# Enhanced intrusion detection through dual reduction and robust mean

**Archi Taha, Mohammed Benattou**
LaRI Computer Science Research Laboratory, Ibn Tofail University, Kenitra, Morocco

## Article Info

## ABSTRACT

The exponential growth of online networks necessitates a paradigm shift in intrusion detection systems (IDS). Traditional methods falter under the massive influx of data, resulting in high false positives and reduced detection accuracy. This research introduces a novel approach combining principal component analysis (PCA) and linear discriminant analysis (LDA), augmented by robust generalized sample mean, to enhance IDS performance. PCA efficiently reduces data dimensionality, while LDA extracts critical features that differentiate normal network traffic from anomalies. The robust generalized sample mean counteracts the influence of outliers, ensuring accurate and reliable analysis. Implemented on the UNSW-NB15 dataset, our method achieves an average 6% reduction in false positives and a 10% increase in detection rate. Additionally, our testing methodology closely mirrors real-world conditions, making the results more representative of practical scenarios compared to existing work. These advancements demonstrate substantial improvements in IDS performance and robustness over existing techniques.

*Corresponding Author:*

Taha Archi
LaRI Computer Science Research Laboratory, Ibn Tofail University
Kenitra, Morocco
Email: architaha1@gmail.com

## 1. INTRODUCTION

In today's digital landscape, the protection of online domains against cyber threats has become increasingly critical. Intrusion detection systems (IDS) serve as a vital line of defense, designed to identify and thwart unauthorized activities within networks. These systems employ a wide array of techniques, from traditional rule-based methods to cutting-edge machine learning algorithms, to safeguard digital infrastructures. The significance of IDS lies in their capacity for continuous monitoring and alerting, which is essential to staying ahead of increasingly sophisticated cyber adversaries [1]. As cyber threats evolve, the integration of artificial intelligence (AI) and advanced data analytics into IDS has become indispensable [2].

However, the rapid expansion of data and the growing complexity of digital networks have intensified the necessity for more robust IDS to counter emerging cyber threats [3]. Adversaries continuously exploit vulnerabilities, necessitating more agile and adaptive IDS solutions. Traditional rule-based systems, while effective in certain scenarios, often rely on static signatures and predefined rules, which are inadequate in the face of dynamic and evolving threats. This limitation can result in increased false positives and missed detections, leaving networks vulnerable to undetected breaches [4].

To address these shortcomings, the incorporation of machine learning techniques into IDS frameworks has ushered in a new era of threat detection. Deep learning, particularly through convolutional neural networks

(CNNs), has emerged as a powerful tool for identifying complex patterns and anomalies within large datasets. These models enable IDS to adapt and learn from evolving attack strategies with greater precision. Hybrid classification approaches, which combine multiple machine learning techniques, have further enhanced IDS performance, showcasing the potential of machine learning to significantly improve network security [5], [6]. Additionally, the application of deep learning in specialized environments, such as mobile ad hoc networks (MANET), has shown promising results [7].

Despite these advancements, the sheer volume and high dimensionality of network traffic data present significant challenges for machine learning-based IDS. Managing such large-scale data can overwhelm these systems, compromising their efficiency and accuracy. Moreover, machine learning-based IDS must contend with issues such as handling diverse data sources, privacy concerns, susceptibility to adversarial attacks, and the complexities of deployment in large-scale environments [8], [9]. The presence of noise and irrelevant features within datasets exacerbates these challenges, leading to decreased performance [10].

To mitigate these issues, dimensionality reduction techniques like principal component analysis (PCA) and linear discriminant analysis (LDA) are essential [11], [12]. PCA helps reduce dimensions by capturing the most significant variance, while LDA improves classification by maximizing class separability. However, both techniques have inherent limitations—PCA is sensitive to outliers, and LDA assumes linear separability, which may not always be applicable in complex datasets [13], [14].

Combining dimensionality reduction with advanced machine learning methods can significantly enhance the performance of IDS [15]. However, outliers and irregularities in datasets remain a persistent challenge. Outliers can distort data structures, negatively impacting both dimensionality reduction and machine learning techniques. To alleviate these effects, preprocessing steps such as outlier removal and the application of robust statistical methods are crucial [16], [17]. Techniques like the generalized sample mean (GSM) can enhance IDS robustness, ensuring more accurate threat detection [18].

Achieving robust network security requires the integration of machine learning, dimensionality reduction, and robust statistical methods. By leveraging the strengths of each approach while addressing their weaknesses, IDS can evolve into formidable defenses against increasingly sophisticated cyber threats [19], [20]. As the threat landscape continues to evolve, the demand for innovative, adaptive IDS solutions to protect digital infrastructures becomes ever more pressing.

Our research tackles these critical challenges by developing an IDS framework that integrates PCA and LDA [21]. This innovative framework combines the strengths of both dimensionality reduction techniques to enhance the effectiveness of IDS in managing large datasets. By incorporating the GSM into both PCA and LDA, we aim to improve IDS robustness against outliers, thereby addressing a significant gap in current IDS methodologies. Our approach is designed to yield results that closely mirror real-life scenarios, achieving a high detection rate (DR) and a lower false positive rate (FPR) in a realistic setting using the UNSW-NB15 dataset. This demonstrates superior performance and robustness compared to traditional methods.

This paper consists of the following: section 2 offers a detailed theoretical foundation for GSM and its integration with PCA-LDA, providing insights into how these techniques can be effectively applied to real-world data. Section 3 presents an overview of our IDS model, detailing the implementation and the specific contributions of each component. In section 4, we rigorously evaluate the performance and robustness of the proposed method through extensive experimentation and analysis. Finally, section 5 summarizes the key findings from our evaluation and suggests potential directions for future research, emphasizing the ongoing need for innovation in IDS methodologies to keep pace with the ever-evolving cyber threat landscape.

## 2. FORMAL MODEL

This section explores the application of the GSM to enhance the robustness of PCA and LDA against outliers. We begin by examining the mathematical foundations of GSM and demonstrating its effectiveness in mitigating the impact of anomalous data points. Subsequently, we outline the overarching strategy for our combined PCA-LDA approach, explaining the rationale behind its integration. By combining the outlier-resistant properties of GSM with the complementary strengths of PCA and LDA, we aim to significantly improve the performance of our IDS.

### 2.1. Generalized sample mean

In various fields, data analysis and optimization techniques often seek robustness to enhance reliability. The generalized mean emerges as a powerful tool in this context, providing flexibility and adaptability through

its parameter p, which controls the emphasis on different parts of the dataset. Unlike traditional means, which may be sensitive to outliers and deviations from the norm, the generalized mean can adjust to diverse data characteristics by varying p. Consider a vector $\mathbf{V}$ of M positive values:

$$\mathbf{V} = (v_1, v_2, \ldots, v_M), \text{ where } v_i > 0 \text{ for } i = 1, 2, \ldots, M$$

The generalized mean $G_m$ of $\mathbf{V}$, for a power $p$ with $p \neq 0$, is defined as (1):

$$G_m = \left( \frac{1}{M} \sum_{i=1}^{M} v_i^p \right)^{\frac{1}{p}} \tag{1}$$

- When $p = 1$, the generalized mean reduces to the arithmetic mean.
- When $p = 0$, it approaches the geometric mean (with a slight modification to account for zero values)
- When $p = -1$, it becomes the harmonic mean.

The research in [22], [23] propose a significant improvement to PCA by leveraging the GSM. This approach replaces the traditional, outlier-sensitive mean with the generalized mean, enhancing the robustness of mean estimation. By adjusting the parameter 'p,' which controls the influence of individual data points, this method effectively mitigates the impact of outliers. The traditional sample mean, is viewed as the centroid in terms of least squares, as given in (2), where m being the vector mean.

$$S_m = \arg\min_m \left[ \frac{1}{M} \sum_{i=1}^{M} \|v_i - m\|_2^2 \right] \tag{2}$$

In (2) highlights a critical issue: the objective function's reliance on squared distances amplifies the influence of outliers in training samples. To address this challenge and achieve more robust sample mean estimation in the presence of outliers, a novel optimization approach is proposed. This approach involves replacing the arithmetic mean used in (2) with the generalized mean.

$$g_{\text{sm}} = \arg\min \left( \frac{1}{M} \sum_{i=1}^{M} \left( \left( \|v_i - m\|_2^2 \right)^p \right)^{\frac{1}{p}} \right) \tag{3}$$

When $p = 1$, this problem becomes equivalent to (2). As the value of $p$ decreases, the influence of large numbers on the objective function diminishes. This demonstrates that setting $p$ greater than 1 ($p > 1$) mitigates the adverse impact of outliers. We can leverage the property that $x^p$ is a monotonically increasing function for positive values of $x$ ($x > 0$) to reformulate this problem for $p > 0$ as (4):

$$g_{\text{sm}} = \arg\min \left( \sum_{i=1}^{M} \left( \|v_j - m\|_2^2 \right)^p \right) \tag{4}$$

To qualify as a local minimum, $g_{\text{sm}}$ requires the gradient of the objective function in (4) with respect to $m$ to reach zero. In simpler terms as (5):

$$\frac{\partial}{\partial m} \sum_{j=1}^{M} \left( \|v_j - m\|_2^2 \right)^p = 0 \tag{5}$$

Unfortunately, finding a closed-form solution for the equation above proves challenging. The research in [22], [23] offers a valuable insight: the generalized mean of positive numbers can be expressed as a non-negative linear combination of the elements in the set. This expression can be further simplified as shown in (6):

$$\sum_{i=1}^{M} a_i^p = b_1 a_1 + \cdots + b_M a_M \quad \text{with} \quad b_i = a_i^{p-1} \quad \text{and} \quad i = 1, \ldots, M \tag{6}$$

In the derivation, we decompose (4) into the form of (6) and consider the weight $b_i$ in (6) as a constant. Then, in (4) can be approximated by a quadratic function of $\|v_i - m\|_2^2$, which can easily be optimized.

$$\sum_{j=1}^{M} \left( \|v_j - m\|_2^2 \right)^p \approx \sum_{j=1}^{M} \gamma_j^{(iter)} \|v_j - m\|_2^2 \tag{7}$$

Where:

$$\gamma_j^{(iter)} = \left( \left| v_j - m^{(iter)} \right|_2^2 \right)^{p-1} \tag{8}$$

The approximation reaches its peak accuracy when $m = m^{(iter)}$. The next step involves finding the value of $g_{sm}^{(iter+1)}$ that minimizes the approximated function. This minimization considers the previously computed values of $\gamma_j^{(iter)}$. The solution is given by (9):

$$g_{\text{sm}}^{(iter+1)} = \frac{1}{\sum_{j=1}^{M} \gamma_j^{(iter)}} \sum_{j=1}^{M} \gamma_j^{(iter)} v_j \tag{9}$$

The following algorithm provides a summary of the steps to find the GSM.

---

**Algorithm 1** Generalized sample mean

---

**Require:** matrix data and $p$
1:  $T \leftarrow 0$
2:  $g_{\text{sm}}^{(T)} \leftarrow m$
3:  **repeat**
4:      **Approximation:** for fixed $g_{\text{sm}}^{(T)}$, compute $\gamma_j^{(T)}$ using (8)
5:      **Minimization:** Using the computed $\gamma_j^{(T)}$, update $g_{\text{sm}}^{(T+1)}$ according to (9)
6:      $T \leftarrow T + 1$
7:  **until** a termination condition is met
**Ensure:** $g_{\text{sm}} = g_{\text{sm}}^{(T)}$

---

To illustrate the effect of the GSM, we generate 100 random data points following a Gaussian distribution. These points reside in a two-dimensional space. Here, we introduce 30 outliers, a considerable portion of the total 100 data points, outliers are marked red on the plot for easy identification.

The script calculates two different means: the arithmetic mean and the GSM (using only one iteration). We observe in Figure 1, that the arithmetic mean is heavily influenced by the outliers. In contrast, the GSM, calculated with power values of p = 0.3 and p = 0.7, shows less sensitivity to outliers. Interestingly, the mean with p = 0.3 is closer to the "normal" points compared to p = 0.7. This demonstrates that the GSM is less affected by outliers, and the power parameter (p) has a significant impact on the calculated mean.
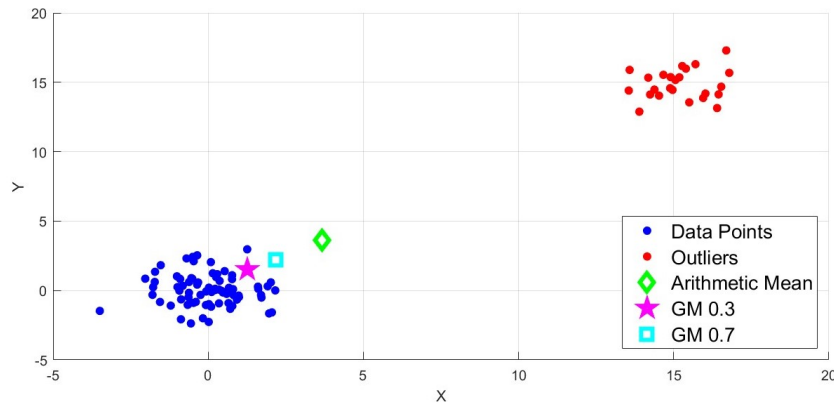


Figure 1. Visualization of data distribution with outliers and mean comparisons

## 2.2. Principal component analysis-linear discriminant analysis (GSM) combination approach

PCA for dimensionality reduction comes with trade-offs. It excels at capturing the most variable aspects of data, but it ignores class labels. This can lead to features that are not ideal for classification, potentially keeping irrelevant information instead of focusing on what separates the classes.

LDA tackles classification by maximizing separation between classes. However, it makes assumptions about the data that may not be true in practice, like normally distributed classes with similar properties. These unrealistic assumptions can hinder LDA's performance. Additionally, outliers can significantly impact LDA's results, reducing its effectiveness.

A combined approach using PCA and LDA can address these limitations. By using PCA first, we can reduce the number of features and lessen the impact of noise. Then, LDA takes over to maximize class separation in the lower-dimensional space. This two-step process ensures we keep the most informative features while making them highly relevant for classifying the data.

Furthermore, incorporating a robust mean, like the GSM, strengthens this combined method. This specific type of mean is less affected by outliers and non-normal data distributions, making the overall approach more reliable and accurate. The following steps outline the GSM PCA-LDA process for dimensionality reduction.

a. Extract key features: Compute the principal components of the data matrix $Y$ using GSM-PCA.

b. Select informative components: Retain the top $m$ principal components with the highest variance, where $m$ defines the target dimensionality.

c. Project data to lower dimension: Obtain a lower-dimensional representation $R$ by projecting the data using the selected components.

d. Calculate class-wise central tendency: Compute the GSM for each class in the projected data $R$, representing the center of each class.

e. Compute class separation: Calculate the between-class and within-class scatter matrices using the generalized means, capturing the distribution of data points within and between classes.

f. Find optimal transformation: Obtain the projection matrix $Z$ via eigenvalue decomposition of the scatter matrices, identifying the best direction to project the data for class separation.

g. Transform data to final representation: Transform the data to the new subspace using $Z$ to obtain the final reduced-dimensional representation $Z_d$.

## 3. CLASSIFICATION METHOD

To evaluate the impact of training data size on IDS performance, we manipulate the volume of training data. A key challenge in IDS training is differentiating between various attack types, often due to inherent similarities among them. We address this challenge by treating each attack type as a separate class, in addition to a class for normal network connections. We then incrementally increase the number of connections for both normal and attack types within each class, applying the same approach to the testing data. This method allows us to systematically assess how data quantity influences the effectiveness of IDS models in identifying diverse network security threats.

The UNSW-NB15 dataset serves as a critical resource for this research. It offers a unique blend of real-world network traffic captured from a controlled environment and synthetically generated contemporary attacks, reflecting the modern threat landscape [24]. This combination makes UNSW-NB15 ideal for training and evaluating machine learning models for effective IDS development [25]. The dataset encompasses a wide range of attack categories, including denial-of-service (DoS), exploits, worms, backdoors, reconnaissance, analysis, shellcode, and generic attacks. This variety allows researchers to build robust IDS models capable of detecting a broad spectrum of network security threats. Despite its strengths, UNSW-NB15 assumes a normal data distribution, potentially affecting results due to outliers.

As the number of features in the UNSW-NB15 dataset can be high (49 features, reaching 211 features after one-hot encoding), we employ a dimensionality reduction process to reduce the feature set to a more manageable size. This helps improve the efficiency and accuracy of machine learning algorithms. We utilize a

two-step approach for this purpose. First, we perform PCA-GSM to reduce the number of features from 211 to 10. Subsequently, we apply LDA-GSM to further reduce the features to 2 dimensions.

Figure 2 gives an overview about the IDS model. The process begins by dividing the data into training and testing matrices, both containing normal traffic and attack samples. These matrices undergo dimensionality reduction through PCA-GSM, creating projected versions of each. They then undergoes a second reduction using LDA-GSM, further refining the features. This dual reduction approach results in compact, information-rich representations of both test and training data. A classifier is then trained on the training reduced matrice, learning to distinguish between normal network activity and various cyber threats such as DoS attacks, exploits, worms, and backdoors.



Figure 2. IDS model based on the combination of PCA-LDA (GSM)

## 4. FINDINGS AND ANALYSIS

This section evaluates the performance of our proposed PCA-LDA (GSM) method compared to traditional PCA. To assess the effectiveness of these approaches, we employ two standard metrics: DR and FPR. These metrics quantitatively measure the system's ability to accurately identify intrusions while minimizing false alarms.

For our experiments, we utilize decision trees, regression trees, and k-nearest neighbors (KNN) as classification models. We partition the dataset into training and testing matrices to evaluate model performance under various conditions. Specifically, we incrementally increase the size of both the training and test sets to assess the impact of data volume on model accuracy. To enhance the reliability and generalizability of our results, each performance metric for a given training set size is the average of 30 independent test runs. This approach mitigates the influence of random fluctuations and provides a more robust performance estimate.

To address the class imbalance in the UNSW-NB15 dataset, we increase the training data for all attack categories while maintaining a fixed representation of the underrepresented attack types. Specifically, we set 400 instances each for analysis and backdoor, 200 instances for shellcode, and 30 instances for worms. The remaining attack categories (fuzzers, DoS, exploits, generic, and reconnaissance) are assigned an equal number of instances to ensure a balanced dataset. This approach allows us to enhance the model's learning of underrepresented attacks while maintaining an equitable distribution across all attack types, thereby improving the model's overall accuracy and robustness.

The test data constitutes 10% of the training set, maintaining a consistent representation of underrepresented attack categories. To ensure a comprehensive evaluation, we iteratively test each scenario 30 times, covering a diverse range of cases. This approach enhances the model's ability to generalize across different network attacks, thereby improving both accuracy and robustness in real-world scenarios.

We investigate the influence of the GSM's power parameter (p) on both PCA and LDA. To achieve robust results, we conduct tests with p-values of 0.3 and 0.7. The DR and FPR are subsequently calculated as averages of thirty randomly selected evaluations for each training data size, demonstrating the consistency of results for large datasets.

### 4.1. Detection rate classification

In this section, we analyze the DR using PCA and a combination of robust PCA and LDA. The DR is a crucial metric for assessing the performance of an IDS, as it precisely quantifies the system's ability to accurately identify malicious activities. By comparing DR values obtained from different methodologies and classification models, we can evaluate their effectiveness in detecting diverse attack types and identify the optimal approach for enhancing intrusion detection capabilities. It is calculated using the following formula:

$$DR = \frac{TP}{TP + FN} \times 100$$

Where:

True positives (TP): Instances where the model correctly identifies positive cases.

False negatives (FN): Cases where the model incorrectly identifies negative cases.

By examining the DR, we can assess the effectiveness of each model in accurately identifying different types of attacks. This analysis helps us understand the strengths of PCA and the robust PCA-LDA combination in detecting each specific attack category.

Across Figures 3 to 5, a consistent trend emerges in the DR% depicted, showcasing similar outcomes overall. Notably, the integration of PCA-LDA (GSM) yields a significant improvement in DR% across all tree classification methods. Particularly striking are the results observed in pure DR% with regression tree classification, where DR% ranges impressively from 80% to 90%, compared to a lower 66% with PCA, peaking at 78% optimally. Furthermore, employing the combination approach with regression trees results in a remarkable 15% increase in DR%.
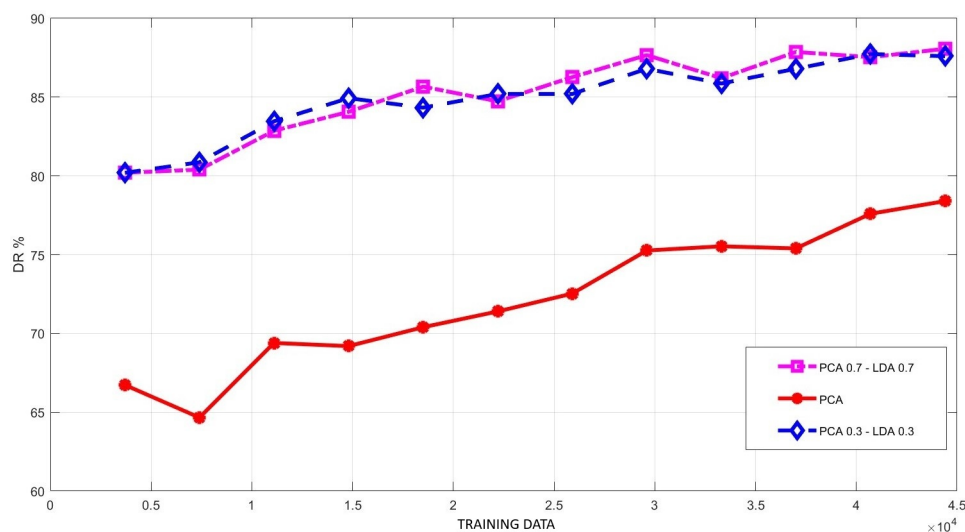


Figure 3. DR with decision tree
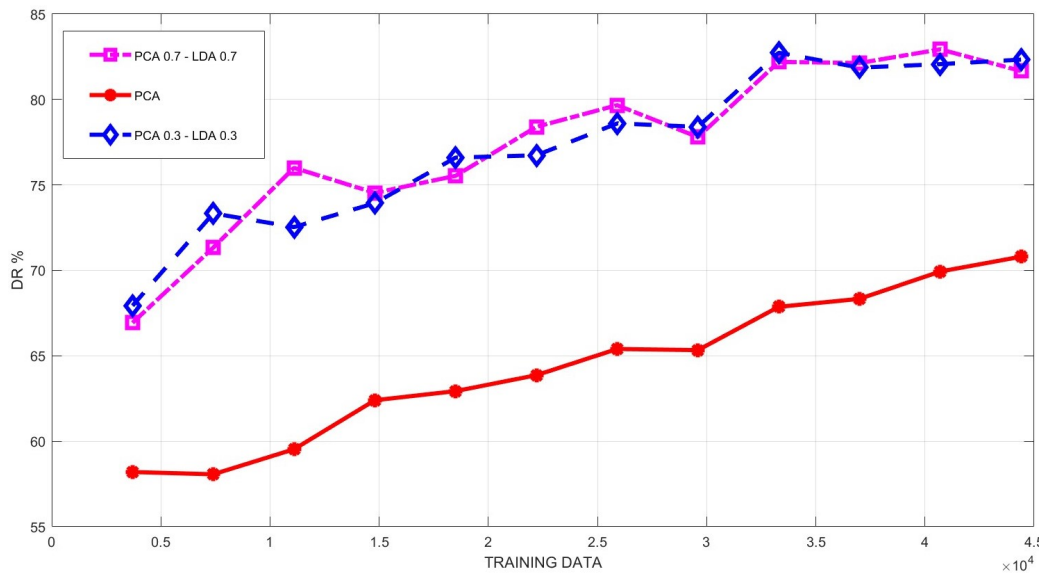


Figure 4. DR with regression tree

Figure 5. DR with KNN

Conversely, KNN exhibits the least favorable performance in both DR% and the overall enhancement achieved through the PCA-LDA combination. Starting at a modest 66%, KNN reaches a maximum of 82%, compared to PCA's initial 58%, which only climbs to 70%. However, despite this, the combined methodology boosts the overall DR% by 10%.

In the case of decision trees, PCA-LDA begins with a 71% DR%, reaching up to 81%, while PCA starts at 55%, ascending only to 62%, marking the lowest DR% among the three cases. Notably, decision trees demonstrate the highest overall boost in DR%, averaging over 20%, emphasizing how the combination of PCA-LDA (GSM) can enhance DR% even in unfavorable conditions. In conclusion, the integration of PCA-LDA (GSM) consistently improves DR% across all tree cases, even with minimal training data, yielding robust and reliable results.

## 4.2. False positive rate classification

In this section, we evaluate the FPR using PCA and the combination of robust PCA and LDA. The FPR is another critical metric that quantifies the proportion of benign instances erroneously classified as malicious by the model. Minimizing the FPR is essential to reduce the number of false alarms and maintain system efficiency while ensuring adequate security. It is calculated using the following formula:

$$FPR = \frac{FP}{FP + TN} \times 100$$

Where:

False positives (FP): Situations where the model incorrectly identifies positive cases.

True negatives (TN): Scenarios where the model correctly identifies negative cases.

Analyzing the FPR is essential to ensure that the IDS does not produce an excessive number of false alarms, which can lead to unnecessary interventions and reduced trust in the system. This section compares the performance of PCA and the robust PCA-LDA combination in minimizing false positives across different attack categories. When examining the FPR across Figures 6 to 8, the combination of PCA-LDA (GSM) consistently yields superior FPR across all classification methods. It's notable that the FPR in the case of regression tree tends to be higher compared to both decision tree and KNN methods. Specifically, when employing PCA, the FPR begins at 16% and decreases to 10%, while the combination method starts at 9% and decreases to less than 4%. Overall PCA-LDA (GSM) reduces the FPR by an average of 6%.
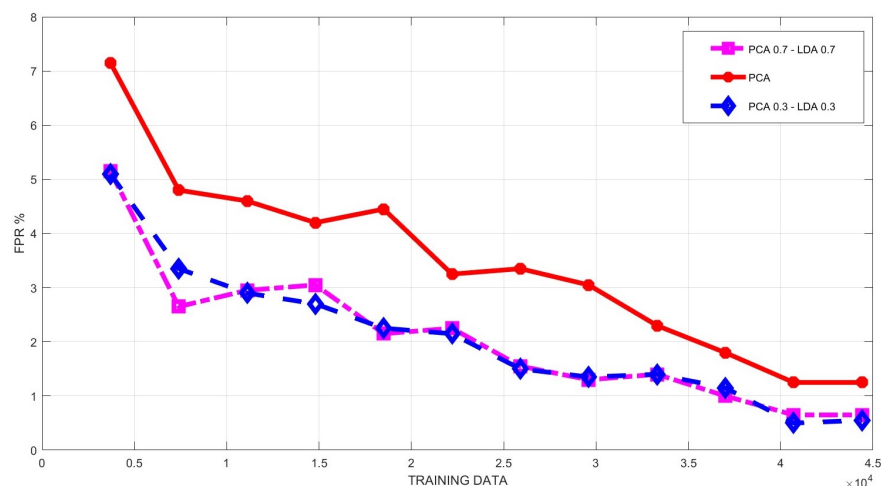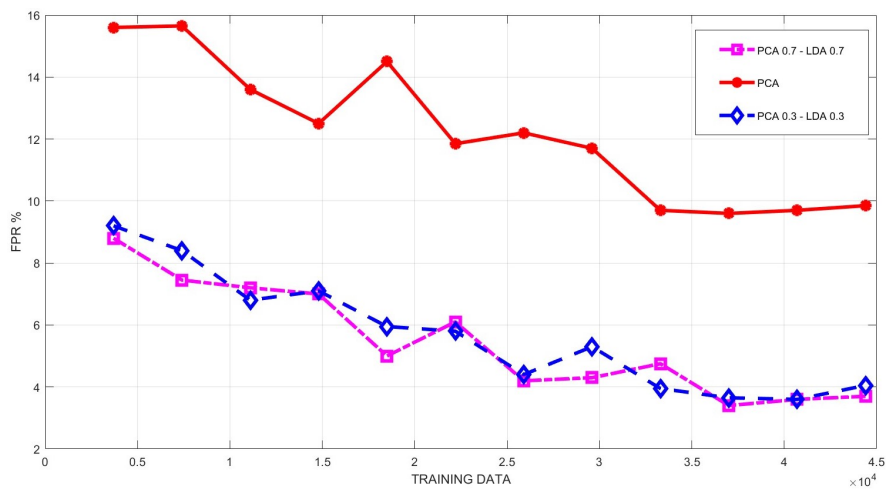
Figure 6. FPR with decision tree
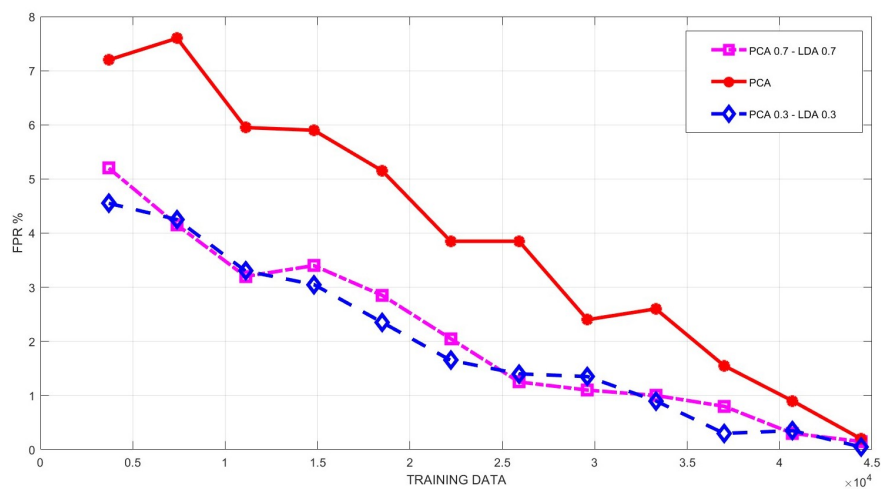


Figure 7. FPR with regression tree



Figure 8. FPR with KNN

Switching to decision tree classification, the FPR starts at 7% with PCA and diminishes to 1%, whereas employing the PCA-LDA (GSM) combination initiates at 5% and declines to 0.5%. It's evident that the PCA-LDA (GSM) combination consistently yields a superior FPR, averaging at 1.5%. In the case of KNN classification, PCA start at a 7% FPR, while the PCA-LDA (GSM) combination starts at 5%. With an increase in training data, both methods converge to nearly identical FPRs of 0.5%. Overall, the combination approach reduces the FPR by 2% in this scenario.

## 5.    CONCLUSION

This study investigated the effectiveness of a combined PCA and LDA for IDS. The goal of combining PCA and LDA was to reduce dimensionality while retaining the most critical information for intrusion detection. The GSM, a robust mean, further improves performance when handling large datasets by mitigating the impact of outliers. The parameter 'p' in the GSM plays a crucial role in achieving this by lessening the influence of outliers in the process. The findings of the study are compelling. The PCA-LDA (GSM) combination consistently led to significant improvements across all three classification methods (regression tree, decision tree, and KNN) compared to using PCA alone. These improvements manifested as increased DR, reduced FPR for training set sizes. Notably, the improvement was most pronounced for regression trees, showcasing a remarkable 15% increase in DR. One of the most significant aspects of this study is the fact that the consistent improvement in performance across various metrics (DR, FPR, and accuracy) was achieved while handling large datasets. The findings held true even with repeated evaluations (30 times for each training data size). This strengthens the validity of the PCA-LDA (GSM) approach for real-world IDS applications, which typically deal with massive amounts of network traffic data.

## REFERENCES

[1]  Amarudin, R. Ferdiana, and Widyawan, "A systematic literature review of intrusion detection system for network security: research trends, datasets and methods," in *2020 4th International Conference on Informatics and Computational Sciences (ICICoS)*, IEEE, 2020, pp. 1–6, doi: 10.1109/ICICoS51170.2020.9299068.

[2]  H. V. Vo, H. P. Du, and H. N. Nguyen, "AI-powered intrusion detection in large-scale traffic networks based on flow sensing strategy and parallel deep analysis," *Journal of Network and Computer Applications*, vol. 220, 2023, doi: 10.1016/j.jnca.2023.103735.

[3]  K. M. Sudar, P. Nagaraj, P. Deepalakshmi, and P. Chinnasamy, "Analysis of intruder detection in big data analytics," in *2021 International Conference on Computer Communication and Informatics (ICCCI)*, IEEE, 2021, pp. 1–5, doi: 10.1109/IC-CCI50826.2021.9402402.

[4]  O. H. Abdulganiyu, T. A. Tchakoucht, and Y. K. Saheed, "A systematic literature review for network intrusion detection system (IDS)," *International Journal of Information Security*, vol. 22, no. 5, pp. 1125–1162, 2023, doi: 10.1007/s10207-023-00682-2.

[5]  A. Khraisat, I. Gondal, P. Vamplew, and J. Kamruzzaman, "Survey of intrusion detection systems: techniques, datasets and challenges," *Cybersecurity*, vol. 2, no. 1, 2019, doi: 10.1186/s42400-019-0038-7.

[6]  S. Choudhary and N. Kesswani, "A hybrid classification approach for intrusion detection in IoT network," *Journal of Scientific and Industrial Research*, vol. 80, no. 9, pp. 809–816, 2021, doi: 10.56042/jsir.v80i09.43878.

[7]  S. Laqtib, K. El Yassini, and M. L. Hasnaoui, "A deep learning methods for intrusion detection systems based machine learning in manet," in *Proceedings of the 4th International Conference on Smart City Applications*, New York, USA: ACM, 2019, pp. 1–8, doi: 10.1145/3368756.3369021.

[8]  M. Sarhan, S. Layeghy, N. Moustafa, M. Gallagher, and M. Portmann, "Feature extraction for machine learning-based intrusion detection in IoT networks," *Digital Communications and Networks*, vol. 10, no. 1, pp. 205–216, 2024, doi: 10.1016/j.dcan.2022.08.012.

[9]  A. Thakkar and R. Lohiya, "A review on challenges and future research directions for machine learning-based intrusion detection systems," *Archives of Computational Methods in Engineering*, vol. 30, no. 7, pp. 4245–4269, 2023, doi: 10.1007/s11831-023-09943-8.

[10]  M. A. Umar, Z. Chen, K. Shuaib, and Y. Liu, "Effects of feature selection and normalization on network intrusion detection," *TechRxiv*, pp. 1–18, 2024, doi: 10.36227/techrxiv.12480425.v3.

[11]  T. Kurita, "Principal component analysis (PCA)," in *Computer Vision*, Cham: Springer International Publishing, 2020, pp. 1–4, doi: 10.1007/978-3-030-03243-2_649-1.

[12]  S. Balakrishnama and A. Ganapathiraju, "Linear discriminant analysis-a brief tutorial," *Institute for Signal and Information Processing*, Department of Electrical and Computer Engineering, Mississippi State University, Mississippi, United States of America, 1998.

[13]  A. Taha and B. Mohammed, "An automatic truncated mean approach for PCA in intrusion detection systems," in *2023 10th International Conference on Wireless Networks and Mobile Communications (WINCOM)*, IEEE, 2023, pp. 1–6, doi: 10.1109/WIN-COM59760.2023.10323011.

[14]  R. Abdulhammed, H. Musafer, A. Alessa, M. Faezipour, and A. Abuzneid, "Features dimensionality reduction approaches for machine learning based network intrusion detection," *Electronics*, vol. 8, no. 3, 2019, doi: 10.3390/electronics8030322.

[15]  G. Lu and X. Tian, "An efficient communication intrusion detection scheme in AMI combining feature dimensionality reduction and improved LSTM," *Security and Communication Networks*, vol. 2021, 2021, doi: 10.1155/2021/6631075.

[16]    G. Meng, B. Wang, Y. Wu, M. Zhou, and T. Meng, "A hybrid dimensionality reduction method for outlier detection in high-dimensional data," *International Journal of Machine Learning and Cybernetics*, vol. 14, no. 11, pp. 3705–3718, 2023, doi: 10.1007/s13042-023-01859-w.

[17]    V. Vaidya and J. Vaidya, "Impact of dimensionality reduction on outlier detection: an empirical study," in *2022 IEEE 4th International Conference on Trust, Privacy and Security in Intelligent Systems, and Applications (TPS-ISA)*, IEEE, 2022, pp. 150–159, doi: 10.1109/TPS-ISA56441.2022.00028.

[18]    A. Taha and B. Mohammed, "A robust intrusion detection model based on a combination of PCA-GM and trunced LDA," in *2023 10th International Conference on Wireless Networks and Mobile Communications (WINCOM)*, IEEE, 2023, pp. 1–6, doi: 10.1109/WINCOM59760.2023.10322896.

[19]    R. A. Al Hasan and E. K. Hamza, "An improved intrusion detection system using machine learning with singular value decom position and principal component analysis," *International Journal of Intelligent Engineering and Systems*, vol. 16, no. 4, pp. 25–38, 2023, doi: 10.22266/ijies2023.0831.03.

[20]    A. Thakkar, N. Kikani, and R. Geddam, "Fusion of linear and non-linear dimensionality reduction techniques for feature reduction in LSTM-based intrusion detection systems," *Applied Soft Computing*, vol. 154, 2024, doi: 10.1016/j.asoc.2024.111378.

[21]    J. Yang and J. Y. Yang, "Why can LDA be performed in PCA transformed space?," *Pattern Recognition*, vol. 36, no. 2, pp. 563–566, 2003, doi: 10.1016/S0031-3203(02)00048-1.

[22]    J. Oh and N. Kwak, "Generalized mean for robust principal component analysis," *Pattern Recognition*, vol. 54, pp. 116–127, 2016, doi: 10.1016/j.patcog.2016.01.002.

[23]    J. Oh, N. Kwak, M. Lee, and C. H. Choi, "Generalized mean for feature extraction in one-class classification problems," *Pattern Recognition*, vol. 46, no. 12, pp. 3328–3340, 2013, doi: 10.1016/j.patcog.2013.06.018.

[24]    N. Moustafa and J. Slay, "UNSW-nb15: a comprehensive data set for network intrusion detection systems (UNSW-NB15 network data set)," in *2015 Military Communications and Information Systems Conference (MilCIS)*, IEEE, 2015, pp. 1–6, doi: 10.1109/MilCIS.2015.7348942.

[25]    N. Moustafa and J. Slay, "The evaluation of network anomaly detection systems: statistical analysis of the UNSW-NB15 data set and the comparison with the KDD99 data set," *Information Security Journal*, vol. 25, no. 1–3, pp. 18–31, 2016, doi: 10.1080/19393555.2015.1125974.

## BIOGRAPHIES OF AUTHORS

**Archi Taha** 🆔 🔗 🔗 🔗 is a Ph.D. student at Ibn Tofail University. He obtained his master's degree from the same university in 2019. His primary research interest lies in the field of cybersecurity, specifically focusing on intrusion detection systems (IDS). He is dedicated to optimizing and developing feature extraction algorithms for these systems. He can be contacted at email: architaha1@gmail.com.

**Mohammed Benattou** 🆔 🔗 🔗 🔗 is a Professor of Computer Science at Ibn Tofail University in Kenitra, where he led the Computer Science and Telecommunication Laboratory. He has held various academic positions in France, including at the University of Pau, University of Orsay Paris XI, 3IL, and Xlim Laboratory. His research interests focus on distributed testing, secure testing, and software testing. He can be contacted at email: mbenattou@yahoo.fr.