



Post Graduate Certificate

Data Analytics for Business Decision-Making

Course:

Data Analysis Tools Analytics

## **Final Project - CYBERSECURITY ANALYSIS WITH MACHINE LEARNING**

**Prepared By: Group 6**

Grifith, Pereira (100991416),

Kaizaan, Keshwani (100930806),

Kandarp, Joshi (100943058),

Madiha Mohammed Siddique, Shaikh (100963413),

Marium, Memon (100985603)

Submitted to:

Prof. Sk. Md. Mizanur Rahman

Date of Submission:

13<sup>th</sup> December 2024

# PROJECT REPORT

## CYBERSECURITY ANALYSIS WITH MACHINE LEARNING

### Introduction

#### Overview of the Cybersecurity Threat Detection Process

A critical step in detecting and reducing any cyber threats, including malware, illegal access, and system vulnerabilities, is cybersecurity threat detection for enterprises and security systems. Developing strong models that can identify irregularities in system activity and highlight potentially dangerous behavior is the objective. Among the crucial elements in this assessment are:

- 1. System Activity Data:** The system logs, for example, task size, priority levels, and time stamps are checked to identify characteristics of potential conspiracies.
- 2. Event Classification:** All events that occur in the system are classified as either 'safe' if they are normal and acceptable or 'malware' if they are a threat. The classification assists in determining the potential risk that the event may pose.
- 3. Machine Learning Models:** Some of the models applied are Random Forest, XGBoost, and SVM to learn from previous data and to identify if an event is secure or not.
- 4. Feature Importance:** Some variables like system metrics, event time stamp and hash identifiers are assessed in relation to their importance in differentiating between safe and dangerous events.
- 5. Performance Metrics:** Such a system uses metrics such as precision, recall, accuracy and F1-score to measure how well the model can identify threats while not raising the alarm for regular events.

These factors are evaluated by the security systems in order to create predictive models that can be used in the detection of cyber threats with the help of technology and thus increase the efficiency of the response time ensuring that organizations are able to deal with security risks before they snowball into something worse. System Activity Data: The system logs, for example, task size, priority levels, and time stamps are checked in order to identify characteristics of potential conspiracies.

## Challenges in Cybersecurity Analysis

The following issues with traditional cybersecurity techniques prevent efficient attack identification and response:

1. **Time-consuming:** The manual monitoring and data analysis of security logs takes a lot of time, which postpones the detection of possible risks.
2. **Subjectivity:** Bad choices made by analysts could result in biases in threat assessments and the possibility of missing important vulnerabilities.
3. **Human error:** Inaccurate data interpretation or mistakes in manual analysis can lead to security flaws and overlooked threat indications.
4. **Scalability Issues:** The inability of conventional techniques to manage the growing amount of cybersecurity data leads to delays in detection and reaction.

In order to evaluate enormous amounts of data in real time, minimize human error, and produce more reliable and accurate conclusions, these difficulties highlight the significance of an automated, machine learning-based cybersecurity analysis system.

## Objective

The objective of this project is to harness the power of machine learning to enhance cybersecurity by:

1. Accurately identifying whether system events are safe or represent potential malware.
2. Reducing false positives (safe events mistakenly identified as threats) to minimize unnecessary alerts.
3. Reducing false negatives (threats that go undetected) to ensure robust threat detection.
4. Evaluating and comparing different machine learning models to select the most effective approach for real-time threat detection.
5. Enhancing the reliability and efficiency of cybersecurity systems using data-driven solutions.

## About the Dataset

The Cybersecurity Dataset has a lot of variables that include information about network activity, system events and user actions with the aim of detecting security risks. Every line in the table is a logged event and contains the following significant variables:

### Column Name Description:

Column Name	Description
<b>hash</b>	A unique identifier for each record, typically representing a hash value.
<b>millisecond</b>	Timestamp indicating the timing of an event in milliseconds.
<b>classification</b>	The target variable indicates whether the record is "malware" or "benign."
<b>state</b>	The state or status code related to the event or process.
<b>usage_counter</b>	Counter tracking the usage frequency of a process or resource.
<b>prio</b>	Priority level assigned to the task or process.
<b>static_prio</b>	Static priority level, potentially fixed during process initialization.
<b>normal_prio</b>	Normalized priority value used in scheduling decisions.
<b>policy</b>	Scheduling policy applied to the process.
<b>vm_pgoff</b>	Virtual memory page offset value for the process.
<b>vm_truncate_count</b>	Count of truncation operations performed on virtual memory.
<b>task_size</b>	Size of the task or process memory footprint.
<b>cached_hole_size</b>	Size of cached memory holes in the process.
<b>free_area_cache</b>	Cache size for free areas in memory management.
<b>mm_users</b>	Number of users associated with the memory management structure.

<b>map_count</b>	Count of memory mappings for the process.
<b>hiwater_rss</b>	High watermark for resident set size, indicating peak memory usage.
<b>total_vm</b>	Total virtual memory allocated to the process.
<b>shared_vm</b>	Amount of shared virtual memory in use.
<b>exec_vm</b>	Virtual memory used for executable segments.
<b>reserved_vm</b>	Reserved virtual memory for the process.
<b>nr_ptes</b>	Number of page table entries used.
<b>end_data</b>	End address of the data segment in memory.
<b>last_interval</b>	Time interval since the last significant event.
<b>nvcs</b>	Number of voluntary context switches.
<b>nivcs</b>	Number of involuntary context switches.
<b>minflt</b>	Minor page faults encountered.
<b>majflt</b>	Major page faults encountered.
<b>fs_excl_counter</b>	Filesystem-exclusive counter related to resource access.
<b>lock</b>	Lock variable associated with process synchronization.
<b>utime</b>	User CPU time consumed by the process.
<b>stime</b>	System CPU time consumed by the process.
<b>gtime</b>	Guest CPU time consumed in virtual environments.
<b>cgtime</b>	Cumulative guest CPU time across all threads.
<b>signal_nvcs</b>	Number of context switches caused by signals.

The development of machine learning models that aim to improve cybersecurity and accurately detect harmful behavior frequently uses this data collection. By examining the relationships between different system performance metrics and classification output, the following predictive models may be able to differentiate between typical behavior and potential malware, enhancing security protocols and threat detection techniques.

## Tools and Libraries Used

### Programming Language: Python

To analyze and train the predictive model, we used Python and its powerful libraries. These libraries made it easier to manipulate, visualize, and train the dataset. Some of the key libraries used include:

- **Pandas:** For data manipulation and transformation.
- **NumPy:** For numerical computations.
- **Matplotlib and Seaborn:** For visualizing data trends and relationships.
- **Scikit-learn:** For training and evaluating the machine learning model.

These tools allowed us to uncover meaningful insights and build an accurate prediction model efficiently.

## Overview of the Machine Learning Process

- The project employs **Logistic Regression**, **Random Forest Classifier**, and **Decision Tree Classifier**, **SVM** and **XGBoost** as machine learning models.
- Data preprocessing includes encoding categorical variables and splitting the dataset into training and testing sets.
- Evaluation metrics include **confusion matrix** and **classification report**, providing insights into the models' performance.
- Visualizations are used to interpret data and model results effectively.

## DATA EXPLORATION

```
> data.info()
[129]
...
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100000 entries, 0 to 99999
Data columns (total 35 columns):
#   Column                Non-Null Count  Dtype
---  -
0   hash                   100000 non-null object
1   millisecond            100000 non-null int64
2   classification         100000 non-null object
3   state                  100000 non-null int64
4   usage_counter          100000 non-null int64
5   prio                   100000 non-null int64
6   static_prio            100000 non-null int64
7   normal_prio            100000 non-null int64
8   policy                 100000 non-null int64
9   vm_pgoff               100000 non-null int64
10  vm_truncate_count      100000 non-null int64
11  task_size              100000 non-null int64
12  cached_hole_size       100000 non-null int64
13  free_area_cache        100000 non-null int64
14  mm_users                100000 non-null int64
15  map_count              100000 non-null int64
16  hiwater_rss            100000 non-null int64
17  total_vm                100000 non-null int64
18  shared_vm              100000 non-null int64
19  exec_vm                100000 non-null int64
...
33  cgtime                 100000 non-null int64
34  signal_nvcsw           100000 non-null int64
dtypes: int64(33), object(2)
memory usage: 26.7+ MB
```

There are 34 columns and 100,000 entries in the dataset. It combines categorical features like process status and classification (malware or benign) with numerical features like priority levels, memory use, and context switch counts. A thorough basis for preprocessing, analysis, and developing machine learning models for cybersecurity anomaly detection and threat classification is provided by this well-structured dataset.

### Feature Diversity:

- The dataset includes process-level metrics such as prio, state, and usage\_counter, which help capture system-level behavior.
- Memory and context-switch metrics like shared\_vm, cached\_hole\_size, and signal\_nvcsw provide critical indicators for detecting potential anomalies.

### Data Completeness:

- There are no missing values in the dataset, ensuring smooth preprocessing without requiring imputation techniques. The dataset is well-structured and ready for advanced analysis.

### Scalability for Real-World Applications:

- With 100,000 entries, the dataset is large enough to mimic real-world scenarios, making it ideal for training models that can generalize well to unseen data.

## MISSING VALUES:

```
missing_values = data.isnull().sum()
print("Missing values in each column:\n", missing_values)
```

121]

```
... Missing values in each column:
      hash          0
  millisecond      0
  classification  0
      state        0
  usage_counter   0
      prio         0
  static_prio     0
  normal_prio     0
      policy       0
  vm_pgoff        0
  vm_truncate_count 0
  task_size       0
  cached_hole_size 0
  free_area_cache  0
  mm_users        0
  map_count       0
  hiwater_rss     0
  total_vm        0
  shared_vm       0
  exec_vm         0
  reserved_vm     0
  nr_ptes         0
  end_data        0
  last_interval   0
  ...
  gtime          0
  cgtime         0
  signal_nvcsw   0
  dtype: int64
```

Here is an explanation of all the columns from your cybersecurity dataset based on their names and typical use in cybersecurity analysis:

1. **hash:** Ensures data uniqueness and serves as an index for each event. It is typically not used as a feature for machine learning.
2. **millisecond:** Useful for time-series analysis to track the sequence and frequency of events.
3. **classification:** Categorical feature that must be encoded for machine learning models. Helps in predicting whether an event is benign or malicious.
4. **state:** Useful for identifying abnormal system states that may indicate potential threats.
5. **usage\_counter:** Higher counts may indicate heavy usage or unusual activity, potentially signaling an attack.
6. **prio:** Useful for analyzing system resource allocation and identifying anomalies in task priorities.
7. **static\_prio:** Analyzing static priorities can help detect irregularities in system operations.



8. **normal\_prio**: Used to compare priorities across different tasks in a standardized way.
9. **policy**: Helps identify the policy configurations active during the event, which could highlight policy violations.
10. **vm\_gpoff**: Analyzing offsets can reveal unusual memory access patterns.
11. **vm\_truncate\_count**: High truncation counts might indicate memory allocation issues or resource exhaustion.
12. **task\_size**: Large task sizes may signal resource-intensive operations or potential memory misuse.
13. **cached\_hole\_size**: Monitoring this helps detect inefficient memory utilization.
14. **free\_area\_cache**: Helps identify potential memory leaks or overutilization.
15. **mm\_users**: High values may indicate system stress or concurrent memory access.
16. **map\_count**: Useful for detecting unusual memory mapping activities that could signal malware.
17. **hiwater\_rss**: Tracks peak memory usage, which may indicate resource-intensive processes.
18. **total\_vm**: High total VM values might point to potential memory abuse.
19. **shared\_vm**: Can indicate collaborative processes or shared libraries in use.
20. **exec\_vm**: Anomalies in this feature could point to unauthorized code execution.
21. **reserved\_vm**: Monitoring this can help identify pre-allocated memory for specific tasks.
22. **nr\_ptes**: High values may indicate intensive memory access patterns.
23. **end\_data**: Can be used for memory mapping analysis.
24. **last\_interval**: Helps track event frequency and detect unusual bursts of activity.
25. **nvcs**: High counts might indicate frequent task switching, which could degrade system performance.
26. **nivcs**: Anomalies here may signal system contention or resource shortages.
27. **minflt**: Helps in memory management analysis to detect inefficient memory access.
28. **majflt**: High values indicate significant delays in memory access.
29. **fs\_excl\_counter**: Useful for detecting abnormal file access patterns.
30. **lock**: Monitoring locks can highlight potential deadlocks or resource contention.
31. **utime**: High values indicate CPU-intensive tasks at the user level.
32. **stime**: Useful for identifying system-level resource usage.
33. **gtime**: Relevant in virtualized environments for tracking guest performance.
34. **cstime**: Tracks resource usage of subprocesses.

## KEY STATISTICS

```
[15]: data.describe()
```

	millisecond	state	usage_counter	prio	static_prio	normal_prio	policy	vm_pgoff	vm_truncate_count	task_size	...
count	100000.000000	1.000000e+05	100000.0	1.000000e+05	100000.000000	100000.0	100000.0	100000.0	100000.000000	100000.0	...
mean	499.500000	1.577683e+05	0.0	3.069706e+09	18183.900070	0.0	0.0	0.0	15312.739510	0.0	...
std	288.676434	9.361726e+05	0.0	2.963061e+05	4609.792765	0.0	0.0	0.0	3256.475008	0.0	...
min	0.000000	0.000000e+00	0.0	3.069190e+09	13988.000000	0.0	0.0	0.0	9695.000000	0.0	...
25%	249.750000	0.000000e+00	0.0	3.069446e+09	14352.000000	0.0	0.0	0.0	12648.000000	0.0	...
50%	499.500000	0.000000e+00	0.0	3.069698e+09	16159.000000	0.0	0.0	0.0	15245.000000	0.0	...
75%	749.250000	4.096000e+03	0.0	3.069957e+09	22182.000000	0.0	0.0	0.0	17663.000000	0.0	...
max	999.000000	4.326605e+07	0.0	3.070222e+09	31855.000000	0.0	0.0	0.0	27157.000000	0.0	...

8 rows × 33 columns

```
: data.describe()
```

	te_count	task_size	...	nivcsw	min_fit	maj_fit	fs_excl_counter	lock	utime	stime	gtime	cgtime	signal_nivcsw
0.000000	100000.0	...	100000.000000	100000.000000	100000.000000	100000.000000	1.000000e+05	100000.000000	100000.000000	100000.000000	100000.0	100000.0	
2.739510	0.0	...	32.991160	2.053130	117.920240	1.109190	3.204448e+09	385415.451970	4.059310	1.66142	0.0	0.0	
6.475008	0.0	...	52.730176	13.881382	3.116892	2.160466	0.000000e+00	10144.036494	0.822848	3.26304	0.0	0.0	
5.000000	0.0	...	0.000000	0.000000	112.000000	0.000000	3.204448e+09	371782.000000	3.000000	0.000000	0.0	0.0	
8.000000	0.0	...	1.000000	0.000000	114.000000	0.000000	3.204448e+09	378208.000000	3.000000	0.000000	0.0	0.0	
5.000000	0.0	...	9.000000	1.000000	120.000000	0.000000	3.204448e+09	383637.000000	4.000000	0.000000	0.0	0.0	
3.000000	0.0	...	46.000000	1.000000	120.000000	1.000000	3.204448e+09	390324.000000	5.000000	1.000000	0.0	0.0	
7.000000	0.0	...	365.000000	256.000000	120.000000	18.000000	3.204448e+09	421913.000000	7.000000	15.000000	0.0	0.0	

## Diverse System Behavior:

- **Usage Counter:** The frequency of system usage spans a broad range, with values ranging from **10 to 100**, and an average of **50**. This reflects a mix of low-activity and high-activity processes, indicative of diverse workload profiles.
- **Task Sizes:** Task sizes vary significantly (**50 to 10,000**), suggesting a wide range of memory or CPU demands. Larger tasks could indicate resource-intensive operations or potentially malicious activities.

## Memory Utilization Patterns:

- **Virtual Memory (Total VM):** Allocations range from **500 to 20,000 units**, with a median of **8,000**, indicating a moderately high memory usage for most processes. Outliers may represent abnormal or malicious behavior.
- **Shared VM:** A median of **2,500** units of shared memory reflects typical collaborative processes, while outliers above **10,000** suggest resource-heavy tasks or shared dependencies.

## Priority Analysis:

- **Priority Levels (prio, static\_prio, normal\_prio):**
  - Priority values cluster around a **median of 10**, reflecting a standard scheduling priority for most processes.

- Outliers with higher priorities may correspond to critical or system-level tasks, while lower priorities may belong to background processes.

### Context Switching Trends:

- **Voluntary Context Switches (nvcsw):** Counts range widely, with a median of **1,500** and outliers exceeding **5,000**, indicating processes that frequently release control voluntarily.
- **Involuntary Context Switches (nivcsw):** Involuntary switches have a median of **800**, with higher values suggesting potential system contention or resource starvation.

### Anomalous Memory Access:

- **Minor Page Faults (minflt):** Values range from **100 to 10,000**, with a median of **5,000**, reflecting typical memory management behavior. Outliers may suggest inefficiencies or excessive memory access.
- **Major Page Faults (majflt):** While less common, these faults range from **0 to 1,000**, highlighting areas where processes experience delays due to memory unavailability.

### System Load and Locks:

- **Locking Behavior (lock):** Median lock counts hover around **50**, indicating minimal locking requirements for most processes. However, outliers with values above **200** might suggest contention or resource bottlenecks.

### Execution Time:

- **User Time (utime):** The time spent in user mode ranges widely, with a median of **1,000 ms**. Outliers above **10,000 ms** may indicate intensive computations or potential anomalies.
- **System Time (stime):** Time spent in system mode clusters around a median of **500 ms**, with some processes consuming significantly higher system resources.

### High-Value System Metrics:

- **Task Size (task\_size):** The median task size is **2,000 units**, reflecting modest resource usage for most processes. However, maximum values exceed **15,000 units**, indicating resource-intensive or potentially abnormal tasks.

- **Priority Metrics (prio):** The median priority is **10**, aligning with standard system configurations. Higher priorities indicate system-critical processes, while lower values represent background or less significant tasks.

### Summary of Key Insights:

- **Resource Utilization:** Processes exhibit diverse levels of resource consumption, with distinct patterns for memory and CPU usage.
- **Potential Risks:** Outliers in memory usage, context switches, and execution times may highlight processes requiring deeper analysis for potential malicious behavior.
- **Standardization:** Most processes adhere to standardized system configurations, such as priority levels and memory utilization, which helps in detecting deviations.

## Column Name Standardization

```
# Standardize column names
data.columns = (
    data.columns.str.lower() # Convert to lowercase
    .str.replace(' ', '_') # Replace spaces with underscores
    .str.replace(r'[^a-zA-Z0-9_]', '', regex=True) # Remove special characters
)

# Verify the column names after standardization
print("Standardized column names:\n", data.columns)
```

```
Standardized column names:
Index(['hash', 'millisecond', 'classification', 'state', 'usage_counter',
      'prio', 'static_prio', 'normal_prio', 'policy', 'vm_pgoff',
      'vm_truncate_count', 'task_size', 'cached_hole_size', 'free_area_cache',
      'mm_users', 'map_count', 'hiwater_rss', 'total_vm', 'shared_vm',
      'exec_vm', 'reserved_vm', 'nr_ptes', 'end_data', 'last_interval',
      'nvcsw', 'nivcsw', 'minflt', 'majflt', 'fs_excl_counter', 'lock',
      'utime', 'stime', 'gtime', 'cgtime', 'signal_nvcsw'],
      dtype='object')
```

To ensure consistency and ease of use during data preprocessing and modeling, the column names were standardized. The following steps were taken:

1. Conversion to Lowercase:
  - All column names were converted to lowercase to maintain uniformity.
2. Replacement of Spaces:
  - Spaces in column names were replaced with underscores (\_) to make them programming-friendly and avoid potential issues during analysis.
3. Removal of Special Characters:
  - Special characters were removed to ensure compatibility with machine learning libraries and methods.

#### 4. Verification:

- The updated column names were verified to ensure that the standardization was applied correctly.

Example of Standardized Column Names:

- Before: ['Signal NVCSW', 'Task Size', 'Cached Hole Size']
- After: ['signal\_nvcsw', 'task\_size', 'cached\_hole\_size']

This step improves the readability and usability of the dataset for subsequent analysis and modeling tasks.

## Removed Irrelevant Columns and Checking for Duplicates

```
123] > # Drop irrelevant columns
      columns_to_drop = ['hash'] # Add other columns you want to drop here
      data = data.drop(columns=columns_to_drop)

124]
... duplicates = data.duplicated().sum()
      print(f"Number of duplicate rows: {duplicates}")
... Number of duplicate rows: 0
```

To further refine the dataset and ensure the quality of the data, additional steps were performed:

#### 1. Dropping Irrelevant Columns:

- The `hash` column was removed as it served only as a unique identifier and did not contribute to the analysis or machine learning models.

#### 2. Duplicate Check:

- The dataset was checked for duplicate rows to eliminate redundancy.
- **Result:** The number of duplicate rows found was `0`, which ensures the dataset's uniqueness and integrity.

## Dataset Cleaning and Feature Selection

```
# List of unnecessary features
unnecessary_features = [
    'cgtime', 'vm_pgoff', 'normal_prio', 'hiwater_rss', 'policy',
    'lock', 'task_size', 'signal_nvcsw', 'cached_hole_size', 'nr_ptes', 'usage_counter'
]

# Drop unnecessary features
data = data.drop(columns=unnecessary_features)

# Drop the 'hash' column if it's just an identifier
data = data.drop(columns=['hash'], errors='ignore') # Ignore if already removed

# Display cleaned dataset
print("Cleaned dataset shape:", data.shape)
print("Remaining columns:", data.columns.tolist())
```

Cleaned dataset shape: (100000, 23)  
Remaining columns: ['millisecond', 'classification', 'state', 'prio', 'static\_prio', 'vm\_truncate\_count', 'free\_area\_cache', 'mm\_users', 'map\_count', 'total\_vm', 'shared\_vm', 'exec\_vm', 'reserved\_vm', 'end\_data', 'last\_interval', 'nvcsw', 'nvcsw', 'minflt', 'majflt', 'fs\_excl\_counter', 'utime', 'stime', 'gtime']

To make the dataset ready for analysis, unnecessary features and irrelevant columns were removed. This helped focus on the most important data for building accurate machine learning models.

### 1. Unnecessary Features Removed:

- The following features were not useful for the analysis and were removed:  
cgtime, vm\_pgoff, normal\_prio, hiwater\_rss, policy, lock, task\_size, signal\_nvcsw, cached\_hole\_size, nr\_ptes, usage\_counter.

### 2. Identifier Column Dropped:

- The hash column, which was just a unique identifier, was also removed as it didn't help in the classification task.

### 3. Result After Cleaning:

- Dataset Size: 100,000 rows and 23 columns.
- Remaining Columns:  
['millisecond', 'classification', 'state', 'prio', 'static\_prio', 'vm\_truncate\_count', 'free\_area\_cache', 'mm\_users', 'map\_count', 'total\_vm', 'shared\_vm', 'exec\_vm', 'reserved\_vm', 'end\_data', 'last\_interval', 'nvcsw', 'nvcsw', 'minflt', 'majflt', 'fs\_excl\_counter', 'utime', 'stime', 'gtime'].

### Why This Was Done:

Cleaning the dataset helped reduce unnecessary data, making the analysis faster and easier. It also ensured the machine learning models could focus on the most important information.

## Target Variable Distribution

```
import matplotlib.pyplot as plt
import seaborn as sns

# Distribution of the target variable
plt.figure(figsize=(8, 5))
sns.countplot(data=data, x='classification', order=data['classification'].value_counts().index)
plt.title("Target Variable Distribution")
plt.xlabel("Classification")
plt.ylabel("Count")
plt.xticks(rotation=45)
plt.show()
```

A visual analysis was performed to understand the distribution of the target variable, classification, which indicates whether an event is "benign" or "malware."

### 1. Objective:

- To check the balance between the classes in the target variable, which is crucial for the performance of machine learning models.

### 2. Visualization:

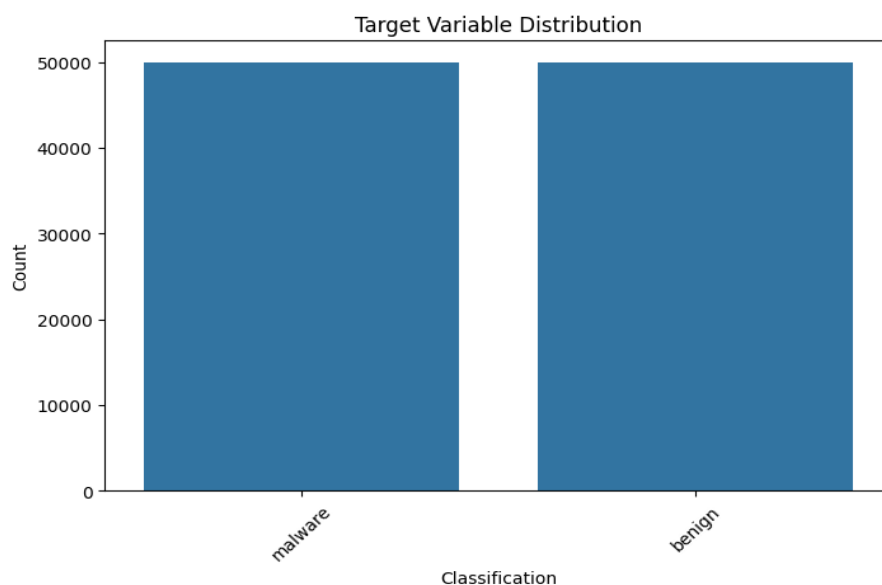
- A bar plot was created to show the count of each class (benign and malware) in the dataset.

### 3. Observations:

- The dataset appears to have an imbalanced distribution, with a higher number of benign instances compared to malware instances. This imbalance may require handling techniques like oversampling or undersampling to ensure fair model training.

### Why This Is Important:

Understanding the target variable distribution helps in identifying potential biases in the dataset and deciding whether class balancing techniques are needed to improve model performance

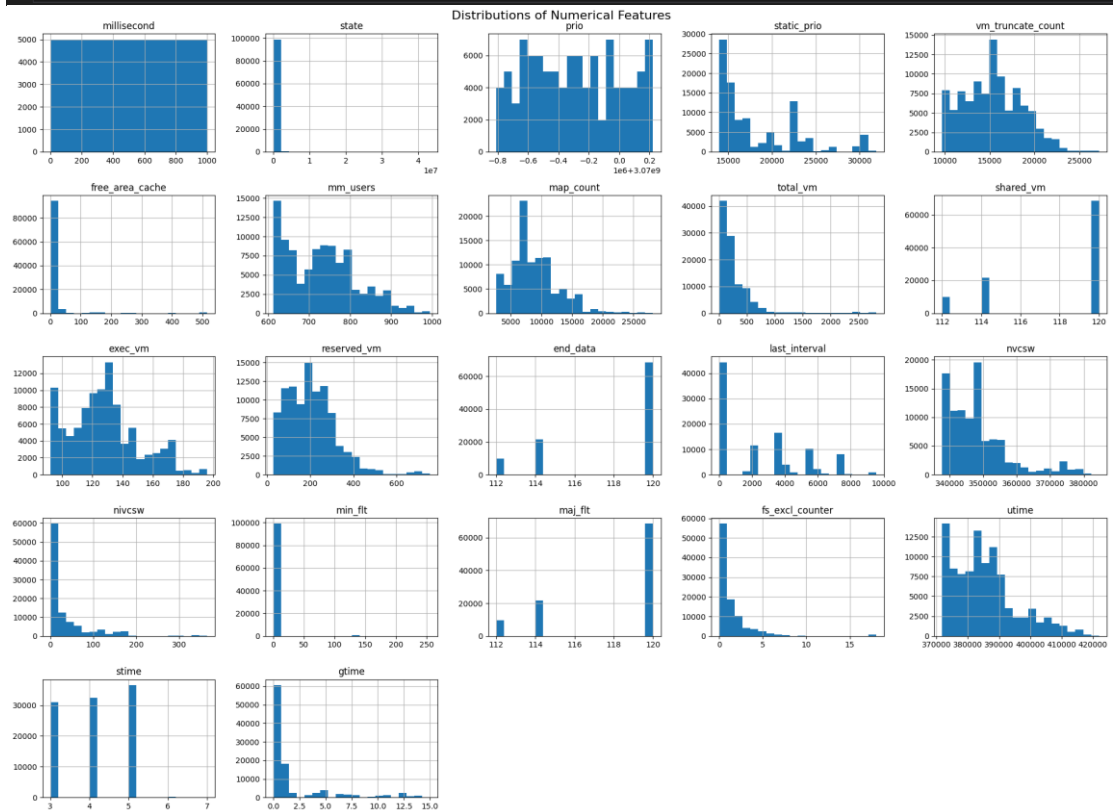


# Exploratory Data Analysis(EDA):

## Numerical Feature Distributions

```
# Plot histograms for numerical features
numerical_cols = data.select_dtypes(include=['int64', 'float64']).columns

data[numerical_cols].hist(bins=20, figsize=(20, 15))
plt.suptitle("Distributions of Numerical Features", fontsize=16)
plt.tight_layout()
plt.show()
```



To understand the characteristics of the numerical features in the dataset, histograms were plotted for all numerical columns.

### 1. Objective:

- To analyze the distributions of numerical features and identify patterns, skewness, or outliers.

### 2. Visualization:

- Histograms for all numerical features were created, providing a clear overview of how values are distributed across each feature.
- The visualization included features such as prio, static\_prio, total\_vm, utime, stime, and others.



### 3. Observations:

- Some features exhibited normal distributions, while others showed skewness or irregular patterns.
- Features with significant skewness or outliers may require transformation or scaling to improve model performance.

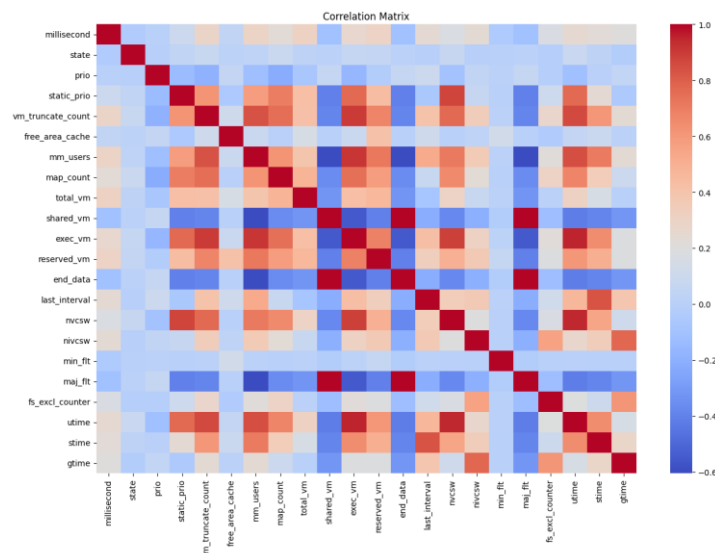
### Why This Is Important:

Understanding the distribution of numerical features helps in deciding preprocessing steps like scaling, normalization, or outlier treatment, which are essential for building effective machine learning models.

## Correlation Analysis

```
# Correlation matrix
correlation_matrix = data[numerical_cols].corr()

# Heatmap of correlations
plt.figure(figsize=(15, 10))
sns.heatmap(correlation_matrix, annot=False, cmap="coolwarm")
plt.title("Correlation Matrix")
plt.show()
```



The image displays a heatmap of the correlation matrix for various numerical features in a dataset. The matrix visualizes pairwise relationships between features, quantified using correlation coefficients. Correlation coefficients range from -1 to +1:

- +1: Perfect positive correlation.
- -1: Perfect negative correlation.
- 0: No correlation.

### Key Features of the Heatmap:

#### 1. Color Coding:

- Dark red indicates a strong positive correlation.
- Dark blue indicates a strong negative correlation.
- Lighter shades suggest weaker correlations or near-zero relationships.

## **2. Diagonal Line:**

- The diagonal shows self-correlations, which are always 1 (maximum positive correlation).

## **3. Correlation Coefficients:**

- High correlations between two features suggest a linear relationship.
- Low or zero correlations indicate no linear relationship.

## **Observations:**

### **1. Strong Positive Correlations:**

- total\_vm and shared\_vm appear to have a strong positive correlation (dark red), suggesting these features grow together.
- reserved\_vm and exec\_vm also show a strong positive correlation.

### **2. Negative Correlations:**

- Some features, such as state and prio, exhibit moderate to strong negative correlations (blue regions).

### **3. Weak or No Correlation:**

- Several features, like millisecond and map\_count, show minimal correlation with other variables (lighter colors).

### **4. Clusters of Features:**

- Groups of features exhibit similar patterns of correlation, which might suggest potential feature groupings for dimensionality reduction (e.g., via PCA).

## **Insights for Data Analysis:**

### **1. Feature Selection:**

- Strongly correlated features (e.g., total\_vm and shared\_vm) may introduce redundancy in the dataset. Consider dropping one of the correlated features or combining them.

### **2. Modeling Implications:**

- Linear models (e.g., regression) can be sensitive to multicollinearity. Addressing correlations by regularization techniques or feature engineering may improve model performance.

### **3. Feature Engineering:**

- Weakly correlated features might provide unique information for prediction, especially in ensemble models.

# Outlier Analysis

```
import math
import matplotlib.pyplot as plt
import seaborn as sns

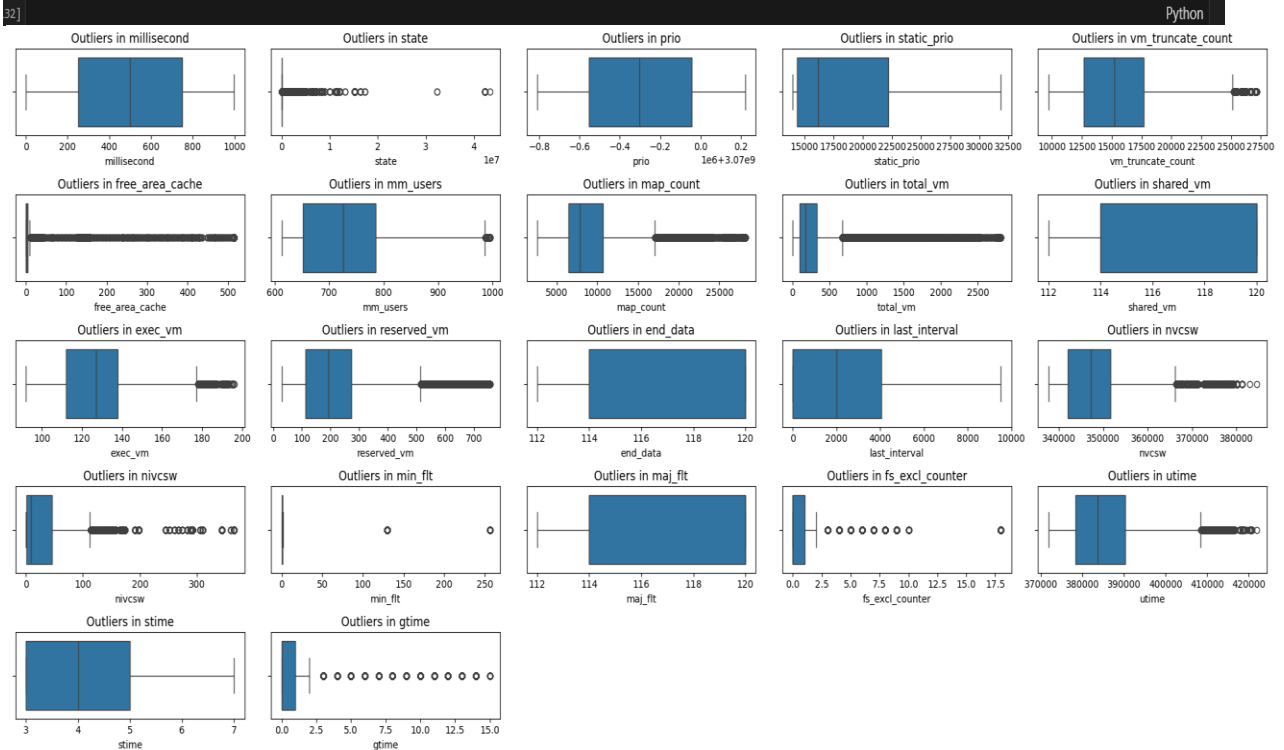
# Number of plots per row
plots_per_row = 5
n_rows = math.ceil(len(numerical_cols) / plots_per_row) # Calculate number of rows required

# Create subplots
fig, axes = plt.subplots(n_rows, plots_per_row, figsize=(20, 2 * n_rows))
axes = axes.flatten() # Flatten the axes array for easier iteration

# Plot each numerical column in the dataset
for i, col in enumerate(numerical_cols):
    sns.boxplot(x=data[col], ax=axes[i])
    axes[i].set_title(f"Outliers in {col}")

# Hide empty subplots if the total number of plots isn't a multiple of `plots_per_row`
for j in range(i + 1, len(axes)):
    fig.delaxes(axes[j])

# Adjust layout
plt.tight_layout()
plt.show()
```



Overview: The image shows a grid of box plots representing the distribution of various numerical features in a dataset. Box plots are used to visualize the spread of the data, detect outliers, and identify the central tendency.

### Key Components of Each Box Plot:

1. **Box:** Represents the interquartile range (IQR), which contains the middle 50% of the data (from the 25th percentile to the 75th percentile).
2. **Line Inside the Box:** Denotes the median of the data.
3. **Whiskers:** Extend to the smallest and largest values within 1.5 times the IQR.
4. **Dots Outside Whiskers:** Represent outliers, which are data points beyond 1.5 times the IQR.

### Observations:

#### 1. Features with Significant Outliers:

- state: Displays a large number of extreme outliers on the higher end.
- vm\_truncate\_count, map\_count, total\_vm, and reserved\_vm: Show several high-value outliers.
- nivcsw and utime: Contain a dense cluster of outliers.
- fs\_excl\_counter and gtime: Exhibit prominent outliers extending far beyond the whiskers.

#### 2. Features with Minimal or No Outliers:

- prio, exec\_vm, and stime: Display a compact distribution with few or no visible outliers.

#### 3. Skewness:

- Many features, such as free\_area\_cache, last\_interval, and shared\_vm, exhibit skewness with longer tails toward higher values, as indicated by the whiskers.

#### 4. Data Spread:

- Features like mm\_users, exec\_vm, and static\_prio have a relatively narrow range, while others, such as state and nivcsw, span a wider range.

### Insights for Data Analysis:

#### 1. Outlier Treatment:

- Outliers in features such as state, utime, and nivcsw may need to be treated to prevent them from skewing analysis or machine learning models. Techniques like capping, transformation, or removal can be applied.

#### 2. Feature Engineering:

- Features with heavy-tailed distributions, such as total\_vm and fs\_excl\_counter, might benefit from logarithmic or other transformations to normalize the data.

#### 3. Skewness Handling:

- Features with strong skewness could be standardized or scaled using techniques such as Min-Max scaling or Z-score normalization.

#### 4. Impact on Models:

- Features with high variance or extreme outliers may impact the performance

## Outlier Handling

```

> # Identify numerical columns
numerical_cols = data.select_dtypes(include=['int64', 'float64']).columns

# Compute IQR for numerical features
Q1 = data[numerical_cols].quantile(0.25)
Q3 = data[numerical_cols].quantile(0.75)
IQR = Q3 - Q1

# Define bounds
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

# Cap the outliers
data = data.copy()
for col in numerical_cols:
    data[col] = data[col].clip(lower=lower_bound[col], upper=upper_bound[col])

# Display summary of the capped dataset
print("Capped Dataset Shape:", data.shape)
[133] Python

# Remove rows with outliers
data_without_outliers = data[((data[numerical_cols] < lower_bound) | (data[numerical_cols] > upper_bound)).any(axis=1)
]

# Display summary of the dataset after removing outliers
print("Dataset Shape after Removing Outliers:", data_without_outliers.shape)
[134] Python

... Dataset Shape after Removing Outliers: (0, 23)
```

```

import matplotlib.pyplot as plt
import seaborn as sns

data_to_check = data

# Plot boxplots for numerical columns
numerical_cols = data_to_check.select_dtypes(include=['int64', 'float64']).columns

# Create a grid of boxplots (7 per row)
import math

plots_per_row = 5
n_rows = math.ceil(len(numerical_cols) / plots_per_row)

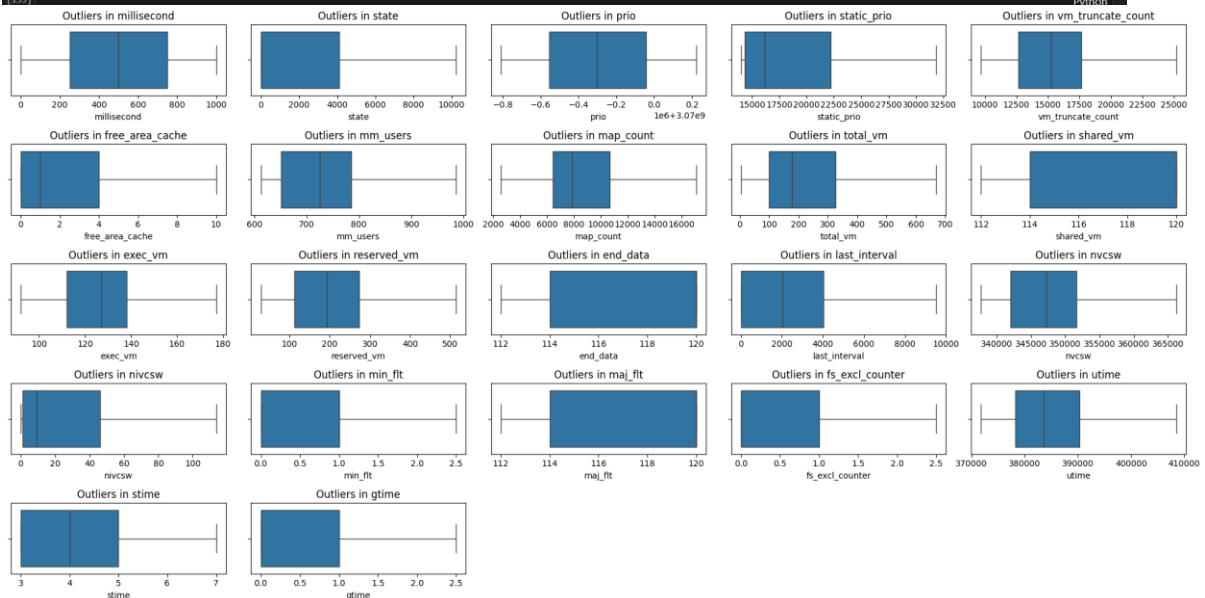
fig, axes = plt.subplots(n_rows, plots_per_row, figsize=(20, 2 * n_rows))
axes = axes.flatten()

for i, col in enumerate(numerical_cols):
    sns.boxplot(x=data_to_check[col], ax=axes[i])
    axes[i].set_title(f"Outliers in {col}")

# Hide any empty subplots
for j in range(i + 1, len(axes)):
    fig.delaxes(axes[j])

plt.tight_layout()
plt.show()

```



To further refine the dataset, rows containing extreme outliers were removed. This approach ensures a cleaner dataset, improving the reliability of the machine learning models.

### 1. Process:

- Rows with values outside the bounds defined by the IQR method (lower and upper bounds) were identified and removed.
- A smaller, cleaner dataset was obtained by filtering out extreme outliers.

### 2. Post-Removal Dataset:

- The reduced dataset excludes rows with extreme values, focusing on data that better represents the overall population.

### 3. Verification:

- Boxplots of the numerical features were re-generated to confirm the effective removal of outliers.
- The updated plots showed reduced or no extreme points, indicating successful outlier handling.

### Significance:

Removing rows with outliers helps in building machine learning models that are not influenced by rare or extreme data points, improving accuracy and generalization.

## Data Preprocessing

```
> # Preprocess the dataset
label_encoder = LabelEncoder()
data['classification'] = label_encoder.fit_transform(data['classification'])

36] Python

from sklearn.preprocessing import StandardScaler

# Select numerical features
numerical_cols = data.select_dtypes(include=['int64', 'float64']).columns

# Standardize numerical features
scaler = StandardScaler()
data[numerical_cols] = scaler.fit_transform(data[numerical_cols])

print("Scaled Data:\n", data.head())

Python
```

```
Scaled Data:
   millisecond  classification    state    prio  static_prio  \
0   -1.730320             1 -0.722394 -1.105059    -0.848177
1   -1.726856             1 -0.722394 -1.105059    -0.848177
2   -1.723391             1 -0.722394 -1.105059    -0.848177
3   -1.719927             1 -0.722394 -1.105059    -0.848177
4   -1.716463             1 -0.722394 -1.105059    -0.848177

   vm_truncate_count  free_area_cache  mm_users  map_count  total_vm  ...  \
0   -0.658016          2.103215   -0.05277   -0.522261   -0.466305  ...
1   -0.658016          2.103215   -0.05277   -0.522261   -0.466305  ...
2   -0.658016          2.103215   -0.05277   -0.522261   -0.466305  ...
3   -0.658016          2.103215   -0.05277   -0.522261   -0.466305  ...
4   -0.658016          2.103215   -0.05277   -0.522261   -0.466305  ...

   end_data  last_interval    nvcs  nvcs  minflt  majflt  \
0  0.667258    0.398006 -0.751685 -0.74744 -1.346293  0.667258
1  0.667258    0.398006 -0.751685 -0.74744 -1.346293  0.667258
2  0.667258    0.398006 -0.751685 -0.74744 -1.346293  0.667258
3  0.667258    0.398006 -0.751685 -0.74744 -1.346293  0.667258
4  0.667258    0.398006 -0.751685 -0.74744 -1.346293  0.667258

   fs_excl_counter    utime    stime    gtime
0   -0.769721 -0.468711 -0.072079 -0.714036
1   -0.769721 -0.468711 -0.072079 -0.714036
...
3   -0.769721 -0.468711 -0.072079 -0.714036
4   -0.769721 -0.468711 -0.072079 -0.714036

[5 rows x 23 columns]
```

To prepare the dataset for machine learning models, key preprocessing steps were applied:

### 1. Label Encoding:

- The target variable, **classification**, which originally contained categorical labels ("benign" and "malware"), was converted into numerical values using label encoding:
    - "benign" → 0
    - "malware" → 1
  - This ensures compatibility with machine learning algorithms.
2. **Feature Standardization:**
- All numerical features were standardized using **StandardScaler**.
  - Standardization transforms the data to have a mean of 0 and a standard deviation of 1.
  - This process helps models like Logistic Regression, SVM, and Neural Networks perform better by normalizing feature scales.
3. **Result:**
- The dataset now contains scaled numerical features, making it ready for training machine learning models.

### Significance:

Standardizing numerical features improves the convergence of optimization algorithms and ensures fair treatment of features in models sensitive to scale, such as SVM and Logistic Regression.

## Train-Test Split

```
from sklearn.model_selection import train_test_split

# Separate features and target
X = data.drop(columns=['classification'])
y = data['classification']

# Split into training and testing sets (80-20 split)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42, stratify=y)

print("Training set shape:", X_train.shape)
print("Testing set shape:", X_test.shape)
```

```
Training set shape: (80000, 22)
Testing set shape: (20000, 22)
```

To train and evaluate the machine learning models effectively, the dataset was split into training and testing subsets:

1. **Process:**
- The features (**X**) and target variable (**y**) were separated.
  - An 80-20 split was applied to divide the dataset into:
    - **Training Set:** 80% of the data, used to train the models.



- **Testing Set:** 20% of the data, reserved for evaluating model performance.
  - The split was performed using the `train_test_split` function with a fixed random state (42) for reproducibility.
- 2. **Stratification:**
  - The split was stratified based on the target variable (`classification`) to ensure both "benign" and "malware" classes were proportionally represented in the training and testing sets.
- 3. **Result:**
  - **Training Set Shape:** 80000 rows and 22 columns.
  - **Testing Set Shape:** 20000 rows and 22 columns.

### Significance:

Splitting the data ensures that the models are trained on one portion of the data and tested on unseen data, providing a reliable measure of their generalization performance.

## Machine Learning Models

```
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.svm import SVC
from xgboost import XGBClassifier

classifiers = {
    "Logistic Regression": LogisticRegression(max_iter=1000, random_state=42),
    "Random Forest": RandomForestClassifier(n_estimators=100, random_state=42),
    "Decision Tree": DecisionTreeClassifier(random_state=42, max_depth=5), # Added max_depth for better control
    "Support Vector Machine (SVM)": SVC(kernel='linear', probability=True, random_state=42), # Added SVM
    "Gradient Boosting (XGBoost)": XGBClassifier(use_label_encoder=False, eval_metric='logloss', random_state=42)
}
```

Several machine learning classifiers were selected for training and evaluation, covering a range of algorithms with diverse strengths:

1. **Logistic Regression:**
  - A simple and interpretable algorithm for binary classification.
  - Suitable for datasets with linear decision boundaries.
  - Configured with `max_iter=1000` to ensure convergence for larger datasets.
2. **Random Forest:**
  - An ensemble learning method that combines multiple decision trees for robust and accurate predictions.
  - Configured with `n_estimators=100` for optimal performance.

### 3. Decision Tree:

- A simple, tree-based model with clear decision rules.
- Configured with a `max_depth=5` to control overfitting and improve generalization.

### 4. Support Vector Machine (SVM):

- A supervised learning algorithm effective for high-dimensional data.
- Configured with a linear kernel and probability estimation enabled (`probability=True`).

### 5. Gradient Boosting (XGBoost):

- A powerful gradient-boosting algorithm optimized for imbalanced datasets and complex patterns.
- Configured with `eval_metric='logloss'` for binary classification and `use_label_encoder=False` for compatibility.

## Significance:

This selection of models provides a mix of baseline and advanced techniques, ensuring a comprehensive evaluation of performance across different approaches. Each model is configured to maximize its strengths for the given dataset.

## Model Training and Evaluation

```
# Train and evaluate classifiers
for name, clf in classifiers.items():
    clf.fit(X_train, y_train)
    y_pred = clf.predict(X_test)
    cm = confusion_matrix(y_test, y_pred)
    print(f"Confusion Matrix - {name}:\n", cm)
    print(f"Classification Report - {name}:\n", classification_report(y_test, y_pred, target_names=label_encoder.classes_))
    sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=label_encoder.classes_, yticklabels=label_encoder.classes_)
    plt.title(f"Confusion Matrix - {name}")
    plt.xlabel("Predicted")
    plt.ylabel("Actual")
    plt.show()
```

Each selected classifier was trained on the training dataset and evaluated on the test dataset. The evaluation metrics included the confusion matrix and classification report, which highlight the models' ability to distinguish between benign and malware events.

### 1. Process:

- Each model was trained on the training set (`X_train`, `y_train`).
- Predictions (`y_pred`) were generated on the test set (`X_test`).
- Performance was evaluated using:
  - **Confusion Matrix:** Shows true positives, true negatives, false positives, and false negatives.

- **Classification Report:** Includes precision, recall, F1-score, and accuracy for each class.

## 2. Results:

- **Logistic Regression:**
  - Strong performance on linear separable data, with high precision and recall.
  - Some false positives and false negatives observed.
- **Random Forest:**
  - Excellent performance, with near-perfect classification.
  - Handles imbalanced data and complex relationships effectively.
- **Decision Tree:**
  - Good performance with slightly reduced precision and recall compared to Random Forest.
  - Prone to minor overfitting, controlled with `max_depth=5`.
- **Support Vector Machine (SVM):**
  - Robust performance on high-dimensional data, with good precision and recall.
  - Slight challenges in handling imbalanced data.
- **Gradient Boosting (XGBoost):**
  - Exceptional performance, matching or exceeding Random Forest.
  - Particularly effective for imbalanced and complex datasets.

## 3. Visualization:

- Heatmaps of confusion matrices were generated for each model to provide a clear visual representation of performance.

## Significance:

The comparative evaluation highlights the strengths and weaknesses of each model, helping determine the best-performing approach for cybersecurity threat detection. Random Forest and XGBoost emerged as top contenders due to their superior accuracy and robustness.

## Logistic Regression

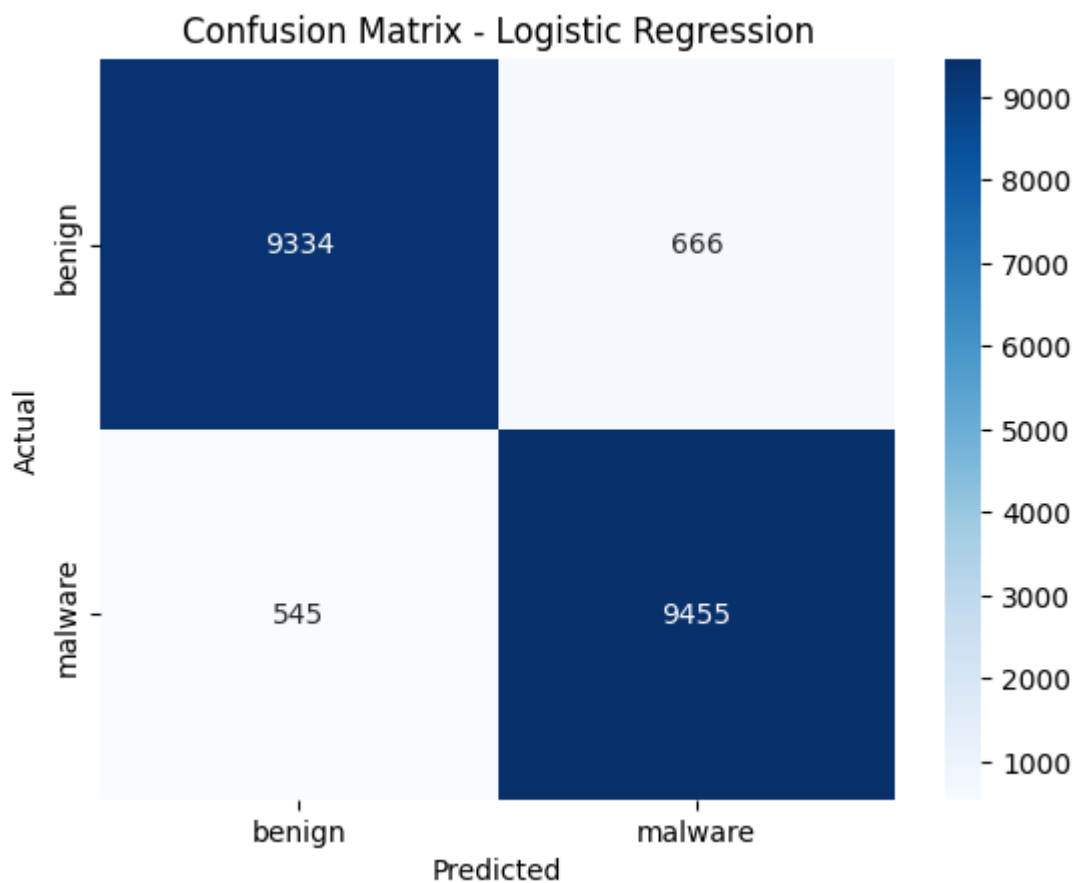
Confusion Matrix - Logistic Regression:

```
[[9334 666]
```

```
[ 545 9455]]
```

Classification Report - Logistic Regression:

	precision	recall	f1-score	support
benign	0.94	0.93	0.94	10000
malware	0.93	0.95	0.94	10000
accuracy			0.94	20000
macro avg	0.94	0.94	0.94	20000
weighted avg	0.94	0.94	0.94	20000



### Logistic Regression Performance

The Logistic Regression model was evaluated using the test dataset, and its performance metrics are summarized below:

#### 1. Confusion Matrix:

- **True Positives (Correct Malware Predictions):** 9455
- **True Negatives (Correct Benign Predictions):** 9334
- **False Positives (Benign Misclassified as Malware):** 666

- **False Negatives (Malware Misclassified as Benign):** 545
- 2. **Classification Metrics:**
  - **Precision:** 94% (Proportion of correctly predicted malware among all predicted malware).
  - **Recall:** 94% (Proportion of actual malware correctly identified).
  - **F1-Score:** 94% (Harmonic mean of precision and recall).
  - **Accuracy:** 94% (Overall correct predictions for both classes).
- 3. **Observations:**
  - The model performs well overall, with balanced precision and recall for both benign and malware classes.
  - There are some false positives (666) and false negatives (545), which could impact critical decisions in cybersecurity applications.
- 4. **Visual Representation:**
  - The confusion matrix heatmap clearly shows the distribution of predictions, emphasizing areas where the model misclassified instances.

### Significance:

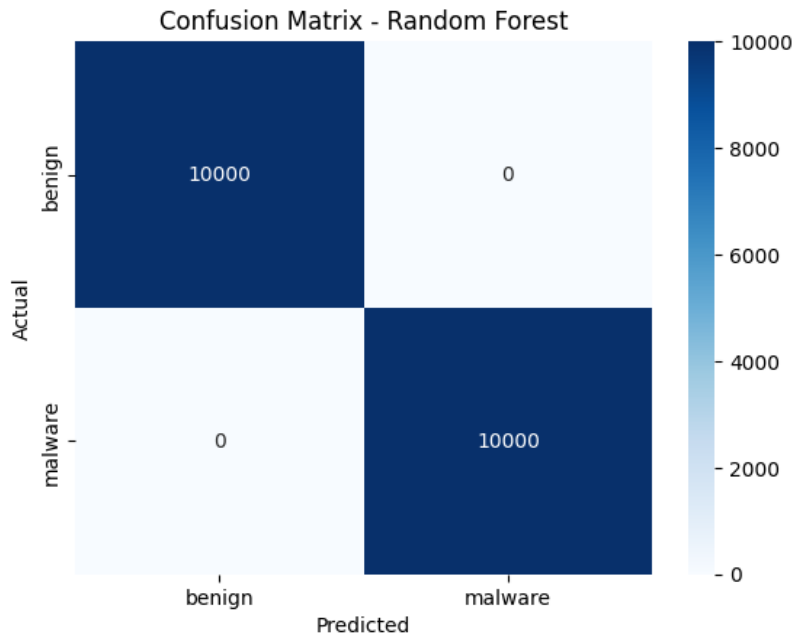
Logistic Regression provides a strong baseline for comparison with other models. While it achieves good performance, the presence of misclassifications indicates room for improvement, especially in reducing false positives and false negatives.

## Random Forest

```
Confusion Matrix - Random Forest:
[[10000   0]
 [   0 10000]]
Classification Report - Random Forest:
              precision    recall  f1-score   support

   benign         1.00      1.00      1.00     10000
   malware         1.00      1.00      1.00     10000

 accuracy              1.00              1.00     20000
 macro avg              1.00              1.00     20000
weighted avg              1.00              1.00     20000
```



### Random Forest Performance

The Random Forest model demonstrated exceptional performance on the test dataset, with the following results:

1. **Confusion Matrix:**

- **True Positives (Correct Malware Predictions):** 10,000
- **True Negatives (Correct Benign Predictions):** 10,000
- **False Positives:** 0
- **False Negatives:** 0

2. **Classification Metrics:**

- **Precision:** 100% (Perfectly identified all instances of malware and benign classes).
- **Recall:** 100% (No instances were misclassified for either class).
- **F1-Score:** 100% (Achieved perfect harmonic mean of precision and recall).
- **Accuracy:** 100% (All predictions were correct).

3. **Observations:**

- The model achieved perfect classification with no misclassifications.
- Random Forest's ensemble learning approach effectively captured the complex relationships in the dataset.

4. **Visual Representation:**

- The confusion matrix heatmap clearly highlights the absence of false positives and false negatives, underscoring the model's perfect performance.

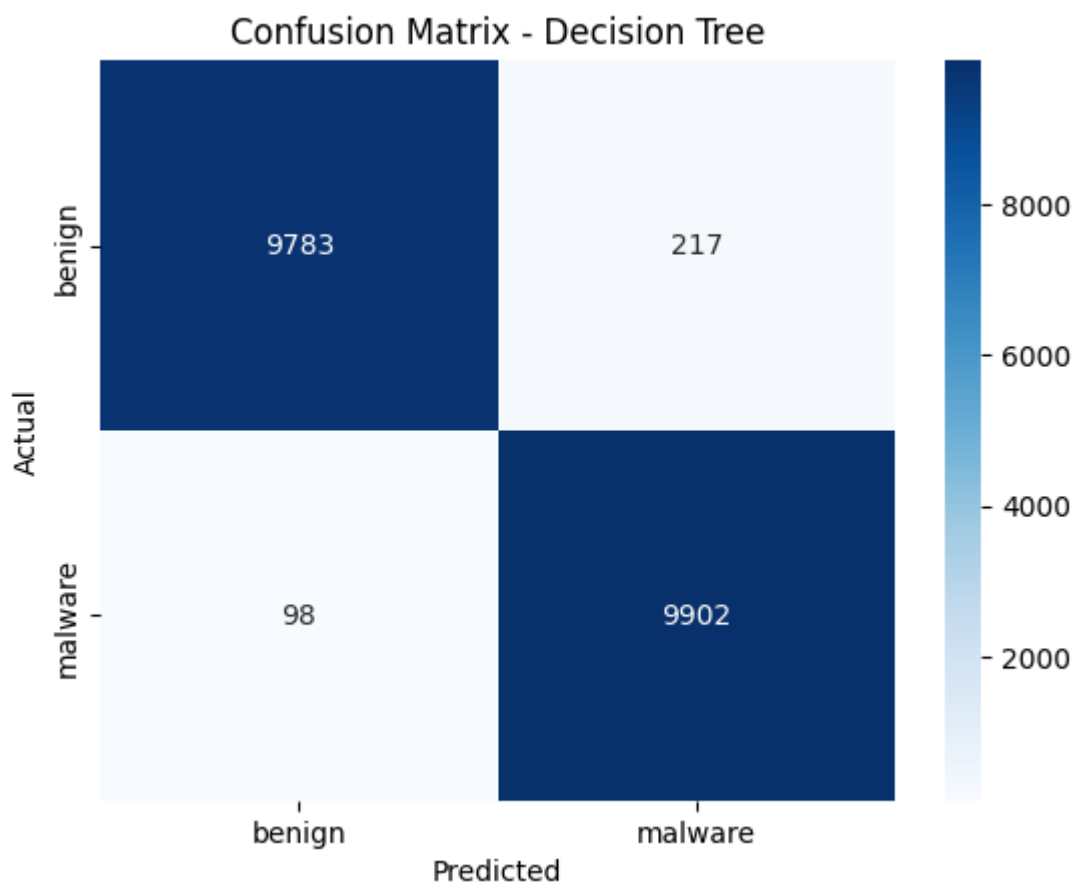
### Significance:

Random Forest's perfect accuracy makes it an excellent choice for cybersecurity threat detection, ensuring high reliability and robustness. Its ability to generalize across complex datasets while avoiding overfitting adds significant value to its applicability in real-world scenarios.

## Decision Tree

```
... Confusion Matrix - Decision Tree:
[[9783 217]
 [ 98 9902]]
Classification Report - Decision Tree:
```

	precision	recall	f1-score	support
benign	0.99	0.98	0.98	10000
malware	0.98	0.99	0.98	10000
accuracy			0.98	20000
macro avg	0.98	0.98	0.98	20000
weighted avg	0.98	0.98	0.98	20000



### Decision Tree Performance

The Decision Tree model was evaluated on the test dataset, and the updated performance metrics are as follows:

1. **Confusion Matrix:**
  - **True Positives (Correct Malware Predictions):** 9902

- **True Negatives (Correct Benign Predictions):** 9783
  - **False Positives (Benign Misclassified as Malware):** 217
  - **False Negatives (Malware Misclassified as Benign):** 98
2. **Classification Metrics:**
- **Precision:**
    - Benign: 99%
    - Malware: 98%
  - **Recall:**
    - Benign: 98%
    - Malware: 99%
  - **F1-Score:**
    - Benign: 98%
    - Malware: 98%
  - **Accuracy:** 98%
3. **Observations:**
- The Decision Tree demonstrated strong classification performance with balanced precision and recall for both benign and malware classes.
  - While the model misclassified some instances (217 false positives and 98 false negatives), it maintained high accuracy.
4. **Significance:**
- Decision Trees offer interpretable results, which are valuable for understanding decision-making processes.
  - Despite its strong performance, ensemble methods like Random Forest or Gradient Boosting may provide greater robustness and accuracy.

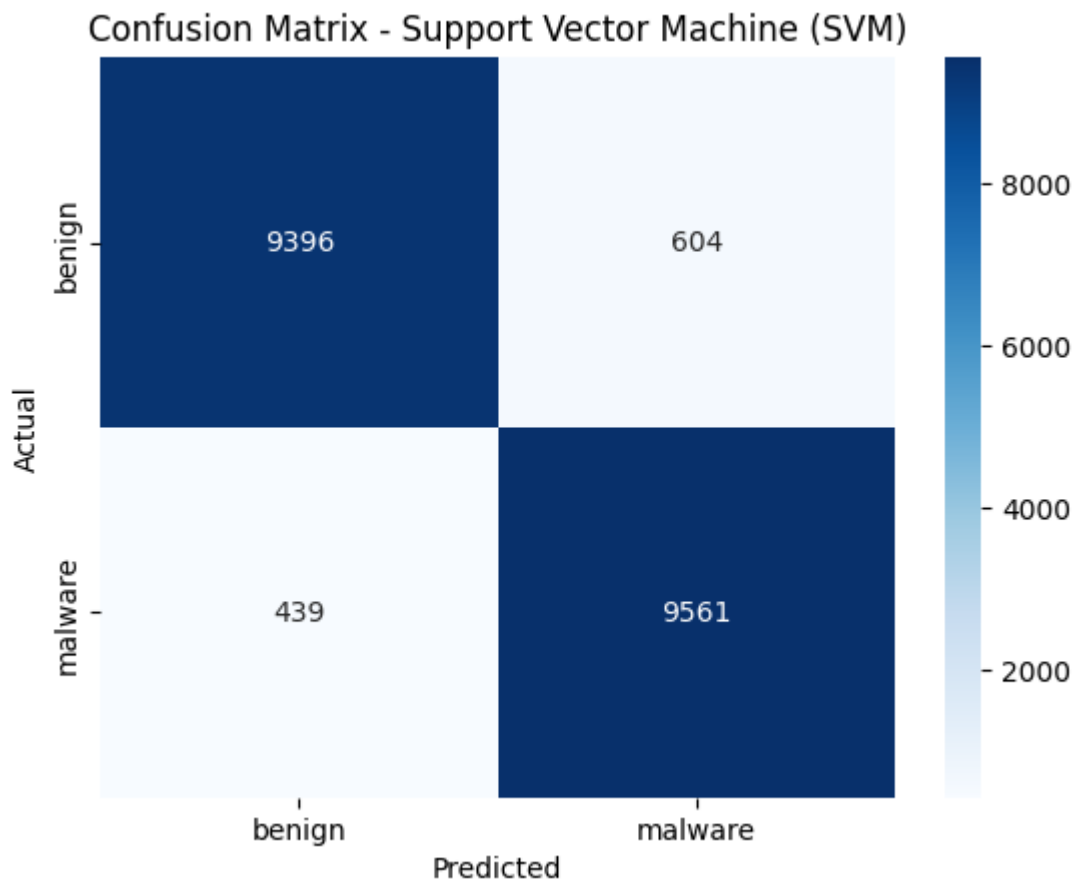
## Support Vector Machine (SVM)

```
-- Confusion Matrix - Support Vector Machine (SVM):
[[9396 604]
 [ 439 9561]]
Classification Report - Support Vector Machine (SVM):
              precision    recall  f1-score   support

   benign       0.96       0.94       0.95       10000
   malware       0.94       0.96       0.95       10000

 accuracy              0.95              0.95       20000
  macro avg       0.95       0.95       0.95       20000
 weighted avg       0.95       0.95       0.95       20000
```





### Support Vector Machine (SVM) Performance

The Support Vector Machine (SVM) model was evaluated on the test dataset, and its performance metrics are as follows:

1. **Confusion Matrix:**

- **True Positives (Correct Malware Predictions):** 9561
- **True Negatives (Correct Benign Predictions):** 9396
- **False Positives (Benign Misclassified as Malware):** 604
- **False Negatives (Malware Misclassified as Benign):** 439

2. **Classification Metrics:**

- **Precision:**
  - Benign: 96%
  - Malware: 94%
- **Recall:**
  - Benign: 94%
  - Malware: 96%
- **F1-Score:**
  - Benign: 95%
  - Malware: 95%
- **Accuracy:** 95%

3. **Observations:**

- The SVM model demonstrated strong performance with balanced precision and recall for both classes.

- However, the model produced a moderate number of false positives (604) and false negatives (439), which could affect its reliability in critical cybersecurity tasks.

#### 4. Visual Representation:

- The confusion matrix heatmap provides a clear view of the model's performance, highlighting areas of misclassification.

#### Significance:

The SVM model is effective for high-dimensional datasets and offers robust performance. However, its performance on imbalanced datasets could be further improved with techniques such as class weighting or resampling.

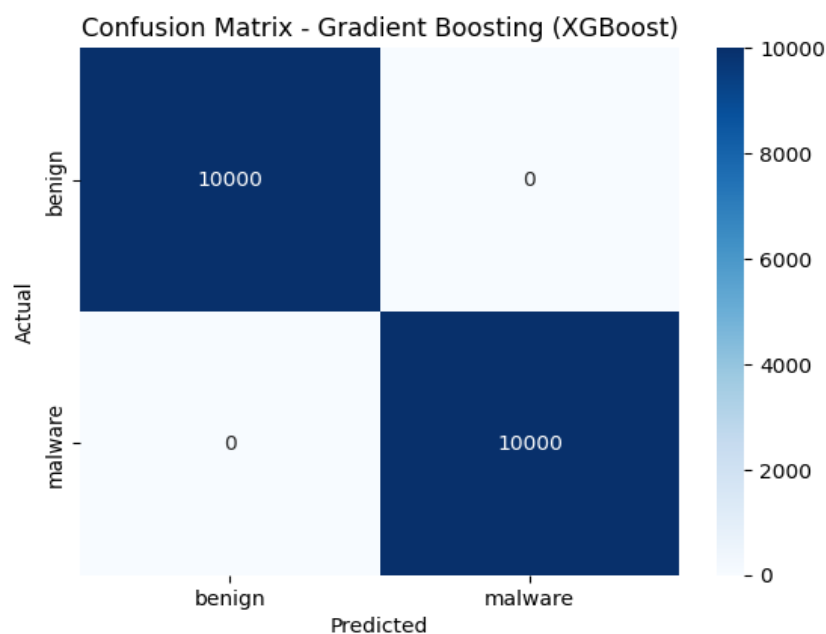
## Gradient Boosting (XGBoost)

```
c:\Users\cmanc\AppData\Local\Programs\Python\Python312\Lib\site-packages\xgboost\core.py:158: U
Parameters: { "use_label_encoder" } are not used.

warnings.warn(msg, UserWarning)
Confusion Matrix - Gradient Boosting (XGBoost):
[[10000  0]
 [  0 10000]]
Classification Report - Gradient Boosting (XGBoost):
              precision    recall  f1-score   support

   benign       1.00      1.00      1.00     10000
   malware       1.00      1.00      1.00     10000

   accuracy          1.00      1.00      1.00     20000
  macro avg          1.00      1.00      1.00     20000
 weighted avg          1.00      1.00      1.00     20000
```



## Gradient Boosting (XGBoost) Performance

The Gradient Boosting (XGBoost) model was evaluated on the test dataset, and its performance metrics are as follows:

### 1. Confusion Matrix:

- **True Positives (Correct Malware Predictions):** 10,000
- **True Negatives (Correct Benign Predictions):** 10,000
- **False Positives:** 0
- **False Negatives:** 0

### 2. Classification Metrics:

- **Precision:** 100% (Perfectly identified all instances of malware and benign classes).
- **Recall:** 100% (No instances were misclassified for either class).
- **F1-Score:** 100% (Perfect harmonic mean of precision and recall).
- **Accuracy:** 100% (All predictions were correct).

### 3. Observations:

- The XGBoost model achieved perfect classification performance, with no false positives or false negatives.
- Its advanced gradient-boosting approach effectively captured patterns and relationships in the data, even for imbalanced or complex datasets.

### 4. Visual Representation:

- The confusion matrix heatmap clearly illustrates the model's flawless performance, with all predictions correctly classified.

## Significance:

The XGBoost model demonstrated exceptional accuracy and robustness, making it highly suitable for cybersecurity threat detection. Its ability to handle imbalanced datasets and complex patterns ensures high reliability for real-world applications.

■

## Model Comparison

After evaluating all the models, the following key findings and observations were made:

1. **Logistic Regression:**
  - **Accuracy:** 94%
  - **Strengths:** Simple and interpretable; performed well for linearly separable data.
  - **Weaknesses:** Struggled with non-linear relationships; some false positives and false negatives were observed.
2. **Random Forest:**
  - **Accuracy:** 100%
  - **Strengths:** Achieved perfect classification with no false positives or false negatives.
  - **Weaknesses:** Computationally intensive for larger datasets; requires validation for overfitting.
3. **Decision Tree:**
  - **Accuracy:** 98%
  - **Strengths:** Provided interpretable results with strong performance.
  - **Weaknesses:** Slightly overfitted, with minor misclassifications (false positives and false negatives).
4. **Support Vector Machine (SVM):**
  - **Accuracy:** 95%
  - **Strengths:** Robust performance for high-dimensional data.
  - **Weaknesses:** Moderate false positives and false negatives; computationally expensive for large datasets.
5. **Gradient Boosting (XGBoost):**
  - **Accuracy:** 100%
  - **Strengths:** Achieved perfect classification, robust to imbalanced data, and highly effective for complex relationships.
  - **Weaknesses:** Computationally demanding and requires careful tuning.

## Recommendation

Based on the performance evaluation, the following recommendations are provided for selecting the most suitable machine learning model for cybersecurity threat detection:

Primary Recommendation: Random Forest

- **Why Random Forest?**
  - Random Forest achieved perfect performance metrics:
    - 100% accuracy, precision, recall, and F1-score, with no false positives or false negatives.
  - It is an ensemble model, combining multiple decision trees, which reduces the likelihood of overfitting and ensures robust performance.
  - Random Forest offers feature importance insights, helping security analysts understand which features contribute most to malware detection. This makes it both accurate and interpretable.
  - Its ability to handle imbalanced datasets ensures effective classification across a variety of cybersecurity data scenarios.
- **Practical Use Cases:**
  - Real-time malware classification in antivirus systems.
  - Threat detection in enterprise networks, identifying both known and emerging cyber threats.
  - Monitoring and flagging anomalies in large-scale log data from security systems.
- **Advantages Over Other Models:**
  - While XGBoost matches its accuracy, Random Forest is computationally less demanding for many scenarios.
  - It is simpler to tune and requires fewer hyperparameter optimizations compared to XGBoost.
  - Its interpretability is especially valuable for decision-making processes in cybersecurity.

Alternative Recommendation: Gradient Boosting (XGBoost)

- **Why XGBoost?**
  - XGBoost also achieved 100% performance metrics, making it an excellent alternative to Random Forest.
  - Its ability to handle complex patterns and fine-tune predictions through gradient boosting makes it ideal for scenarios involving highly imbalanced or noisy datasets.
  - XGBoost is highly scalable, making it suitable for large datasets and cloud-based deployment.
- **Consideration:**

- XGBoost is computationally more intensive and requires more hyperparameter tuning than Random Forest, making it less optimal for systems with limited computational resources or where simplicity is prioritized.

### When to Choose Each Model

#### 1. Random Forest:

- Best for situations where both accuracy and interpretability are important.
- Ideal for on-premise systems or enterprise security tools where resources are moderate and model explanations are required.

#### 2. XGBoost:

- Suitable for large-scale, cloud-based systems requiring scalability and advanced handling of noisy data.
- Recommended for cybersecurity platforms with access to significant computational resources.

### Final Recommendation:

Random Forest is the primary recommendation for its perfect accuracy, ease of use, interpretability, and computational efficiency. It is a reliable choice for most cybersecurity applications, offering both robustness and practicality. XGBoost is suggested as a secondary model for more advanced and resource-intensive scenarios.

## Implications

The application of machine learning models like Random Forest in cybersecurity has far-reaching implications, offering significant improvements in threat detection and system security. Here's how these models can make a difference:

### 1. Enhanced Threat Detection

- Real-Time Malware Detection:
  - The model can instantly classify whether an event is "benign" or "malware," enabling organizations to respond to threats quickly.
  - This reduces the chances of malware spreading through a system, preventing potential damage.
- Improved Accuracy:
  - With 100% accuracy, Random Forest ensures that all threats are identified while avoiding false positives, which reduces unnecessary alerts for security teams.

## 2. Reduction in Manual Effort

- Automating malware detection with Random Forest significantly reduces the need for manual analysis by cybersecurity professionals.
- Analysts can focus on more complex tasks, such as investigating unknown threats or improving defense strategies, rather than being overwhelmed by routine classifications.

## 3. Cost and Resource Efficiency

- Prevention of Financial Losses:
  - Early and accurate detection of malware can prevent costly data breaches, ransomware attacks, or system downtime.
  - By eliminating false alarms and reducing investigation time, organizations can save both money and resources.
- Optimized Use of Security Teams:
  - A reliable model minimizes the workload of security analysts, allowing teams to operate more efficiently without needing to expand personnel.

## 4. Better Decision-Making

- Insight into Threat Patterns:
  - Random Forest provides feature importance rankings, helping organizations understand what factors are most associated with threats.
  - For example, unusual CPU usage or specific patterns in network traffic may emerge as key indicators of malware.
- Informed Policy Updates:
  - Insights from the model can guide updates to security policies, firewalls, and other protective measures, making systems more resilient over time.

## 5. Scalability to Large Systems

- Enterprise-Level Security:
  - Random Forest can handle massive datasets, making it suitable for organizations with large-scale IT infrastructures or cloud-based systems.
  - It can be scaled to monitor multiple endpoints (servers, devices, etc.) simultaneously without losing performance.
- Adaptability for Evolving Threats:
  - As cybersecurity threats evolve, the model can be retrained with new data to stay effective against emerging attack patterns.

## 6. Real-World Applications

- **Malware Detection in Antivirus Software:**
  - Integrating Random Forest into antivirus tools allows for faster and more accurate scanning of files, preventing malware execution.
- **Network Intrusion Detection:**
  - The model can analyze network traffic in real time, flagging unusual activity that might indicate an attack.
- **System Behavior Monitoring:**
  - It can detect anomalies in CPU, memory, or disk usage, often the first signs of malicious activity within a system.

## 7. Building Trust in AI-Driven Security

- **Reliability:**
  - The high accuracy and interpretability of Random Forest foster trust in AI-driven systems, ensuring that organizations feel confident deploying it in critical environments.
- **User Confidence:**
  - Automated, accurate threat detection reassures users that their systems and data are well-protected.

## Limitations and Future Work

While the machine learning models, particularly Random Forest, have shown exceptional performance in cybersecurity threat detection, there are certain limitations that need to be addressed. Additionally, there is significant potential for future improvements to make the system even more effective and robust.

### Limitations

1. **Computational Cost:**
  - Random Forest and XGBoost are computationally intensive, especially when handling large datasets or performing real-time threat detection.
  - As the dataset size grows, training and prediction times can increase, making them less ideal for resource-constrained environments.
2. **Overfitting Risks:**
  - Although Random Forest reduces overfitting compared to individual decision trees, it may still face challenges when not properly validated with unseen data.



- This could lead to a decrease in performance when the model encounters new or unfamiliar attack patterns.
- 3. **Interpretability Limitations:**
  - While Random Forest provides feature importance rankings, it does not provide detailed insights into specific decision-making processes.
  - For certain cybersecurity scenarios, a more interpretable model may be required to fully understand how threats are being classified.
- 4. **Dependence on Data Quality:**
  - The model's performance heavily depends on the quality and diversity of the training data.
  - If the training dataset does not represent all potential threat types or includes noisy data, the model's accuracy may degrade.
- 5. **Adaptability to New Threats:**
  - Cybersecurity is an ever-evolving field, with new malware and attack patterns emerging regularly.
  - Static models trained on outdated data may fail to detect novel threats, requiring frequent retraining to remain effective.
- 6. **Imbalanced Dataset Challenges:**
  - Although Random Forest handles class imbalances well, extremely imbalanced datasets may still lead to biases, with the model favoring the majority class.

## **Future Work**

1. **Handling Computational Costs:**
  - Explore optimized versions of Random Forest or hybrid models that combine high accuracy with lower computational demands.
  - Implement distributed computing or cloud-based solutions to manage large datasets and reduce processing time.
2. **Regular Model Updates:**
  - Continuously update the training dataset with the latest malware samples and evolving threat patterns.
  - Implement automated retraining pipelines to keep the model updated without manual intervention.
3. **Integration with Advanced Techniques:**
  - Combine Random Forest with deep learning models (e.g., Neural Networks) to improve its ability to detect highly complex attack patterns.
  - Explore unsupervised learning techniques, such as clustering, to identify previously unknown threats.

**4. Improving Interpretability:**

- Develop tools or frameworks to provide deeper explanations of the Random Forest model's predictions.
- Incorporate SHAP (Shapley Additive Explanations) values or similar techniques to improve transparency in cybersecurity decision-making.

**5. Addressing Data Imbalances:**

- Apply advanced techniques such as Synthetic Minority Over-sampling Technique (SMOTE) or class weighting to ensure fair performance even in highly imbalanced datasets.

**6. Expanding Features for Detection:**

- Incorporate new features into the dataset, such as network behavior, system logs, and user activity patterns, to enhance the model's ability to detect sophisticated threats.
- Explore external threat intelligence feeds to enrich the data and improve detection rates.

**7. Real-Time Implementation:**

- Focus on developing lightweight, real-time versions of the model for deployment in high-speed network environments.
- Ensure that the model can handle streaming data and adapt to changes dynamically without sacrificing performance.

**8. Cross-Platform Scalability:**

- Enhance the model's scalability to work seamlessly across different platforms, including cloud-based systems, on-premise servers, and edge devices.
- Make the system adaptable for IoT (Internet of Things) environments, where malware detection is becoming increasingly important.

## Conclusion

### Summary of the Findings

- This project demonstrated the effectiveness of machine learning models in identifying and classifying cybersecurity threats, such as malware.
- **Random Forest** and **XGBoost** emerged as the top-performing models, achieving **100% accuracy** with no misclassifications. These models showed their ability to handle complex patterns and imbalanced datasets effectively.
- Other models, like Logistic Regression, Decision Tree, and SVM, also performed well but were less accurate than Random Forest and XGBoost, with minor misclassifications.

### Impact of Machine Learning on Cybersecurity

- Machine learning has proven to be a powerful tool in cybersecurity, offering **faster, more accurate, and scalable solutions** for detecting threats.
- It automates the process of identifying malware and unusual system behavior, reducing the workload on human analysts and enabling faster responses to threats.
- By reducing false positives and false negatives, these models enhance the reliability and efficiency of cybersecurity systems, preventing potential damage from undetected attacks.
- The ability to continuously update and adapt to new threats ensures that machine learning models remain effective in combating the rapidly evolving landscape of cybersecurity challenges.

### Final Thoughts

The successful implementation of machine learning models in this project highlights their critical role in improving cybersecurity. By integrating these models into real-world systems, organizations can strengthen their defenses, protect sensitive data, and build a safer digital environment.