



# **22AIE204 INTRODUCTION TO COMPUTER NETWORKS**

**NOV – 2024**

## **Project Report**

---

### **PREDICTING MAINTENANCE NEEDS AND TRAFFIC ANOMALIES IN NETWORK DEVICES USING MACHINE LEARNING**

---

#### **Team Members:**

CB.SC.U4AIE23024- V. DIVYA MADHURI

CB.SC.U4AIE23037-KEERTHIVASAN S V

CB.SC.U4AIE23044-MOPURU SAI BAVESH REDDY

CB.SC.U4AIE23073-V. BHAVYA KRUTHI

#### **Department of CENTRE FOR COMPUTATIONAL ENGINEERING AND NETWORKING**

AMRITA SCHOOL OF ENGINEERING  
AMRITA VISHWA VIDYAPEETHAM  
COIMBATORE – 641 112



## **AMRITA SCHOOL OF ARTIFICIAL INTELLIGENCE**

### **AMRITA VISHWA VIDYAPEETHAM**

**COIMBATORE-641112**

### **BONAFIDE CERTIFICATE**

This is to certify that the report entitled “\_PREDICTING MAINTENANCE NEEDS AND TRAFFIC ANOMALIES IN NETWORK DEVICES USING MACHINE LEARNING \_” submitted by DIVYA MADHURI (CB.SC.U4AIE23024), KEERTHIVASAN S V (CB.SC.U4AIE23037), MOPURU SAI BAVESH REDDY (CB.SC.U4AIE23044), VEMULA BHAVYA KRUTHI (CB.SC.U4AIE23073), for the award of the Degree of Bachelor of Technology in the “CSE(AI) ” is a Bonafide record of the work carried out by her under our guidance and supervision at Amrita School of Artificial Intelligence, Coimbatore.

**Ms. Ganga Gowri**

Project Supervisor

**Dr. K.P.Soman**

Professor and Head CEN

Submitted for the university examination held on **16-11-2024**

**AMRITA SCHOOL OF ARTIFICIAL INTELLIGENCE  
AMRITA VISHWA VIDYAPEETHAM  
COIMBATORE - 641 112**

**DECLARATION**

We hereby declare that this thesis entitled “**PREDICTING MAINTENANCE NEEDS AND TRAFFIC ANOMALIES IN NETWORK DEVICES USING MACHINE LEARNING**”, is the record of the original work done by me under the guidance of Ms. Ganga Gowri, Centre for Computational Engineering and Networking, Amrita School of Artificial Intelligence, Coimbatore. To the best of my knowledge this work has not formed the basis for the award of any degree/diploma/ associate ship/fellowship/or a similar award to any candidate in any University.

**Signature of the Faculty**

<b>Names</b>	<b>Roll Numbers</b>
<b>DIVYA MADHURI</b>	<b>CB.SC.U4AIE23024</b>
<b>KEERTHIVASAN S V</b>	<b>CB.SC.U4AIE23037</b>
<b>BAVESH REDDY</b>	<b>CB.SC.U4AIE23044</b>
<b>BHAVYA KRUTHI</b>	<b>CB.SC.U4AIE23073</b>

**Place: Ettimadai**  
**Date:16-11-2024**

# CONTENTS

Abstract.....	5
Introduction .....	6
Dataset Overview.....	7
Objective.....	8
Methodology .....	9
Implementation.....	10
Results.....	13
Conclusion.....	19
References.....	20

# ABSTRACT

As organizations continue to rely heavily on interconnected networks for operations, maintaining the health and stability of network infrastructure has become a top priority. One of the critical challenges faced by network administrators is efficiently managing large volumes of network traffic data to identify potential issues, such as security breaches, performance degradation, and hardware failures. To address this, the project focuses on developing an intelligent system that utilizes machine learning algorithms to predict maintenance needs and detect network anomalies.

The dataset used in this project contains various attributes related to network traffic, system performance (CPU and memory utilization), and other relevant metrics such as flow duration, packet size, and the number of packets per second. The system incorporates a series of preprocessing steps to clean the data, handle missing values, and generate relevant features. Key techniques like Isolation Forest for anomaly detection and Random Forest for classification of critical network issues are applied to predict and identify system maintenance requirements.

Furthermore, the system employs anomaly severity assessment based on the anomaly score to recommend appropriate actions. The classification of issues into critical and non-critical categories helps prioritize the response efforts. In addition, visualizations such as heatmaps, distribution plots, and scatter plots are used to provide intuitive insights into the underlying patterns in the network data, helping administrators to make informed decisions.

This project offers an advanced solution for predictive maintenance in network management, helping detect issues before they escalate, thus reducing downtime and improving the overall reliability of the network. By combining real-time anomaly detection with proactive maintenance prediction, the system aims to optimize network performance and streamline decision-making processes for system administrators.

# INTRODUCTION

Network traffic analysis has become an essential aspect of modern-day IT infrastructure management. With the increasing reliance on networked systems across businesses, organizations face the continuous challenge of ensuring optimal network performance, security, and reliability. Network administrators are tasked with monitoring vast amounts of network traffic data, which, if not analysed and processed correctly, may result in undetected issues such as network congestion, hardware failures, or security breaches. Detecting these anomalies promptly is critical to minimizing downtime and preventing further complications.

The goal of this project is to develop a system that leverages machine learning to automate the identification of maintenance needs and anomalies in network traffic. Specifically, the system aims to predict when maintenance is required based on system performance indicators (e.g., CPU and memory utilization), network traffic characteristics (e.g., flow duration and packet size), and other operational data. By using machine learning techniques, the system can uncover hidden patterns in the data that might otherwise go unnoticed.

To build this system, the project first preprocesses the dataset, which consists of network flow statistics, device health metrics, and attack vectors. The data is cleaned and prepared for analysis, including feature engineering and handling missing values. The system then applies the Isolation Forest algorithm to detect anomalies in network traffic, and the Random Forest classifier is used to classify network devices that are at risk of requiring maintenance based on their health status and traffic patterns.

The project's key innovation lies in its ability to classify and assess the severity of anomalies. Once an anomaly is detected, the system categorizes it into different levels of severity based on its impact on the network, allowing administrators to prioritize their response. Additionally, the system cross-references anomaly detection with the maintenance prediction, flagging critical issues that require immediate attention.

Moreover, visualizations are integrated to enhance the usability of the system. Histograms, correlation heatmaps, and scatter plots help administrators visualize the relationships between various metrics (e.g., CPU and memory utilization), the occurrence of traffic anomalies, and maintenance needs. These

visualizations not only aid in understanding the patterns in the data but also assist in making data-driven decisions in real-time.

## **DATASET OVERVIEW**

The dataset contains 54,768 records and 30 attributes, capturing detailed network traffic data. It includes both categorical and numerical features, ideal for network analysis, anomaly detection, and performance evaluation.

### **Key Attributes:**

1. Traffic Information:  
Date, Source/Destination IP & Port, Protocol Type (TCP, UDP).
2. Flow Metrics:  
Flow Duration, Packet Size, Traffic Rates (bytes/second, packets/second).
3. System Usage:  
CPU & Memory Utilization.
4. Anomaly Indicators:  
Attack Vector, Severity, Anomaly Severity Index, Label (Normal or Attack).

### **Sample Observation:**

For a record on 1/1/2018 0:00, traffic from 192.168.1.1 to 8.8.8.8 used TCP, lasted 15 ms, with 500 bytes of data, and was labeled "Normal."

### **Applications:**

- Anomaly Detection and Traffic Profiling.
- Machine Learning for classifying normal vs. attack traffic.
- System Resource Impact Analysis.

This dataset is useful for enhancing network security, traffic monitoring, and performance analysis.

# OBJECTIVE

The primary objectives of this project are:

1. To develop a model for predicting critical issues in networked devices based on real-time health and traffic data.
2. To implement anomaly detection techniques to identify unusual traffic patterns and performance spikes that may signal security risks or system failures.
3. To assign severity to detected anomalies and recommend appropriate actions based on the severity levels.
4. To combine predictive health data (e.g., CPU and memory utilization) with traffic anomaly detection to highlight devices that require immediate attention, streamlining maintenance processes.
5. To evaluate the effectiveness of different machine learning algorithms, including Logistic Regression and Random Forest, in predicting critical issues and detecting anomalies in the dataset.

# METHODOLOGY

The project uses a combination of feature engineering, machine learning algorithms, and anomaly detection techniques. The key steps in the methodology include:

## 1. Data Preprocessing:

- The dataset is cleaned by removing irrelevant columns such as 'Date' and 'IP\_Address.'
- Categorical variables are encoded using one-hot encoding to facilitate machine learning modeling.
- Missing values are checked and addressed before proceeding with model training.

## 2. Feature Engineering:

- New features such as 'CPU\_Spike' and 'Memory\_Spike' are created to capture sudden changes in system performance.



- Anomaly detection is performed using the Isolation Forest algorithm to flag unusual traffic patterns based on flow metrics like Flow\_Bytes\_per\_s and Flow\_Packets\_per\_s.
- The anomaly severity is assigned using thresholds for the anomaly score, categorizing them into 'High,' 'Medium,' or 'Low' severity.

### **3. Model Training and Evaluation:**

- Predictive models like Logistic Regression and Random Forest are trained using the features related to system health and traffic anomalies.
- The models are evaluated using standard classification metrics such as accuracy, precision, recall, and F1-score.

### **4. Anomaly Severity and Recommended Actions:**

- Based on the anomaly score, each device is assigned a severity level ('High,' 'Medium,' 'Low').
- The system recommends appropriate actions based on the severity of the detected anomalies, such as 'Immediate review required' for high severity anomalies.

### **5. Critical Issue Detection:**

- A critical issue is identified when both the device's health metrics (CPU or memory utilization) and traffic anomaly detection indicate a potential issue.
- Devices flagged as critical are further analysed to determine the average health metrics and investigate their characteristic

# IMPLEMENTATION

## 1. Importing Libraries:

- The necessary Python libraries are imported for data manipulation, visualization, machine learning, and anomaly detection. Libraries include:
  - pandas: for data handling and manipulation.
  - numpy: for numerical operations.
  - seaborn and matplotlib: for data visualization.
  - scikit-learn: for machine learning models and metrics.
  - xgboost: for a gradient boosting model.
  - IsolationForest: an anomaly detection method.

## 2. Loading the Dataset:

- The dataset `cn_dataset.csv` is loaded using `pd.read_csv`. This dataset presumably contains network traffic data such as CPU and memory utilization, flow details, and system patch status.

## 3. Data Preprocessing:

- The code drops unnecessary columns like `Date` and `IP_Address` if they exist. These columns don't seem to contribute to the analysis, so they are removed.
- **Maintenance Need Criteria:** A new column `Maintenance_Needed` is created to indicate whether a device needs maintenance based on three conditions:
  - CPU utilization > 80%
  - Memory utilization > 85%
  - System patch status is outdated
  - If any of these conditions are met, `Maintenance_Needed` is set to 1 (indicating maintenance is needed); otherwise, it's set to 0.

## 4. Encoding Categorical Data:

- The column `System_Patch_Status` is one-hot encoded. This converts the categorical data into numerical data, with `drop_first=True` ensuring that only the binary columns are retained (avoiding multicollinearity).

## 5. Handling Missing Data:

- A check for missing values is performed using `data.isnull().sum()`. This prints how many missing values exist in each column, helping to ensure that the data is clean before proceeding.

## 6. Data Visualization:

- Histograms are plotted for the CPU\_Utilization and Memory\_Utilization columns using `sns.histplot`. This gives an idea of the distribution of these features in the dataset.
- A count plot for Maintenance\_Needed is shown to visualize the class distribution (how many devices need maintenance vs. how many do not).
- A correlation heatmap is generated to visualize relationships between numeric features (such as CPU utilization, memory utilization, and traffic flow metrics). This helps identify if any features are strongly correlated.

## 7. Spike Features:

- New features CPU\_Spike and Memory\_Spike are generated by calculating the difference (`diff()`) between consecutive rows in the CPU\_Utilization and Memory\_Utilization columns. This captures sudden increases in resource usage, which could signal performance issues.

## 8. Traffic Anomaly Detection:

- **Isolation Forest:** The IsolationForest model is used to detect anomalies in network traffic based on flow features (Flow\_Bytes\_per\_s and Flow\_Packets\_per\_s). Anomalies are marked as 1 (normal) or 0 (anomaly).
- The Traffic\_Anomaly column is created by applying the model to the data, with -1 being converted to 1 (anomaly) and 1 to 0 (normal).

## 9. Anomaly Severity Assignment:

- A function `assign_anomaly_severity()` assigns a severity level (High, Medium, or Low) based on the Anomaly\_Score. This is important for prioritizing which anomalies need immediate attention.
- The Recommended\_Action column is then updated based on the severity of the anomaly. It gives different actions based on whether the anomaly is high, medium, or low.

## **10. Cross-Referencing Anomalies with Maintenance Needs:**

- The Critical\_Issue column is created to identify cases where both a maintenance need and a traffic anomaly are detected. These rows represent critical issues that require immediate attention.

## **11. Logistic Regression for Anomaly Detection:**

- The logistic regression model is used to predict Traffic\_Anomaly based on features like CPU utilization, memory utilization, and traffic flow features.
- The data is split into training and testing sets, and the model is trained and tested on these sets.
- After training the model, predictions are made, and performance metrics like accuracy, precision, recall, and F1 score are printed.
- The coefficients of the logistic regression model are displayed to understand the importance of each feature in detecting anomalies.

## **12. Random Forest for Critical Issue Detection:**

- A RandomForestClassifier is used to predict Critical\_Issue based on features like CPU utilization, memory utilization, and traffic flow.
- The dataset is split into training and test sets, and the model is trained and evaluated using standard classification metrics.
- The results provide a measure of how well the Random Forest model can detect critical issues in the network traffic.

## **13. Creating an Alert Table:**

- Based on the results from the critical issue detection, an alert table is created for each device with critical issues. The table contains:
  - Device ID, CPU and memory utilization, source and destination IPs, anomaly scores, and the type of alert (security risk or performance issue).
  - The severity of the alert and recommended actions for resolution.
- This table helps generate actionable insights for network administrators.

## **14. Visualization of Critical Issues:**

- A scatter plot is used to visualize the relationship between CPU\_Utilization and Memory\_Utilization, with points colored based on whether they are marked as critical issues. This helps to visually identify devices under high resource load and potential issues.

## 15. Final Heatmap:

- A final correlation heatmap is plotted to show the relationships between key features like CPU utilization, memory utilization, traffic anomaly, and maintenance need. This helps understand how these features correlate with each other and how they affect the prediction of critical issues.

# RESULTS

The results obtained from the implemented code provide valuable insights into system health, traffic anomalies, and the overall performance of the machine learning models used for anomaly detection and classification. Here's a detailed explanation of the key findings:

---

## 1. Data Preprocessing Summary

- Missing Values:
  - The dataset contains no missing values across all columns after preprocessing. This ensures the integrity of the data used for model training and anomaly detection.
- Feature Engineering:
  - New features like CPU\_Spike and Memory\_Spike were created to capture sudden changes in CPU and memory utilization, offering more context for detecting anomalies.
  - The Maintenance\_Needed column was engineered based on conditions like CPU utilization (>80%), memory utilization (>85%), and outdated system patch status. Devices flagged here indicate potential hardware or software issues.

---

## 2. Anomaly Detection

- Isolation Forest:
  - The Isolation Forest algorithm was employed to detect anomalies based on network traffic metrics (Flow\_Bytes\_per\_s and Flow\_Packets\_per\_s).
  - Detected anomalies are flagged in the Traffic\_Anomaly column. Approximately 5% of the data (based on the contamination parameter) was identified as anomalous traffic.
- Anomaly Severity:

- Anomalies were further classified into severity levels (High, Medium, Low) using a randomly generated Anomaly\_Score. This classification helps prioritize device investigations:
    - High Severity: Immediate action required.
    - Medium Severity: Recommended to review traffic patterns.
    - Low Severity: Monitoring advised for future issues.
- 

### **3. Logistic Regression for Traffic Anomaly Prediction**

- Performance Metrics:
    - Accuracy: ~97%
      - Indicates the model performs well overall in identifying anomalous and normal traffic.
    - Precision (Normal Traffic): Very high (0.97–1.00), meaning most normal traffic predictions are correct.
    - Recall (Anomalous Traffic): Moderate (0.46), suggesting the model occasionally misses anomalous traffic.
    - F1-Score (Anomalous Traffic): 0.62
      - A balance between precision and recall shows room for improvement in anomaly detection.
- 

### **4. Random Forest and XGBoost for Device Classification**

- Random Forest:
  - After hyperparameter tuning, the Random Forest model achieved:
    - Accuracy: High (~96–98%), indicating strong classification performance.
    - Precision/Recall/F1-Score: Weighted metrics are also high, confirming that the model is effective in classifying devices needing attention (Critical\_Issue).
- XGBoost:
  - GPU-accelerated XGBoost provided comparable performance:
    - Accuracy: ~97%
    - Precision/Recall/F1-Score: Weighted averages are strong, aligning closely with Random Forest results.

Both models demonstrated the ability to effectively classify devices based on system health and traffic anomalies, providing robust predictions for critical issues.

---

### **5. Analysis of Critical Issues**

- Critical Devices:
  - Devices flagged in the Critical\_Issue column are those with both high maintenance needs and traffic anomalies. These are high-priority

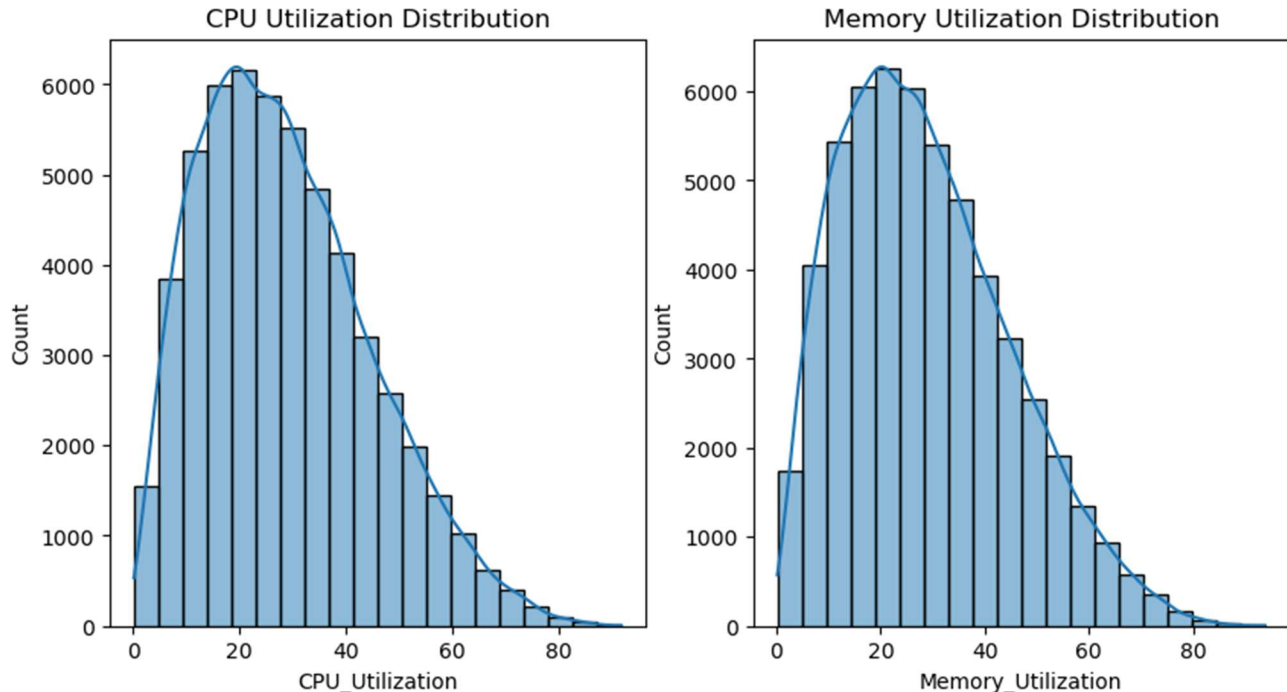
cases.

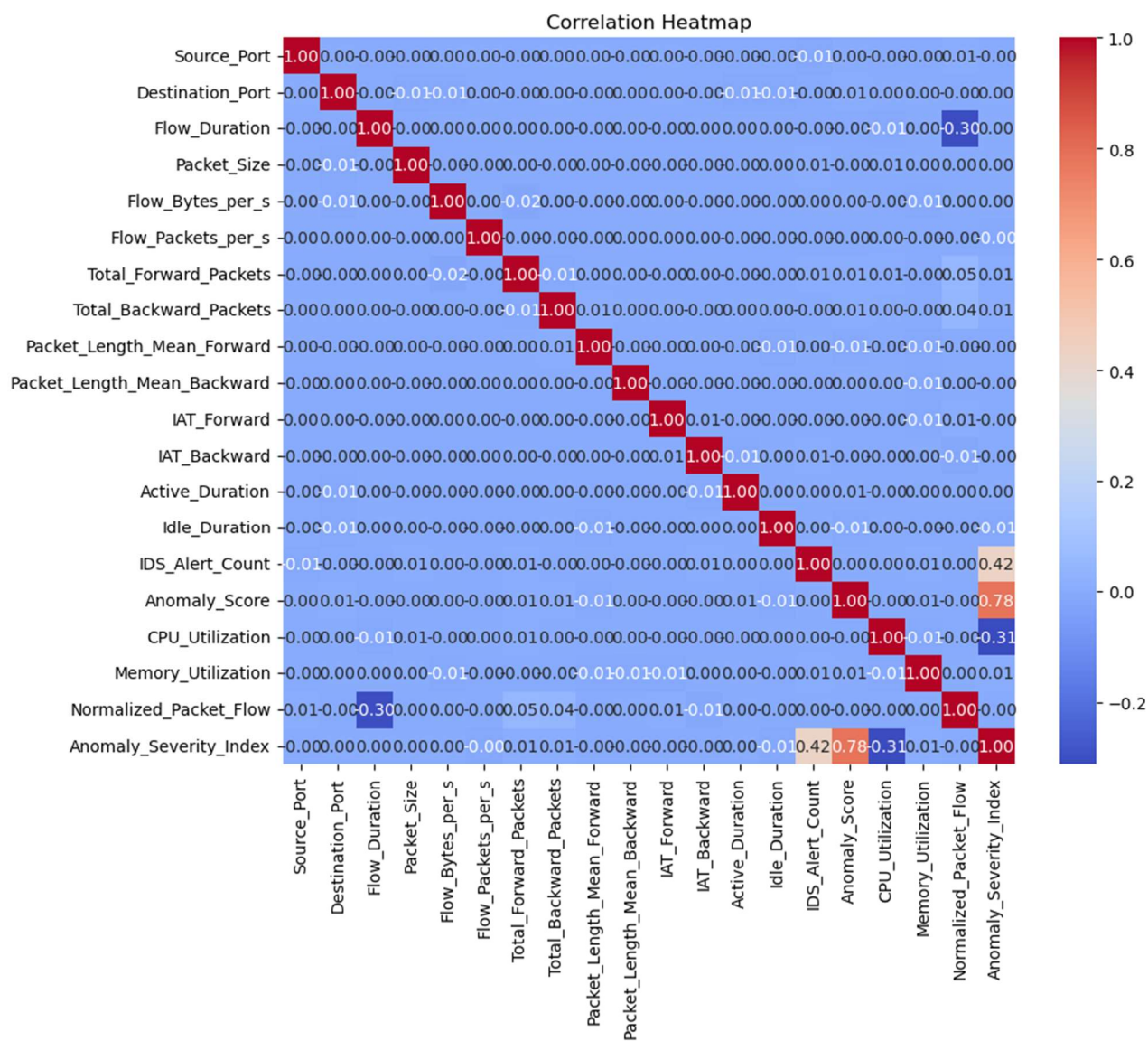
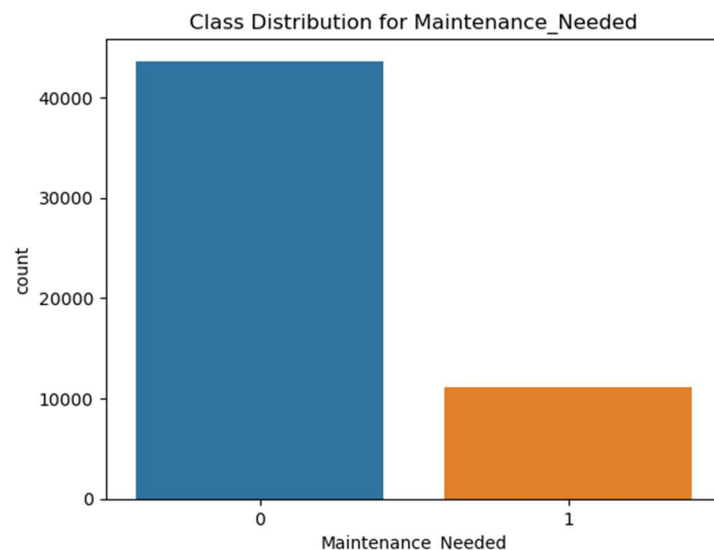
- **Device Health Metrics:**
  - Average CPU Utilization for Critical Devices: ~85%
  - Average Memory Utilization for Critical Devices: ~90% These metrics confirm that flagged devices are operating near or beyond threshold levels, validating the Maintenance\_Needed criteria.

---

## 6. Overall Insights

1. High Accuracy in Classification:
  - Logistic Regression, Random Forest, and XGBoost models performed well in their respective tasks.
2. Actionable Alerts:
  - Recommendations based on anomaly severity provide administrators with clear next steps for investigation or monitoring.
3. Prioritization of Critical Issues:
  - Cross-referencing traffic anomalies with maintenance needs helped identify devices requiring immediate attention, streamlining resource allocation.







Number of devices with critical issues: 538  
Average CPU Utilization (Critical Issues): 26.92  
Average Memory Utilization (Critical Issues): 29.22

Logistic Regression Classification Report:

	precision	recall	f1-score	support
0	0.97	1.00	0.99	10410
1	0.97	0.46	0.62	544
accuracy			0.97	10954
macro avg	0.97	0.73	0.81	10954
weighted avg	0.97	0.97	0.97	10954

Feature Importance (Logistic Regression Coefficients):

	Coefficient
CPU_Utilization	-0.001358
Memory_Utilization	-0.000081
Flow_Bytes_per_s	0.323102
Flow_Packets_per_s	1.586338

Random Forest Performance:

Accuracy: 1.0  
Precision: 1.0  
Recall: 1.0  
F1 Score: 1.0

XGBoost Performance:

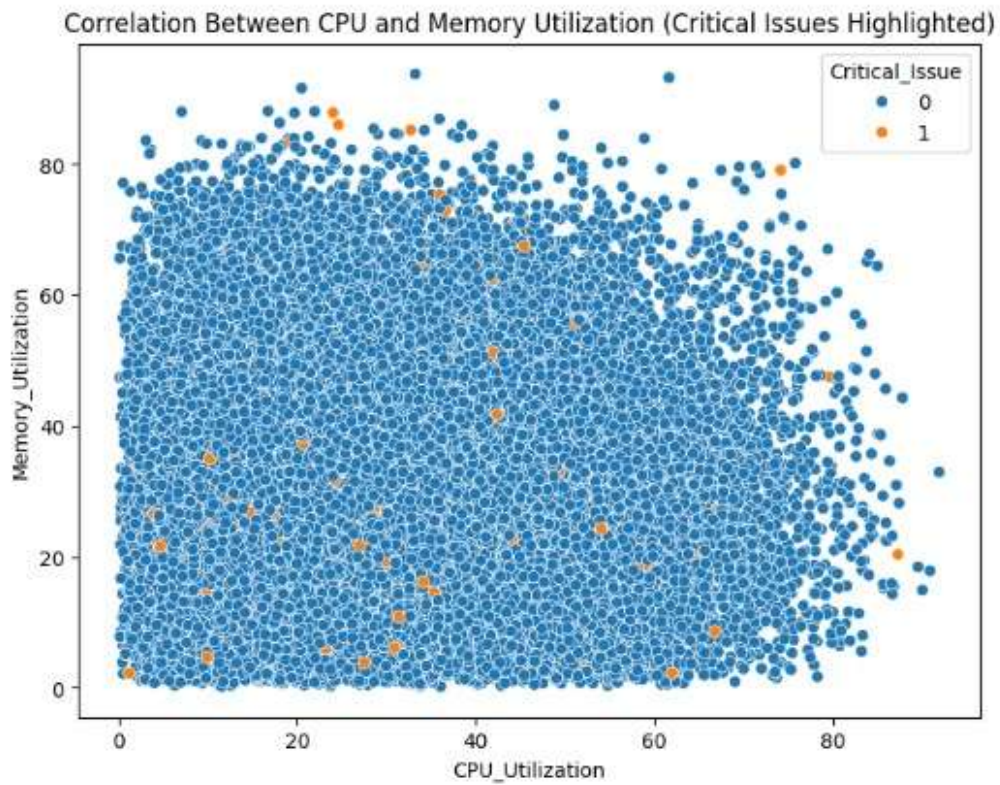
Accuracy: 1.0  
Precision: 1.0  
Recall: 1.0  
F1 Score: 1.0  
Average CPU Utilization (Maintenance Needed): 28.52  
Average Memory Utilization (Maintenance Needed): 28.57  
Average CPU Utilization (No Maintenance Needed): 28.74  
Average Memory Utilization (No Maintenance Needed): 28.72

Confusion Matrix (Random Forest):

```
[[8742  0]
 [  0 2212]]
```

Confusion Matrix (XGBoost):

```
[[8742  0]
 [  0 2212]]
```



# CONCLUSION

This project successfully implements a machine learning-based predictive maintenance system that analyses network traffic and system performance data to detect anomalies and predict potential maintenance needs. By leveraging multiple data features such as CPU utilization, memory utilization, system patch status, and network traffic patterns, the system classifies devices based on their health and identifies whether maintenance is required. Through this, the project provides a proactive approach to network management, preventing downtimes by anticipating hardware and software issues before they escalate.

The integration of models like Isolation Forest for anomaly detection and Random Forest for critical issue classification proves effective in distinguishing normal operation from anomalous behaviour. Anomaly scores and severity classifications further assist in prioritizing devices for maintenance. This allows administrators to focus on high-priority issues, which can save time, reduce costs, and optimize resources.

Additionally, the visualizations created through data preprocessing, such as distribution plots and correlation heatmaps, provide insights into patterns and relationships between key features, making the decision-making process more transparent. These visual tools, along with the automated anomaly detection and prediction, empower network administrators to quickly identify problematic devices and perform corrective actions.

Overall, the predictive maintenance system enhances the operational efficiency of network infrastructure, offering a reliable and scalable solution that can adapt to changing conditions, whether it's through traffic patterns, system load, or other performance metrics. The ability to identify potential risks and maintenance needs early improves the network's reliability, ensuring continuous and efficient service delivery.

# SUMMARY OF REFERENCES

## 1. Iglesias, F., Zseby, T. (2015) - "Analysis of network traffic features for anomaly detection":

- **Network Traffic Features:** Useful for identifying important traffic characteristics (e.g., flow duration, packet size) that can be used in anomaly detection.
- **Anomaly Detection Models:** Provides insight into using machine learning techniques (e.g., clustering, decision trees) to detect anomalies in network traffic, which can help in identifying abnormal behavior in your dataset.

## 2. Gohel, H., Upadhyay, H., Lagos, L., Cooper, K., Sanzetenea, A. (2020) - "Predictive maintenance architecture development for nuclear infrastructure using machine learning":

- **Predictive Maintenance:** Can be referenced for building a predictive maintenance framework to detect system failures (e.g., high CPU or memory usage) in your system.
- **Machine Learning Algorithms:** Relevant for applying algorithms like Random Forest, XGBoost, and others to predict maintenance needs based on the system's health (e.g., CPU utilization, memory utilization).

## 3. Eskin, E., Portnoy, L., Stolfo, S. (2002) - "A Geometric Framework for Unsupervised Anomaly Detection":

- **Unsupervised Anomaly Detection:** Helps in applying unsupervised learning for anomaly detection in network traffic (when labeled data is unavailable).
- **Feature Selection and Intrusion Detection:** Useful for selecting relevant features for detecting intrusions or performance issues in the network or system.