

Big Data and Cybersecurity: Using Analytics to Predict and Prevent Cyber Threats

Dinesh Kumar Budagam¹, Naresh Kumar Miryala²

¹Visa Inc, Foster City, CA, USA

Email: [dbudagam\[at\]gmail.com](mailto:dbudagam[at]gmail.com)

²Meta Platforms Inc. CA, USA

Abstract: *In the period of digital transformation, Cyber threats have become more difficult, which simplifies data of easy and the defenses are more necessary in active state. This study proposes a novel ensemble-based analytical framework using Random Forest (RF) combined with Bayesian Inference for predicting and preventing cyber threats. The approach leverages big data analytics to process massive volumes of real-time network logs, system behavior data and user activity patterns to identify anomalies and potential breaches. Through detailed analysis, the model demonstrates high detection accuracy and low false positive rates by effectively capturing non-linear threat patterns and incorporating probabilistic reasoning. Interpretation of the results shows that integrating with probabilistic models enhances prediction reliability across varied threat scenarios, including malware propagation and insider attacks. The model's robustness and effectiveness are validated by the comparative results on standard benchmark datasets, which demonstrate an accuracy of 96.7% and specificity of 97%.*

Keywords: Cybersecurity, Big Data, Bayesian Inference, Random Forest.

1. Introduction

In recent years, digital connectivity growth has driven over 4.57 billion users, generating vast data volumes [1]. People and devices linked to internet generate around 2.5 quintillion bytes of data every second [2]. This massive amount of data has led to rapid development and adoption of big data technology across many industries, with cybersecurity being one of the most important uses. Big data technologies have become a crucial tool for controlling and reducing cyber threats by handle and analyze enormous volumes of data in real-time [3–4]. One of the primary uses of big data in cybersecurity is Security Information and Event Management (SIEM) systems. To identify and address any security threats, SIEM systems collect, store, and examine logs, messages and security events. Traditionally, SIEM systems relied on relational databases, but with the increase in data volume, non-relational databases, such as NoSQL, have become the standard for handling these large data sets efficiently [5].

To improve their capacity to recognize and stop cyberthreats, SIEM systems are anticipated to use predictive analytics models and techniques in their upcoming evolution. Prominent cybersecurity firms like Red Lambda and Palantir are already utilizing big data technology to create comprehensive user profiles and identify anomalous activity through the analysis of data from several sources [6]. Corporate communications, personnel systems, Access Control Systems (ACS), Customer Relationship Management (CRM) systems, Web and intranet data and Internet of Things (IoT) and Industrial IoT (IIoT) systems are some examples of these sources [7]. Additional sources of useful real-time data that are utilized to identify possible risks include social networks, external news feeds, and other aggregators. Through the integration and analysis of these disparate information sources, organizations uncover security threats and acquire a greater understanding of user behavior.

The potential of big data technologies to do real-time online analysis of streaming and packet data highlights its importance in cybersecurity. This makes it possible for businesses to identify important cybersecurity incidents as they occur and take prompt action in response [8]. Furthermore, use of big data in isolation and processing of complex security incidents allows for the discovery of hidden patterns and potential threats. Cybersecurity teams have a distinct advantage in the combating against cyberthreats because real-time analysis of vast amounts of data.

The capacity of big data to offer proactive security measures is among its most important benefits in cybersecurity. Big Data technologies enable enterprises to anticipate and stop security incidents before they have a major impact on vital infrastructure or cause harm by continuously monitoring data and analyzing trends [9]. The use of big data technology has the ability to detect new threats and weaknesses before they materialize. These tools forecast the probability of a cyberattack by analyzing data trends and abnormalities, enabling organizations to take precautionary steps to reduce the risk [10]. Because it shifts the focus from reacting to threats after they occur to preventing harm, this predictive capability is groundbreaking in field of cybersecurity. Threat analysis, response deployment, data gathering and data storage are the four primary parts of a big data-powered cybersecurity analytics platform.

These are described in more detail below:

- **Collection of Data** is the process of obtaining unprocessed, real-time data from a variety of sources, including network traffic, user behavior logs, endpoints, cloud services and IoT devices.
- Large amounts of organized and unstructured data are handled effectively and securely by data storage, frequently with the use of scalable Big Data systems.
- **Threat Analysis** is the process of looking at the gathered data with Machine Learning (ML), advanced analytics

Volume 14 Issue 5, May 2025

Fully Refereed | Open Access | Double Blind Peer Reviewed Journal

www.ijsr.net

and anomaly detection.

- **Response Deployment** helps security teams prevent or mitigate cyber incidents by enabling prompt action based on insights.

2. Different types of Cybersecurity Threats

2.1 Malware

Malware, short for malicious software, is any software intentionally designed to harm, exploit, or otherwise compromise devices, networks, or data. Cybercriminals use this malware to steal sensitive information, disrupt operations, gain unauthorized access, or in few cases demand ransoms from individuals or organizations. Large amount of data is exfiltrated in recent malware attacks. Malware is developed as harmful software that invades or corrupts your computer network. The goal of malware is to cause havoc and steal information or resources for monetary gain or sheer sabotage intent. Malicious software designed to damage equipment and snip data from people or companies, including viruses, worms and spyware. These apps often cause serious problems for productions and people.

Few of the examples of Malware are Virus, worms, Spyware and Adware.

2.2 Ransomware

Ransomware is a damaging malware that converts important data on a computer in which the attacker difficulties payment. In general, Ransomware is a malware that employs encryption to hold a victim's information at ransom. A user or organization's critical data is encrypted so that they cannot access files, databases, or applications. A ransom is then demanded to provide access. Ransomware is often designed to spread across a network and target database and file servers and can thus quickly paralyze an entire organization. Few of the examples of Ransomware are readily available malware kits to generate new malicious code, driven by significant financial incentives. This malware kits are used by cybercriminals. Use of known good generic interpreters to create cross-platform ransomware such as Ransom32 which incorporates Node.js with a JavaScript payload is another well-known example of Ransomware. Use of new techniques, such as encrypting the complete disk instead of selected files.

2.3 Denial-of-Service (DoS) Attacks

This kind of attack excesses a network or website with traffic flow, reducing its requirements and ensuing a delay or crash. A denial-of-service (DoS) attack is a type of cyber-attack in which a malicious actor aims to render a computer or other device unavailable to its intended users by interrupting the device's normal functioning. DoS attacks typically function by overwhelming or flooding a targeted machine with requests until normal traffic is unable to be processed, resulting in denial-of-service to addition users. A DoS attack is characterized by using a single computer to launch the attack. The primary goal of a DoS attack is to flood the capacity of a targeted machine, resulting in denial-of-service to additional requests. We can group multiple attack vectors of DoS attacks by identifying their similarities

2.4 Man in the Middle Attacks (MitM)

When two individuals or systems connect, this risk involves illegally listening to or changing information by exchange without the other person's knowledge. A man-in-the-middle (MITM) attack is a cyberattack in which a hacker steals sensitive information by eavesdropping on communications between two online targets such as a user and a web application.

After discreetly placing themselves in the middle of two-party communications, MITM attackers intercept sensitive data such as credit card numbers, account information and login credentials. Attackers use the sensitive information such as credit card numbers, account information and login credentials to commit cybercrimes such as making unauthorized purchases, hijacking financial accounts and identity theft. MITM attacker might also intercept on private communications between two people besides exchanges between a user and an application. This is achieved by attacker by diverting and relaying messages between the two people which also involves altering or replacing messages to control the conversation

2.5 Insider Threats

An insider threat is a security risk that emerges from within the premises of an organization. Such threats quite often come from a current or former employee, contractor, vendor or partner who has appropriate user credentials and uses them inappropriately to harm the organization's networks, systems and data. Such insider threats, or risks, are the result of the abuse by current or former employees of their authorized access in the system, for the purpose of their own gain or to retaliate against others, by causing damage (e.g., stealing sensitive information, leaking information, or denying the access to the service to other users). An insider threat can be intentional or unintentional. No matter the intent, the result is compromised confidentiality, availability, and/or integrity of enterprise systems and data.

Insider threats often neglected are major source of data breaches. Traditionally existing cybersecurity strategies, policies, procedures and systems often focus on external threats, leaving the organization vulnerable to attacks from within. Because the insider already has valid authorization to data and systems, it's difficult for security professionals and applications to distinguish between normal and harmful activity.

Malicious insiders have a distinct advantage over other categories of malicious attackers because of their familiarity with enterprise systems, processes, procedures, policies and users. They are keenly aware of system versions and the vulnerabilities therein. Organizations must therefore tackle insider threats with at least as much rigor as they do external threats.

2.6 Zero-day Exploits

These kinds of attacks prey on undiscovered defects or problems in software that creators are unaware of. These flaws allow attackers to hurt or steal data because there is

currently no update or remedy. A zero-day exploit is a cyberattack vector that takes advantage of an unknown or unaddressed security flaw in computer software, hardware or firmware. "Zero day" refers to the fact that the software or device vendor has zero days to fix the flaw because malicious actors can already use it to access vulnerable systems. The unknown or unaddressed vulnerability is referred to as a *zero-day vulnerability* or *zero-day threat*. A *zero-day attack* is when a malicious actor uses a zero-day exploit to plant malware, steal data or otherwise cause damage to users, organizations or systems.

3. Cybersecurity Defenses

Cyber defense is the practice of protecting networks, devices and data from unauthorized access or criminal use. It encompasses a range of technologies and practices, such as vulnerability management, network security, endpoint security, data security and identity and access management (IAM). Managed cyber defense is important because it helps protect your organization from cyber threats, for example, data breaches, malware infections, ransomware attacks, denial-of-service attacks and social engineering like phishing attacks. We will discuss few of cybersecurity defenses as in below section.

3.1 Threat Detection and Prevention

Threat detection is considered as an essential and critical element of cybersecurity today. As cyberattacks grow more sophisticated and devastating, threat detection systems help quickly identify attacks in progress and block vulnerabilities or mitigate them before they result in a security breach leading lot of financial losses. Sophisticated analytics and Artificial Intelligence (AI) analyses enormous data sets in real-time to search for odd patterns or behaviours that could indicate a cyberattack. Companies prevent or minimize damage by acting quickly. There are different models and approaches for building a threat detection and response tool. This mechanism involves incorporating Zero Trust, where all users need frequent authorization. Irrespective of the model and threat detection method, threat detection and response should and must meet the needs of your business. By utilizing the effective threat detection and response, we can protect the applications and sensitive data against advanced attacks.

3.2 Patch Management

Patch management is an important process in cybersecurity which involves updating and applying software patches to fix vulnerabilities, enhance functionality, and ensure compliance. Software vendors release patches to address security issues such as security flaws, bugs, and performance issues. This patch provides critical updates to maintain a secure and efficient IT environment. By implementing the effective patch management minimizes the risk of cyberattacks by ensuring that systems are always up to date with the latest protections. By patching known security flaws in a timely manner, regular software and system updates lower the likelihood that attackers would take benefit of them regular threats and new zero-day exploits that target unpatched software faults are both defended against by this

proactive patch management.

3.3 Network Security Measures

Network security measures can be referred as security controls which we add to your networks to protect confidentiality, integrity, and availability. These controls continue to evolve and emerge, but there is a lot of fundamental knowledge that readily available. It takes consistent efforts to keep attackers out of network. Firewalls, proxies, and gateways serves the purpose. Secure Virtual Private Networks (VPNs), firewalls and Intrusion Detection and Prevention Systems (IDS/IPS) cooperate to monitor and manage network traffic in order to stop unwanted access and identify questionable activity. Few of the network security measures are access Management, security monitoring, Firewalls, Anti-Malware Software, Application Security, Data Risk Management, Email Security, Security Information and Event Management (SIEM), Backup and Disaster Recovery (BDR), Endpoint Security, Virtual Private Network (VPN), Web Security

3.4 Data Backup and Recovery

Data backup and recovery is considered as the last line of defense against data loss, corruption, and cyber threats. This critical cybersecurity measure ensures that organizations can recover their critical information even when there is failure in other security layers. Data backup is the process of creating and storing copies of important data in secure locations. Businesses quickly reestablish significant data to regular, safe backups in case ransomware attacks cause it to be lost, corrupted or encoded. This indicate that business continuousness by avoiding permanent data loss and decreasing downtime.

3.5 Multi-Factor Authentication (MFA)

MFA and other identity verification layers prevent unwanted access even if credentials are compromised. When additional identification is required, such a security token, a code texted to a phone, or a fingerprint and security is enhanced. Multi-factor authentication (MFA) is a multi-layered security access management process which grants users access to a network, system, or application only after confirming their identity with more than one credential or authentication factor.

This is typically achieved by configuring a combination of a username and password, along with an additional factor such as verification code or one-time password (OTP) sent via text or email, a security token from an authenticator app, or a biometric Identifier). The security posture of organizations can be notably enhanced by requiring more than one factor of authentication. "This is because even if the initial authentication was compromised or disabled, access is denied unless the user also has possession or control of the second factor of authentication.

3.6 Incident Response Planning

By establishing explicit incident response protocols, companies quickly identify and mitigate the effects of

cyberattacks. This timely strategy minimizes operating downtime, limits damage, and encourages a quicker recovery.

4. Role of Big Data in Cybersecurity Defense

Before we dive about how big data is used in Cybersecurity, we need to understand what a bigdata is. Big data is referred to large volume of structured and unstructured of data produced by the multiple sources. The sources for these examples include (but are not restricted to) social media, e-commerce, mobile sensors, etc. These types of data are characterized by its volume, variety, velocity, and veracity. The capability of streaming, processing and analyzing this data in real-time has tremendous impact across many domains, including cyber security. Big data enables companies to detect, assess and respond to threats faster and more effectively. In terms of cyber security, big data involves the gathering and inspection of large amounts of data sets to identify patterns, trends and synonymous anomalies as to potential security risks. Examples of this datasets is network logs, user activity data, and external threat intelligence. In cybersecurity domain, one of the primary benefits of big data is its ability to enhance threat detection. Big data enables to analyze the large datasets from captured from different sources in real time. This mechanism benefits organizations to identify unusual patterns and behaviors that may indicate a cyber-attack

4.1 Big Data Enhances Threat Visibility and Situational Awareness

Cybersecurity teams may collect and analyse network traffic, log files, user behaviour, and application usage across distributed platforms to big data technologies. weird trends like traffic spikes, weird logins, or unauthorized access attempts are easily detected using advanced analytics.

4.2 Integration of Real-Time Analytics in Security Operations Centers (SOCs)

SOCs are in charge of continuously observing threats and taking appropriate action. SOCs use Big Data and real-time analytics to automatically gather and process data from cloud environments, firewalls, IDS and endpoint devices. Analysts take immediate action in response to threats uncover to real-time dashboards and alerts, which greatly speeds up response times.

4.3 Threat Intelligence Aggregation and Correlation

To improve detection capabilities, big data platforms gather, compile and correlate threat intelligence feeds from many internal and external sources.

5. ML Techniques used for predicting Cybersecurity

Machine learning is critical in enhancing cybersecurity by enabling organizations to detect and respond to threats more effectively and efficiently. Below are few of the examples of ML techniques used for cybersecurity such as

Supervised learning algorithms, Unsupervised learning algorithms and Hybrid approaches.

Supervised learning algorithms such as support vector machines (SVMs) and neural networks for threat detection

Unsupervised learning algorithms like clustering and principal component analysis (PCA) for anomaly detection

Hybrid approaches, combining multiple machine learning techniques, are often employed to enhance the accuracy and robustness of cybersecurity systems in detecting and mitigating evolving threats.

5.1 Decision Tree

Decision trees are used in cybersecurity to build rule-based models that categorize user or network activity into categories that are either benign or malevolent. Decision trees are hierarchical structures that recursively split data into smaller subsets based on the most significant attributes. In cybersecurity, decision trees are used for intrusion detection, malware classification, and risk assessment tasks. Major advantage of Decision tree is transparency and interpretability which makes them suitable for generating rules to detect security threats.

Advantages

It aids security analysts in comprehending the logic underlying warnings by offering unambiguous, interpretable guidelines for identifying threats.

Disadvantages

Overfitting to noise in cyber data is easy, particularly when the dataset contains a large number of unrelated attributes or uncommon attack types.

5.2 Support Vector Machine

In cybersecurity, SVMs create hyperplanes in high-dimensional feature space to categorize network traffic, emails, or user activity as either normal or suspect. SVM is a supervised learning algorithm that separates data points by finding the hyperplane that maximally separates different classes in the feature space.

It's possible to apply SVM in cybersecurity for tasks such as malware detection, intrusion detection, and spam filtering.

Advantages

Proficient at identifying minute irregularities or cyberattacks, including zero-day exploits, in high-dimensional spaces.

Disadvantages

High training and prediction timeframes make it unscalable for big, real-time cybersecurity datasets.

5.3 Random Forest

A group of decision trees cooperating to increase threat detection accuracy, like spotting malware or intrusion patterns.

Random forest is a commonly used machine learning algorithm that combines the output of multiple decision trees to reach a single result. Its ease of use and flexibility have fueled its adoption, as it handles both classification and regression problems. **Random Forest** can be defined as a **supervised machine learning algorithm** that belongs to the ensemble learning family. It works by creating a "**forest**" of **decision trees** during training and outputs the **mode (classification)** or **mean (regression)** of the individual trees' predictions. Important notes to be considered are a **Decision Tree** which is a model that splits data based on feature values to make decisions, Ensemble which can be referred as Random Forest combines multiple trees to improve accuracy and reduce overfitting.

Advantages

Its ensemble nature makes it extremely accurate at identifying intricate cybersecurity risks.

Disadvantages

It is more computationally intensive and challenging to comprehend in real-time threat monitoring systems.

5.4 K-Nearest Neighbor (KNN)

KNN classifies a new instance based on the behaviour of its neighbours and compares it with previous instances to identify dangers.

Advantages

The availability of previous labelled data makes it simple to implement and adapt to new kinds of cyberattacks.

Disadvantages

Computationally costly and slow for detecting threats in real time, particularly when dealing with extensive network data.

5.5 Naïve bayes

Naive Bayes is used in cybersecurity to identify suspect user activity based on probabilities, filter spam, and detect URLs that are fraudulent.

Advantages

For bulk text-based threats, such as spam email classification, it is incredibly quick and efficient.

Disadvantages

In complex cyberattacks where aspects are interdependent, the assumption of feature independence frequently proves to be incorrect.

5.6 Logistics Regression

Using input features such as IP reputation or access patterns, logistic regression is utilized to forecast the probability of a cybersecurity occurrence.

Advantages

Login anomaly detection and other basic binary classification tasks benefit greatly from its simplicity, interpretability, and speed.

Disadvantages

It limited in its ability to detect non-linear patterns present in covert activities or complex cyberattacks.

5.7 Isolation Forest

An anomaly detection method called Isolation Forest is used in cybersecurity to find uncommon and uncommon network behaviours or zero-day threats.

Advantages

Effective for identifying outliers in unlabelled data; efficient for high-dimensional data.

Disadvantages

In extremely thick or overlapping data distributions, this is less effective since it incorrectly labels common but uncommon actions as dangers.

5.8 Extreme Gradient Boosting (XGBoost)

In cybersecurity, XGBoost is a scalable and effective gradient boosting implementation used for fraud detection, malware detection, and network intrusion classification. Extreme Gradient Boosting, also known as XGBoost, is a scalable and optimized algorithm in computer science that improves the speed and prediction performance of Gradient Boosting Machines (GBM). It achieves this by using a new tree learning algorithm and leveraging parallel and distributed computing to accelerate model discovery. XGBoost is recognized for its high prediction success and has found applications in various domains.

Advantages

High efficiency and accuracy, good handling of imbalanced datasets and missing values, and integrated regularization to reduce overfitting.

Disadvantages

It requires meticulous hyperparameter adjustment; it is rather resource-intensive and difficult.

5.9 Deep Belief Networks (DBN)

In cybersecurity, DBNs are layered generative models that are used to identify complex risks like behavioural anomalies or Advanced Persistent Threats (APTs).

Advantages

Effectively manages intricate and non-linear patterns; able to extract deep, abstract features from unprocessed data.

Disadvantages

Training takes a lot of time and requires a lot of computing power; the results are difficult to interpret due to their black-box nature.

5.10 Fuzzy Logic

Fuzzy logic systems detect fuzzy attack fingerprints or soft thresholds in intrusion detection by handling ambiguous or imprecise input data.

Advantages

It handles ambiguity and partial truth well; it simulates human decision-making in ambiguous situations.

Disadvantages

With high-dimensional raw data, rule development is less effective and are subjective and difficult.

5.11 Genetic Algorithm (GA)

GA are optimization methods for feature selection, changing attack patterns, and fine-tuning rule sets in IDS.

Advantages

In adaptive cybersecurity systems, this ability to search large search regions for the most effective detection settings is essential.

Disadvantages

It is slow to converge, computationally costly and findings change from run to run.

5.12 Bayesian Networks

Bayesian networks are used to evaluate the likelihood of cyber events such as insider threats or system breaches by representing relationships among variables. Bayesian Networks are one of the most widely used types of probabilistic graphical models. Offering effective solutions for decision making and inference under uncertainty, these networks play a critical role in artificial intelligence, machine learning and data analysis. Bayesian networks represent dependency relationships between variables in complex problems and are capable of probabilistic inference.

Bayesian Networks are directed acyclic graphs that show the conditional dependencies between a set of random variables. These networks describe the dependence of each variable on the other variables using Bayes' Theorem. Bayes' Theorem is a probability theorem for updating the probability of an event with observations of other events.

A **Bayesian Network (BN)**, also known as a **Bayes Net** or **Belief Network**, is a **probabilistic graphical model** that represents a set of variables and their conditional dependencies via a **directed acyclic graph (DAG)**.

- **Key Concepts:**
- **Nodes:** Represent random variables.
- **Edges:** Represent conditional dependencies.
- **Conditional Probability Tables (CPTs):** Each node has a CPT that quantifies the effect of the parent nodes on the node.

Bayesian Networks apply **Bayes' Theorem** to update the probability estimate for a hypothesis as more evidence is available.

Bayesian Networks are used in cybersecurity to model uncertainty and infer potential threats from incomplete or noisy data. Their ability to reason probabilistically makes them ideal for: **Intrusion Detection Systems (IDS), Threat Intelligence and Risk Assessment, Malware Detection and Decision Support for Security Analysts**

Advantages

It combines probabilistic reasoning with domain knowledge.

Disadvantages

High-dimensional data makes it difficult to construct and train, and it is dependent on assumptions.

6. Challenges in implementing ML techniques

- High costs and infrastructure complexity
- Difficult to interpret in real-time threat monitoring systems
- Not scalable for large, real-time cybersecurity datasets

The use of traditional cybersecurity techniques is summarized to highlight their limitations in handling modern threats. The main contributions of this paper:

- To integrate big data analytics into cybersecurity for real time detection and prediction of cyber threats using high-volume, high-dimensional data.
- To apply Kernel PCA as a pre-processing method to reduce dimensionality, remove noise and retain significant features for accurate classification.
- To develop a RF combined with Bayesian Inference that enhances classification performance in detecting cybersecurity threats.
- To improve intrusion detection accuracy and minimize false positives through optimized feature selection and intelligent classification.

6.1 Threat Detection using Kernel PCA

PCA aims to preserve as much variance as feasible while reducing a dataset's dimensionality. This is achieved by highlighting and identifying the most effective PCs, which are arranged from highest to lowest in terms of the amount of variance that capture.

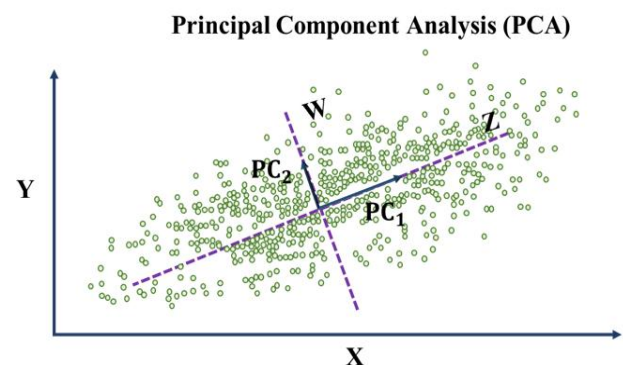


Figure 1: Graphical Representation of PCA

PCA is a powerful statistical technique used for dimensionality reduction while retaining the essential structure of the data. A common way to visualize the outcome of PCA is through a 2D graphical representation using the first two PCs are PC1 and PC2 plotted on the X and Y axes, respectively.

PC1, shown on the X-axis, represents the direction of maximum variance in the data. It captures the greatest amount of variability, meaning it explains largest portion of total dataset variance. PC2, displayed on Y-axis, is

orthogonal to PC1 and represents the second-most significant direction of variance, capturing additional, uncorrelated patterns in the data. The graphical representation allows us to observe patterns, clusters, or separations among data points that are not have been apparent in the higher-dimensional space.

Kernel techniques are used in kernel PCA, an improved PCA methodology that permits nonlinear dimensionality

reduction. Kernel PCA implicitly projects input data into a higher-dimensional feature space where complicated, nonlinear patterns become linearly separable, rather than altering it directly. The most informative components are subsequently extracted by applying PCA inside this modified space. The original dataset's complex nonlinear patterns are captured using this method.

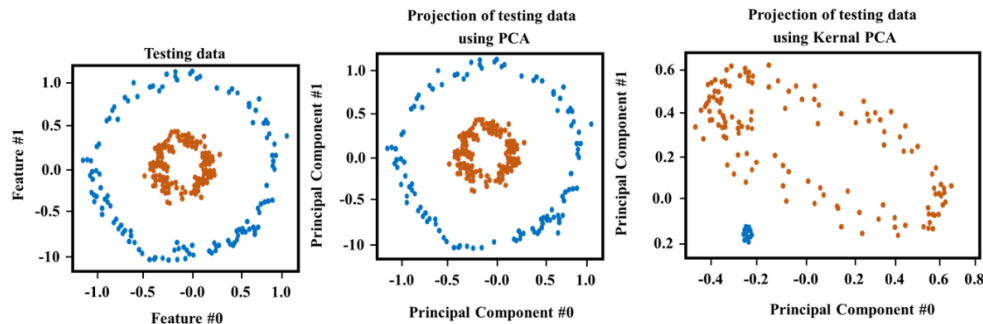


Figure 2: Graphical representation of Kernel PCA

Advantages

- Kernel PCA extract more informative features from non-linear data, enhancing the performance of downstream tasks like classification, anomaly detection or clustering.
- It supports various kernels giving users the flexibility to tailor the algorithm to specific data characteristics or domains.
- In cybersecurity, Kernel PCA is particularly effective for detecting anomalies in high-dimensional network traffic

4.2 Classification using RF combined with Bayesian Inference approach

RF combined with Bayesian Inference in cybersecurity is an advanced hybrid approach that leverages the high predictive power of Random Forest with the probabilistic reasoning and uncertainty estimation of Bayesian Inference. The RF model is centered around the optimal selection of hyper-parameters, which significantly influence its predictive accuracy. To enhance this process, a Bayesian Optimization-based RF model is employed to predict the refractivity parameters of atmospheric channels. The range of hyper-parameters in an RF model is represented as each column with certain feature selection criterion, minimum samples per data, and hyper-parameter.

Bayesian Optimization is integrated into the RF model to systematically search for the most effective hyper-parameter configurations, allowing the model to achieve better generalization and performance. Additionally, K-fold Cross-Validation (CV) is applied to the BO-RF model to rigorously evaluate its training process and reduce the risk of overfitting.

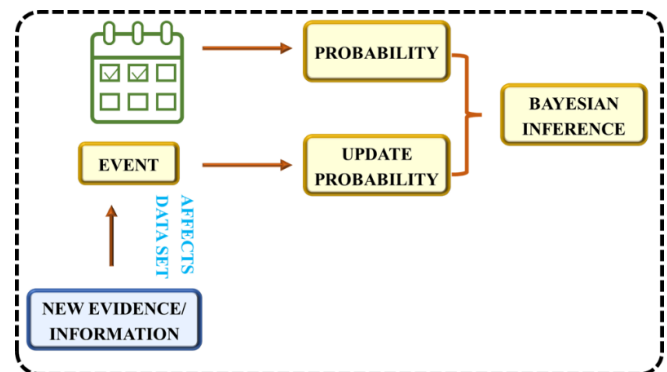


Figure 3: Architecture of Bayesian Inference

Figure 3 illustrates the process of Bayesian Inference, which involves updating the probability of an event based on new evidence or information. Initially, a probability is assigned to an event. As new data or evidence becomes available, it affects the dataset and leads to an update in the event's probability. This process of revising the prior probability based on observed evidence results in a posterior probability. The RF's hyper-parameters are then optimized through BO, which builds a probabilistic surrogate model to estimate and fine-tune the atmospheric duct parameters.

Advantages

- RF provides strong performance even on noisy or imbalanced cybersecurity data.
- Bayesian inference allows security systems to make decisions with confidence scores.
- Uncertainty modeling helps filter low-confidence classifications.

7. Proposed System Description

With today's digital world, cybersecurity has become an essential concern. We use big data analytics to successfully anticipate and stop possible cyberthreats. Conventional cybersecurity approaches frequently rely on rule-based detection methods and isolated systems, which are inadequate for spotting complex and dynamic threats.

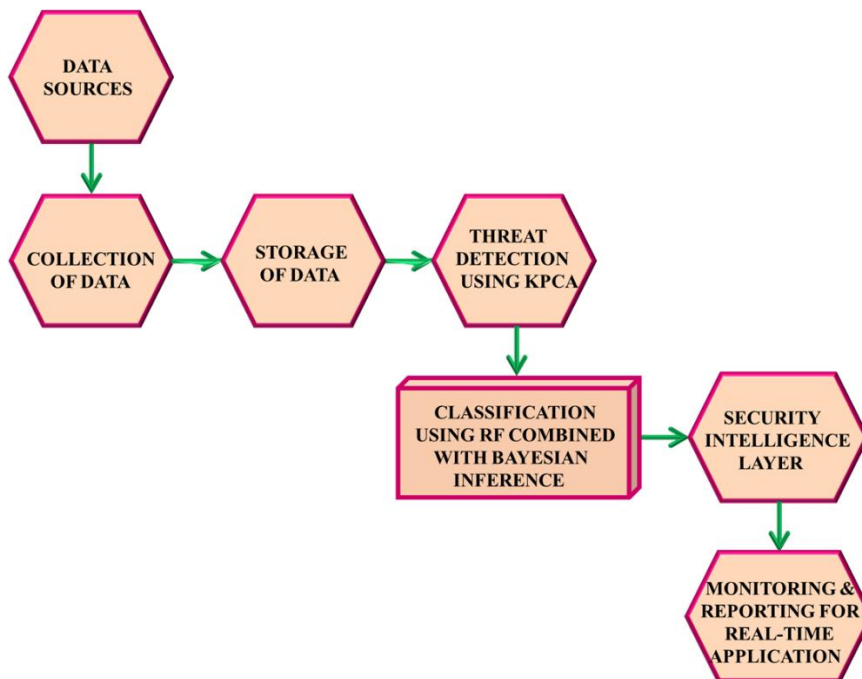


Figure 4: Proposed Block Diagram

Figure 4 depicted a security framework for effective threat identification and categorization through the use of cutting-edge intelligence approaches. The first step in the process is collecting data from various sources, which is then directed into the data collection stage. The data storage step is where this gathered data is methodically arranged and kept. After storage, Kernel PCA, a dimensionality reduction technique, is utilized to discover anomalies and possible risks in the data. Following the detection of threats, the data is sent for classification utilizing a Bayesian Inference and RF method, which allows for improved accuracy using optimization strategies inspired by biology. Important security insights are informed by the classified data, which is a component of the security intelligence layer. Finally, the framework supports monitoring and reporting for real-time application, ensuring proactive threat response and situational awareness.

8. Result and Discussion

This research investigates the proposed RF in conjunction with Bayesian inference. Swarm intelligence concepts are used to assess the hybrid algorithm's convergence speed, flexibility and resilience in handling enormous volumes of cybersecurity data. Additionally, to demonstrate advancements in predictive capacity and system reliability, the efficacy of the suggested model is contrasted with current cybersecurity methodologies.

Table 1: Comparison of Accuracy

Sl. No	Techniques	Accuracy
1.	PCA [11]	95%
2.	ANN [12]	77%
3.	LR [13]	94.2%
4.	KNN [14]	74%
5.	Naïve Bayes [15]	70.79%
6.	XGBOOST [16]	94.2%
7.	Kernal PCA	96.7%

Table 1 presents a comparative analysis of the accuracy of various machine learning techniques applied in cybersecurity. PCA and XGBoost follow closely with 95% and 94.2% respectively, showing their effectiveness in dimensionality reduction and ensemble learning. Logistic Regression also performs well with an accuracy of 94.2%, while ANN and KNN exhibit moderate accuracies of 77% and 74%. Naïve Bayes, while fast and simple, records the lowest accuracy at 70.79%, Kernel PCA demonstrates the highest accuracy at 96.7%, indicating its superior ability to capture complex.

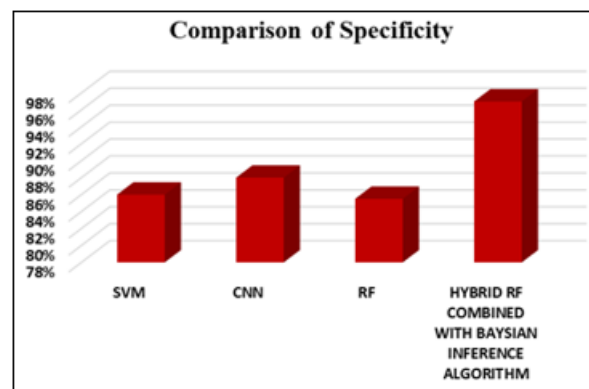


Figure 5: Comparison of Specificity

Figure 5 presents a comparison of the specificity achieved by different methods. The RF combined with Bayesian Inference approach demonstrates the highest specificity of 97%, outperforming SVM [17], CNN [18] and RF [19].

9. Related Works

A cybersecurity awareness focused on software and email security is effectively analyzed using statistical methods such as the normal distribution. Using normal distribution analysis, this methodology gathers user response data from surveys or system logs and uses it to measure awareness levels, detect outliers, and comprehend average behavior.

The simplicity and efficacy of this method lie on its ability to visualize trends in user behavior and gauge the level of security awareness among a population. However, the accuracy of inferences is limited because not all data naturally follow a normal distribution [20].

An AI driven ERP systems to increase cybersecurity resilience in real time threat environments combines large-scale data processing with intelligent decision-making. Using AI algorithms, the procedure analyzes massive, real-time datasets to identify anomalies and potential threats in Enterprise Resource Planning (ERP) systems. Organizations are able to preserve operational continuity while promptly responding to security breaches because to this connection. Automated, real-time threat identification and mitigation across organizational systems is possible with this method. However, there face a challenge on implementation's complexity and expense [21].

A cybersecurity risk assessment in smart city infrastructures using probability theory has emerged as a robust approach for quantifying uncertainties in potential threats. By using threat intelligence and historical data to model the likelihood of cyber incidents, this methodology allows for predictive analysis for proactive protection. It offers a data-driven, structured framework for effectively allocating resources and setting risk priorities. This depends on precise and thorough data, though, because inaccurate or out-of-date data might produce inaccurate risk assessments [22].

A big data framework for cybersecurity utilizing an optimized deep learning algorithm based on a genetic algorithm presents a novel approach to detecting and responding to threats in real time. This approach makes use of constant flow of diverse, large-scale data from user systems, networks, and Internet of Things devices. By simulating natural selection, the GA enhances model parameters and architecture, guaranteeing strong performance against changing cyberthreats. High detection accuracy and automated learning from dynamic data streams are two advantages of this method. However, by facilitating scalable, intelligent threat detection and mitigation in dynamic digital environments, the framework greatly enhances proactive cybersecurity [23].

A cybersecurity prevention framework using the IoT integrated with Particle Swarm Optimization (PSO) offers a proactive solution for mitigating threats across connected environments. The methodology involves modelling the behavior of IoT devices and identifying vulnerabilities through data collected in real time. PSO mimics the social behavior of particles in a swarm to search for optimal solutions, such as detecting intrusions or reducing false alarms in threat detection systems. This technique is adaptability and fast convergence, which is essential for resource-constrained IoT systems. Additionally, PSO's lightweight computation makes it suitable for decentralized security in edge devices. However, its tendency to get trapped in local optima if not properly tuned by this technique [24].

An enhanced cybersecurity framework that integrates Ant Colony Optimization (ACO) with an Artificial Neural

Network (ANN), in the TensorFlow platform for adaptive learning and real-time threat detection. ACO is employed to optimize key parameters of the neural network, such as weight initialization and learning rates, by mimicking the intelligent pathfinding behavior of ants. The methodology involves using ACO to optimize critical hyperparameters of the ANN such as weights, biases, and learning rates to improve learning efficiency and classification accuracy for anomaly detection in IoT-generated data. It provides its high adaptability and improved accuracy, making it effective for resource-constrained, heterogeneous IoT environments that demand intelligent, lightweight, and scalable security solutions. However, it is computational overhead introduced by ACO, particularly when dealing with extremely large and dynamic datasets, which affect the speed and energy efficiency of real-time IoT systems [25].

A study introduces a cybersecurity enhancement approach for Industry 4.0 wireless communication networks using a Fuzzy Harmony Search (FHS) technique. The methodology combines the Harmony Search Algorithm a metaheuristic inspired by musical improvisation with fuzzy logic to address the uncertainties and dynamic nature of cyber risks in complex industrial environments. This approach is its robustness in uncertain environments and adaptability to evolving threats, making it well-suited for the interconnected and real-time communication systems found in Industry 4.0. However, its potential increase in computational complexity due to the integration of fuzzy rule sets, which require more processing time for large-scale systems [26].

10. Conclusion

The proposed framework effectively addresses the growing cybersecurity challenges in the Big Data and post-5G era by integrating advanced data analytics and intelligent optimization techniques. By employing Kernel PCA for dimensionality reduction, the model ensures efficient handling of high-dimensional network traffic data, enhancing the relevance of features for accurate threat identification. The incorporation of a RF combined with Bayesian Inference further refines the classification process, leading to improved detection rates, faster convergence and greater adaptability to evolving cyber-attack patterns. The comparative results on standard benchmark datasets shows the accuracy of 96.7% and specificity of 97% to validate the model's robustness and efficacy, showcasing its potential to outperform traditional methods in terms of precision, detection accuracy and computational efficiency.

References

- [1] Chowdhury, Salman Mohammad Abdullah, Nayem Uddin Prince, Rakibul Hasan and L. A. Mim. "The role of predictive analytics in cybersecurity: Detecting and preventing threats." *World Journal of Advanced Research and Reviews* 23, no. 2 (2024): 1615-1623.
- [2] Lekkala, Seshagirirao, Raghavaiah Avula, and Priyanka Gurijala. "Big Data and AI/ML in Threat Detection: A New Era of Cybersecurity." *Journal of Artificial Intelligence and Big Data* 2, no. 1 (2022): 32-48.

- [3] Sarker, Iqbal H. "Machine learning for intelligent data analysis and automation in cybersecurity: current and future prospects." *Annals of Data Science* 10, no. 6 (2023): 1473-1498.
- [4] Duary, Shomili, Pratyusha Choudhury, Sushruta Mishra, Vandana Sharma, Deepak Dasaratha Rao, and Adedapo Paul Aderemi. "Cybersecurity threats detection in intelligent networks using predictive analytics approaches." In *2024 4th International Conference on Innovative Practices in Technology and Management (ICIPTM)*, pp. 1-5. IEEE, 2024.
- [5] Rizvi, Mohammed. "Enhancing cybersecurity: The power of artificial intelligence in threat detection and prevention." *International Journal of Advanced Engineering Research and Science* 10, no. 5 (2023): 055-060.
- [6] Wang, Lidong, and Randy Jones. "Big data analytics in cyber security: network traffic and attacks." *Journal of Computer Information Systems* 61, no. 5 (2021): 410-417.
- [7] Dandamudi, Jaideep Sajja, Sai Ratna Prasad and Amit Khanna. "Advancing Cybersecurity and Data Networking Through Machine Learning-Driven Prediction Models." *International Journal of Innovative Research in Computer Science and Technology* 13, no. 1 (2025): 26-33.
- [8] AL-Hawamleh, Ahmad Mtair. "Predictions of cybersecurity experts on future cyber-attacks and related cybersecurity measures." *momentum* 3 (2023): 15.
- [9] Yedalla, Jayasudha. "Building cyber-Resilient Smart Cities: The role of AI and big data in urban security." *International Journal of Science and Research (IJSR)* 14, no. 2 (2025): 648-652.
- [10] Yeboah-Ofori, Abel, Shareeful Islam, Sin Wee Lee, Zia Ush Shamszaman, Khan Muhammad, Meteb Altaf, and Mabrook S. Al-Rakhami. "Cyber threat predictive analytics for improving cyber supply chain security." *IEEE Access* 9 (2021): 94318-94337.
- [11] Al-Fawa'reh, Mohammad, Mustafa Al-Fayoumi, Shadi Nashwan, and Salam Fraihat. "Cyber threat intelligence using PCA-DNN model to detect abnormal network behavior." *Egyptian Informatics Journal* 23, no. 2 (2022): 173-185.
- [12] AL-Ghamdi, Abdullah Saad AL, Mahmoud Ragab, Maha Farouk S. Sabir, Ahmed Elhassanein, and Ashraf A. Gouda. "Optimized Artificial Neural Network Techniques to Improve Cybersecurity of Higher Education Institution." *Computers, Materials & Continua* 72, no. 2 (2022): 4.
- [13] Adejumo, A., and C. Ogburie. "Strengthening finance with cybersecurity: Ensuring safer digital transactions." *World Journal of Advanced Research and Reviews* 25 (2025).
- [14] Kabanda, Gabriel. "Performance of machine learning and big data analytics paradigms in cybersecurity and cloud computing platforms." *Global Journal of Computer Science and Technology: G Interdisciplinary* 21, no. 2 (2021): 1-25.
- [15] B. Kranthikumar and R. L. Velusamy, "SQL injection detection using REGEX classifier," *J. Xi'an Univ. Archit. Technol.*, vol. Volume XII, no. VI, pp. 800–809, 2020.
- [16] Farhan, Ammar Hatem, Omar Salah F. Shareef, and Rehab Flaih Hasan. "The Effect of False Predictions of Machine Learning on the Security of the Big Data Environment." *Iraqi Journal of Science* (2025): 361-374.
- [17] Saikin, Saikin, Sofiansyah Fadli, and Maulana Ashari. "Optimization of Support Vector Machine Method Using Feature Selection to Improve Classification Results." *JISA (Jurnal Informatika dan Sains)* 4, no. 1 (2021): 22-27.
- [18] Patasik, Eva Sapan, and Sri Yulianto. "Classification of Regional Languages Using Methods Gradient Boots and Random Forest." *Jurnal Teknik Informatika (JUTIF)* 4, no. 5 (2023): 1249-1255.
- [19] PASSIAS, PETER G., COLE A. BORTZ, KATHERINE E. PIERCE, HADDY ALAS, AVERY BROWN, DENNIS VASQUEZ-MONTES, SARA NAESSIG ET AL. "A SIMPLER, MODIFIED FRAILITY INDEX WEIGHTED BY COMPLICATION OCCURRENCE CORRELATES TO PAIN AND DISABILITY FOR ADULT SPINAL DEFORMITY PATIENTS." *INTERNATIONAL JOURNAL OF SPINE SURGERY* 14, NO. 6 (2020): 1031-1036.
- [20] Alqahtani, Mohammed A. "Cybersecurity Awareness Based on Software and E-mail Security with Statistical Analysis." *Computational Intelligence and Neuroscience* 2022, no. 1 (2022): 6775980.
- [21] Chinta, Purna Chandra Rao, Krishna Madhav Jha, Vasu Velaga, Chethan Moore, Kishankumar Routhu, and GANGADHAR SADARAM. "Harnessing Big Data and AI-Driven ERP Systems to Enhance Cybersecurity Resilience in Real-Time Threat Environments." Available at SSRN 5151788 (2024).
- [22] Kalinin, Maxim, Vasilii Krundyshev and Peter Zegzhda. "Cybersecurity risk assessment in smart city infrastructures." *Machines* 9, no. 4 (2021): 78.
- [23] Hussien, Noha, Sally M. Elghamrawy, Mofreh Salem, and Ali I. El-Desouky. "A fully streaming big data framework for cyber security based on optimized deep learning algorithm." *IEEE Access* 11 (2023): 65675-65688.
- [24] Alterazi, Hassan A., Pravin R. Kshirsagar, Hariprasad Manoharan, Shitharth Selvarajan, Nawaf Alhebaishi, Gautam Srivastava, and Jerry Chun-Wei Lin. "Prevention of cyber security with the internet of things using particle swarm optimization." *Sensors* 22, no. 16 (2022): 6117.
- [25] Sadu, Vijaya Bhaskar, Kumar Abhishek, Omaia Mohammed Al-Omari, Sandhya Rani Nallola, Rajeev Kumar Sharma, and Mohammad Shadab Khan. "Enhancement of cyber security in IoT based on ant colony optimized artificial neural adaptive Tensor flow." *Network: Computation in Neural Systems* (2024): 1-17.
- [26] Diao, Zhifeng, and Fanglei Sun. "Fuzzy Harmony Search Technique for Cyber Risks in Industry 4.0 Wireless Communication Networks." *Processes* 11, no. 3 (2023): 951.

Author Profile

¹**Dinesh Kumar Budagam** is a seasoned IT professional with over 15 years of diverse experience in the technology sector. He has extensive experience working in Microsoft technologies, Hadoop Big Data Systems, Data Engineering, Implementing Security in Cloud & Big Data ecosystem, Privacy & Security, and Cybersecurity. He worked in IBM and has also served as a consultant for Microsoft and Meta, beginning his career as a developer and progressing to roles as a Technical Lead and Senior Technical Manager. Currently, Dinesh Budagam serves as a Senior Manager and Senior Cybersecurity Consultant at VISA, specializing in cybersecurity. In his role as a Security Architect, Dinesh is at the forefront of driving key initiatives aimed at safeguarding the organization's most critical assets. He leads efforts to design and implement robust security frameworks and strategies specifically tailored for Big Data and Cloud environments. His work focuses on ensuring that these platforms maintain high levels of security, compliance, and resilience against evolving cyber threats. By collaborating with cross-functional teams, he helps align security policies with business objectives, ensuring the protection of sensitive data and the integrity of VISA's digital infrastructure. Dinesh Kumar Budagam holds a Bachelor of Technology in Electrical and Electronics Engineering and a master's degree in software engineering, providing a solid foundation for his technical expertise. In addition to his academic credentials, he is an IBM Certified Solution Architect for Big Data Analytics, Cloud Certified Architect, Certified System Architect, Project Management Professional (PMP) and Senior Member of IEEE. Furthermore, he has completed an advanced cybersecurity course at Stanford University, demonstrating a commitment to staying at the forefront of industry advancements.