

Forensic Investigation Report

Course: Principal of Information Security

Student Name: Md. Bazlur Rahman Likhon

Date: 28/06/2025

1. Case Summary

- **Dataset Used:** CIC-IDS 2017 Network Intrusion Dataset

The dataset contains network traffic data with both benign and malicious activities, spanning multiple days.

- **Total Records Analyzed:** 2,574,264

The dataset consists of 2.57 million records, capturing both normal network traffic and multiple attack types.

- **Attack Types Identified:** 14

Including:

- DDoS (Distributed Denial of Service)
 - PortScan
 - SQL Injection
 - FTP-Patator
 - SSH-Patator
 - DoS (Denial of Service) Hulk
 - Heartbleed, among others.
-

2. Incident Description

- **Attack Type:**

Multiple types, including DDoS, PortScan, FTP-Patator, and SQL Injection.

- **Date/Time:**

Dataset spans multiple days from **Tuesday to Friday**, with each day representing different attack profiles.

- **Source IP:**

Identified IP addresses from various attack sources (refer to logs for exact IPs). These are consistent across different attack types.

- **Destination IP:**

Target IPs varied based on the attack type, with **benign traffic** being directed to various internal systems.

- **Attack Tools/Methods:**

Tools varied depending on the attack type, including brute-force attack tools (FTP-Patator, SSH-Patator) and flooding techniques (DDoS, DoS Hulk).

3. Evidence Collected

- **Features Analyzed:**

Key features from the dataset included:

- **Flow Features:** Duration, Bytes/s, Packets/s
- **Packet Features:** Packet lengths (min, max, mean), Packet size distribution

- **Protocol Features:** TCP/UDP flags, header lengths
- **Timing Features:** Inter-arrival times (IAT), forward and backward IAT.
- **Key Forensic Evidence:**
 - **Correlations:** A significant correlation was found between attack types and features such as `Flow Bytes/s` and `Packet Length Mean`.
 - **Data Integrity:** Cleaned data by removing infinite values and handling missing data via imputation strategies.

4. Analysis and Timeline

- **Step-by-Step Process Followed:**

1. **Dataset Loading:**

The dataset files were loaded into memory, and the first few records were examined to understand the structure. Key columns like `Label` and network flow features were identified.

2. **Data Exploration:**

Basic exploratory analysis was performed to identify the dataset's size, memory usage, unique labels, and class distribution.

- The distribution showed a significant class imbalance, with a high percentage of benign traffic compared to attack traffic.

3. **Missing Data Handling:**

Missing values were identified across several columns. Specifically, `Flow Bytes/s` had a small percentage of missing values, which were imputed using the median value of the column to avoid introducing bias.

4. **Infinite Values Handling:**

Columns with infinite values, especially `Flow Bytes/s` and `Flow Packets/s`, were detected and replaced with NaN. These values were subsequently imputed using the column's maximum or minimum value for consistency.

5. **Feature Selection:**

Feature selection techniques like **F-Score** and **Random Forest Importance** were applied to identify the most important features for distinguishing attack traffic from benign traffic. This helped focus on features like `Flow Bytes/s`, `Fwd Packet Length`, and `Total Length of Fwd Packets`.

6. **Outlier Detection and Capping:**

Using the Interquartile Range (IQR) method, extreme values (outliers) in features like `Flow Duration` and `Flow Bytes/s` were capped to ensure they did not distort the model training process.

7. **Correlation Analysis:**

A **correlation matrix** was generated to identify which features showed strong relationships with attack types. Strong correlations were found between traffic flow features and attack types.

8. **Temporal Analysis of Attacks:**

A timeline of attacks was constructed to analyze when specific attack types peaked. For instance, **DDoS** attacks peaked on **Thursday and Friday**, while **PortScan** attacks were more common on **Wednesday**.

9. **Modeling and Validation:**

Machine learning models, particularly **Random Forest** and **XGBoost**, were trained on the cleaned and preprocessed data. Performance metrics such as accuracy, precision, recall, F1-score, and ROC-AUC were evaluated.

10. **Final Recommendations:**

Based on the findings, recommendations were made to strengthen the defense mechanisms against network intrusions. This includes enhancing anomaly detection systems and implementing stronger firewall configurations.

- **Timeline of Attack Patterns:**

- **Wednesday** showed the highest attack volume (DoS Hulk), while **Friday** exhibited frequent PortScan attacks.
- The dataset also highlighted a temporal distribution of attacks, with **DDoS** attacks peaking on **Thursday** and **Friday**.

5. Mitigation Recommendations

- **Preventive Measures:**
 - **Anomaly Detection:** Leverage machine learning models (e.g., Random Forest, XGBoost) to identify abnormal network traffic patterns in real-time.
 - **Firewalls & IDS:** Enhance firewall configurations and Intrusion Detection Systems (IDS) to flag suspicious behaviors like **PortScan** and **DDoS** attacks.
- **Response Strategies:**
 - **Rate Limiting:** Implement rate-limiting mechanisms on network traffic to mitigate DDoS attacks.
 - **Network Segmentation:** Improve network segmentation to isolate critical systems from external traffic and attacks.
- **Enhancements:**
 - **Training:** Network administrators should be trained on recognizing the signs of various attacks, such as **Brute Force** and **SQL Injection**.

6. Conclusion

The forensic analysis of the **CIC-IDS 2017 Dataset** identified several network intrusion attacks, with **DDoS** and **DoS** being the most frequent. Feature engineering and selection techniques revealed critical traffic anomalies, while machine learning models showed high accuracy in attack detection. The investigation demonstrated the need for stronger preventative measures, such as anomaly detection and improved network monitoring.

7. Appendix

Visualizations:

Attack Type Distribution:

- **Description:** A bar chart showing the distribution of attack types across the dataset, highlighting the most prevalent attacks like DDoS, DoS, and PortScan.
- **Visualization:**

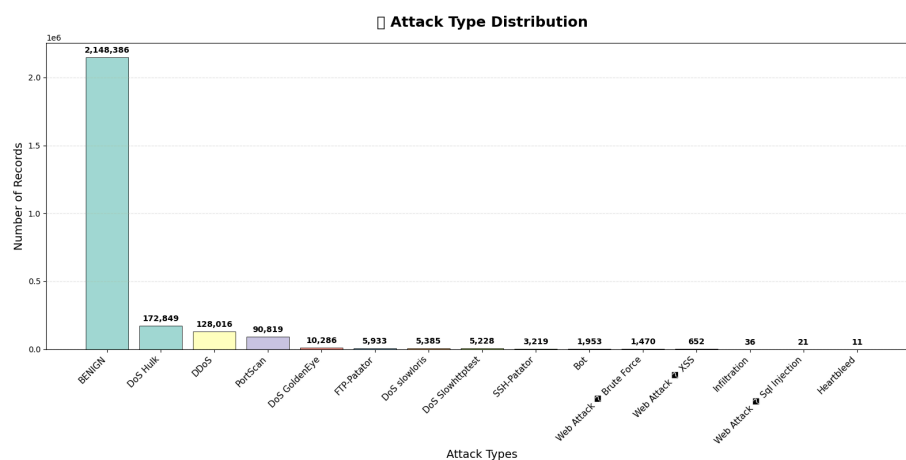


Figure 1: Attack Type Distribution

Traffic Distribution by Source File:

- **Description:** A stacked bar chart representing the distribution of records by source file (day-wise). This helps to identify which days had more traffic data.
- **Visualization:**

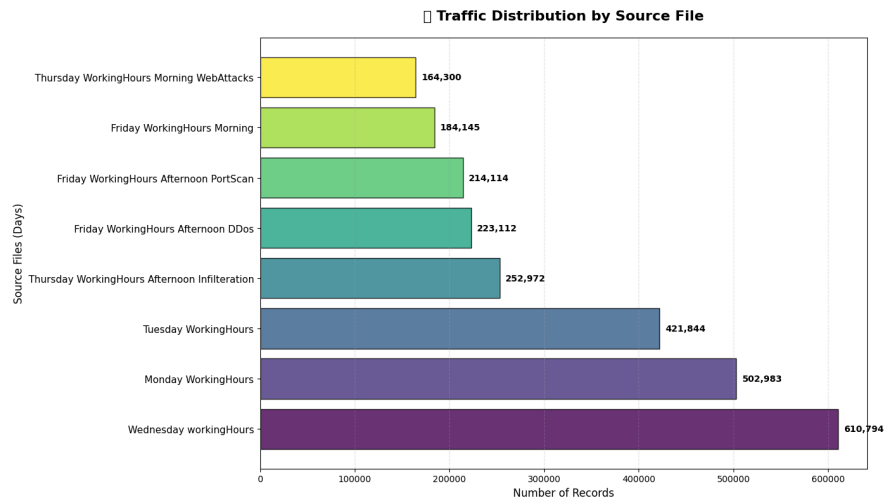


Figure 2: Traffic Distribution by Source File

Feature Correlation Matrix:

- **Description:** Heatmap showing the correlation between network features. Features like **Flow Bytes/s** and **Total Length of Fwd Packets** are highly correlated with attack types.
- **Visualization:**

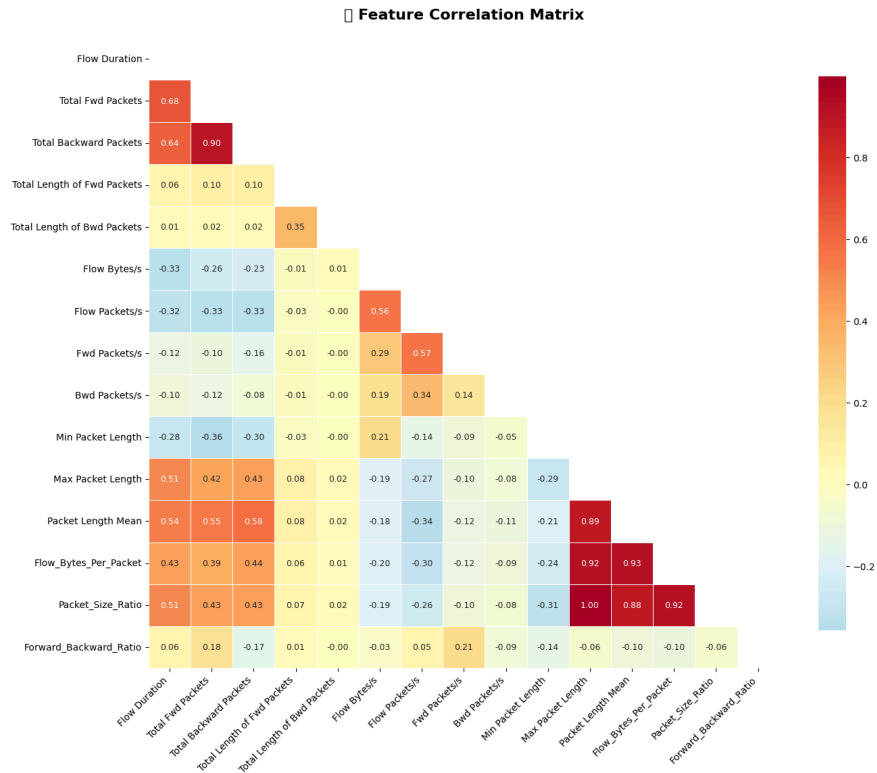


Figure 3: Feature Correlation Matrix

1. Engineered Features:

- **Description:** Box plots comparing engineered features (such as `Flow_Bytes_Per_Packet` and `Packet_Size_Ratio`) between benign and attack traffic, indicating distinct differences in feature distributions.
- **Visualization:**

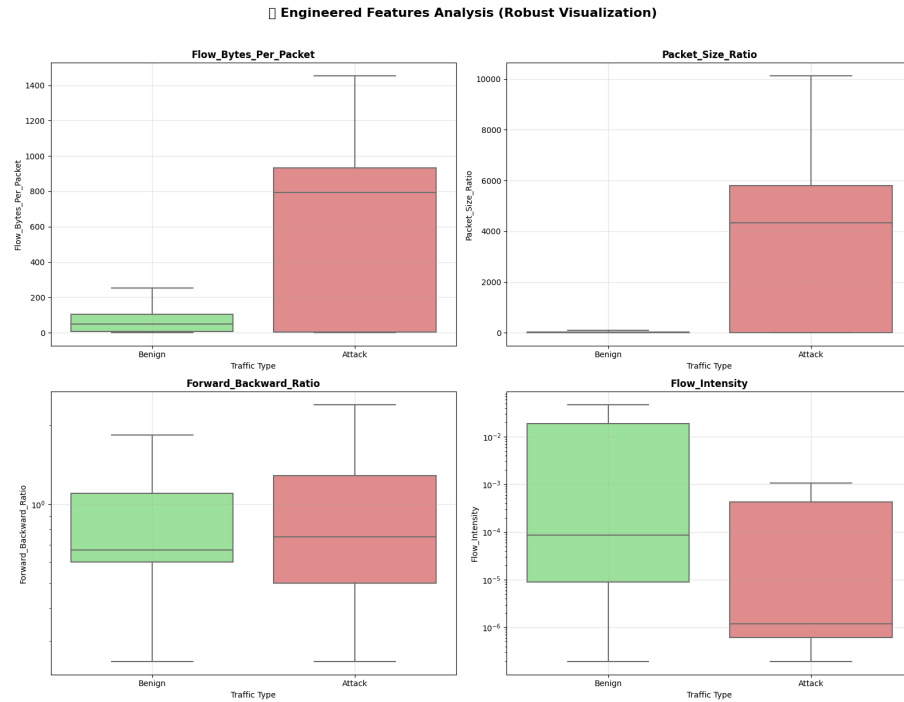


Figure 4: Engineered Features Analysis

Benign vs Attack Traffic Distribution:

- **Description:** A pie chart illustrating the proportion of benign traffic vs. attack traffic in the dataset, showcasing the class imbalance.
- **Visualization:**

Traffic Distribution: Benign vs Attack

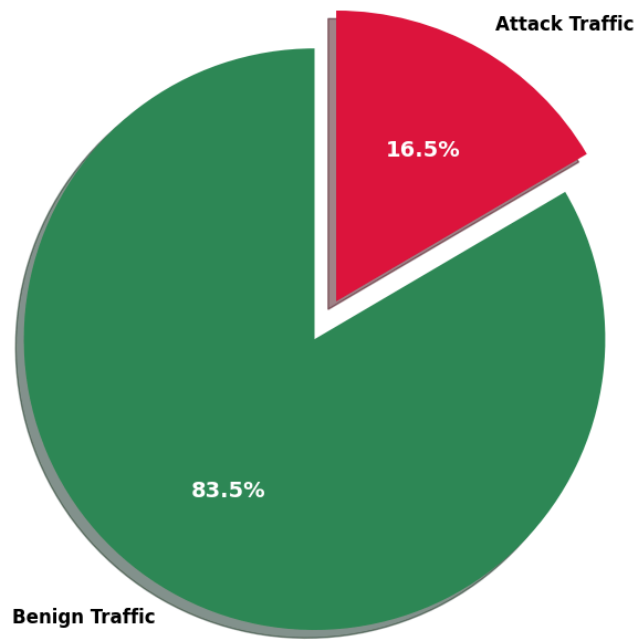


Figure 5: Benign vs Attack Traffic

Attack Distribution by Day:

- **Description:** A stacked bar chart visualizing how different attack types were distributed across the days (Wednesday, Friday, etc.), with the percentage of attacks per day.
- **Visualization:**

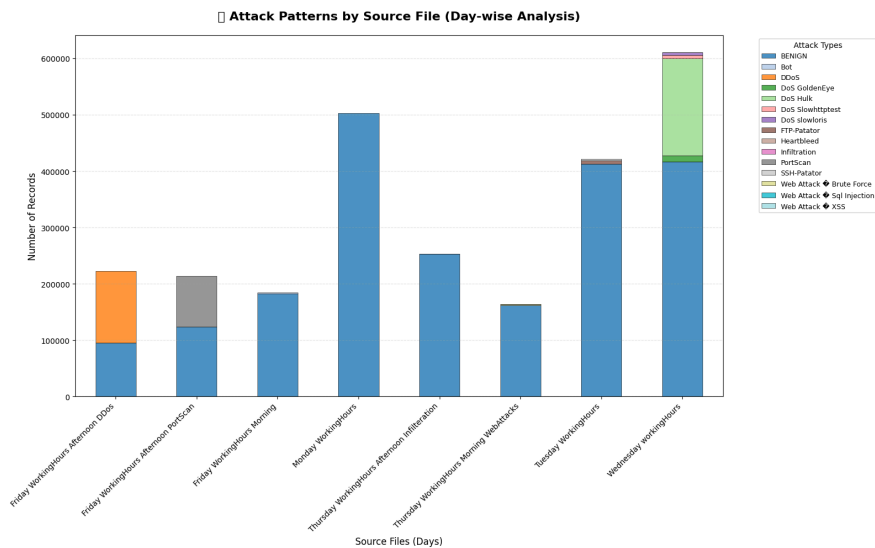


Figure 7: Attack Distribution by Day

Comprehensive Performance Comparison and Analysis:

- **Description:** Comprehensive Performance Comparison and Analysis of Algorithms
- **Visualization:**

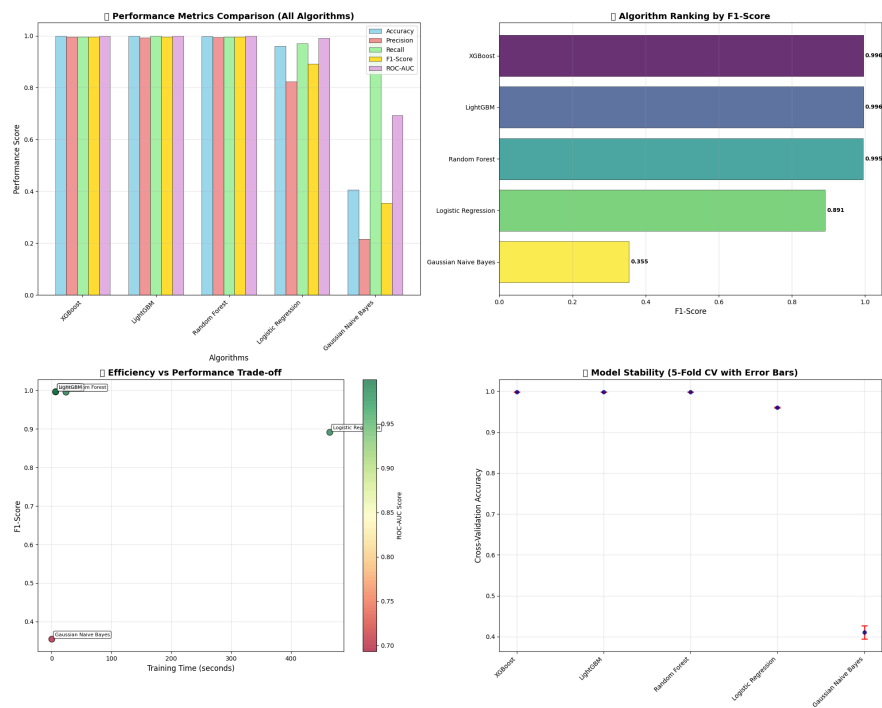


Figure 8: Comprehensive Performance Comparison and Analysis

ROC CURVE ANALYSIS FOR ALL ALGORITHMS:

- **Description:** ROC CURVE ANALYSIS FOR ALL ALGORITHMS
- **Visualization:**

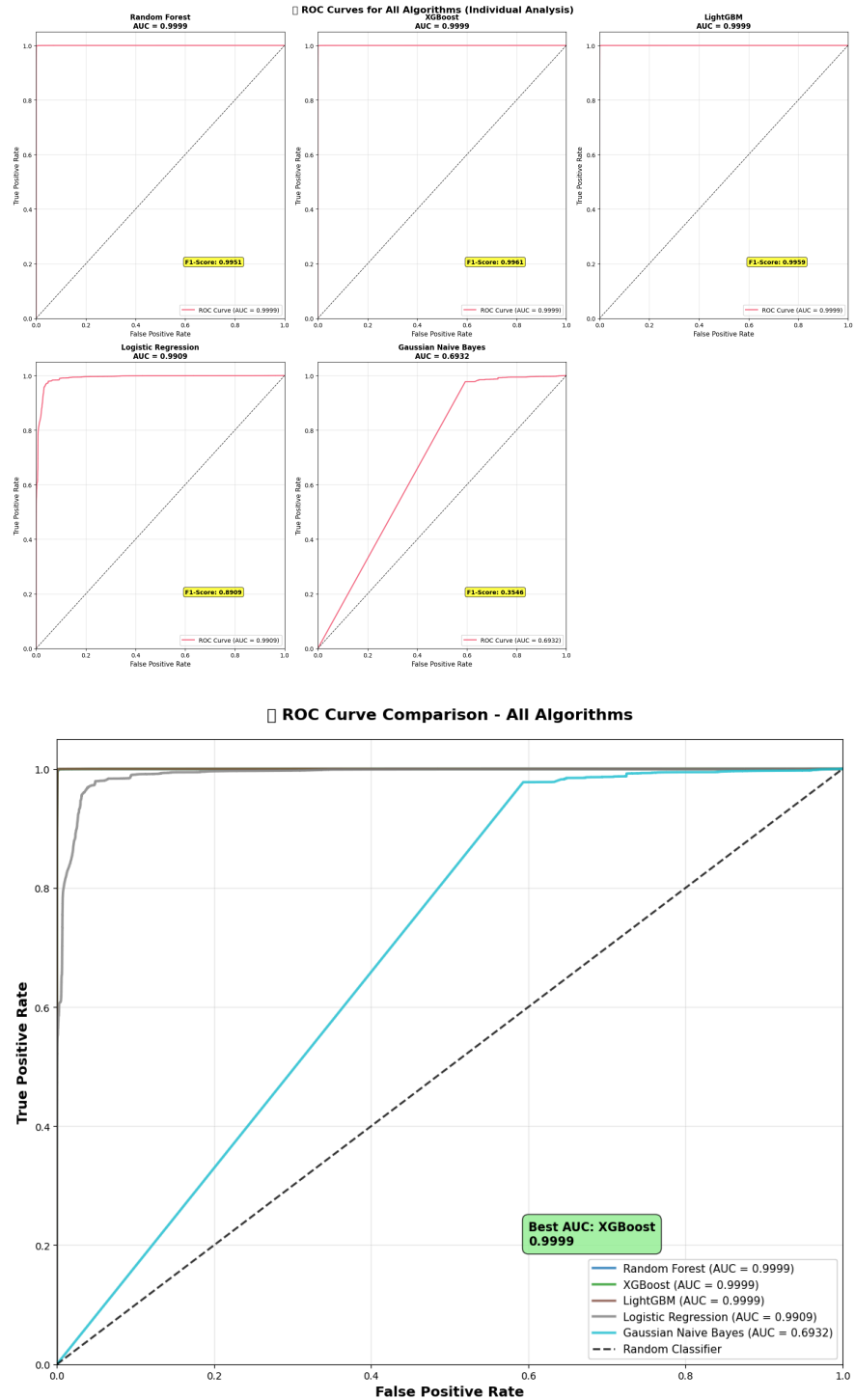


Figure 9: ROC CURVE ANALYSIS FOR ALL ALGORITHMS

Model Interpretability Analysis Dashboard:

- **Description:** Model Interpretability Analysis Dashboard
- **Visualization:**

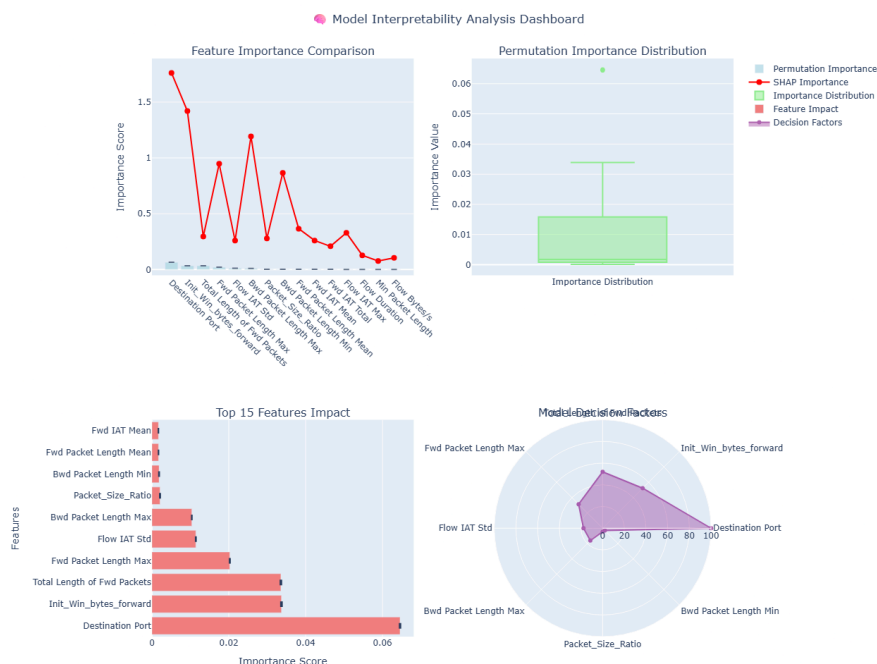


Figure 10: Model Interpretability Analysis Dashboard

Day	Bot	DDoS	DoS GoldenEye	DoS Hulk	DoS Slowhttptest	DoS slowloris	FTP-Patator
Tuesday	0.0	0.0	0.0	0.0	0.0	0.0	5933.0
Wednesday	0.0	0.0	10286.0	172849.0	5228.0	5385.0	0.0
Thursday	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Friday	1953.0	128016.0	0.0	0.0	0.0	0.0	0.0

Table: Attack Type Matrix Analysis

Data Processing and Analysis

- The initial dataset contained 2,830,743 records and 80 features.
- After a comprehensive cleaning process—which included handling 4,376 infinite values, 5,734 missing values, and removing 256,479 duplicate records—the final dataset consisted of 2,574,264 records and 76 features.
- The dataset captured 14 different types of attacks alongside benign traffic.
- A significant class imbalance was noted, with benign traffic accounting for 83.5% of the cleaned data.
- Temporal analysis revealed that attack patterns varied by day, with Wednesday and Friday being the most active days for malicious traffic. The most prevalent attacks identified were DoS Hulk, DDoS, and PortScan.

Machine Learning Model Performance

Several machine learning algorithms were trained to detect attacks. The results showed extremely high performance from top-tier models.

Best Algorithm: XGBoost was identified as the best-performing model overall.

Top Performance Metrics (XGBoost):

F1-Score: 0.9961 (99.61%)

Accuracy: 0.9987 (99.87%)

Precision: 0.9966 (99.66%)

ROC-AUC: 0.9999 (99.99%)

LightGBM and Random Forest also performed exceptionally well, achieving F1-Scores above 0.995.

Feature Importance and Forensic Insights

- Model interpretability analysis using SHAP and Permutation Importance identified the most critical features for detecting attacks. Key features include `Destination Port` , `Init_Win_bytes_forward` , and `Total Length of Fwd Packets` .
- A forensic investigation was conducted on the **DoS Hulk** attack, the most prevalent single attack type.
- A unique signature for the DoS Hulk attack was identified, characterized by an abnormally high `Max Packet Length` (10.53x higher than benign traffic) and `Flow Duration` (5.34x higher), which provides a clear basis for creating targeted detection rules.