cat
ORDER BY total_bytes DESC
"""

result4 = hive_context.sql(query4).toPandas()

```

Figure 12 Query: Average duration and total byte count per category of attack

The bar and line plots are used together to visualise the results, as shown in Figure 13. The graph shows the overall byte count for each assault category, with the average duration represented by a line plot. The entire byte count is shown on the left y-axis, and the average time is shown on the right y-axes, respectively. According to the findings, the attack categories with the highest total byte counts are "Exploits," "DoS," and "Generic." "Fuzzers," "Exploits," and "DoS" have the longest average durations. These findings have significant implications for designing and implementing security measures in each category to detect and mitigate attacks.

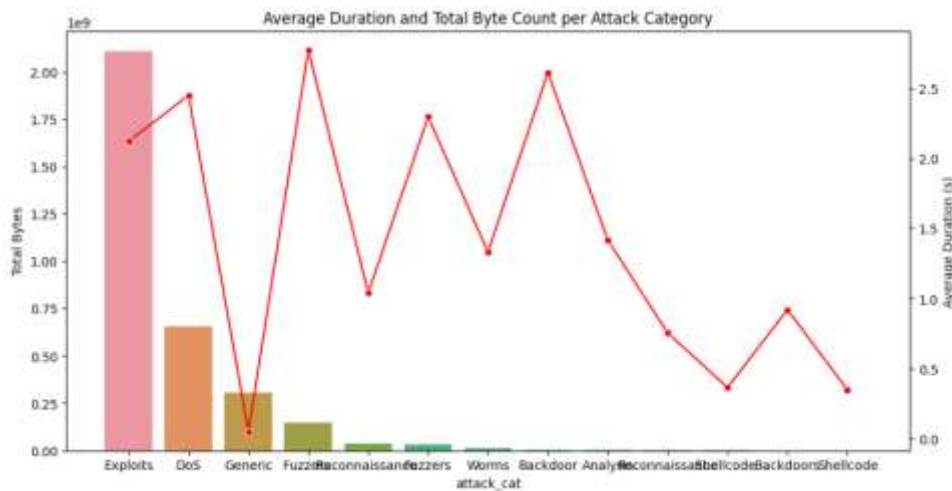


Figure 13 Average duration and total byte count per attack category

3.5 Query 5: Proportion of attacks with a non-zero response body length, grouped by the attack category

The following query was created to analyse the proportion of attacks in the UNSW-NB15 dataset with non-zero response body length for each attack category. The query selects the attack category and computes the total number of assaults, number of attacks with a non-zero response body length, and percentage of attacks with a non-zero response body length for each category, omitting the 'benign' traffic category. The findings are then tabulated and organised by attack category and ranked in descending order according to the percentage of attacks with a non-zero response body length.

Query 5: Proportion of attacks with non-zero response body length, grouped by attack category

```
[20] query5 = """
      SELECT attack_cat,
             COUNT(*) as total_attacks,
             SUM(CASE WHEN res_bdy_len > 0 THEN 1 ELSE 0 END) as attacks_with_body,
             (SUM(CASE WHEN res_bdy_len > 0 THEN 1 ELSE 0 END) / COUNT(*)) * 100 as percentage_with_body
      FROM unsw_hive
      WHERE attack_cat != 'benign'
      GROUP BY attack_cat
      ORDER BY percentage_with_body DESC
      """

      result5 = hive_context.sql(query5).toPandas()
```

Figure 14 Query: Proportion of attacks with a non-zero response body length, grouped by the attack category

The analysis attack type had the highest percentage of attacks with a non-zero response body length, followed by “Worms” and “Exploits.” This implies that the presence of a non-zero response body length is more connected with certain types of attacks than others. On the other hand, the “Fuzzers,” “Reconnaissance,” “Shellcode,” and “Backdoors” attack categories have a very low percentage of attacks with a non-zero response body length. These findings can be used to create more effective intrusion detection and prevention systems that consider each attack category’s unique characteristics, allowing for improved cyber threat defence.

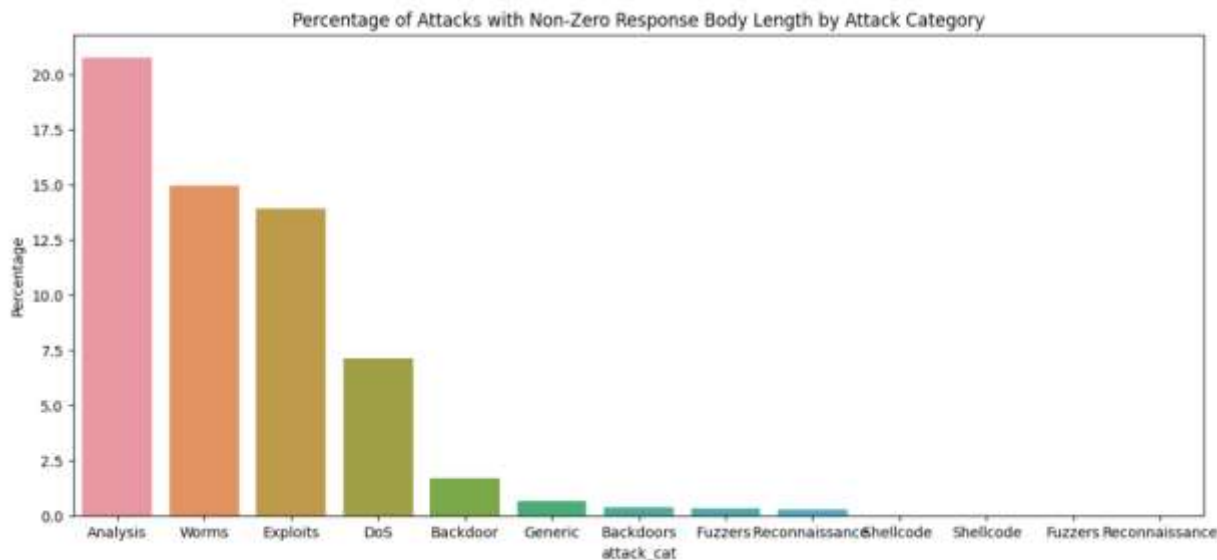


Figure 15 Proportion of attacks with non-zero response body length, grouped by attack category

4. Advanced Analytics Using PySpark

Big data analysis and interpretation are critical components of data science and are used to derive important insights from massive amounts of data. To gain a better understanding of the data, various analytical techniques, such as descriptive statistics, correlation, hypothesis testing, and density estimation, must be employed. PySpark was used to analyse and interpret the UNSW-NB15 dataset using at least four analytical methodologies.

4.1 Descriptive Summary

A descriptive analysis is a statistical analysis in which the main features of a dataset are summarised and described. The purpose of descriptive analysis is to provide a high-level overview of the dataset, highlight potential flaws or patterns, and inform subsequent studies. The initial step in descriptive analysis is to identify which of the dataset's numeric columns may be studied. In this task, 39 columns were chosen, comprising information such as the source and destination ports, connection length, and amount of bytes exchanged. Summary statistics, such as the count, mean, standard deviation, minimum, and maximum values, have been computed using these columns. The summary statistics were produced in PySpark using the describe() method, which generates a DataFrame containing the summary statistics for each numeric column. The toPandas() method was used to convert the Spark DataFrame to a Pandas DataFrame for easier display and analysis. The results showed that the 'dur' variable has a mean of 0.6588 s and a maximum of 8786.6377 s, indicating that the dataset has some unusually long durations. The "sbytes" variable has a mean of 4340.07 bytes and a maximum of 14355774 bytes, indicating various packet sizes. The "ct_srv_src" and "ct_srv_dst" variables have a maximum value of 67, which suggests that a few sources and destinations are responsible for a large number of connections. Overall, by analysing the summary statistics, we can obtain an idea of the range of values for each numeric column and the distribution of values across the dataset.

	summary	sport	dport	dur	sbytes	dbytes	att1	att2	class	dclass	...	is_ftp_login
0	count	2539739	2539739	2539739	2539739	2539739	2539739	2539739	2539739	2539739	...	2539739
1	mean	30536.53104503206	11235.096789085767	1.6588634003800393	4340.072263330996	36432.01132478574	62.78149802097214	30.779443978384	5.184547223159545	16.331423819534212	...	0.01735335796319878
2	stddev	20441.216752342224	18438.20083577195	13.925767633289283	56489.39812206232	181105.38488024756	74.62689956848524	42.851922109518734	22.51836838779722	56.597886163285116	...	0.12348518417617913
3	min	0	0	0.0	0	0	0	0	0	0	...	0
4	max	65535	65535	8786.6377	14355774	14657531	255	254	5319	5507	...	4

5 rows x 13 columns

Figure 16 Descriptive analysis results

4.2 Hypothesis testing

Hypothesis testing is a statistical tool for determining whether observed differences in sample data are real or random variation. In this task, we investigate the possibility of a statistically significant difference in the mean value of the "dur" characteristic between normal and attack traffic. Following the null and alternative hypotheses for this task:

- The null hypothesis (H0): The mean "dur" value is the same for both normal and attack traffic.
- Alternative hypothesis (HA): The mean "dur" value is different for normal traffic and attack traffic.

The "dur" values for both regular and intrusive network traffic are collected from the dataset to test the hypothesis. The t-statistic and p-value were then computed using a two-sample t-test. The t-test assumes that the data follows a normal distribution and that both groups' variances are comparable.

Hypothesis Testing

- Applying Hypothesis Testing To check if the mean dur value is different between normal traffic and attack traffic
- Null hypothesis (H0): The mean 'dur' value is the same for both normal traffic and attack traffic.
- Alternative hypothesis (HA): The mean 'dur' value is different for normal traffic and attack traffic.

```
[ ] normal_traffic = data.filter(data.Label == 0).select('dur').rdd.flatMap(lambda x: x).collect()
    attack_traffic = data.filter(data.Label == 1).select('dur').rdd.flatMap(lambda x: x).collect()

    t_statistic, p_value = ttest_ind(normal_traffic, attack_traffic)

[ ] print("P Value is:",p_value)

P Value is: 0.0018154559048992132
```

Figure 17 Hypothesis testing code

The null hypothesis is rejected when the p-value is less than the preset significance level, commonly set at 0.05. Boxplot of 'dur' attribute for both normal and attack traffic was created to visualise the results of the hypothesis. Figure 18 illustrates the median, interquartile range, and outliers for each group, as well as the p-value. In this case, the p-value was 0.0018, which is less than the significance level. As a result, the null hypothesis was rejected, indicating that normal and attack traffic have different means, as corroborated by the median 'dur' value of the plot, which is greater for normal traffic than for attack traffic.

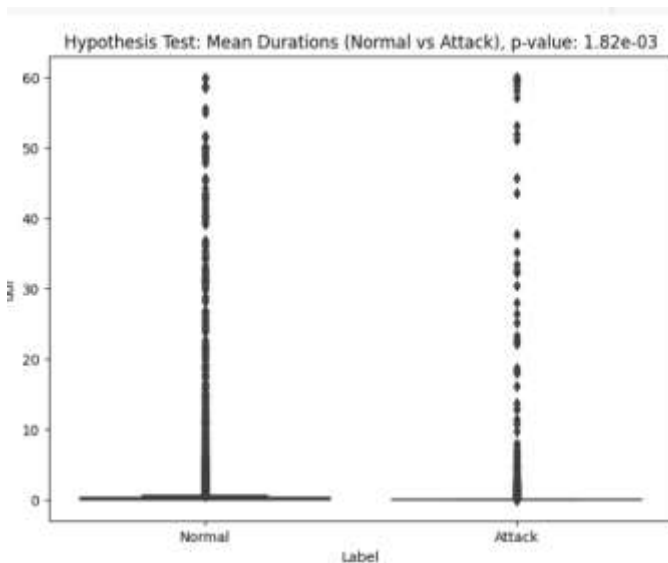


Figure 18 Hypothesis test boxplot

4.3 Correlation

The Pearson correlation was implemented to understand the relationship between numerical features in the UNSW-NB15 dataset. Pearson's correlation coefficient is a statistical measure that quantifies the magnitude and direction of a linear association between two variables (Turney , 2022). The values for Pearson correlation coefficient ranges from -1 to 1, where a value of -1 signifies a complete negative correlation, indicating that as one variable increases, the other decreases. A value of 1 indicates a perfect positive correlation, indicating that both variables increase or decrease together, and a value of 0 indicates no correlation. The categorical variables were eliminated from the dataset, and the correlation between the numerical features was assessed to identify

highly correlated features. In the code, we used the drop() method to eliminate categorical variables from the dataset and the VectorAssembler() function to create a correlation assembler. The Correlation.corr() function was then used to calculate the Pearson correlation coefficient and the results were recorded in a data frame. Finally, we printed out highly correlated pairs of characteristics (correlation coefficient > 0.8 or -0.8) as shown in Figure 22.

```
columns_to_drop = ['dstip', 'proto', 'state', 'service', 'attack_cat', 'srcip']
correlation_matrix = data.drop(*columns_to_drop)

## Creating Correlation Assembler
vector_col = "corr_features"
assembler = VectorAssembler(inputCols=correlation_matrix.columns, outputCol=vector_col)
vect_cor = assembler.transform(correlation_matrix).select(vector_col)

## Correlation Matrix
matrix = Correlation.corr(vect_cor, vector_col)

## Applying Pearson Correlation
correlation = Correlation.corr(vect_cor, vector_col, "pearson").collect()[0][0]

rows = correlation.toArray().tolist()
cor_df = spark.createDataFrame(rows, correlation_matrix.columns)

## Visualize
cor_df_pd = cor_df.toPandas()
cor_df.show()
```

Figure 19 Correlation coefficient

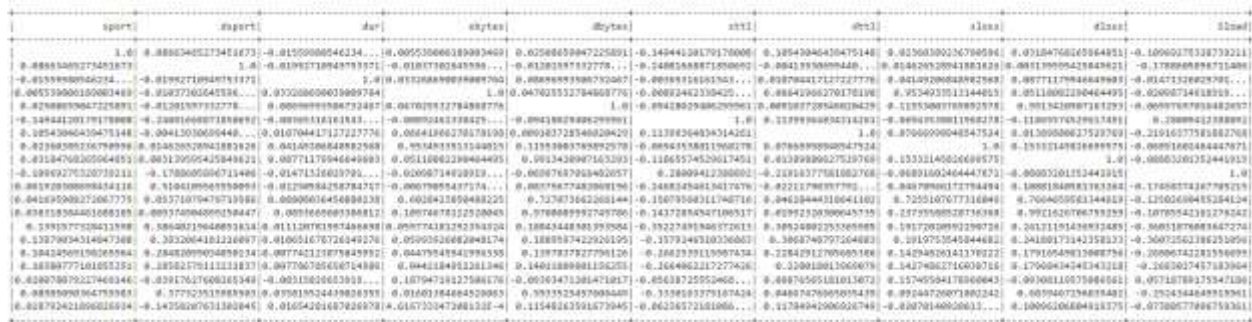


Figure 20 Correlation Graph of the

```

Highly correlated pairs (correlation coefficient > 0.8 or < -0.8):
sloss - sbytes: 0.9534833513144815
dloss - dbytes: 0.9913428987163293
Dpkts - dbytes: 0.9788889992745786
Dpkts - dloss: 0.9921626786759259
Dpkts - Spkts: 0.8225776151624578
dwin - swin: 0.9972869954888567
itime - stime: 0.9999999997797463
synack - tcpett: 0.9313053022033396
ackdat - tcpett: 0.9188987922483882
is_sm_ips_ports - Sintpkt: 0.8866298486949893
ct_state_ttl - sttl: 0.9868746751965445
ct_ftp_cwd - is_ftp_login: 0.8496322369558772
ct_srv_dst - ct_srv_src: 0.9567312281508591
ct_dst_ltm - ct_srv_src: 0.8366557815433892
ct_dst_ltm - ct_srv_dst: 0.8598876413790725
ct_src_ltm - ct_srv_src: 0.8405144171558931
ct_src_ltm - ct_srv_dst: 0.8246727417893495
ct_src_ltm - ct_dst_ltm: 0.9385165238013475
ct_src_dport_ltm - ct_srv_src: 0.8697481965388258
ct_src_dport_ltm - ct_srv_dst: 0.875945140285581
ct_src_dport_ltm - ct_dst_ltm: 0.9681552558311022
ct_src_dport_ltm - ct_src_ltm: 0.9453356577337514
ct_dst_sport_ltm - ct_srv_src: 0.8485524119738095
ct_dst_sport_ltm - ct_srv_dst: 0.848826348393225
ct_dst_sport_ltm - ct_dst_ltm: 0.8801999225817902
ct_dst_sport_ltm - ct_src_ltm: 0.8802688290121267
ct_dst_sport_ltm - ct_src_dport_ltm: 0.9214888967474458
ct_dst_src_ltm - ct_srv_src: 0.9421756410442572
ct_dst_src_ltm - ct_srv_dst: 0.9518264424997848
ct_dst_src_ltm - ct_dst_ltm: 0.876487558129229
ct_dst_src_ltm - ct_src_ltm: 0.8567978501003428
ct_dst_src_ltm - ct_src_dport_ltm: 0.9188925284471405
ct_dst_src_ltm - ct_dst_sport_ltm: 0.8803824737828199
Label - sttl: 0.9044186464509674
Label - ct_state_ttl: 0.873894824214898

```

Figure 21 Highly Correlated Features

The results show that numerous numerical features are highly associated with one another. For example, the correlation coefficient between “sloss” and “sbytes” is 0.95, indicating a strong positive correlation between these two features. Similarly, the correlation coefficient between “dloss” and “dbytes” is 0.99, whereas it is 0.96 for “ct_srv_dst” and “ct_srv_src.” Furthermore, some features are significantly connected with two or more other features, such as “Dpkts,” which is highly correlated with both “dloss” and “Spkts,” with correlation values of 0.99 and 0.82, respectively.

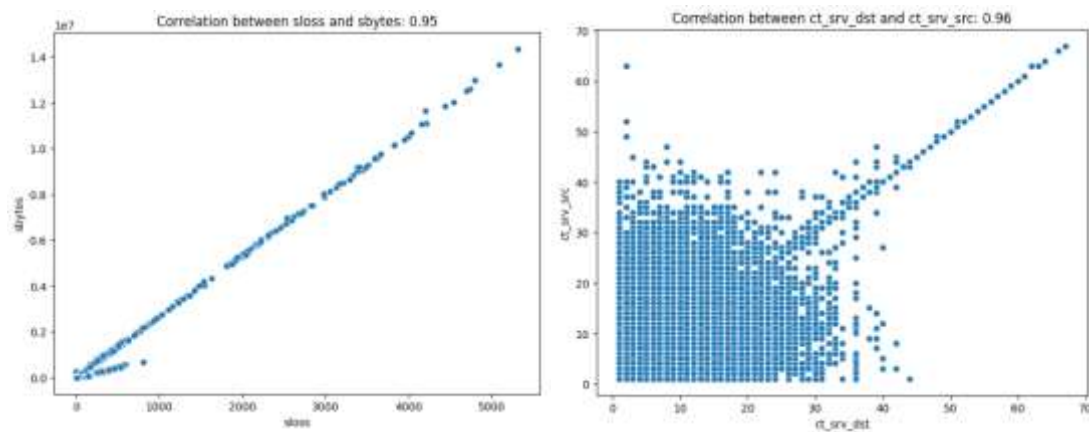


Figure 22 Correlation Plot

4.4 Principal component analysis (PCA)

PCA is a technique used to reduce the dimensionality of data. This technique involves identifying patterns in data and then translating the data into a space with fewer dimensions (Volpi, 2020). The goal is to maintain as much diversity in the original data as possible. Simply put, this assists in identifying the major factors that account for the dataset's variability. For this task, the PCA was used on the UNSW dataset. The Vector Assembler was created to assemble all data columns into a single feature column. The data are then reduced to two dimensions using PCA with $k = 2$. The resulting two dimensions are saved in the `pca_features` column, as shown in Figure 23.

Principal Component Analysis

```
[11] assembler = VectorAssembler(inputCol="numeric_columns", outputCol="features")
data_vector = assembler.transform(data).select("features")
```

Apply PCA

```
[12] pca = PCA(k=2, inputCol="features", outputCol="pca_features")
model = pca.fit(data_vector)
pca_result = model.transform(data_vector).select("pca_features")
```

```
[14] # Extract explained variance (PCA scores)
explained_variance = model.explainedVariance.toArray()
```

Figure 23 PCA code

The explained variance, which is a measure of how much variation in the original dataset is retained by each primary component, is calculated. The first principal component (PCA 1) covers 99.88% of the variation in the UNSW dataset, whereas the second principal component (PCA 2) captures only 0.12%. This shows that the first principal component can explain the majority of the variation in the dataset and that the second principal component is less important in explaining the variation in the data, as shown in Figure 24 by the scatterplot which used to visualise the data in the new two-dimensional space created by PCA.

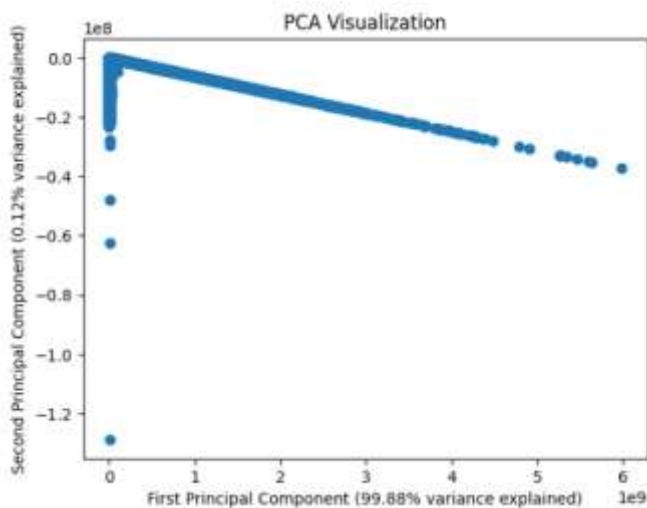


Figure 24 PCA Component Scatter Graph

5. Classification Model

Classification is a supervised learning strategy that uses a categorical variable as the goal variable. It involves using a labelled dataset to train a model and predicting the class label for future cases. The RF algorithm is an ensemble learning strategy that builds numerous decision trees and aggregates their output to enhance overall prediction accuracy and avoid overfitting. This method has several advantages, including noise resistance, overfitting resistance, and the capacity to handle large datasets with many features. As a result, the RF technique is commonly utilised for classification tasks such as network intrusion detection. The utilisation of classification algorithms is critical in enhancing the efficacy of IDS by successfully differentiating between distinct categories of activity. For this task, our primary goal is to create classifiers for binary and multi-class classification using the RF method, a powerful ensemble learning approach that has demonstrated amazing precision and durability in a variety of ML scenarios.

5.1 Binary classification

The classification issue in network intrusion detection can be treated as a binary or multi-class problem. The algorithm differentiates between two classes: normal traffic and attack traffic. The initial step is to generate an index for the categorical target variable, which is then used to create a feature vector with numeric features. Then, the dataset is partitioned into an 80% training set and a 20% testing set using a randomised seed. The RandomForestClassifier is created using features and indexed target columns. A parameter grid is formulated to assess various permutations of the number of trees (numTrees) and the maximum tree depth (maxDepth) (RandomForest—PySpark 3.4.0 Documentation, n.d.) to enhance the model's performance. Cross-validation is used to train the model and assess its efficacy on data that has not been previously observed. The results were evaluated using different evaluation metrics, such as accuracy, precision, recall, f1-score, and confusion matrix. Table 1 summarises the results of the binary classification using the Random Forest classifier, with values expressed as percentages.

Table 1. Binary Classification Results

Class	Precision	Re-call	F1Score	Support
Normal (0.0)	1.00	1.00	1.00	443,683
Attack (1.0)	1.00	1.00	1.00	64,623
Accuracy			1.00	508,306
Macro Avg	1.00	1.00	1.00	508,306
Weighted Avg	1.00	1.00	1.00	508,306

By comparing the random forest classifier outputs to the evaluation criteria, we achieved a binary classification accuracy of 99.99%, demonstrating the usefulness of the algorithm for intrusion detection. The RF classifier had a precision score of 99.99%, recall score, and F1-score of 99.99%. These findings demonstrate the classifier's excellent performance in discriminating between normal and attack traffic in the UNSW-NB15 dataset. Figure 25 depicts the confusion matrix for binary classification, which provides more insight into the performance of the random forest classifier. The diagonal elements of the matrix indicate that all the data points are correctly

categorised and there was not a single misclassified data. Overall, the confusion matrix indicates the capacity of the RF classifier to correctly categorise the normal and attack labels on the UNSW-NB15 dataset.

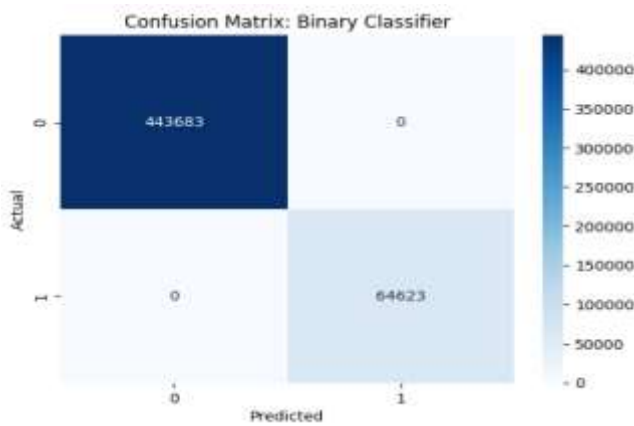


Figure 25. Binary classification confusion matrix.

```

idxC_label = StringIndexer(inputCol="label", outputCol="label_idx")
scaler = VectorAssembler(inputCols=["num1","num2","num3","num4","num5","num6","num7","num8","num9","num10","num11","num12","num13","num14","num15","num16","num17","num18","num19","num20","num21","num22","num23","num24","num25","num26","num27","num28","num29","num30","num31","num32","num33","num34","num35","num36","num37","num38","num39","num40","num41","num42","num43","num44","num45","num46","num47","num48","num49","num50","num51","num52","num53","num54","num55","num56","num57","num58","num59","num60","num61","num62","num63","num64","num65","num66","num67","num68","num69","num70","num71","num72","num73","num74","num75","num76","num77","num78","num79","num80","num81","num82","num83","num84","num85","num86","num87","num88","num89","num90","num91","num92","num93","num94","num95","num96","num97","num98","num99","num100","num101","num102","num103","num104","num105","num106","num107","num108","num109","num110","num111","num112","num113","num114","num115","num116","num117","num118","num119","num120","num121","num122","num123","num124","num125","num126","num127","num128","num129","num130","num131","num132","num133","num134","num135","num136","num137","num138","num139","num140","num141","num142","num143","num144","num145","num146","num147","num148","num149","num150","num151","num152","num153","num154","num155","num156","num157","num158","num159","num160","num161","num162","num163","num164","num165","num166","num167","num168","num169","num170","num171","num172","num173","num174","num175","num176","num177","num178","num179","num180","num181","num182","num183","num184","num185","num186","num187","num188","num189","num190","num191","num192","num193","num194","num195","num196","num197","num198","num199","num200","num201","num202","num203","num204","num205","num206","num207","num208","num209","num210","num211","num212","num213","num214","num215","num216","num217","num218","num219","num220","num221","num222","num223","num224","num225","num226","num227","num228","num229","num230","num231","num232","num233","num234","num235","num236","num237","num238","num239","num240","num241","num242","num243","num244","num245","num246","num247","num248","num249","num250","num251","num252","num253","num254","num255","num256","num257","num258","num259","num260","num261","num262","num263","num264","num265","num266","num267","num268","num269","num270","num271","num272","num273","num274","num275","num276","num277","num278","num279","num280","num281","num282","num283","num284","num285","num286","num287","num288","num289","num290","num291","num292","num293","num294","num295","num296","num297","num298","num299","num300","num301","num302","num303","num304","num305","num306","num307","num308","num309","num310","num311","num312","num313","num314","num315","num316","num317","num318","num319","num320","num321","num322","num323","num324","num325","num326","num327","num328","num329","num330","num331","num332","num333","num334","num335","num336","num337","num338","num339","num340","num341","num342","num343","num344","num345","num346","num347","num348","num349","num350","num351","num352","num353","num354","num355","num356","num357","num358","num359","num360","num361","num362","num363","num364","num365","num366","num367","num368","num369","num370","num371","num372","num373","num374","num375","num376","num377","num378","num379","num380","num381","num382","num383","num384","num385","num386","num387","num388","num389","num390","num391","num392","num393","num394","num395","num396","num397","num398","num399","num400","num401","num402","num403","num404","num405","num406","num407","num408","num409","num410","num411","num412","num413","num414","num415","num416","num417","num418","num419","num420","num421","num422","num423","num424","num425","num426","num427","num428","num429","num430","num431","num432","num433","num434","num435","num436","num437","num438","num439","num440","num441","num442","num443","num444","num445","num446","num447","num448","num449","num450","num451","num452","num453","num454","num455","num456","num457","num458","num459","num460","num461","num462","num463","num464","num465","num466","num467","num468","num469","num470","num471","num472","num473","num474","num475","num476","num477","num478","num479","num480","num481","num482","num483","num484","num485","num486","num487","num488","num489","num490","num491","num492","num493","num494","num495","num496","num497","num498","num499","num500","num501","num502","num503","num504","num505","num506","num507","num508","num509","num510","num511","num512","num513","num514","num515","num516","num517","num518","num519","num520","num521","num522","num523","num524","num525","num526","num527","num528","num529","num530","num531","num532","num533","num534","num535","num536","num537","num538","num539","num540","num541","num542","num543","num544","num545","num546","num547","num548","num549","num550","num551","num552","num553","num554","num555","num556","num557","num558","num559","num560","num561","num562","num563","num564","num565","num566","num567","num568","num569","num570","num571","num572","num573","num574","num575","num576","num577","num578","num579","num580","num581","num582","num583","num584","num585","num586","num587","num588","num589","num590","num591","num592","num593","num594","num595","num596","num597","num598","num599","num600","num601","num602","num603","num604","num605","num606","num607","num608","num609","num610","num611","num612","num613","num614","num615","num616","num617","num618","num619","num620","num621","num622","num623","num624","num625","num626","num627","num628","num629","num630","num631","num632","num633","num634","num635","num636","num637","num638","num639","num640","num641","num642","num643","num644","num645","num646","num647","num648","num649","num650","num651","num652","num653","num654","num655","num656","num657","num658","num659","num660","num661","num662","num663","num664","num665","num666","num667","num668","num669","num670","num671","num672","num673","num674","num675","num676","num677","num678","num679","num680","num681","num682","num683","num684","num685","num686","num687","num688","num689","num690","num691","num692","num693","num694","num695","num696","num697","num698","num699","num700","num701","num702","num703","num704","num705","num706","num707","num708","num709","num710","num711","num712","num713","num714","num715","num716","num717","num718","num719","num720","num721","num722","num723","num724","num725","num726","num727","num728","num729","num730","num731","num732","num733","num734","num735","num736","num737","num738","num739","num740","num741","num742","num743","num744","num745","num746","num747","num748","num749","num750","num751","num752","num753","num754","num755","num756","num757","num758","num759","num760","num761","num762","num763","num764","num765","num766","num767","num768","num769","num770","num771","num772","num773","num774","num775","num776","num777","num778","num779","num780","num781","num782","num783","num784","num785","num786","num787","num788","num789","num790","num791","num792","num793","num794","num795","num796","num797","num798","num799","num800","num801","num802","num803","num804","num805","num806","num807","num808","num809","num810","num811","num812","num813","num814","num815","num816","num817","num818","num819","num820","num821","num822","num823","num824","num825","num826","num827","num828","num829","num830","num831","num832","num833","num834","num835","num836","num837","num838","num839","num840","num841","num842","num843","num844","num845","num846","num847","num848","num849","num850","num851","num852","num853","num854","num855","num856","num857","num858","num859","num860","num861","num862","num863","num864","num865","num866","num867","num868","num869","num870","num871","num872","num873","num874","num875","num876","num877","num878","num879","num880","num881","num882","num883","num884","num885","num886","num887","num888","num889","num890","num891","num892","num893","num894","num895","num896","num897","num898","num899","num900","num901","num902","num903","num904","num905","num906","num907","num908","num909","num910","num911","num912","num913","num914","num915","num916","num917","num918","num919","num920","num921","num922","num923","num924","num925","num926","num927","num928","num929","num930","num931","num932","num933","num934","num935","num936","num937","num938","num939","num940","num941","num942","num943","num944","num945","num946","num947","num948","num949","num950","num951","num952","num953","num954","num955","num956","num957","num958","num959","num960","num961","num962","num963","num964","num965","num966","num967","num968","num969","num970","num971","num972","num973","num974","num975","num976","num977","num978","num979","num980","num981","num982","num983","num984","num985","num986","num987","num988","num989","num990","num991","num992","num993","num994","num995","num996","num997","num998","num999"]
scaler.transform(trainData).write("data/scaler_output")

# Training
print("Training started")
model = RandomForestClassifier()
model.fit(trainData)
print("Training completed")

```

Figure 26. Binary classification code.

5.2 Multiclassification

Similar to binary classification, the RF classifier was used for multi-classification problems, except that the StringIndexer used to index the category column was changed from “label” to “attack_cat” and VectorAssembler was used to assemble all numerical features into a single “features” column. Using an 80-20 split ratio, the data was then divided into training and test sets. The performance of the model was evaluated using the same evaluation metrics. Table 2 shows that the model achieved an overall accuracy of 0.98, suggesting that it correctly classified most cases in the dataset. Furthermore, the table shows that the model had differing degrees of success in recognising traffic from various classes. While the model performed exceptionally well in identifying classes such as benign, generic, fuzzy, and DoS, it struggled to correctly classify instances of several other classes, such as backdoors and worms, which can also be seen through the confusion matrix, as shown in Figure 27. Overall, the model may not be adequate for identifying specific types of network attacks, and additional study and fine-tuning may be required to improve its performance. It is also worth noting that the class imbalance of the dataset may have contributed to the poorer precision and recall scores for some classes, as some of these classes had on average only 30–500 data points. Overall, the results indicate that the RF can be a useful tool for detecting specific

types of network attacks, but its performance may be limited by the nature of the data and the specific attack types being targeted.

Table 2. Multiclassification Results

Class	Precision	Re-call	F1Score	Support
Benign	1	1	1	443683
Generic	1	0.98	0.99	43458
Exploits	0.57	0.97	0.72	8887
Fuzzers	0.78	0.87	0.83	3839
DoS	0.81	0.02	0.04	3268
Reconnaissance	0.84	0.68	0.75	2445
Fuzzers	0.94	0.18	0.3	1043
Analysis	0.81	0.13	0.22	567
Backdoor	0.7	0.02	0.04	339
Reconnaissance	0	0	0	341
Shellcode	0.76	0.56	0.65	261
Backdoors	0	0	0	95
Shellcode	0	0	0	37
Worms	1	0.02	0.05	43
Macro Average	0.66	0.39	0.4	508306
Weighted Average	0.99	0.98	0.98	508306

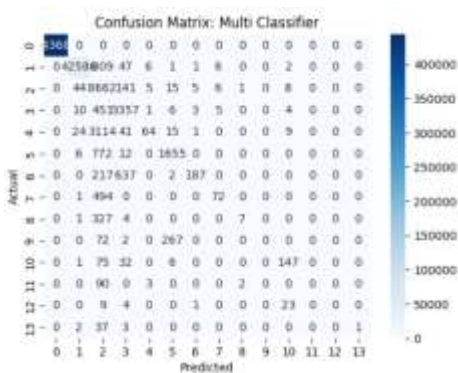


Figure 27 Multiclassification confusion matrix

Discussion

This study reinforces the role of big data analytics in addressing the scalability and accuracy challenges of IDSs. Traditional IDS solutions fail to efficiently manage large-scale traffic; however, the integration of Hive and PySpark provides both computational speed and analytical depth. The RF model delivered near-perfect binary classification accuracy, confirming its robustness for distinguishing normal from malicious traffic. Multi-class results, while strong overall, revealed limitations in detecting minority attack classes such as backdoors and worms. This suggests the need for further work on class imbalance, ensemble methods, or hybrid DL models. The analysis of traffic by source IPs, protocols, and attack categories provides practical insights for security practitioners in prioritising defence strategies.

6. Conclusion

This study demonstrates that combining Apache Hive and PySpark significantly enhances the effectiveness of NID systems. The proposed framework achieved state-of-the-art accuracy in both binary and multi-class scenarios by leveraging the UNSW-NB15 dataset. The findings confirm that big data tools can scale intrusion detection to meet modern network demands while providing interpretable analytical insights. Future research should focus on addressing class imbalance, incorporating real-time detection, and exploring deep learning-based extensions to further strengthen IDS capabilities.

References

- Alrawashdeh T, Alhamid M. 2020. Intrusion detection system using machine learning International Journal of Advanced Computer Science and Applications, 11(6), 7-14. <https://doi.org/10.14569/IJACSA.2020.0110602>
- Aminanto, E., Wibisono, H., & Adi, K. (2017). Intrusion Detection System Using Cloud Computing Data Mining Techniques Journal of Telecommunication, Electronic and Computer Engineering, vol. 9, no. 3–8, pp. 43–47.
- Chen, X., Li, D., Chen, M., and Zou, D. (2019). Cybersecurity and privacy protection: Survey, taxonomy, and open issues. IEEE Communications Surveys & Tutorials, 21(3), 2333-2370. <https://doi.org/10.1109/COMST.2019.2914962>
- Federal Bureau of Investigation (FBI). (2021). Internet Crime Complaint Centre (IC3) Report 2020. [Retrieved from https://www.ic3.gov/Media/PDF/AnnualReport/2020_IC3Report.pdf](https://www.ic3.gov/Media/PDF/AnnualReport/2020_IC3Report.pdf)
- Wang, Y., Wang, J., Huang, L., & Yao, X. (2018). Intrusion detection system based on improved GMM algorithm Journal of Physics: Conference Series, 1096, 032023. <https://doi.org/10.1088/1742-6596/1096/3/032023>
- Chowdhury, M., Zaharia, M., Ma, J., Jordan, M. I., and Stoica, I. (2011). Managing data transfers in computer clusters with orchestra. ACM SIGCOMM Computer Communication Review, 41(4), 98-109. <https://doi.org/10.1145/2043164.2018448>
- Moustafa, N. (2015). The UNSW-NB15 dataset. Research Data Australia. [Available from: https://researchdata.edu.au/the-unsw-nb15-dataset/1957529](https://researchdata.edu.au/the-unsw-nb15-dataset/1957529)
- Turney, S. (2022). Pearson correlation coefficient (r) | Guide and examples. Scribbr. <https://www.scribbr.com/statistics/pearson-correlation-coefficient/>

- Volpi, G. F. (2020). The most gentle introduction to PCA. Towards Data Science. <https://towardsdatascience.com/the-most-gentle-introduction-to-principal-component-analysis-9ffae371e93b>
- Li, J., Wu, Y., & Zhang, H. (2021). Deep learning methods for network intrusion detection: A survey. Computers & Security, 102, 102153. <https://doi.org/10.1016/j.cose.2020.102153>
- Shone, N., Ngoc, T. N., Phai, V. D., & Shi, Q. (2022). Hybrid deep learning approach for network intrusion detection Journal of Information Security and Applications, 67, 103182. <https://doi.org/10.1016/j.jisa.2022.103182>
- Zhang, Y., Sun, Y., & Lin, X. (2023). Scalable ML for big data intrusion detection in cloud environments Future Generation Computer Systems, 144, 85-97. <https://doi.org/10.1016/j.future.2023.01.005>
- Alqahtani, A., & Wang, H. (2024). A survey on big data analytics for cybersecurity: Challenges and opportunities. IEEE Access, 12, 5573-5590. <https://doi.org/10.1109/ACCESS.2024.3349557>