# Enhancing Intrusion Detection with Principal Component Analysis and Random Forest

**A T Devi [1], Dr. Levina Tukaram[2], Rashmi Purad[3]**

[1]Professor, Department of Information Science and Engineering, KNSIT Bangalore-64, Karnataka, India
Email: *devi.vijay1905[at]gmail.com*

[2]Department of Computer Science and Engineering, KNSIT Bangalore-64, Karnataka, India
Email: *levinad[at]knsit.com*

[3] Professor, Department of Information Science and Engineering, KNSIT Bangalore-64, Karnataka, India
Email: *rashmipurad0811[at]gmail.com*

**Abstract:** *With the rising number of cybersecurity attacks, intrusion detection systems (IDS) play a vital role in spotting unauthorized access and lowering security risks. This study presents a framework that integrates principal component analysis (PCA) for dimensionality reduction with random forest (RF) for classification tasks. The proposed model is evaluated using the NSL-KDD dataset, showing notable gains in detection accuracy, reduced error rates, and faster processing compared to common models like SVM and Naïve Bayes. Results show that the proposed approach achieves 96.78% accuracy while keeping the error rate at 0.21%.*

**Keywords:** Intrusion detection system (IDS), Principal Component Analysis (PCA), Random Forest (RF), Machine Learning, Cybersecurity, NSL-KDD dataset

## 1. Introduction

An IDS is a protective system built to observe network activity and detect any unusual or malicious behavior. It acts as a safeguard by alerting administrators to possible breaches. IDS can be broadly divided into two categories: network-based intrusion detection systems (NIDS) and host-based intrusion detection systems (HIDS). NIDS inspects network packets to find suspicious patterns, whereas HIDS focuses on monitoring the activities of a specific device or host. Detection methods include signature-based detection, which uses known attack patterns, and anomaly-based detection, which identifies deviations from normal behavior using machine learning approaches.
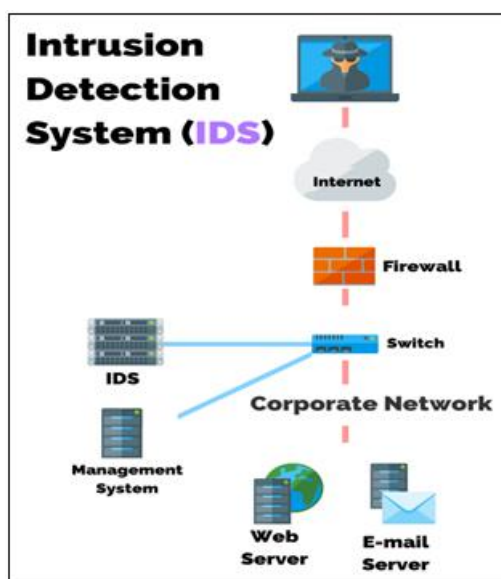


**Figure 1:** System Diagram for IDS

### 1.1 Motivation for the Project

The increasing sophistication of cyber threats presents a growing challenge for organizations seeking to protect their digital assets. As businesses gather and process more data, the potential vulnerabilities also increase, making it harder for traditional security measures to keep pace with emerging risks. This creates a critical need for advanced systems that can detect malicious activity in real-time, with high precision and minimal human oversight.

The aim of this paper is to investigate the effectiveness of various techniques in identifying cyber threats, while also developing a system that not only improves detection performance but also assists in real-time decision-making for cybersecurity teams.

### 1.2 Limitations in Current IDS Systems

Despite advancements in machine learning for IDS, several limitations persist:

- High Computational Overhead:
- Complex models like deep learning require more time and hardware resources.
- False Alarms: Even optimized IDS models may occasionally flag normal behavior as malicious, requiring manual investigation.
- Difficulty in Feature Interpretation: PCA transforms features into principal components, which are harder to map back to real network parameters.
- Data Imbalance: IDS datasets often have more normal records than attack ones, which affects classifier learning.
- Generalization to Real-Time Scenarios: Some models trained on offline data perform poorly in real-time deployments due to unseen patterns.

## 1.2 Principal Component Analysis (PCA)

PCA is a widely used method for reducing data dimensions in machine learning, converting high-dimensional datasets into fewer dimensions' while preserving important trends and variations. In IDS, the main aim of PCA is to remove irrelevant features, speeding up computations while keeping detection accuracy high. Selecting the most valuable features improves IDS performance and reduces overfitting, leading to better adaptability to new data.

The PCA process involves:
1) Start with the dataset: The dataset contains d dimensions, representing the features.
2) Calculate the mean vector: Compute the average vector for each of the d dimensions.
3) Covariance matrix: Building the covariance matrix to find relationships among features
4) Eigen decomposition: Determine the eigenvectors ($e_1$, $e_2$, $e_3$... $e^d$) and eigenvalues ($\lambda_1$, $\lambda_2$, $\lambda_3$... $\lambda^d$) of the covariance matrix.
5) Sort eigenvalues: Rank the eigenvalues in descending order and select the top n eigenvectors corresponding to the largest eigenvalues, forming a matrix of reduced dimensions ($d_n = m$).
6) New sample space: Apply this new eigenvector matrix to transform the original data into a lower-dimensional space.

The resulting "principal components" hold the most important data patterns while lowering dimensionality.

## 1.3 Random Forest (RF) Classifier

The Random Forest method builds multiple decision trees and merges their predictions to boost accuracy and stability. Unlike many traditional classifiers, RF reduces overfitting and performs well on large datasets with noisy or irrelevant data.

For IDS, RF effectively distinguishes between normal and malicious network activity, making it a strong choice for cybersecurity applications.
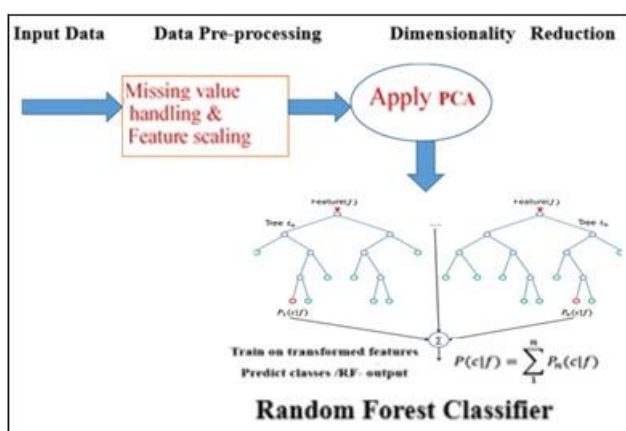


**Figure 1:** working of Random Forest with PCA

This Work Proposes Combining PCA for Feature Selection with RF For Classification to Achieve High Detection Accuracy with Less Computation. Similar Approaches Have Been used for Iot Device Identification, where Deep learning models categorize network traffic into binary and multi-class formats with strong accuracy [7].

## 2. Literature Review

The authors in [9] presented an IDS solution using SVM and Naïve Bayes, finding that SVM provided better results. Experiments were conducted using the KDD dataset to compare detection and false alarm rates. Another study performed three experiments with feature selection, Naïve Bayes, adaptive boosting, and partial decision trees, showing their usefulness in detecting intrusions [9].

Further research in [10] showed that combining artificial neural networks with feature selection can outperform SVM. Using the NSL-KDD dataset, these approaches delivered strong outcomes [10].

A detailed review of machine learning-based IDS techniques compared multiple algorithms by accuracy and false alarm rates, highlighting their ability to improve detection performance. A logistic regression and belief propagation-based approach was introduced in [12], improving detection speed over earlier methods. Work in [13] applied advanced machine learning for feature extraction, improving dataset quality and IDS performance.

Research in [14] identified key contributions by Chowdhury et al. and models from Shenfield, Ayesh, and Day. A comparative analysis in [15] found that extreme learning machine (ELM) could outperform SVM and random forest in certain classification scenarios.

A fuzzy rule-based feature enhancement technique in [16] boosted dataset quality, resulting in improved detection accuracy.

## 3. Problem Domain

Handling large volumes of network traffic data with many features is challenging due to time and resource demands. Redundant and irrelevant data slow processing and reduce classification accuracy. A method is needed that can reduce the number of features while preserving intrusion-related patterns.

## 4. Proposed Solution

The proposed framework targets real-time intrusion detection by combining PCA and RF to improve speed and accuracy. Traditional models often face high computation costs and reduced accuracy due to redundant features; this approach processes live network traffic efficiently.
1) Dataset: NSL-KDD dataset is used, similar to the implementation in related work. PCA reduces redundant features, boosting detection speed and reducing overfitting.
2) Classification: RF classifier identifies intrusions, implemented using Spark ML for distributed computing.
3) Deployment: Django is used to build a web dashboard for real-time intrusion alerts and visualization.

4) Performance Evaluation: Accuracy, precision, recall, and F1-score are calculated.

**Table 1:** Performance Evaluation

| Approach Used | Processing Duration (minutes) | Accuracy (%) | Error (%) |
|---|---|---|---|
| Support Vector Machine | 4.57 | 84.34 | 2.67 |
| Naïve Bayes Classifier | 9.12 | 80.85 | 3.49 |
| Decision Tree Model | 12.36 | 89.91 | 0.78 |
| PCA + Random Forest | 3.42 | 96.78 | 0.21 |

# 5. Outcome and Next Steps

Findings show that PCA significantly reduces the number of features, resulting in faster model training and prediction. RF achieves high accuracy, confirming its suitability for intrusion detection. The trade-off between dimensionality reduction and accuracy is studied, identifying an optimal point where computational cost is minimized without losing detection quality.
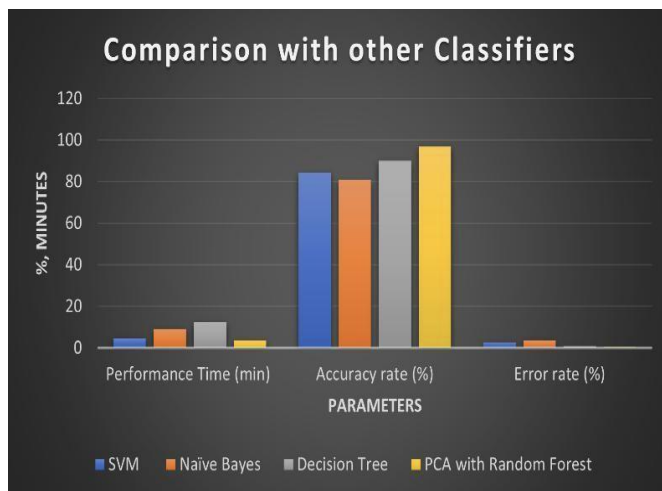


**Figure 2:** Comparisons of the Classifiers

**Tools and Technologies Used**

The successful implementation of an Cyber Threat Detection demands a combination of robust tools and technologies. In this project, Python was selected as the primary programming language due to its simplicity, versatility, and extensive support for machine learning libraries.

The major libraries and frameworks utilized include:
- Scikit-learn: A comprehensive library for machine learning providing essential functions for preprocessing, dimensionality reduction (PCA), and Classification (Random Forest).
- Pandas and NumPy: These libraries were essential for efficient data handling, manipulation, and numerical computations.
- Matplotlib and Seaborn: Used for data visualization, plotting feature distributions, correlation heatmaps, PCA scree plots, and model performance metrics like confusion matrices and ROC curves.
- Flask: A lightweight web framework to create a dashboard interface for real-time intrusion alerts visualization.

- Dataset: The NSL-KDD dataset was chosen due to its widespread acceptance in IDS research and its improvements over the original KDD'99 dataset, addressing redundancy and imbalance issues.

# 6. Conclusion

This work presents an IDS framework combining PCA and RF to enhance speed and accuracy. It handles high-dimensional data well and outperforms SVM, Naïve Bayes, and Decision Tree models. On the NSL-KDD dataset, it achieved 96.78% accuracy, 0.21% error rate, and a processing time of 3.24 minutes. These results lead to better detection rates and fewer false positives.

With the rapid growth of internet-connected systems, security threats are increasing. This approach offers a strong and efficient intrusion detection method. Future work can look at combining PCA with deep learning techniques to further improve scalability and accuracy.

# References

[1] A. Bonnaccorsi, "On the Relationship between Firm Size and Export Intensity," Journal of International Business Studies, XXIII (4), pp. 605-635, 1992. (journal style)

[2] R. Caves, Multinational Enterprise and Economic Analysis, Cambridge University Press, Cambridge, 1982. (book style)

[3] M. Clerc, "The Swarm and the Queen: Towards a Deterministic and Adaptive Particle Swarm Optimization," In Proceedings of the IEEE Congress on Evolutionary Computation (CEC), pp. 1951-1957, 1999. (conference style)

[4] H.H. Crokell, "Specialization and International Competitiveness," in Managing the Multinational Subsidiary, H. Etemad and L. S, Sulude (eds.), Croom-Helm, London, 1986. (book chapter style)

[5] K. Deb, S. Agrawal, A. Pratab, T. Meyarivan, "A Fast Elitist Non-dominated Sorting Genetic Algorithms for Multiobjective Optimization: NSGA II," KanGAL report 200001, Indian Institute of Technology, Kanpur, India, 2000. (technical report style)

[6] J. Geralds, "Sega Ends Production of Dreamcast," vnunet.com, para. 2, Jan. 31, 2001. [Online]. Available: http://nl1.vnunet.com/news/1116995. [Accessed: Sept. 12, 2004]. (General Internet site)

[7] Rohit Kumar Singh Gautam, Er. Amit Doegar; 2018 8th International Conference on Cloud Computing, Data Science & Engineering (Confluence) "An Ensemble Approach for Intrusion Detection System Using Machine Learning Algorithms."

[8] Kazi Abu Taher, Billal Mohammed Yasin Jisan, Md. Mahbubur Rahma, 2019 International Conference on Robotics, Electrical and Signal Processing Techniques (ICREST) "Network Intrusion Detection using Supervised Machine Learning Technique with Feature Selection."

[9] L. Haripriya, M.A. Jabbar, 2018 Second International Conference on Electronics, Communication and Aerospace Technology (ICECA)" Role of Machine Learning in Intrusion Detection System: Review"

[10] Nimmy Krishnan, A. Salim, 2018 International CET Conference on Control, Communication, and Computing (IC4) "Machine Learning-Based Intrusion Detection for Virtualized Infrastructures"

[11] Mohammed Ishaque, Ladislav Hudec, 2019 2nd International Conference on Computer Applications & Information Security (ICCAIS) "Feature extraction using Deep Learning for Intrusion Detection System."

[12] Aditya Phadke, Mohit Kulkarni, Pranav Bhawalkar, Rashmi Bhattad, 2019 3rd International Conference on Computing Methodologies and Communication (ICCMC)"A Review of Machine Learning Methodologies for Network Intrusion Detection."

[13] Iftikhar Ahmad, Mohammad Basheri, Muhammad Javed Iqbal, Aneel Rahim, IEEE Access (Volume: 6) Page(s): 33789 – 33795 "Performance Comparison of Support Vector Machine, Random Forest, and Extreme Learning Machine for Intrusion Detection."

[14] B. Riyaz, S. Ganapathy, 2018 International Conference on Recent Trends in Advanced Computing (ICRTAC)" An Intelligent Fuzzy Rule-based Feature Selection for Effective Intrusion Detection."

## Author Profile

**A T Devi** received the B.E and M.Tech. degrees in Electronics & Communication Engineering from Anna University Chennai and MS University Tirunelveli 2006 and 2009, respectively. Pursuing Ph D from VTU Belagavi Karnataka, working with KNS Institute of Technology with 07 years of Teaching experience, area of Interest is Automata Theory, Computer Graphics, Machine Learning &Deep Learning.

**Dr. Levina T** is an Associate Professor in the Department of Computer Science and Engineering at K.N.S. Institute of Technology, Bengaluru, India. She holds a Ph.D. in Computer Science and an M. Tech in Computer Networking Engineering, with over 20 years of academic and research experience. Her areas of interest include Data Science, Artificial Intelligence, Cloud Computing, Internet of Things (IoT), and Network Security.

**Rashmi Purad** is an Assistant Professor in the Department of Information Science and Engineering at K.N.S. Institute of Technology, Bengaluru, India. Received the B.E in Information Science & Engineering and M.Tech. Computer Science & Engineering from VTU Belgavi. Karnataka,2010 and 2013, respectively. with 07 years of Teaching experience, area of Interest Data Structure, Clouud Computing & IoT.