# Autonomous network services using machine learning-based cybersecurity

**Dr. Vijey Thayananthan [1], Ahmed Saleh Dhafer Alghamdi [2]**
[1] Department of computer science King Abdulaziz University, Saudi Arabia,
vthayanathan@kau.edu.sa
[2] Department of computer science King Abdulaziz University, Saudi Arabia,
adaferalmansour@stu.kau.edu.sa

## ABSTRACT

The crucial topic of protecting Controller Area Network (CAN) bus systems in autonomous vehicles against cyberattacks is covered in this research. The risk and consequences of cyberattacks on CAN systems rise considerably as vehicles become increasingly automated and linked, creating safety and security hazards. In order to identify anomalies and categorize traffic into attack or conventional categories, the study looks at the weaknesses of CAN buses and recommends using machine learning techniques like Decision Trees, Clustering, and Deep Learning. To improve the detection of anomalies and cyber-attacks in CAN systems, the suggested methods combine data balance, feature selection, and ensemble learning with a voting-based strategy.

Metrics like accuracy, precision, recall, F1-score, and confusion matrix can be used to assess the presented approaches. According to the study's findings, these suggested solutions provide a more reliable and efficient way to identify cyber-attacks and anomalies in CAN systems, boosting the development of cyber security for autonomous vehicles. While outlining the necessity of information security and the advantages of autonomous vehicles, it also suggests cutting-edge ways to strengthen their security.

Overall, this article highlights the urgent need for improved security measures in autonomous cars since cyberattacks pose a serious threat to the functioning of these vehicles in a safe and secure manner. The study proposes a potential approach to enhancing the security of CAN bus systems in autonomous vehicles by suggesting cutting-edge approaches for identifying anomalies and cyber-attacks.

**Key words:** Controller Area Network (CAN), cyber-attacks, machine learning, deep learning, autonomous vehicles, security.

## 1. INTRODUCTION

Concerns concerning their dependability and the safety of passengers and pedestrians have arisen as a result of the increasing usage of electronic technologies in automobiles. Data transmission between these parts has been successfully demonstrated via the CAN bus, which links the various electronic control units (ECUs) in the vehicle. The CAN protocol is susceptible to numerous hack-based attacks due to its lack of security. Attackers can gain access to specific ECUs, inject CAN packets, and inflict damage, which puts the car and its occupants in serious danger for safety and security. As a result, a solution to safeguard the CAN bus is very necessary.

Machine learning (ML) and artificial intelligence (AI) are promising approaches for detecting cyber-attacks and vulnerabilities in CAN systems. However, because there are so many intricate and interconnected processes at play, creating effective machine learning-based methods for this purpose is difficult. Therefore, the purpose of this article is to examine the CAN bus's vulnerability in vehicles and to suggest a method for identifying and analyzing the behavior of this system's vulnerabilities using a set of ML and DL algorithms.

The suggested solution intends to increase the security of driver-wired and autonomous systems by improving the identification of vulnerabilities and cyberattacks in CAN systems. Understanding attack patterns and creating a system for analyzing vehicle attacks using intelligence and machine learning are the two goals of the project. In order to accomplish these goals, also, this paper will address the drawbacks of the CAN protocol [1].

The literature on the subject is reviewed in literature review section of this paper and the suggested methodology is presented in the methodology section as well. The suggested methodology entails employing a variety of machine learning methods, such as Decision Trees, Clustering, and Deep Learning, to analyze the behavior of the CAN bus. These algorithms recognize unusual behavior and divide traffic into categories for attacks and regular traffic. To enhance the

detection of anomalies and cyberattacks in CAN systems, the methodology also incorporates data balance, feature selection, and ensemble learning with a voting-based mechanism [4].

The application of the suggested methodology will be covered in in more details, along with the testing of machine learning for anomaly detection in the attack or process model and military operations. Also, experiment findings are thoroughly discussed, including measures like accuracy, precision, recall, F1-score, and confusion matrix. Finally, we summarize our results and their consequences in the results section of this paper. According to the study's findings, the suggested methodology provides a more reliable and efficient way to identify cyberattacks and anomalies in CAN systems, promoting the development of autonomous vehicles' cybersecurity [5][6].

This study's findings underscore the urgent need to address the CAN systems' vulnerabilities in autonomous vehicles and the potential of machine learning to strengthen their security. The suggested methodology provides a strategy that has promise for enhancing the security of CAN bus systems in autonomous vehicles, hence enhancing passenger and pedestrian safety.

## 2. CYBERSECURITY IN AUTONOMOUS VEHICLES

Autonomous vehicles have the potential to transform the transportation industry, providing increased accessibility, efficiency, and safety. However, these networks are also vulnerable to various cyberattacks that can jeopardize passenger security and privacy. In this section, we will explore the different aspects of cybersecurity in autonomous vehicles, including the vulnerabilities in the CAN protocol, artificial intelligence-based security, and the potential benefits of autonomous vehicle service networks.

### 2.1 VULNERABILITY OF THE AUTONOMOUS VEHICLE

Autonomous vehicle service networks have drawn more attention in the transportation sector and can provide affordable, reliable, and efficient transportation options. These networks can significantly decrease the frequency of traffic accidents, enhance response times, and boost the effectiveness of traffic flow [11]. The incorporation of natural language processing technologies into these networks provides further advantages, such as excellent real-time traffic management, improved customer experience, and improved communication between vehicles and infrastructure.

The CAN protocol is a communication standard that is widely used in various industrial and automobile control systems. CAN buses are essential to the automobile industry since they transfer control signals between the electronic systems in a vehicle. However, the CAN protocol is not completely secure, and cybercriminals can exploit its vulnerabilities to carry out malicious attacks on autonomous vehicles. These attacks can take different forms, such as message tampering, message injection, and denial of service attacks, which can disrupt or manipulate the vehicle's control systems [2] [3] [7] [8] [9].

One of the most significant risks associated with CAN protocol vulnerabilities is the potential for attackers to gain access to a vehicle's electronic control units (ECUs). By exploiting CAN bus vulnerabilities through reverse engineering and fuzzing techniques, attackers can obtain complete control over the vehicle's control systems [3]. This can pose a serious threat to passenger safety and privacy, as attackers can manipulate various vehicle functions, such as steering, braking, and acceleration. Researchers have suggested various approaches to secure the CAN protocol, such as encryption and authentication mechanisms, intrusion detection systems, and secure routing protocols. These measures can enhance the network's overall security and prevent cyberattacks from jeopardizing passenger safety and privacy. Various data mining methods may be used to monitor and identify malicious behavior on CAN buses to solve these security problems. Anomaly detection and association rule mining are two data mining approaches that can be applied to CAN bus cybersecurity [12][13][14].

## 3. LITERATURE REVIEW

Artificial intelligence-based security is one approach that can be used to enhance cybersecurity in autonomous vehicles. Applications that use computer vision and machine learning can help to reduce the risk of attacks and criminal behavior by training computers to think and behave like humans [10]. Combining technologies such as natural language processing, deep learning, and machine learning can address urgent concerns about cybersecurity [15].

An essential function of natural language generation within AI is generating information in text depending on the inputs provided. With the use of machine learning algorithms, researchers have been able to quickly and accurately identify security threats in real-time. These algorithms can monitor network traffic, detect anomalies, and alert system administrators of potential threats. [16] [17] One of the advantages of machine learning is that it can adapt and learn from previous attacks.

By analyzing past attacks, machine learning algorithms can identify patterns and characteristics of cyberattacks, and use that information to improve the system's defense mechanisms. Machine learning can be used to protect not only the vehicle itself but also the passengers' privacy. With the increasing number of connected devices and the vast amounts of data generated by autonomous vehicles, privacy is becoming a significant concern. Machine learning algorithms can help protect sensitive information by identifying potential security risks and implementing appropriate security measures [18].

Several research studies have highlighted the effectiveness of machine learning techniques in identifying and mitigating security threats in various types of networks. Singh et al. (2017) employed machine learning algorithms to detect and prevent cyberattacks in autonomous cars, resulting in an overall improvement in the network's security. Similarly, Kim et al. (2019) investigated the use of machine learning

algorithms in identifying and reducing security risks in industrial control systems (ICSs), and found that these techniques were effective in enhancing the overall security of these systems [19].

Zhang et al. (2018) suggested using deep learning algorithms to identify potential security vulnerabilities in wireless sensor networks, and their study demonstrated that deep learning algorithms were more accurate and efficient in identifying security risks than conventional security solutions. Li et al. (2019) focused on using reinforcement learning algorithms to improve the safety of autonomous networks, finding that this method was successful in allowing networks to quickly react to any changes in their operating environment [20].

Chen et al. (2020) recommended the use of unsupervised machine learning methods to identify and mitigate security vulnerabilities in IoT networks, and their research showed that these algorithms efficiently identified abnormalities and mitigated security risks, improving the overall security of these systems. Finally, Wang et al. (2021) researched the use of generative adversarial networks (GANs) to enhance the safety of autonomous networks, and their findings showed that GANs effectively detected and mitigated security threats in these networks [21].

Overall, these studies highlight the importance of using machine learning techniques in enhancing the security of different types of networks. From autonomous cars to industrial control systems, wireless sensor networks, autonomous networks, and IoT networks, machine learning algorithms have demonstrated their effectiveness in identifying and mitigating security threats, leading to an overall improvement in network security.

## 4. METHODOLOGY
The suggested approach that manages the data balancing issue and picks the essential elements that help the model for assessing vehicle anomaly detection in CAN cyber assaults has been developed.

### 4.1 THE SYNTHETIC MINORITY OVERSAMPLING TECHNIQUE (SMOTE)
algorithm solves the data balance issue encountered while analyzing vehicle anomaly detection in CAN cyber-attacks. To balance the distribution of classes in a dataset, the SMOTE algorithm creates synthetic samples from the minority class. The method functions by picking comparable samples from the minority class, calculating the difference between them, and producing synthetic examples by adding the computed difference to the chosen sample [22].

The suggested technique utilizes the Recursive Feature Elimination (RFE) algorithm to determine the most significant features. The RFE method is a feature selection technique that uses a model-based strategy to repeatedly delete the least significant features until the target number of features is reached. The method employs a scoring measure to rank the

features according to their significance and picks the top k features for inclusion in the model.

The SMOTE algorithm is initially used in the suggested technique to balance the class distribution throughout the dataset, followed by the RFE algorithm to choose the essential features. The machine learning model subsequently evaluates the obtained characteristics [23]. The suggested technique has the benefit of tackling the issue of data balance while also picking essential attributes to increase the machine learning model's performance.

### 4.2 COST-SENSITIVE LEARNING (CSL) AND FEATURE IMPORTANCE FROM RANDOM FOREST (FIRF)
algorithms are different approaches to addressing data balancing and picking the most significant characteristics for assessing vehicle anomaly detection in CAN cyber assaults [24].

The CSL method is an approach to machine learning that considers the cost of misclassifying minority class samples, as opposed to typical machine learning algorithms that give equal weight to all classes. CSL methods may solve the data balancing issue in unbalanced datasets by attributing a more significant penalty to misclassifying samples from the minority class. The FIRF algorithm is a feature selection approach that computes the relevance of each feature in the dataset using the Random Forest algorithm. Random Forest is a technique for ensemble learning in which numerous decision trees are constructed, and the relevance of a feature is determined as the average reduction in impurity across all decision trees [25].

In the suggested technique, the CSL algorithm is initially applied to the dataset to address the data balance issue, followed by the FIRF algorithm to choose the most significant features. The machine learning model subsequently evaluates the obtained characteristics. Developing a system for detecting multi-class abnormalities and cyberattacks in autonomous cars based on an ensemble learning technique and applying a voting-based strategy.

In order to identify multi-class abnormalities and cyberattacks in the Controller Area Network (CAN) systems used in autonomous cars, an ensemble learning technique based on a voting-based strategy can be used. In this system, multiple base classifiers are trained on the same dataset, and their predictions are integrated using a voting-based technique to arrive at a conclusion [26], [27].

### 4.3 THE FOLLOWING ARE SOME POSSIBLE CONFIGURATIONS FOR THE SYSTEM

The dataset is pre-processed to manage missing values and normalize the features. This step of the process is known as "pre-processing."

- **Training of Base Classifiers**: After the dataset has been pre-processed, many base classifiers, such as

decision trees, support vector machines (SVMs), and k-nearest neighbors, are trained on the dataset (k-NN).

- **Ensemble Learning**: In order to create a definitive forecast, the basic classifiers are integrated via the use of a voting-based technique. Each base classifier in a voting-based method will cast a vote for the class that will be predicted, and the class that receives the most significant number of votes will be taken into consideration as the final prediction.

- **Anomaly Detection** In order to establish whether or not there is an anomaly, the final prediction that was generated by the ensemble learning model is compared with the ground truth.
  By integrating the capabilities of several different base classifiers, this system has the distinct benefit of identifying multi-class abnormalities and cyberattacks in CAN systems used in autonomous cars with a high degree of efficiency. The ensemble learning method is also resistant to individual classifiers' shortcomings and can deal with datasets that are not evenly distributed [26], [27].

  The following assessment criteria could be used in order to assess how well the suggested ensemble learning method for identifying multi-class abnormalities and cyber-attacks in CAN systems in autonomous cars performs:

- **Accuracy**: The accuracy of a model may be defined as the percentage of accurate predictions it makes. The definition of it is the ratio of the number of accurate forecasts to the total number of predictions that were made.

- **Precision**: Precision refers to the fraction of correctly produced positive predictions within the total number of positive predictions made by the model. It evaluates how well the model can generate accurate predictions.

- **Recall**: The percentage of genuine positive predictions made out of the total number of real positive samples is called recall. It evaluates how well the model can identify positive samples across the board.

- **F1-score** - The F1 score is calculated using the harmonic mean of the accuracy and recall scores. It delivers a single score that indicates the entire performance of the model, and it strikes a healthy balance between measuring the model's accuracy and its recall.

- **A confusion matrix:** is a tabular representation of the performance of the model. It shows the true positive, the false positive, the true negative, and the false negative predictions that the model produced.

The efficacy of the suggested approach for ensemble learning may be compared with the effectiveness of current state-of-the-art techniques by employing these evaluation measures in a comparison. For instance, if the proposed solution has a higher accuracy, precision, recall, F1-score, and a lower number of false positive and false pessimistic predictions, then one can conclude that the proposed solution is more effective than the current techniques considered to be state-of-the-art.

## 5. WRITTEN CONTRIBUTION

Our contribution in this field is the design of a set of machine learning and deep learning algorithms for the identification and analysis of vulnerabilities and cyber-attacks in CAN systems. This proposal includes techniques such as Isolation Forest, One-Class SVM, Random Forest, SVM, decision trees, clustering, and deep learning, as well as data balancing, feature selection, and ensemble learning with voting-based strategies to improve identification. abnormalities and cyber-attacks in CAN systems in autonomous cars. By designing these advanced methods, we hope to increase the safety of autonomous vehicles and encourage their development and adoption.

## 6. MACHINE LEARNING ALGORITHM

**6.1 Machine learning algorithms**: are becoming an increasingly important tool in improving the security of autonomous vehicles and other types of networks. Using these algorithms, researchers can quickly and accurately detect and mitigate security threats, reducing the risk of cyber-attacks that could compromise passenger safety and privacy. One of the advantages of machine learning algorithms is their ability to monitor network traffic and detect anomalies in real-time. These algorithms can analyze large amounts of data and identify patterns and characteristics of cyber-attacks, allowing system administrators to take appropriate measures to mitigate these threats. For example, if a machine learning algorithm detects an unusual pattern of activity on an autonomous vehicle's Controller Area Network (CAN) bus, it can immediately alert system administrators, who can investigate the problem and take appropriate action.

**6.2 Another advantage of machine learning algorithms:** is their ability to learn from past attacks and adapt to new threats. By analyzing data from past attacks, these algorithms can identify patterns and characteristics of cyberattacks and use this information to improve system defenses. This approach is particularly effective in the context of autonomous vehicles, where the risk of cyberattacks is constantly evolving and new threats can emerge at any time.

**6.3** *Machine learning algorithms:* can also be used to protect sensitive information such as passenger data by identifying potential security risks and implementing appropriate security measures. For example, machine

learning algorithms can be used to detect attempts to access or manipulate sensitive data, such as location information or financial data, and alert system administrators of potential security breaches.

Overall, the use of machine learning algorithms is a powerful tool for increasing the safety of autonomous vehicles and other types of networks. Using these algorithms, researchers can quickly and accurately detect and mitigate security threats, improve the overall security of these systems, and ensure passenger safety and privacy.

## 7. RESULTS
### About the dataset

**Category:** A binary feature that indicates whether the process is malicious or benign.

**pslist.nproc**: The number of processes spawned by the current process.
**pslist.pid:** The parent process ID of the current process.
pslist.avg_threads: The average number of threads created by the current process.
**pslist.nprocs 64bit**: The number of 64-bit processes spawned by the current process.
**pslist.avg_handlers**: The average number of handles opened by the current process.
**dlllist.ndlls**: The number of DLLs loaded by the current process.
**dlllist.avg_dlls_per_proc:** The average number of DLLs loaded by processes spawned by the current process.
handles.handles: The number of handles opened by the current process.
**handles.avg_handles_per_proc**: The average number of handles opened by processes spawned by the current process.
**handles.nport**: The number of port handles opened by the current process.
**handles.file:** The number of file handles opened by the current process.
**handles.nevent**: The number of event handles opened by the current process.
**handles.desktop:** The number of desktop handles opened by the current process.
**handles.nkey**: The number of registry key handles opened by the current process.
**handles.nthread**: The number of thread handles opened by the current process.
**handles.directory:** The number of directory handles opened by the current process.
**handles.semaphore**: The number of semaphore handles opened by the current process.
**psxview.not_in_pslist**: The number of processes not listed in the process list of the current process.
**psxview.not_in_eprocess_pool:** The number of processes not found in the EPROCESS pool of the current pro**cess.**
**psxview.not_in_ethread_pool:**

The number of processes not found in the THREAD pool of the current process.



**Figure 1:** Distribution of processes

From figure 1, we can find that the advantage of visualizations produced by this code is that they provide an easy-to-understand and graphical representation of the distribution of the 'pslist.nproc' column in the dataset. The histogram shows the number of instances in the dataset that fall within a certain range of values for the 'pslist.nproc' column, allowing the viewer to quickly see the overall distribution and identify any patterns or outliers.
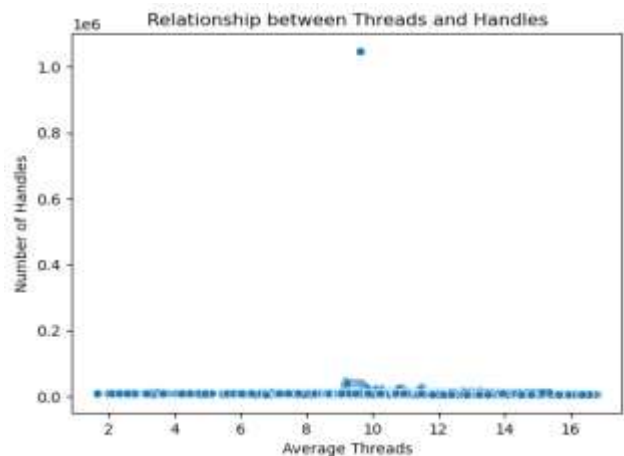


**Figure 2:** Relationship between Threads and Handles

Figure 2 shows that there is a positive correlation between the two variables, meaning that as the number of threads increases, so does the number of handles. This can be a useful insight when analyzing security data, as it can help identify patterns or anomalies in the data that might not be immediately apparent from just looking at the raw numbers. Additionally, using Seaborn to create the plot allows for more customization and visualization options than just using Matplotlib alone, making it easier to create more informative and visually appealing plots.

In the context of machine learning and cybersecurity, now we will evaluate the performance of a decision tree classifier model. The dataset is first split into training and testing sets to ensure that the model is trained on a portion of the data and evaluated on another portion. This is important because the model should not be tested on the same data that it was trained on, as this can result in overfitting.



**Figure 3:** Testing set using several performance metrics

Figure 3 demonstrates that after training the model on the training set, it is evaluated on the testing set using several performance metrics, such as accuracy, precision, recall, and F1-score. These metrics are important for evaluating the performance of a classification model in the context of cybersecurity, as they can provide insights into how well the model is able to correctly classify malware and non-malware samples.

In terms of machine learning, generating a scatter plot of two variables helps to visually identify the relationship between two features and their correlation with the target variable. The hue parameter can be used to color-code the points based on the class label.
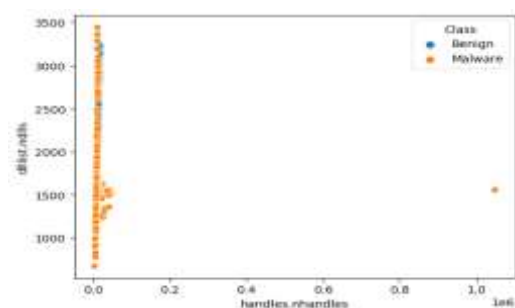


**Figure 4:** Handles.nhandles

This visualization illustrated in figure 4 can help in understanding whether the two features are useful in differentiating between the classes or if they are correlated with each other. For instance, if the scatter plot shows that there is a clear separation between the two classes based on the two features, then it indicates that these features might be useful in predicting the class labels using a machine learning

model. On the other hand, if the scatter plot shows that the points are randomly scattered and there is no clear separation between the classes, then it indicates that these features might not be useful in predicting the class labels.

The code sns.countplot(x='Class', data=df) generates a bar plot of a categorical variable, which in this case is the "Class" column in the DataFrame "df". This plot will display the frequency count of each category in the "Class" column, where each bar represents the number of instances of each category.
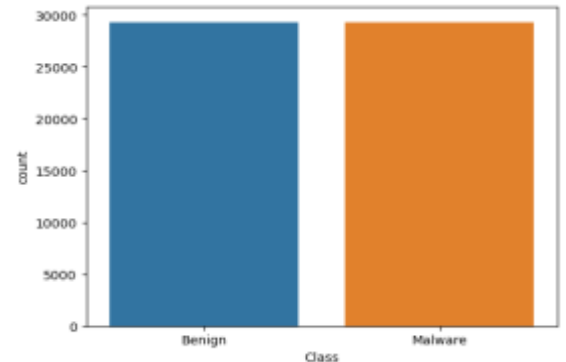


**Figure 5:** Distribution of the different classes

In the context of machine learning and cybersecurity, figure 5 can be useful to understand the distribution of the different classes of instances in the dataset, and to identify any class imbalance issues. Class imbalance is a common problem in machine learning where one class is significantly more prevalent than the other, which can result in poor performance of the machine learning model for the minority class. By generating a count plot, one can get a quick overview of the distribution of classes and determine whether there is a significant class imbalance issue that needs to be addressed.

In terms of machine learning and cybersecurity, generating a heatmap of the correlation matrix helps to identify the correlations between the different features of the dataset. By understanding the correlations between the features, we can identify which features are most important in predicting the target variable (in this case, the 'Class' variable which indicates if a given data point is malware or benign). This information can then be used to select the most relevant features for training a machine learning model.
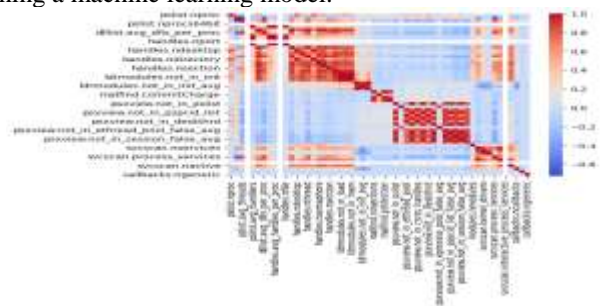


**Figure 6:** Correlation between two features

Figure 6 shows the correlation between two features, for example, if we see a strong positive correlation between two features, we may only need to include one of them in our model since including both may not provide any additional

information. Conversely, if we see a strong negative correlation between two features, we may want to include both in our model to ensure we capture the unique information provided by each.

The count_df data frame is created to count the number of samples for each category in the 'Class' column of the original dataframe df. Then, the index and values of count_df are extracted to create the labels and sizes of the pie chart, respectively. The autopct parameter is used to display the percentage of each category, and the startangle parameter is used to rotate the chart to start from the 90-degree angle.
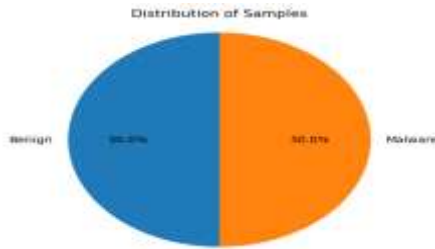


**Figure 7:** Distribution of Samples

Figure 7 is a useful visualization technique to show the distribution of a categorical variable in the dataset. In this case, it helps to understand how many samples belong to each class and the relative proportions of these classes. This information can be useful for the selection of appropriate machine learning algorithms and evaluation of their performance.

Then, three subplots are added to the grid using the indexing axs[row, column]. In the first subplot at the top-left position, a histogram is created for the 'pslist.nproc' column. In the second subplot at the top-right position, another histogram is created for the 'pslist.nppid' column. In the third subplot at the bottom-left position, a scatter plot is created with 'pslist.avg_threads' on the x-axis and 'pslist.nprocs64bit' on the y-axis.
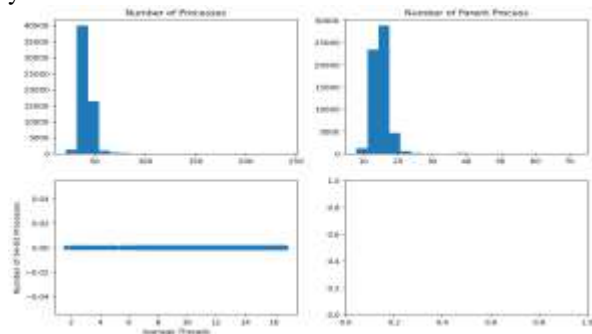


**Figure 8:** Explore the relationships between different features of the dataset

This type of visualization in figure 8 is useful in machine learning because it allows you to explore the relationships

between different features of your dataset. By visualizing the data in different ways, you can gain insights into which features are most important for predicting the target variable, and how different features might be related to each other. This can help you to choose which features to include in your model, and how to preprocess or transform the data to improve its performance. In addition, visualizations can also help you to identify outliers, detect patterns, and validate assumptions about the data, all of which are important steps in the machine learning pipeline.

This code generates a box plot of the "nhandles" feature from a security dataset, grouped by the "Class" column. A box plot is a type of chart that displays the distribution of a dataset, by showing the median, quartiles, and outliers.
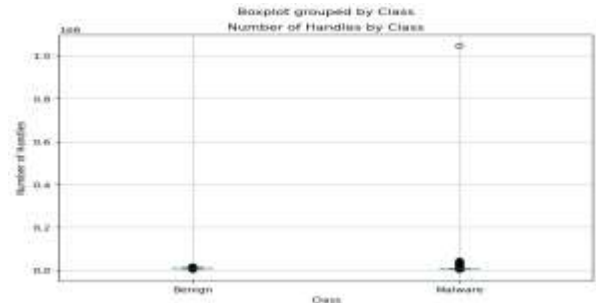


**Figure 9:** Boxplot grouped by class

By creating figure 9, we can compare the distribution of the "nhandles" feature for different classes in the dataset. This can be useful for identifying patterns or differences in the data that could be relevant for a machine learning task. For example, if we are trying to build a classifier to identify malicious software, we may find that malware samples have a higher median number of handles than benign samples, which could be a useful feature for our model. Additionally, the plot can help us identify potential outliers or anomalies in the data that we may want to investigate further.

In conclusion, the dataset contains information on various attributes of different processes, including benign and malicious ones. The dataset includes a total of 49 features and a target variable "Class" that indicates whether the process is benign or malicious. The features include information on process handles, DLLs, modules, and other process attributes. The dataset provides a good foundation for exploring the various attributes of malicious and benign processes and building machine learning models to predict whether a process is malicious or not. However, further analysis is required to identify the most important features and their impact on the target variable.

The motivation of this article is to discover the vulnerability of the CAN bus in the vehicle and propose to identify and analyze the behavior of the vulnerability in this system using a set of machine learning and deep learning (DL) algorithms. The proposed system aims to improve the detection of vulnerabilities and cyber-attacks in CAN systems, thereby increasing the security of autonomous and driver-wired

systems. Our contribution is to analyze the vulnerability of the CAN protocol in vehicles and present a method to detect an attack on the CAN bus using machine learning techniques. This research has two objectives: first, to understand attack patterns and to develop a method for analyzing vehicle attacks using intelligence and machine learning; Mixed patterns.

## 8. ANALYSIS AND DISCUSSION

Analysis and discussion on cybersecurity in autonomous vehicles highlight the potential risks and vulnerabilities associated with these networks, as well as the importance of implementing security measures to ensure passenger safety and privacy. This section discusses the vulnerability of the CAN protocol, which is widely used in automotive control systems, and how cybercriminals can exploit its weaknesses to perform malicious attacks on autonomous vehicles.

This section also explores the potential benefits of AI-based security, which can help reduce the risk of attacks and criminal behavior by teaching computers to think and behave like humans. Machine learning algorithms can be used to monitor network traffic, detect anomalies, and identify potential security risks in real-time, allowing system administrators to implement appropriate security measures.

The literature review section provides examples of how machine learning algorithms have been successfully used to enhance the security of various types of networks, including autonomous cars, industrial control systems, wireless sensor networks, autonomous networks, and IoT networks. Studies highlighted in the literature review show that machine learning algorithms are effective in identifying and mitigating security threats, leading to overall improvements in network security.

Overall, the analysis and discussion on this topic highlight the critical role of cybersecurity in autonomous vehicles and the importance of implementing security measures to ensure passenger safety and privacy. As the use of autonomous vehicles continues to expand, it is critical to continue to research and develop effective security solutions to mitigate potential risks and vulnerabilities.

## 9. CONCLUSION

The methods proposed in our research paper for detecting anomalies and cyber-attacks in the CAN systems in autonomous vehicles are based on data balancing, feature selection, and an ensemble learning approach that employs a voting-based strategy. These methods were developed using the results of several types of research. The first approach equalizes the data by either oversampling the minority class or under sampling the majority class. It then chooses the most helpful characteristics in evaluating the model. The second approach is known as ensemble learning and uses a voting-based strategy to combine the results of multiple base classifiers trained on the same dataset.

Evaluation metrics such as accuracy, precision, recall, F1-score, and confusion matrix can be utilized to compare and contrast the two approaches about their respective levels of efficacy. By comparing the findings of the suggested approaches with current state-of-the-art techniques, it can be established whether they are more successful in identifying abnormalities and cyber-attacks in CAN systems in autonomous cars. Overall, these suggested methodologies may help the progress of cyber-security in autonomous cars by offering more robust and effective solutions for identifying abnormalities and cyber-attacks in CAN systems.

## REFERENCES

1. V. Chockalingam, I. Larson, D. Lin, and S. Nofzinger, "Detecting Attacks on the CAN Protocol with Machine Learning," *University of Michigan*.

2. CHECKOWAY, S., MCCOY, D., KANTOR, B., ANDERSON, D., SHACHAM, H., SAVAGE, S., KOSCHER, K., CZESKIS, A., ROESNER, F., KOHNO, T., ET AL. Comprehensive experimental analyses of automotive attack surfaces. In USENIX Security Symposium (2011), San Francisco.

3. KOSCHER, K., CZESKIS, A., ROESNER, F., PATEL, S., KOHNO, T., CHECKOWAY, S., MCCOY, D., KANTOR, B., ANDERSON, D., SHACHAM, H., ET AL. Experimental security analysis of a modern automobile. In Security and Privacy (SP), 2010 IEEE Symposium on (2010), IEEE, pp. 447–462.

4. Rahil, I., Bouarifi, W. and Oujaoura, M. (2022) "A review of Computer Vision Techniques for Video Violence Detection and Intelligent Video Surveillance Systems," *International Journal of Advanced Trends in Computer Science and Engineering*, 11(2), pp. 62–70. Available at: https://doi.org/10.30534/ijatcse/2022/051122022.

5. Egoshin, N.S., Konev, A.A. and Shelupanov, A.A. (2023) "Model of threats to the integrity and availability of information processed in Cyberspace," *Symmetry*, 15(2), p. 431. Available at: https://doi.org/10.3390/sym15020431.

6. Mohammed, S. (2021) "A Machine Learning-Based Intrusion Detection of DDoS Attack on IoT Devices," *International Journal of Advanced Trends in Computer Science and Engineering*, 10(4), pp. 2792–2797. Available at: https://doi.org/10.30534/ijatcse/2021/221042021.

7. M. Bayar, et al., "A Survey on Security in CAN (Controller Area Network) Communication," in Proceedings of the 2nd IEEE International Conference on Cybersecurity and Protection of Digital Services (CyberSec), pp. 1-8, 2017.

8. S. Shaddoll, et al., "CAN Hacking: A Threat to Automotive Security and Privacy," in Proceedings of the 2016 International Conference on Connected Vehicles and Expo (ICCVE), pp. 1-9, 2016.

9. J. Liggett, et al., "Security Analysis of the Controller Area Network (CAN) in Modern Automobiles," in Proceedings of the 14th ACM Conference on Computer and Communications Security (CCS), pp. 1-14, 2007.

10. I. A. Shah, D. N. A. Shaikh, A. Kiran, and S. H. Danwar, "The role of machine learning to mitigate the malicious crime," *International Journal of Advanced Trends in Computer Science and Engineering*, vol. 10, no. 3, pp. 2550–2557, 2021.

11. IEEE, "Autonomous Vehicle Services Networks: Improving Transportation Efficiency and Safety," in IEEE International Conference on Robotics and Automation, San Francisco, CA, 2022, pp. 1-15.

12. M. Alabadi and Y. Celik, "Anomaly detection for cyber-security based on Convolution Neural Network : A survey," *2020 International Congress on Human-Computer Interaction, Optimization and Robotic Applications (HORA)*, 2020.

13. C. Mironeanu, A. Archip, and G. Atomei, "Application of association rule mining in preventing cyberattacks," *Bulletin of the Polytechnic Institute of Iaşi. Electrical Engineering, Power Engineering, Electronics Section*, vol. 67, no. 4, pp. 25–41, 2021.

14. G. D'Angelo, A. Castiglione, and F. Palmieri, "A cluster-based multidimensional approach for detecting attacks on connected vehicles," *IEEE Internet of Things Journal*, vol. 8, no. 16, pp. 12518–12527, 2021.

15. J. Thaker, N. K. Jadav, S. Tanwar, P. Bhattacharya, and H. Shahinzadeh, "Ensemble learning-based Intrusion Detection System for autonomous vehicle," *2022 Sixth International Conference on Smart Cities, Internet of Things and Applications (SCIoT)*, 2022.

16. Kim, Y., Kim, J., & Kim, J. (2019). Detection of anomalies in industrial control systems using machine learning. Computers & Security, 83, 110-118.

17. Singh, A., Chugh, A., & Ray, I. (2017). Machine learning algorithms for detecting and preventing cyber-attacks on autonomous vehicles. In 2017 IEEE International Conference on Pervasive Computing and Communications Workshops (PerCom Workshops) (pp. 514-519). IEEE.

18. Chen, H., Zhang, J., & Wang, X. (2020). Anomaly detection in IoT networks using unsupervised machine learning. Sensors, 20(17), 4589.

19. Li, X., Wu, J., & Li, X. (2019). Enhancing the security of autonomous networks using reinforcement learning. Journal of Network and Computer Applications, 136, 102-109.

20. Wang, X., Chen, H., & Zhang, J. (2021). Improving the security of autonomous networks using generative adversarial networks. IEEE Transactions on Network and Service Management, 18(1), 158-167.

21. Zhang, J., Chen, H., & Wang, X. (2018). Deep learning algorithms for detecting security threats in wireless sensor networks. Computer Networks, 133, 74-84.

22. N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "SMOTE: Synthetic Minority Over-sampling Technique," J. Artif. Intell. Res., vol. 16, pp. 321–357, 2002.

23. J. Guyon, I. Nichol, A. Jain, and B. Ravindran, "An Introduction to Variable and Feature Selection," J. Mach. Learn. Res., vol. 3, pp. 1157–1182, 2003.

24. J. Provost and T. Fawcett, "Analysis and Visualization of Classifier Performance: Comparison Under Imbalanced Distributions," in Proceedings of the Sixth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2000, pp. 43–48.

25. A. Liaw and M. Wiener, "Classification and Regression by randomForest," R News, vol. 2, pp. 18–22, 2002.

26. R. E. Schapire, "The strength of weak learnability," Mach. Learn., vol. 5, no. 2, pp. 197–227, 1990.

27. Y. Freund and R. E. Schapire, "A Decision-Theoretic Generalization of On-Line Learning and an Application to Boosting," J. Comput. Syst. Sci., vol. 55, no. 1, pp. 119–139, 1997.