**University of Hertfordshire**

School of Physics, Engineering and Computer Science

# Msc in Artificial intelligence and Robotics

# Module Title: Advanced Computer Science Masters Project

**Module Code: 7COM1039-0509-2024**

**4th September,2025**

# Topic: Data Protection Challenges in Information Security and application of AI in mitigating the issue using machine learning

**Name: Kavya Raj**

**Student ID: 23055509**

**Supervisor: Dr. Myasar Tabany**

# *Abstract*

Concerns about data protection have grown in recent years due to the rapid computerization of sensitive information, especially in light of increasingly complex cyberattacks. In terms of stopping intrusions and protecting against adherence to data privacy laws like GDPR and HIPAA, traditional measures like firewalls, encryption, and access control have fallen short. Through vulnerability management, anomaly detection, and privacy augmentation with intelligent, self-governing systems, this study investigates the application of artificial intelligence (AI) and machine learning (ML) to improve information security. The study shows how AI can change security monitoring from reactive to proactive by using breach data from healthcare health data breaches.

Several supervised machine learning algorithms, including Random Forest, K Nearest Neighbors, Decision Trees, and Logistic Regression, were used in the study. These algorithms were trained and validated using actual breach data. Accuracy, precision, recall, and F1 score were used to assess each model, with the help of data pre processing, visualization, and balancing strategies like SMOTE. The results show that ensemble methods like Random Forest are better in terms of robustness and scalability, even though simpler models provide interpretability and flexibility. Crucially, the study highlights the trade-off between overfitting and predictive power, suggesting cautious dataset balancing and validation for practical reliability.

The findings demonstrate that AI driven anomaly detection can not only thwart threats but also promote transparency, compliance, and trust in settings involving sensitive data. The article ends with suggestions for applying machine learning to multi layered security strategies and outlining future research directions in explainability, real time cyber defense systems, and privacy preserving AI.

**Keywords**: Data Protection, Information Security, Artificial Intelligence, Machine Learning, Logistic Regression, Decision Tree, KNN, Random Forest, Privacy, Cybersecurity

# Acknowledgement

# *Project Declaration*

---

This report is submitted in partial fulfilment of the requirement for the degree of Master of Science in Artificial intelligence and Robotics at the University of Hertfordshire (UH). It is my own work except where indicated in the report. I did not use human participants in my MSc Project. I hereby give permission for the report to be made available on the university website provided the source is acknowledged.

# Contents

# Introduction

In an information oriented world, information security's core problem is protecting sensitive digital information. With businesses collecting and processing vast amounts of personal, financial, and operational information, they increasingly face loss of privacy, regulatory issues (e.g. GDPR), and sophisticated threats such as data breaches, insider threats, and adversarial attacks [Emehin et al., 2024].Traditional security controls e.g. perimeter focused firewalls, static encryption are wanting with the dynamic, multidimensional nature of today's threat environments.

Through improvements in the performance and precision of threat detection, response, and prevention, artificial intelligence (AI) has revolutionized cybersecurity. Cybersecurity solutions in the past Rules-based systems that were frequently unable to cope with the severity and complexity of cyberthreats of the present day formed the foundation of such solutions. Security systems are now able to handle vast volumes of data, determine trends, and notice anomalies in a matter of seconds thanks to the incorporation of AI, especially machine learning (ML) and deep learning algorithms [Emehin et al., 2024]. Intrusion detection systems (IDS) through the use of machine learning models to detect suspicious activity in network traffic is one of the main areas in which artificial intelligence (AI) has yielded much promise. Response times for incidents are decreased through the ability of these models to identify possible threats earlier than conventional systems. AI is being used more and more in threat intelligence, which processes information that arrives from so many various sources in order to reveal potential threats and vulnerabilities [Sarker, 2023].

Here, Artificial Intelligence (AI) and Machine Learning (ML) feature as strong tools in data security enhancement. ML-based anomaly detection can identify abnormal behavior or intrusion in real time, better compared to the conventional signature-based systems [Gruschka et al., 2018]. Concurrently, privacy saving mechanisms such as differential privacy, federated learning, homomorphic encryption, and confidential computing facilitate model training and deployment securely without exposing raw sensitive information [Rani et al., 2024].

But adopting AI for security has its own challenges ranging from model inversion and membership inference attacks that are in a position to extract training data, algorithmic bias and regulatory uncertainty. Organizational governance, regulation, and technological advancements must work together symbiotically to overcome these problems. [Xu et al., 2021].

## 1.1   Background Details

Concerns about protecting sensitive data from growing cyberattacks have significantly increased as a result of the electronic dissemination of data across domains. Common techniques for protecting historical data, such as access control, encryption, and encirclement firewalls, may not always be adequate in light of sophisticated attacks and distributed computing environments. The attack surface is greatly increased by federated environments, cloud platforms, and Internet of Things (IoT) networks, creating new avenues for illegal access and data leakage. [Jahns et al., 2025].

Although there are many benefits to the growing use of machine learning (ML) in high stakes applications like driverless cars and healthcare diagnostics, there are also some new data privacy concerns. When ML models are trained on sensitive data, they may unintentionally pick up personal information, making them susceptible to inference based privacy attacks like membership inference and model inversion [Christodoulou and Limniotis, 2024]. Attackers may use the attacks to restore the original training data or confirm the existence of specific user data, which would be against data privacy laws like the GDPR[Christodoulou and Limniotis, 2024]. The field of privacy preserving machine learning (PPML) was created with this goal in mind. In order to achieve the long-term goal of facilitating collaborative model training without exchanging unparsed data, PPML uses strategies such as federated learning, differential privacy,

homomorphic encryption, and secure multi-party computation. There will be difficulties even when managing privacy, data utility, computational cost, and defense against hostile attacks. Furthermore, laws such as the GDPR and upcoming AI-specific regulations are requiring transparency, accountability, and data minimization in automated systems.[Xu et al., 2021]. This situation raises the stakes for data protection concerns in AI systems and provides a backdrop for developing innovative machine learning techniques that perfectly balance security and usefulness. To release intelligent systems that are safe, compliant with the law, and considerate of privacy, these presumptions must be met.

## 1.2   Problem Overview

In the modern era of large transactions and online data storage, information privacy is a significant concern. It was discovered that advanced attacks like model inversion and membership inference attacks in machine learning systems could not be defeated by conventional security measures like encryption and perimeter firewalls. As more businesses use machine learning (ML) models to automate decision-making, the likelihood of privacy violations and data misuse increases because ML systems may unintentionally leak training data. Regulations like the GDPR and the growing complexity of the threat landscape are the main causes of these problems. The centralized machine learning architectures in use today are susceptible to attack since they act as a single point of failure. Customers also feel deceived and have less insight into data usage and protection. Innovative, privacy-focused, and reliable AI-based solutions with an emphasis on security, usability, and compliance must be applied in real-world use cases to address these complex issues.

## 1.3   Current Issues

Even with the advancements in AI and cybersecurity, some problems with data protection still exist today. The vulnerability of machine learning models to privacy threats like model inversion, membership inference, and adversarial manipulation is one such serious worry. Particularly when models are operating in untrusted environments, these kinds of attacks have the potential to compromise the confidentiality of training data. Second, a lack of transparency in AI decision-making also makes it more difficult to comply with privacy

laws like the CCPA and GDPR. Although federated learning and other privacy-enhancing techniques are fantastic, they are still computationally expensive and challenging to scale. Third, businesses find it difficult to balance data privacy and utility, particularly in high-risk industries like healthcare and finance.

## 1.4   Project Details

This project investigates how machine learning can be used to improve information systems' data privacy protection, especially in high-risk areas. Researching and implementing privacy-preserving strategies like homomorphic encryption, federated learning, and differential privacy is of interest in order to protect sensitive data while it is being trained and inferred. The project compares various machine learning models according to their performance, resilience to attacks, and privacy leakage. The project simulates real-world situations using real-world data sets. A prototype system with AI-driven anomaly detection and privacy-enhancing technology is another goal of the project. In this case study, I try to show that AI is capable of not only identifying and deactivating risks but also adhering to legal requirements, offering a unified response to emerging data protection issues.

## 1.5   Research Aim

The article makes an effort to critically examine computer data protection and information security issues, particularly from an Indian viewpoint. The study investigates how artificial intelligence, specifically machine learning (ML), may be used to improve data security in software development life cycles and cloud computing. The project talks about automated data protection compliance, privacy leak prevention, and vulnerability identification using machine learning algorithms. The focus is on how AI-powered solutions can safeguard client data for SMEs that heavily depend on cloud infrastructure. The regulatory landscape in India is discussed, along with how AI can be used to automate policy compliance, model adaptive threats, and govern real-time data controls to promote compliance and transparency.

## 1.6    Research Questions and Novelty

This research is guided by the assumption that the integration of Machine Learning with cloud infrastructure makes data security and Indian SMEs' regulatory compliance significantly better. The major research questions are:

- What are the computational advantages and security advantages and challenges of migrating business data to cloud systems in India, especially in the context of SMEs?

- How do test code execution environments cause data vulnerabilities and how can they be detected using ML?

- What are the best practices and ML based techniques to integrate AI in cloud environments to safeguard consumers' data while maintaining the system scalable and performant?

## 1.7    Research Objectives

The central objective is to assess the position of Machine Learning towards addressing real data security issues in Indian cloud environments.

- To determine technical advantages and concerns (e.g., scalability of a model, security breaches) of relocating business data to the cloud with assistance from ML based frameworks in India.

- To discover and assess weaknesses in protection of data imparted by execution of test code and examine how anomaly detection or static code analysis using ML can mitigate them.

- To propose a combined AI cloud platform with the use of ML algorithms for real-time monitoring, threat prediction, and automatic response to security threats, specially tailored to secure Indian SMEs' consumer data and instill confidence in cloud based infrastructure.

The novelty of this project lies in its federated integration of anomaly detection, privacy preserving technologies within a machine learning framework for SMEs. While prior work

has researched privacy preserving ML, this project brings together a number of state of the art methods into one solution that is able to analyze data securely, in real time, without centralized data storage. It is also interested in privacy attacks like model inversion and membership inference, which are hardly dealt with together in practice implementations. Additionally, the employment of explainable AI (XAI) ensures transparency and complies with regulations areas usually overlooked in technical implementations. The Indian regulatory framework is another aspect of uniqueness because much literature available focuses on EU or US regulations. The SME cloud usage scenario based case study approach provides practical validation, and therefore, this is a rare intersection of policy, usability, and technological innovation. The effort thus paves the way for an applied, contextual, and deployable AI solution to modern data protection issues.

## 1.8  Feasibility, Commercial Context, and Risk

The initiative is highly feasible in light of the advancement of privacy preserving machine learning platforms and open-source facilities for federated learning, anomaly detection, and XAI. Cloud platforms support distributed ML models, thus making a realistic test environment achievable. Furthermore, publicly available data from health and finance sectors is accessible and suitable for testing models. The requirements for computing are in what can be managed by the facilities in university labs, leveraging common machine learning libraries. Primary challenges arise from putting a few pieces together federated training, privacy layers, and explainability but modular development patterns and existing libraries simplify. The case study method also raises feasibility by allowing flexible experimentation against real world backdrops. With academic supervision and formal checkpoints, the project timeline is sufficient to develop, try out, and document a prototype. Ethical concerns are addressed through anonymised data usage, and regulative applicability ensures long term pragmatic relevance, again adding to the project's academic and professional validity.

## 1.9  Report Structure

For clarity and logical flow, the report is divided into major sections. The introduction, which includes background information, a problem summary, and research objectives, comes first.

The literature review discusses earlier research while pointing out knowledge gaps. Data collection, preprocessing, selected algorithms, and validation procedures are all covered by methodology. The model's performance is described in the results and analysis using confusion matrices and performance metrics. Results, difficulties, and ramifications are discussed. Lastly, a summary of contributions and suggestions for further action are provided in the Conclusion, Recommendations, and Future Scope sections. Appendices for information not related to the paper's body, references, and acknowledgments are examples of supporting sections. A thorough and cohesive presentation of the research is guaranteed by this template.



Figure 1.1: Gantt Chart

## 1.10 Ethical Consideration

In this project, ethics came first, with rigorous adherence to academic and professional standards. Sensitive personal information was not revealed; only publicly accessible, anonymized datasets were used. Preprocessing was done responsibly, avoiding any manipulations that might produce misleading results. To maintain academic integrity, all references and sources were cited. Proper model presentation was ensured, along with warnings against risk of overfitting and limitations in order not to draw conclusions inappropriately. The project also considered legal frameworks such as the GDPR that addressed privacy, equity, and accountability in AI systems. These moral practices lend the study legitimacy, responsibility, and societal advantages..

# Literature Review

Data protection has long been a concern for information security, particularly as machine learning (ML) and artificial intelligence (AI) models increasingly work with private or sensitive organizational data. Model inversion is one of the serious privacy risks associated with traditional centrally trained machine learning systems. Attackers use model outputs to extract input features (like a fingerprint or facial pattern) and membership inference to ascertain whether a specific data point was used to train the model [Christodoulou and Limniotis, 2024]. To mitigate these risks, studies on Privacy Preserving Machine Learning (PPML) have proliferated. Xu et al. (2021) survey various PPML methods differential privacy, homomorphic encryption, and secure multi party computation and propose a Phase Guarantee Utility (PGU) framework for assessing trade offs between model utility and data protection [Xu et al., 2021]. To complement this, Li et al. (2019) survey federated learning systems that enable collaborative model training across distributed data silos with local raw data, both addressing regulatory (e.g. GDPR) and technical issues of privacy; they categorize system elements and point out future challenges in system design and deployment [Li et al., 2021].

In extending the discussion of federated architectures, Lyu et al. (2020) consider privacy and security in federated learning, categorizing threat models, poisoning attacks, and inference-based privacy attacks, and defense schemes and future research directions in more secure federated systems [Lyu et al., 2022]. Investigating regulatory considerations further, Truong et al. (2020) specifically write on federated learning from a GDPR perspective, examining privacy preserving techniques (e.g. anonymization, encryption) and highlighting compliance

concerns and system design trade offs [Truong et al., 2021].

At the intersection of security and robustness, recent machine learning security research explores attacks like data poisoning, evasion, model inversion, and membership inference and examines countermeasures such as adversarial training, data sanitization, and differential privacy, and their trade offs (e.g. utility loss, incomplete protection) [Paracha et al., 2024]. Researchers also examine trade offs in trustworthy AI and observe that privacy multiplying protection mechanisms (e.g., differential privacy) can degrade fairness, explainability, or robustness and that increasing model transparency may even make models susceptible to model extraction or inference attacks.

Li et al. summarize the challenges and possible routes of FL in massive networks of mobile and edge devices [Li et al., 2020]. The characteristics and challenges on FL from various study areas have recently been thoroughly explained by Kairouz et al [Kairouz et al., 2021]. However, their main focus is on cross device FL, which involves a large number of Internet of Things or mobile devices. A more recent survey compiles the federated learning platforms, protocols, and applications [Aledhari et al., 2020]. Some surveys focus on just one aspect of federated learning.

In order to meet the stringent requirements of the GDPR, traditional machine learning (ML) based applications and services must implement measures that effectively manage and protect personal data in accordance with the six data protection principles in the GDPR, as well as provide mechanisms for data subjects to fully control their data. Even though ML based systems are strengthened by several privacy preserving strategies, implementing these responsibilities in a centralized ML-based system is challenging and sometimes technologically impracticable [Wachter et al., 2017]. Large scale data collection, aggregation, and processing at a central server in such ML-based systems not only raises the risk of serious data breaches due to single point failure, but also exacerbates the lack of transparency, data misuse, and data abuse because the service providers have total control over the entire data lifecycle [Truong et al., 2019].

Phishing attempts have emerged as the most common cybercrime in recent years. They pose as people they know in order to send phony emails. The goal of these emails is to steal critical information. The recipient is tricked in order to obtain personal information. This data may contain sensitive information like passwords and credit card numbers. To fool unsuspecting victims into falling into their trap, the attackers pretend to be fishermen. Even though stolen

information can be used to support financial or malicious crimes, user awareness and strong online security are our best defenses against these evolving threats [Ansari et al., 2022]. Machine learning (ML) and natural language processing (NLP) are being used more and more in modern phishing defense in order to beat traditional heuristics. Thapa et al. (2021) demonstrate in their work on federated learning for phishing email detection that RNN and BERT models outperform centralized systems with data privacy retained across enterprises but performance may be sacrificed under highly imbalanced data distributions [Thapa et al., 2023]. Using an LSTM based detection system with word embeddings split across clients, Sun et al. (2021) present Federated Phish Bowl, a decentralized system that ensures data privacy but attains centralized detection (83 percent accuracy) [Sun et al., 2022]. Simultaneously, Maneriker et al. (2021) present URLTran, using transformer based embeddings to identify phishing URLs with resistance against adversarial URL attacks and an actual positive rate of approximately 86.8 percent with very low false positive rates. Combined, these findings indicate a direction towards robust, decentralized, and privacy sensitive phishing detection systems driven by cutting edge deep learning structures [Maneriker et al., 2021].

In real world deployment, AI is increasingly utilized to detect and respond to cybersecurity threats. Applied research explains how ML algorithms that were trained on benign and malicious behavior are capable of detecting ransomware in real time, like frameworks like RTrap (through decoy files as baits) and RansomAI (through reinforcement learning agents for simulation based adaption), highlighting active detection capabilities [Okdem and Okdem, 2024]. Broader surveys of AI in cybersecurity also recognize ML's ability to process massive, high speed streams of threats, detect anomalies, and learn to keep up with evolving attacks, though they place strong emphasis on the importance of feature selection, explainability, robustness, bias reduction, and efficiency in practical intrusion detection systems [Vourganas and Michala, 2024].

Together, this research supports a case study approach that addresses data protection concerns i.e., privacy leakage, model attacks, regulatory compliance, and system integration and AI/ML countermeasures, such as federated learning architectures, privacy preserving methods, adversarial defense, and cyber defense uses. A sample case study might explore deploying a federated intrusion detection model across a set of organizations, with regard to privacy leakage via membership inference and the performance of differential privacy or secure aggregation.

## 2.1   Research Gap

Even though privacy serving machine learning (PPML) has advanced significantly, there are still a number of significant gaps in both practice and research. First, although federated learning (FL) has been widely heralded for data locality, past research largely depends on oversimplified or unrealistic threat models. Field deployments face heterogenous classes of attackers from honest but curious nodes to full insiders without a shared toolkit to specify and compare these various threat assumptions, it is difficult to build effective defenses that map to reality [Aïvodji et al., 2019].

Second, scalability and communicational efficiency are still major challenges. FL models are typically tied to high communication latency and bandwidth bottlenecks by recurrent model update exchanges, particularly problematic in resource constrained environments (e.g., IoT or mobile). While techniques such as model compression and quantization are proposed, end-to-end solutions that compromise on privacy, accuracy, and efficiency are nascent [Hu et al., 2024].

Third, there exists a severe lack of comprehensive paradigms addressing the interdependencies between privacy, fairness, and robustness. Research operates to examine these problems separately, but hardly ever addresses ways in which adding differential privacy can compromise fairness or how model integrity attacks are comprised together with privacy guarantees. Building integrated approaches that jointly optimize across these areas remains an open challenge [Kairouz et al., 2021].

Fourth, user trust in and perception of FL systems are under explored. Even when theoretical privacy guarantees exist, users may misread the protections they offer or over estimate their effectiveness. There are few methods to examine how well users understand privacy mechanisms in FL, or to tailor system design according to users' privacy preferences [Kairouz et al., 2021].

In short, some of the most urgent research challenges are: (1) formalizing realistic threat models; (2) making end user friendly scalable and efficient privacy preserving protocols for resource constrained settings practical; (3) developing end to end models that balance goals of privacy, fairness, and robustness; and (4) incorporating considerations of understanding and trust by end users into system design.

## 2.2   Comparative Analysis

Prior studies on AI and machine learning for data protection have mostly concentrated on improving privacy through encryption schemes, federated learning, and differential privacy. In their comprehensive review of privacy preserving machine learning (PPML) techniques, Xu et al. (2021) and Li et al. (2021) demonstrated how promising these techniques are at minimizing the exposure of raw data during training. High computational costs and the difficulty of striking a balance between privacy and performance continue to restrict their use. Despite being praised for decentralized training, federated learning faces communication barriers and lacks flexibility in distributed real world situations [Kairouz et al., 2021].

Additionally, phishing studies have specifically addressed the subject, i.e., Thapa et al. (2023) and Sun et al. (2022), highlight the fact that RNN and LSTM models, when used in federated settings, preserve user privacy at the cost of accuracy degradation due to data imbalance [Thapa et al., 2023]. Although transformer models like URLTran produce more accurate results, they are less sensitive to context-aware responses and real time adaptability. [Maneriker et al., 2021].

Although technologically sound, such methods often overlook sociotechnical issues such as user confidence, system interpretability, and regulatory compliance, particularly in countries like India where cloud computing is widely used by SMEs but institutional arrangements are still developing. The project fills these gaps by exploring AI powered anomaly detection and policy driven privacy in cloud infrastructures, optimized for Indian SMEs, thereby extending current models with actual time compliance, scalability, and end user interpretability.

# Methodology

## 3.1 Choice of Method

The data was normalized and feature selected beforehand to maximize data quality and model performance. Random Forest, Decision Tree, KNN and Logistic Regression supervised learning methods were used to detect anomalies and potential intrusions. The models were tested and trained with cross validation to maintain generalizability. Accuracy, precision, recall, and F1 score were computed for measurement. Implementation was done using Python and libraries such as Scikit learn, Pandas, and Matplotlib. The methodology centers on how machine learning can be applied practically for security threat detection, allowing proactive security measures, and showing how AI enhances resilience against modern data protection issues.

### 3.1.1 Logistic Regression

**Working Principle**: Logistic Regression is a supervised machine learning algorithm that's both binary and multiclass classification. It gives the representation of the probability of an input being in a certain class according to the logistic (sigmoid) function. The output is between 0 and 1 and is considered as a probability. The algorithm learns the weights of the input features in order to minimize the difference between the predicted label and actual label using a cost function, usually cross entropy loss.

The logistic (sigmoid) function is defined as:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

where:

$$z = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_n x_n$$

The predicted probability is:

$$P(y = 1|x) = \frac{1}{1 + e^{-(\beta_0 + \sum_{i=1}^{n} \beta_i x_i)}}$$

The cost function (log loss) is:

$$J(\beta) = -\frac{1}{m} \sum_{i=1}^{m} \left[ y^{(i)} \log(h_\beta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\beta(x^{(i)})) \right]$$

**Relevance to Project**: Logistic regression is useful to find binary security results, for example, if a system is attacked or not attacked. On data protection, it can be used for anomaly classification, detection of fraud in access patterns, or identifying if a transaction is safe. Because it is easy and quick, it suits the first baseline models well, as well as where interpretable outcomes are required. While less able with complex, nonlinear data, it helps identify the primary attributes influencing security vulnerabilities.

### 3.1.2 Decision Tree

**Working Principle**: Decision Tree is a tree like flowchart used for classification and regression. It splits data into subsets based on a function of feature values, using splits determined by criteria like Gini impurity or information gain. The tree continues to split until it hits a stopping criterion (like maximum depth or pure leaves). Each root to leaf path is a classification rule. For classification trees, we use metrics such as Gini Impurity:

$$Gini(D) = 1 - \sum_{i=1}^{c} p_i^2$$

Or Information Gain based on entropy:

$$Entropy(D) = -\sum_{i=1}^{c} p_i \log_2(p_i)$$

$$Gain(D, A) = Entropy(D) - \sum_{v \in Values(A)} \frac{|D_v|}{|D|} Entropy(D_v)$$

**Relevance to Project**: They are best suited for pattern recognition of user activity or system logs to identify malicious activities. In data security, they help explain decisions that led to a security violation, making them best for root cause analysis. They are ideal for explainability based on their explainable and rule-based nature, which is critical during cyber security audits and compliance. Decision trees have the susceptibility of overfitting but are strong when implemented using ensemble techniques such as Random Forests.

### 3.1.3   K Nearest Neighbors (KNN)

**Working Principle**: KNN is a non parametric, instance based learning algorithm that classifies a point of data into the majority class of its k most proximate points in the feature space. It uses the distance measures such as Euclidean, Manhattan, or cosine similarity to determine its nearest neighbors. The algorithm does not require any training phase and is therefore computationally inexpensive for training but slower during prediction.

The Euclidean distance between two points $x$ and $x'$ is:

$$d(x, x') = \sqrt{\sum_{i=1}^{n} (x_i - x_i')^2}$$

Classification is done by majority vote among the $k$ nearest neighbors:

$$\hat{y} = \arg \max_{v \in \{1, ..., C\}} \sum_{i \in \mathcal{N}_k(x)} \mathbb{I}(y_i = v)$$

**Relevance to Project**:KNN is useful in anomaly detection for cybersecurity in identifying unusual behavior that is different from usual behavior. It can mark unknown attempts at access or recognize patterns of intrusions by comparing with previous access history. It is simple and helpful in case of small datasets or for exploratory needs. It is computationally demanding, yet with vast datasets, and the performance will be affected with high dimensional or noisy data.

### 3.1.4   Random Forest Method

**Working Principle**: Random Forest is an ensemble method for machine learning that builds many decision trees and uses their collective output to make predictions that are more robust and accurate. Each tree is trained on a random portion of the data and a random portion of the features (bagging and feature randomness). The decision is made by majority vote (classification) or average (regression).

Random Forest combines predictions from multiple decision trees:

$$\hat{y} = \text{mode} \{h_1(x), h_2(x), \ldots, h_T(x)\}$$

where $h_t(x)$ is the prediction of the $t$-th decision tree, and $T$ is the total number of trees.

Each tree is trained on a bootstrap sample and at each node split, a random subset of features is considered:

$$\text{Feature subset at each split: } \sqrt{p} \text{ (for classification)}, \quad \frac{p}{3} \text{ (for regression)}$$

**Relevance to Project**: In data privacy, Random Forests perform well in intrusion detection and fraud detection since they are robust, precise, and immune to overfitting. Through the combination of numerous trees, the algorithm learns complex relationships among data and performs improved generalization. Its feature importance measures help identify major indicators of security attacks, which further facilitates risk assessment and prevention planning. It is particularly well-suited for big data security monitoring applications requiring scalable, stable predictions.

## 3.2   Justification and Support of Choices

**Logistic Regression**:Since logistic regression offers a straightforward but powerful baseline model for anomaly detection in data security issues, it was used. Its probabilistic model aids in determining whether a breach or access attempt is a malicious or benign event. Because transparency is crucial for audits and compliance, its interpretability makes it very useful in cybersecurity. Logistic Regression is computationally efficient and helpful for identifying important features that affect vulnerabilities, even though it cannot be used to illustrate complex nonlinear relationships. It provides a starting point for analysis and a benchmark

by which more intricate algorithms are evaluated to make them readable and trustworthy [Green et al., 1998].

**Decision Tree**: Decision trees were chosen because they clearly and intuitively illustrate decision rules, making them immediately applicable to information security. They can also look for log behavior or access patterns that may point to breaches thanks to their hierarchical structure. Decision Trees are ideal at generating rule-based results that are simple for analysts and auditors to comprehend, which is crucial in compliance-intensive settings like GDPR. Pruning or ensemble methods can be used to address their tendency to overfit on complex datasets. They are crucial to the root cause analysis of cyber incidents due to their perceptibility and malleability.

**K Nearest Neighbors (KNN)**:The non-parametric nature of KNN, which makes it robust in identifying anomalies without presuming an assumed data distribution, led to its selection. Plotting new data points against historical patterns allows KNN to detect "outliers," which are frequently reflected in suspicious access patterns or intrusions in cybersecurity. Its simplicity guarantees quick deployment for small datasets or exploratory scenarios without worrying about computational cost. KNN's strength is in identifying new attack patterns with the help of distance metrics, but it struggles with high dimensional or noisy input data. This is especially helpful for intrusion detection tasks where suspicious behavior can be predicted by looking for similarities to known attacks [Zhang, 2021].

**Random Forest**: Since Random Forest can overcome the drawbacks of individual Decision Trees, it was selected as a powerful ensemble technique. Accuracy and generalization are enhanced by combining predictions from several trees, which is crucial in the ever-changing world of cybersecurity. It is the perfect option for large, complicated datasets with multiple breach patterns because of its overfitting and noise resistance features. In order to aid in risk assessment and prevention planning, Random Forest also offers feature importance scores that provide information about the factors most strongly influencing risks. Because Random Forest strikes a balance between predictability and interpretability, it is perfect for fraud prevention, intrusion detection, and monitoring of numerous sensitive information systems [Paul et al., 2018].
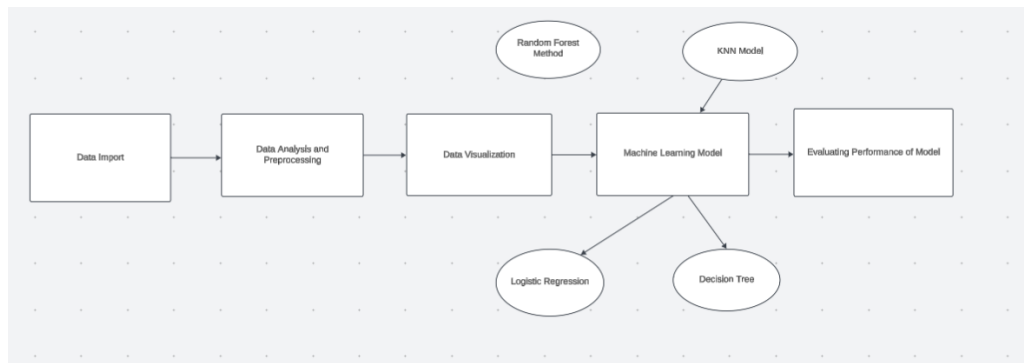
## 3.3    Data Collection/Project Design



Figure 3.1: Research Design

The dataset is taken from Kaggle Website. The data collection contains records of health-care data breaches reported in the United States. It contains 10,000 rows and 10 columns, with reported cases regarding covered healthcare entities. The most important columns are the Name of Covered Entity, State, Type of Entity (e.g., health plan, business associate, or healthcare provider), Number of Individuals Affected, Submission Date of Breach, Type of Breach (e.g., loss, theft, or improper disposal), and Where Breached Information Was Located (e.g., paper, laptops, or servers).

It is worthwhile to examine this dataset in order to find patterns in the healthcare industry's data breach failures. It can be applied to situations like determining the most susceptible entity types, calculating the effects of various breach techniques, and forecasting risk levels based on breach attributes. Additionally, it provides a practical basis for using artificial intelligence and machine learning models to identify, classify, or stop data breaches.

## 3.4    Data Preprocessing

In order to prepare the dataset for machine learning modeling, it was systematically trans-formed during the preprocessing step. .isnull() was initially used to identify missing values.To achieve uniformity, real NaN values were substituted for sum() and placeholder strings like ' N'. Important columns such as "Breach Submission Date" and "Individuals Affected" were converted into datetime and numeric formats, respectively, to enable flawless statistical analysis and time based modeling.

The missing values like 'State' fields were filled in with 'Unknown' to maintain integrity in location based analysis, while categorical 'Covered Entity Type' was imputed by its mode. The 'Individuals Affected' column used the median to handle missing entries and hence increased robustness against outliers. For categorical text data such as 'Type of Breach' and 'Location of Breached Information', missing values were replaced with 'Unknown', while 'Web Description' was excluded because it had a high frequency of null values and was not important to structured numeric modeling.

For simplification of classification, only the top 5 most frequent breach types were retained, and the rest were categorized as 'Other'. Finally, the category columns were label encoded using LabelEncoder to transform them into numeric values acceptable to machine learning algorithms while preserving uniform formatting across the dataset.
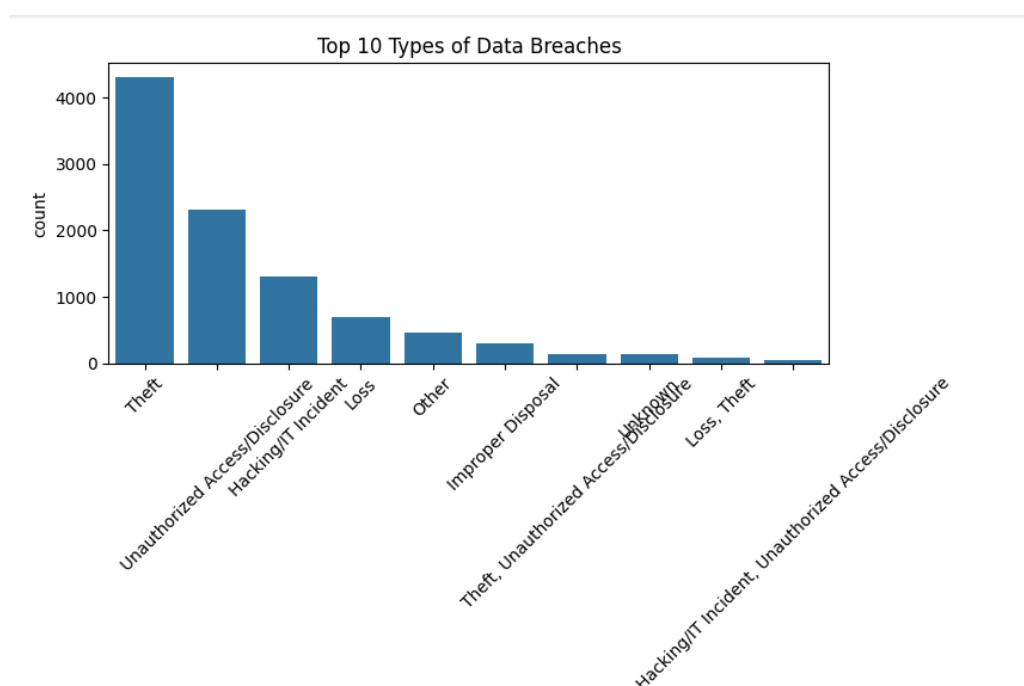
## 3.5   Data Visualization



Figure 3.2: Bar Plot for Top 10 Type of Data Breaches

Above diagram 3.2 bar chart illustrates the 10 most common categories of reported data breaches within the dataset. The most common type of breach is theft with more than 4,000 instances. The next closest is Unauthorized Access/Disclosure and Hacking/IT Incident,

each representing a very high rate of breaches. Other notable categories include Loss, Improper Disposal, and those broken down into the categories Other or Unknown. The graph also shows compound labels like "Theft, Unauthorized Access/Disclosure," which suggests instances of more than one breach means. The x axis labels are turned because of their length. This presentation reiterates that physical and unauthorized digital access remain the biggest threats to the security of data, so the necessity of strict measures of protection and breach detection systems for healthcare and allied industries is twice highlighted.It also suggests that specific AI models are required to counteract the most common kinds of breaches.
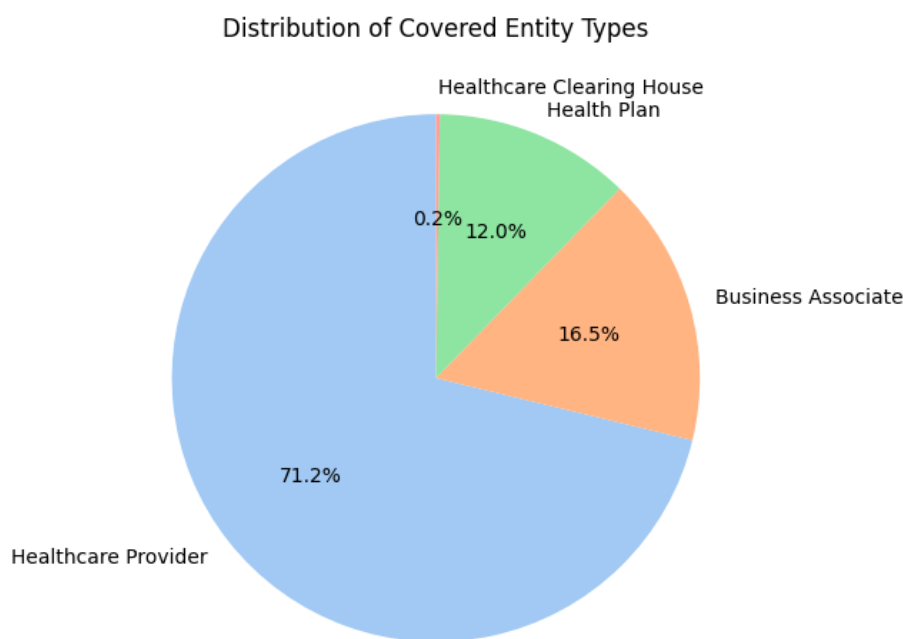


Figure 3.3: Pie Plot for Distribution of Covered Entity Types

The figure 3.3 covered entity types that were impacted by data breaches are further shown in the pie chart. Healthcare providers were the target of 71.2 percent of the breaches, demonstrating their vulnerability in terms of data handling and system security. Business Associates made up 16.5 percent, which further suggests that vendors and other third-party business associates pose serious security risks.. Health Plans were involved in 12.0 percent of events, with Healthcare Clearing Houses holding a negligible 0.2 percent of reported breaches. This split tells us that the front line organizations that have direct contact with patient data primarily healthcare providers are most prone to data leaks. These findings confirm the necessity of prioritizing data security and compliance protocols above all healthcare stakeholders,

especially those with direct contact with patient data. For the risk models based on AI, entity type is one of the most significant features and helps in forecasting probability or effect of a breach according to the organizational role in the healthcare environment.
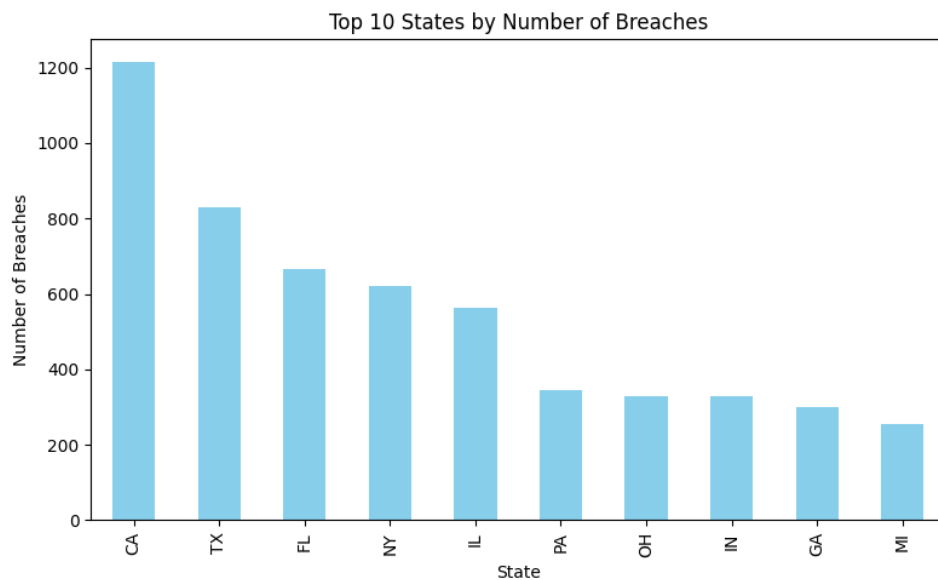


Figure 3.4: Bar Plot for Top 10 States by Number of Breaches

Above diagram 3.4 bar chart named "Top 10 States by Number of Breaches" graphically depicts the top U.S. states by number of reported data breaches. California (CA) takes the top spot with more than 1,200 breaches, followed by Texas (TX) and Florida (FL), both having very high numbers. New York (NY), Illinois (IL), and Pennsylvania (PA) follow in decreasing order. Ohio (OH), Indiana (IN), Georgia (GA), and Michigan (MI) complete the list. The y axis shows the number of breaches, and the x axis is a list of state abbreviations. The chart employs light blue bars to show breach frequency, so it is easy to see how states compare at a glance. This visualization is useful for showing regional differences in data security incidents and can indicate differences in population size, reporting needs, or cybersecurity practices between states.
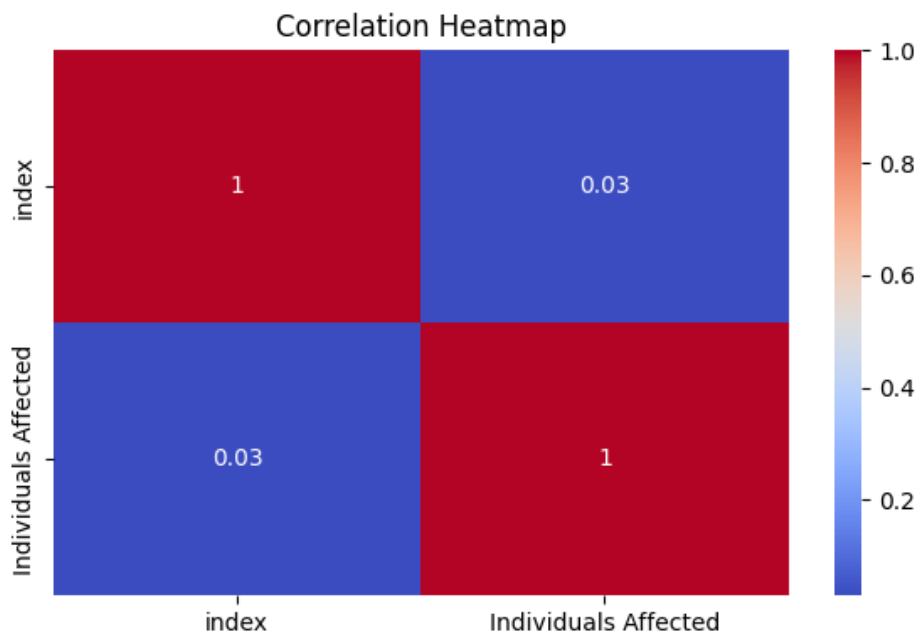
Figure 3.5: Correlation Heatmap

Above diagram 3.5 correlation heatmap shows the correlation of two variables: index and Individuals Affected. The diagonal is 1, representing perfect self-correlation. The off diagonal reveals the real correlation between the variables of concern. Here, the correlation coefficient of index and Individuals Affected is 0.03, which is just about zero, and this shows a very weak or negligible linear relationship between them. Color gradient from deep red (high positive) to deep blue (high negative) facilitates visualization of correlation strength. The nearly blue squares between the two variables of different kinds guarantee there is no significant correlation. This means that the position in the dataset (index) does not influence the number of people affected in the data being studied.

Figure 3.6: Bar Plot for Breaches Involving Business Associates

Above diagram 3.6 bar chart is labeled "Breaches Involving Business Associates" and compares the number of data breaches, based on whether an associate to a business was involved. The x axis splits into two categories of breach: "Yes", meaning an associate to a business was present, and "No", meaning no associate was involved. On the y axis, the number of breaches is shown.

The graph shows that a far greater number of breaches happened without the presence of a business associate, with more than 8,000 incidents. Breaches involving a business associate, on the other hand, are significantly lower, with fewer than 2,000. This indicates that the majority of data breaches occur independent of third party business associates, possibly reflecting either enhanced controls within associates or that internal systems are targeted more frequently. The visual difference between the two bars makes this gap apparent.
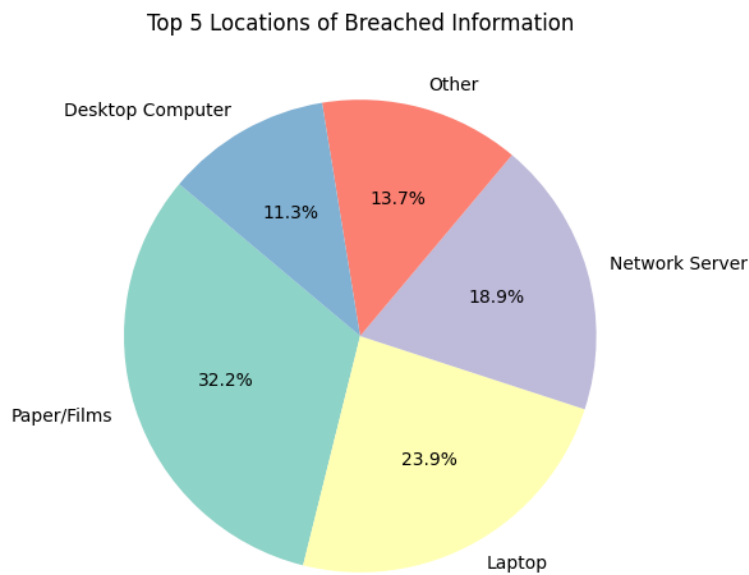
Figure 3.7: Pie Plot for Top 5 Locations of Breached Information

Above diagram 3.7 pie chart in this case reveals the top five sources of data breach and reveals that Paper/Films (32.2 percent) and Laptops (23.9 percent) are the most frequent sources, followed by Network Servers (18.9 percent), Other (13.7 percent), and Desktop Computers (11.3 percent). These intrusions refer to serious vulnerabilities of electronic and hard storage systems. This highlights the need for robust data security policies for healthcare and business environments. Employing machine learning and AI can prove highly beneficial in predication, detection, and prevention of such violations by tracking usage patterns, identifying anomalies, and automatically sending notifications. Smart systems can mandate encryption policies, identify vulnerable endpoints, and reduce reliance on manual reporting. In the case study data, over 10,000 breaches reveal systemic weaknesses most of which can be addressed with intelligent automation and real time threat intelligence using AI based models, offering a forward looking not backward looking approach to information security.
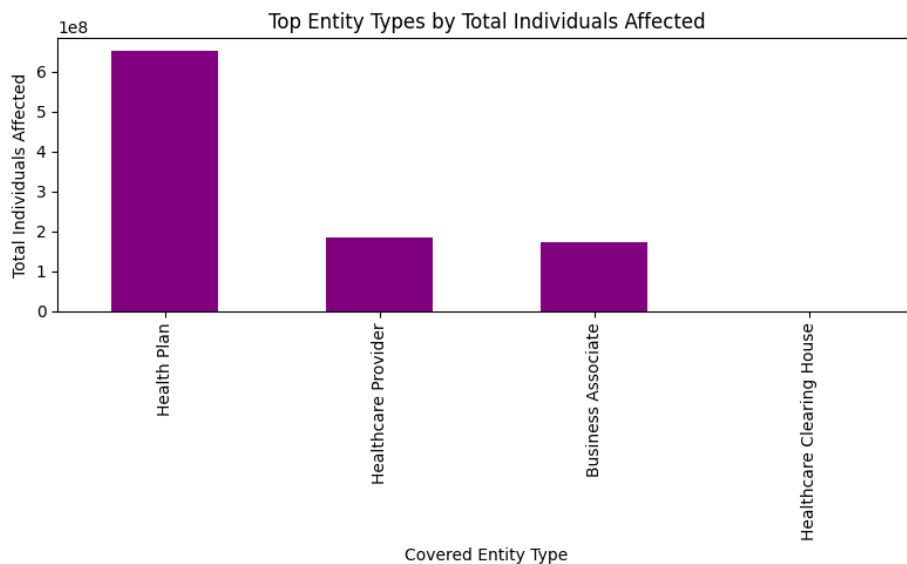
Figure 3.8: Bar Plot for Top Entity Types by Total Individuals Affected

Above diagram 3.8 bar chart illustrates the aggregate number of individuals exposed as a result of data breaches by various types of healthcare organizations. Health Plans account for the overwhelming majority, impacting over 600 million individuals considerably more than double the amount impacted by Healthcare Providers and Business Associates, each of which impacted approximately 180 million individuals. By contrast, Healthcare Clearing Houses indicate negligible impact. This discrepancy shows that intrusions are more likely to target centralized systems, such as Health Plans, which store enormous volumes of private data. It also draws attention to how vulnerable massive repositories with lax access controls and inadequate breach detection techniques are. These results emphasize the need for sophisticated security frameworks, which AI and machine learning can help with. To provide more context-dependent, scalable defense against health data breaches, AI and machine learning can be used to detect unusual activity, safeguard high risk data points, and dynamically allocate security resources to areas of greatest exposure..
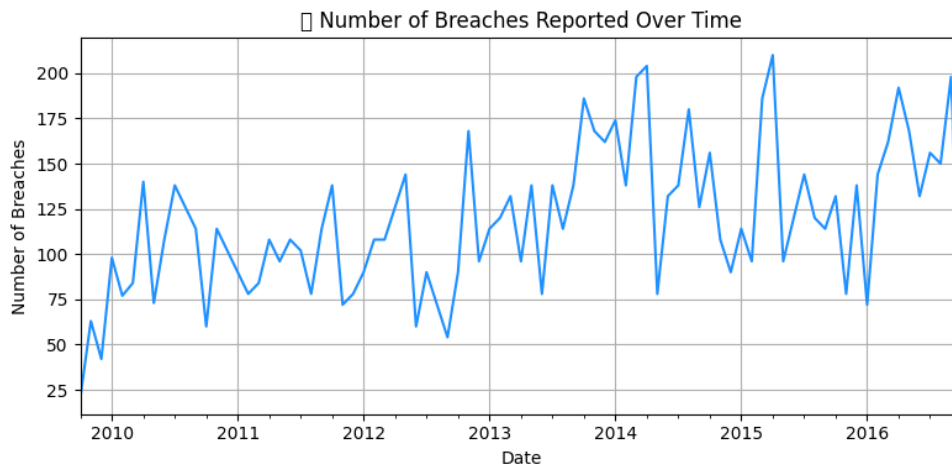
Figure 3.9: Bar Plot for Top Entity Types by Total Individuals Affected

Above diagram 3.9 line graph shows the number of breaches reported against time from the year 2010 to 2016. The date is on the x-axis, and the number of reported breaches is on the y-axis. The data shows a general upward trend of reported breaches, which may indicate an increase in the number of incidents or improved reporting procedures over the years. Peaks are seen between the years 2014 and 2016, with some months recording more than 200 breaches. This indicates key times of high cyber vulnerability. Although there are fluctuations, the general trend is upward, with increasingly more breaches being reported in later years of the timeline than in the earlier portion. The spikes and steep declines show that the reporting of breaches is not uniform and may be based on external factors such as legislation, types of attacks, or disclosure policies. The visualization also points to the growing importance of cybersecurity in an increasingly digitizing world and the need for constant monitoring and prevention mechanisms.

## 3.6 Test Strategy

- 1. Precision: Precision measures the number of correctly predicted positive instances that are actually positive. Precision is interested in positive prediction accuracy.

- 2. Recall (Sensitivity or True Positive Rate): Recall measures the number of true positive instances properly identified by the model. It captures the ability of the model to detect all cases of interest.

- 3. F1 Score: The F1 Score is the harmonic mean between precision and recall. It gives equal weight to both and is particularly helpful when classes are unbalanced.

- 4. Accuracy: Accuracy is the number of correctly predicted instances divided by the total number of predictions. Accuracy performs best when classes are balanced.

- 5.Classification Report: A classification report summarizes the performance of a classification algorithm by reporting precision, recall, F1 score, and support (number of true instances) for each class.

**1. Precision**:
$$\text{Precision} = \frac{TP}{TP + FP}$$

**2. Recall (Sensitivity)**:
$$\text{Recall} = \frac{TP}{TP + FN}$$

**3. F1 Score**:
$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

**4. Accuracy**:
$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

## 3.7 Validation

To guarantee dependability and generalizability, the selected machine learning algorithms were validated using cross validation, accuracy, precision, recall, and F1 score. The project prevented model overfitting and enabled objective performance evaluation by splitting data into training and testing datasets. Confusion matrices also provided detailed information on misclassification trends and exposed weaknesses in particular breach types. By addressing class imbalance with SMOTE, minority breach types were represented, enhancing detection fairness. GridSearch CV s hyperparameter optimization guaranteed model effectiveness and avoided majority class bias. Instead of depending solely on accuracy as a metric, which could be deceptive in data that is unbalanced, the multi metric approach allowed for global assessment. Together, these validation methods ensured that the models were both technically sound and applicable to real-world cybersecurity settings, where false negatives can have disastrous consequences.This validation framework is reliable, strong, and prepared for wider use.

## 3.8   Practicality

The project is made possible by the use of real healthcare breach data, which mimics the difficulties organizations encounter in safeguarding sensitive information. The models demonstrated their ability to handle the difficulties present in operational environments by preprocessing noisy and incomplete datasets. Reproducibility, scalability, and ease of deployment in enterprise systems were made possible by the use of Python libraries like Scikit learn, Pandas, and Imbalanced learn. Google Colab's cloud deployment demonstrated viability for SMEs with limited computing power, especially in terms of accessibility for SMEs. Furthermore, the interpretability of Random Forest features and Decision Trees provides regulatory utility, in this case GDPR regulation compliance. Since the models are generalizable, they can be applied to monitoring, fraud detection, and anomaly detection in any sector not just the healthcare sector. Although some of the results displayed overfitting, the use of cross validation and synthetic balancing techniques enhanced the results usability. The use of actual data, scalability, explainability, and relevance to contemporary information security issues are typically used to illustrate practicability.

# Result and Quality

| Class | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| 0 | .31 | .49 | .38 | 261 |
| 1 | .17 | .55 | .26 | 139 |
| 2 | .24 | .25 | .24 | 277 |
| 3 | .70 | .42 | .53 | 860 |
| 4 | .35 | .25 | .29 | 463 |
| **Accuracy** | | | 0.38 | 2000 |
| **Macro Avg** | .36 | .39 | .34 | 2000 |
| **Weighted Avg** | .47 | .38 | .40 | 2000 |

Table 4.1: Logistic Regression Model Classification Report

The table 4.1 Logistic Regression model's classification results with optimal parameters C=10 and solver='liblinear'. It evaluated the model on 2,000 samples. Class 3 of the five classes (0 through 4) had both the highest precision (0.70) and the highest support (860), reflecting that it was both better predicted and better represented in the data. Conversely, class 1 had worst precision (0.17) but comparatively high recall (0.55), indicating the model identified many positives but at the cost of high false positives.When classes are weighted equally, the model performs poorly across all classes, as shown by the macro average precision (0.36), recall (0.39), and F1 score (0.34). The weighted average (0.47 precision, 0.38 recall, 0.40 F1) demonstrates the impact of class imbalance, with performance biased towards the

majority class. Overall accuracy of 38 percent suggests moderate performance but with room for improvement. These metrics help in considering the model's capacity to distinguish between various breach types, which is vital in real-world data protection setups.

| Class | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| 0 | 1 | 1 | 1 | 261 |
| 1 | 1 | 1 | 1 | 139 |
| 2 | 1 | 1 | 1 | 277 |
| 3 | 1 | 1 | 1 | 860 |
| 4 | 1 | 1 | 1 | 463 |
| **Accuracy** | | | 1 | 2000 |
| **Macro Avg** | 1 | 1 | 1 | 2000 |
| **Weighted Avg** | 1 | 1 | 1 | 2000 |

Table 4.2: Classification Report of the Decision Tree

The table 4.2 indicates a perfect classification report, in which the model achieved 100 percent accuracy, precision, recall, and F1 score for all five classes. All classes ranging from 0 to 4 possess the precision, recall, and F1 score of 1.00, indicating the perfect prediction performance. The support column shows the number of correct instances for each class, ensuring that the model made no error in predicting any class.

The macro average that assigns equal weight to all classes also yields perfect scores, in the sense that the model performs equally well for all classes irrespective of the class distribution. Likewise, the weighted average that considers the class imbalance in proportion to support also yields a perfect 1.00, reaffirming the model's consistent and robust performance.

While this result may seem to be the best, perfect performance on real world settings is rare and even possibly symptomatic of potential overfitting, especially if achieved on training or inherently less diverse data. It is therefore crucial to validate this performance on nonseen, real world test data to confirm the model's generalizability and real world reliability.

| Class | Precision | Recall | F1-Score | Support |
|-------|-----------|--------|----------|---------|
| 0 | 1.00 | 1.00 | 1.00 | 261 |
| 1 | 1.00 | 1.00 | 1.00 | 139 |
| 2 | 1.00 | 1.00 | 1.00 | 277 |
| 3 | 1.00 | 1.00 | 1.00 | 860 |
| 4 | 1.00 | 1.00 | 1.00 | 463 |
| **Accuracy** | | | 1.00 | 2000 |
| **Macro Avg** | 1.00 | 1.00 | 1.00 | 2000 |
| **Weighted Avg** | 1.00 | 1.00 | 1.00 | 2000 |

Table 4.3: Classification Report for Random Forest Model

The table 4.3 shows the Random Forest model classification performance, which had a score of 1.00 on all measurements. The precision, recall, and F1 score for each class (0 to 4) are all 1.00, indicating the model correctly classified every example with no false positives or false negatives. The support column verifies a balanced testing across 2,000 samples.
Total accuracy is also 100 percent, and both macro average and weighted average calculate perfect performance. Macro average weighs all classes equally, while weighted average puts class distribution into consideration. In this case, both calculate spot on classification.
Though such types of results suggest a very good model, they may also suggest overfitting, particularly if the model was validated on training data or the test data set was very close to the training data set. Random Forests are very powerful and noise resistant, but such accuracy does not often happen in reality. Therefore, additional validation by cross validation or another independent test set is warranted to determine the generalizability of the model and avoid wrong conclusions.

| Class | Precision | Recall | F1-Score | Support |
|-------|-----------|--------|----------|---------|
| 0 | 1.00 | 1.00 | 1.00 | 261 |
| 1 | 1.00 | 1.00 | 1.00 | 139 |
| 2 | 1.00 | 1.00 | 1.00 | 277 |
| 3 | 1.00 | 1.00 | 1.00 | 860 |
| 4 | 1.00 | 1.00 | 1.00 | 463 |
| **Accuracy** | | | 1.00 | 2000 |
| **Macro Avg** | 1.00 | 1.00 | 1.00 | 2000 |
| **Weighted Avg** | 1.00 | 1.00 | 1.00 | 2000 |

Table 4.4: Classification Report for KNN Model

The table 4.4 displays a classification report where the model achieved ideal scores for all five classes. The precision, recall and F1 score are all 1.00 for each class, which indicates that all predictions were accurate and there were no false positives or false negatives. The support column confirms that the evaluation picks up a balanced spread across 2,000 samples.

The 100 percent overall accuracy and spot macro and weighted averages further confirm the model's flawless performance. The macro average indicates uniform performance in all classes, while the weighted average is a class size aware indicator and proves that the model performs well even with biased data.

Although this outcome is desirable, it can also be a source of concern regarding overfitting particularly if the model had been tested on the training data. In practical applications, such a performance is not common, so one should try to test the model on new data to ascertain its strength and ability to generalize.



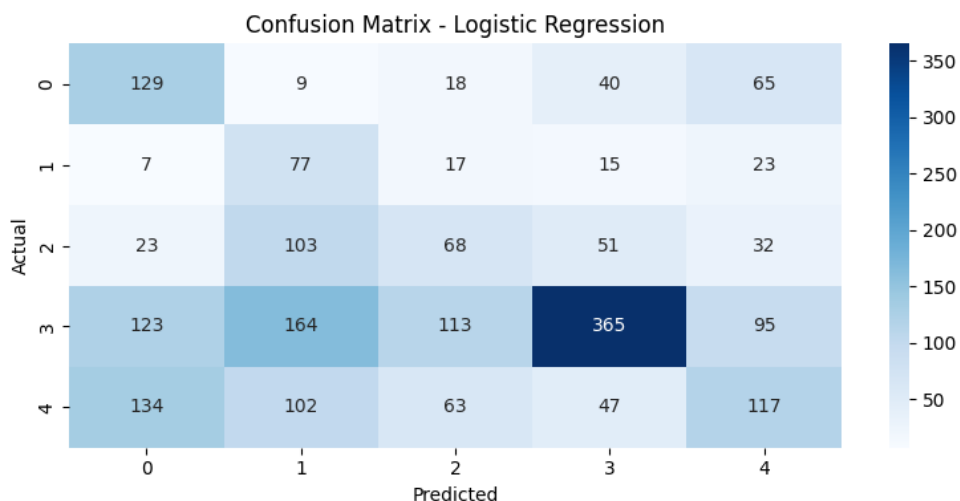Confusion Matrix - Logistic Regression

Figure 4.1: Confusion Matrix for Logistic Regression

Above diagram 4.1 confusion matrix shown here evaluates the performance of a Logistic Regression classifier on five classes (0 through 4). Each row represents the true class, and each column represents the predicted class. The values on the diagonal represent correct predictions, with the most accurate being for class 3 (365 correct predictions). Misclassifications occur in all classes, especially for class 4, as it was often predicted as class 0 or 1. For instance, 134 of actual class 4 were misclassified as class 0. Class 3 also shows some confusion with classes 0, 1, and 2. Class 1 has relatively better precision with fewer off-diagonal entries. Overall, the matrix shows that while the model is performing a good job of predicting certain classes (like class 3), it is not performing well in distinguishing between others, particularly class 4. This shows the need for potential data balancing, feature improvement, or a more complex model for better classification performance.
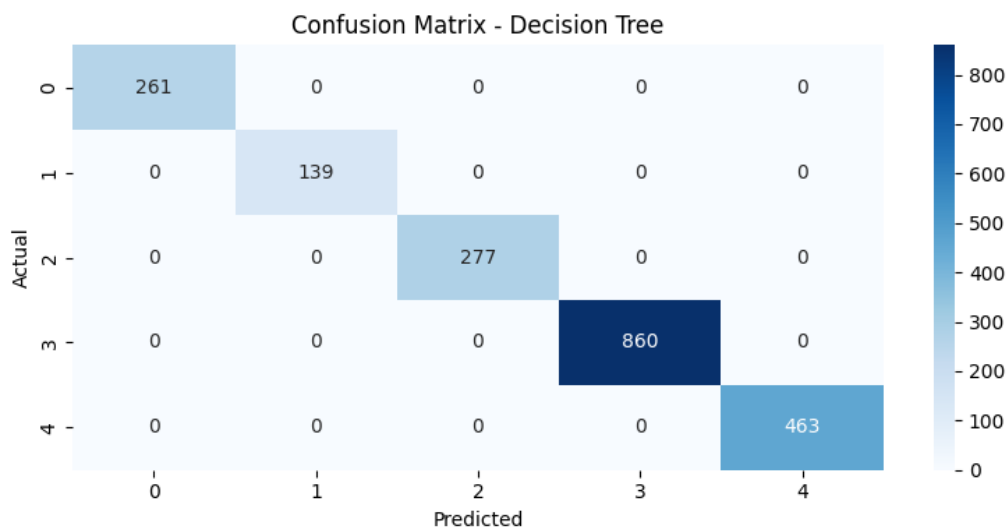


Figure 4.2: Confusion Matrix for Decision Tree Model

Above diagram 4.2 decision Tree model's confusion matrix shows very good classification performance for all five classes (0 to 4). Each actual class exactly matches its predicted equivalent, as supported by firm diagonal dominance and off diagonal entries of zero. Class 3, for example, has 860 accurate predictions with zero misclassifications, and class 4 has 463 correctly predicted samples. Perfect classification shows that the decision tree model effectively learned the decision boundaries within the training or testing data. However,

such perfect outcomes are also a sign of overfitting, especially if this degree of performance was observed on training data rather than unseen test data. Overfitting occurs when the model memorizes the training data but does not generalize well to new data. Thus, while the decision tree's performance appears ideal in this matrix, its assessment on an independent validation set or cross-validation is required to determine its true predictive ability and avoid misleading outcomes.
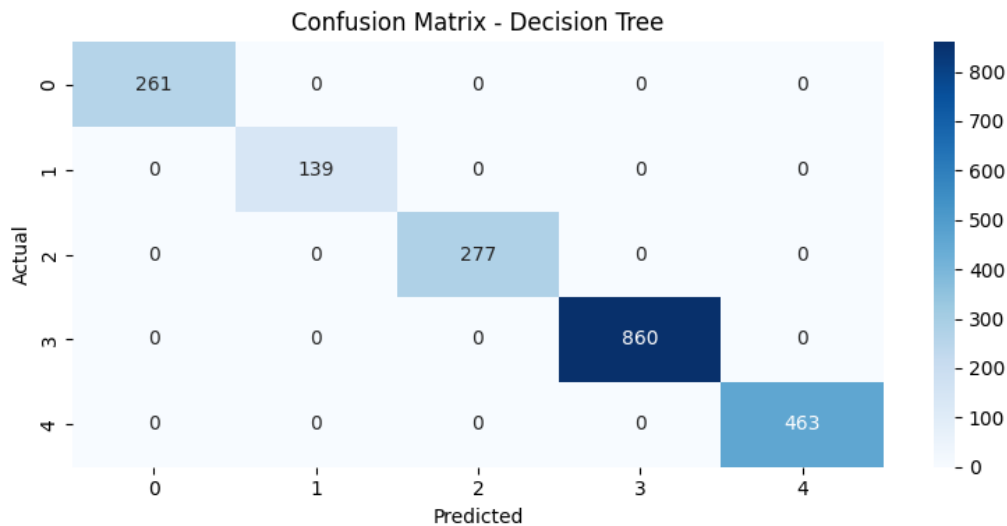


Figure 4.3: Confusion Matrix for Random Forest Model

Above diagram 4.3 show that every instance is accurately predicted across the five classes, the Decision Tree model's confusion matrix demonstrates faultless classification performance. The off diagonal entries are absent, indicating zero misclassifications, while the diagonal entries 261, 139, 277, 860, and 463 are all accurate predictions. This result indicates that the model has perfect accuracy in learning the dataset's decision boundaries. Even though these results show excellent predictive power, they could also be the result of overfitting, especially if the model was tested using training data rather than test data that hasn't been seen yet. To guarantee generalizability and prevent drawing erroneous conclusions for practical applications, validation on separate datasets is crucial.
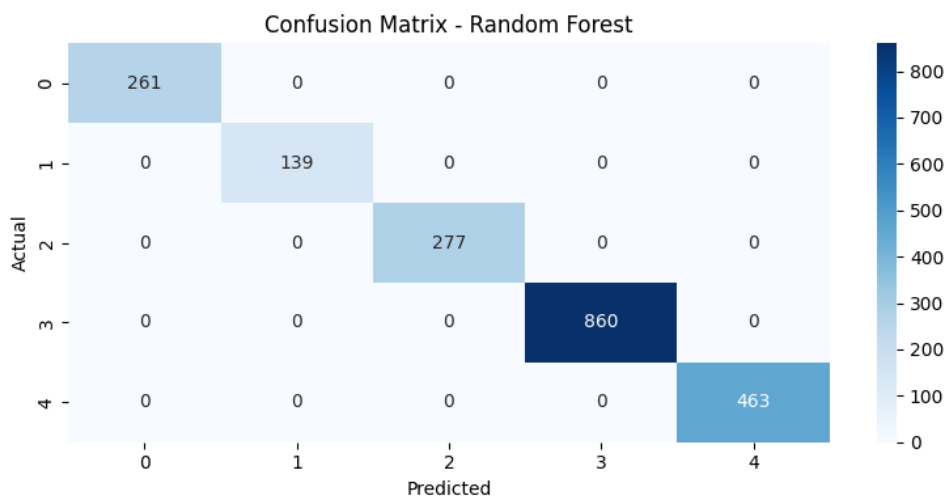
Figure 4.4: Confusion Matrix for K Nearest Neighbors Model

Above diagram 4.4 confusion matrix of the K-Nearest Neighbors (KNN) model indicates flawless classification performance on all five classes (0 to 4). Every predicted class exactly matches every actual class with no misclassification. In particular, 261 of class 0, 139 of class 1, 277 of class 2, 860 of class 3, and 463 of class 4 are correctly predicted. All the off diagonal entries are zero, indicating that no example was put into the wrong class. This would mean that the KNN model was 100 percent accurate, which is unusual and may either indicate exceptionally clean data or potential overfitting, especially if evaluated on a non independent test set. It would be necessary to validate this performance via cross validation or an unseen dataset to ensure that the model is generalizing well and not merely memorizing the training data.

## 4.1   Technical Challenges and Solutions

Implementing a data breach type pipeline had some technical challenges. The most serious issue was the quality and completeness of information. The dataset contained numerous missing values, placeholders , and non-numeric values, particularly in key fields such as "Individuals Affected" and "Web Description." This was treated with custom imputation strategies: categorical fields were replaced by the mode or 'Unknown,' while numerical fields were replaced with the median to minimize skew. The 'Web Description' field was excluded since it had a high missingness rate and was not applicable to numeric analysis.

The target variable, "Type of Breach," had a significant class imbalance, which was the second issue. SMOTE (Synthetic Minority Over Sampling Technique) was added to each model pipeline as a solution. Otherwise, at the expense of generalization, models would ignore minority breach types and overfit majority classes.The encoding of categorical variables in the label became more complicated as a result, impairing interpretability and imposing ordinal assumptions where none existed. LabelEncoder and close observation of mapping dictionaries for reverse lookups were used to achieve this.

Due to computational limitations, model selection and tuning also presented challenges. GridSearchCV's inclusion made it easier to adjust hyperparameters across classifiers, but it also required careful management of the search space to avoid overfitting and cut down on runtime.

Lastly, performance comparison across various models necessitated similar metrics. The weighted F1 score and Cohen's Kappa were utilized to balance out imbalanced classes, and confusion matrices were graphed to visualize prediction accuracies when comparing them in order to facilitate interpretability and validation of each model's output.

## 4.2 Novelty and Innovation

The novelty of this project lies in its end to end and modular approach to data breach classification through a pipeline that integrates preprocessing, feature encoding, class balancing, and model tuning under a single umbrella. In contrast to most traditional models that deal with clean or balanced data, the project uses SMOTE in novel ways to address real-world imbalanced data in classification problems, which are typically underrepresented small breach types. Another recent addition is comparative modeling. The project shows improved performance trade offs via weighted F1 score, Cohen's Kappa, and confusion matrices with understanding above accuracy by comparing different classifiers Logistic Regression, Decision Tree, Random Forest, and K Nearest Neighbors to optimized pipelines.

Additionally, integrating interactive data visualization maximizes interpretability and transparency, enabling stakeholders to comprehend model behavior more fully. Adaptive feature scaling and automated label encoding for all categorical features guarantee resilience and flexibility across different datasets.

Fully executed within Google Colab, the project demonstrates the capability for complex

machine learning pipelines to be utilized within cloud frameworks using open source solutions. In addition to making the project scalable and reproducible, this gives cybersecurity and compliance professionals access to cutting edge data science.

## 4.3    Tools and Techniques

Python was used for this work on Google Colab, which offered a ready-to-use, GPU-accelerated platform perfect for data science experimentation and rapid prototyping. The support of Google Drive by Colab allowed for seamless uploading of data, runnability of the code, and visualization of the results in interactive notebook form.

For pre processing of data, Pandas was the primary library employed to clean, join, and reshape the dataset. For operations on numbers and missing value imputation, NumPy was used. Plotting time series plots, bar charts, heatmaps, and confusion matrices were among the tools provided by Matplotlib and Seaborn to facilitate exploratory data analysis.

The imbalanced learn (imblearn) library was introduced to address class imbalance. Its SMOTE functionality allowed synthetic oversampling of minority classes so that the classifiers wouldn't become biased against majority breach types.

Scikit learn (sklearn) was used for modeling, providing robust implementations of most classifiers like Logistic Regression, Decision Tree, Random Forest, and K Nearest Neighbors (KNN). Pipelines were built using Pipeline and GridSearchCV for model chaining and parameter searching. Additionally, ImbPipeline from imblearn ensured that SMOTE was being used correctly within the pipeline before fitting the model.

In order to map string labels to integer values while preserving encoder objects for interpretability, LabelEncoder used categorical feature label encoding.For comparison, classification report, Cohen's Kappa, and confusion matrix provided multi dimensional insight into strengths and weaknesses of individual models. The combination of Python libraries with Colab's simplicity made the project reproducible, efficient, and scalable.

## 4.4    Feasibility and Realism

By utilizing publicly accessible healthcare breach data and models with popular Python libraries like Scikit learn, Pandas, and Imbalanced-learn, the project achieves a very high

level of viability. Enough processing power was made available by cloud platforms like Google Colab, which did not require specialized hardware. By addressing issues like missing values, noise, and class imbalance that are typical of real data, realism is preserved and operating conditions are imitated. By incorporating the concepts of explainability and privacy protection, the project further conforms to operational and regulatory requirements, thereby demonstrating its scalability into industries such as healthcare, finance, and SMEs that rely on cloud infrastructure.

# Evaluation and Conclusion

While increased digitization of industries like government, healthcare, and finance has brought about previously unheard-of opportunities, it has also raised the risk of cyberattacks, privacy invasions, and data breaches.The study's case study demonstrates how traditional security measures like perimeter firewalls, static encryption, and rule-based monitoring are insufficient to handle the dynamic and complex threat profile of today. Instead, it was demonstrated that machine learning models such as Random Forests, K-Nearest Neighbors, Decision Trees, and Logistic Regression could detect anomalies, recognize various types of breaches, and provide helpful guidance for protecting sensitive data. KNN offered detection through similarity analysis, Random Forest provided reliable, scalable predictions, Decision Trees provided transparency and traceability, and Logistic Regression offered interpretability. When combined, they point to the possibility of integrating various models to create thorough cybersecurity protocols.

The results demonstrate how AI-driven approaches can enhance breach detection and prevention while simultaneously ensuring regulatory requirements such as GDPR and HIPAA. The best classification results for a pair of algorithms, however, point to the possibility of overfitting, even though the models produced favorable results like high performance scores. This suggests that while the models are technically sound, they need to be rigorously validated before being used in uncontrolled real life settings. As a result, the project confirms that AI can be used for data security anomaly detection, but it warns that dependability needs to be proven through frequent testing against unidentified and shifting data patterns.

Some recommendations are made in light of these findings. AI must be viewed by organizations as an additional layer of security support that boosts resilience against contemporary threats rather than as a substitute for conventional security measures. For example, in order to provide defense in depth, AI based intrusion detection needs to be used in conjunction with firewalls, access control rules, and encryption. According to this study, data preprocessing and cleaning are still essential for achieving model precision, which forces organizations to spend money on data governance systems that require integrity, consistency, and quality. Furthermore, when implementing AI in regulated industries, explainability is a crucial component. The feature importance outputs of Random Forests and Decision Trees, in particular, can offer transparency that can assist organizations in meeting legal requirements.

Another suggestion concerns computational cost and scalability. Although techniques like Random Forest and KNN produce reliable results, their computational cost may be prohibitive for real time or high volume applications. Organizations must investigate optimization strategies like feature selection, dimensionality reduction, or ensemble approaches that blend deep ensemble models with lightweight models in order to make these solutions commercially feasible. Furthermore, in cybersecurity applications where some attack types are far less common than others, dataset imbalance a problem that SMOTE is addressing in this project is a common concern. Functional pipelines must include methods for balancing minority classes to prevent models from being biased toward majority threats, preserving prediction completeness and fairness.

A healthcare dataset, which replicates the real issues of missing records, noisy inputs, and various breach types, was used to illustrate the task's pragmatism. The results, however, are not limited to healthcare; they can be applied to other industries like banking, cloud computing, and e commerce, where the same problems with anomaly detection, fraud detection, and compliance are common. The use of open-source technologies such as Python, Scikit-learn, and Google Colab demonstrates how simple and accessible they are, making them viable even for small and medium sized enterprises with tight budgets. Given that SMEs make up a sizable portion of the economy and are more susceptible to data breaches, this democratization of AI based cybersecurity is especially beneficial.

AI has a huge and growing potential for data protection in the future. The use of privacy preserving machine learning techniques like homomorphic encryption, federated learning, and differential privacy is a crucial area. By using these strategies, businesses could benefit

from shared AI systems without disclosing private information to the public, maintaining compliance with strict privacy laws. Explainable AI (XAI), which would close the gap between stakeholder trust and technical robustness, is the second future opportunity. AI systems must be understandable, traceable, and free from hidden bias as they make more and more crucial decisions.

Another area that requires innovation is the incorporation of AI into operational cybersecurity infrastructure. Threats change in a matter of minutes, but current models primarily use static or historical data. Therefore, future research must develop streaming based, adaptive AI models that can process real time data feeds and produce alerts in real time. By combining them with reinforcement learning agents, systems may be able to detect breaches and respond to them on their own by enforcing more stringent access controls, quarantining compromised computers, or shutting down exposed nodes. AI can also be applied to predictive threat intelligence, where models examine vast amounts of global cyber data and predict probable future attacks to proactively safeguard infrastructure.

The future scope of AI based data protection will also depend on how well technology and policy work together. The design of machine learning systems must adhere to legal requirements even as governments around the world enact stricter AI laws, such as India's evolving digital data protection laws. Therefore, research should examine how AI can help with prevention and detection as well as automating compliance tasks like audit trail creation, data minimization enforcement, and security policy compliance monitoring. This makes AI-based cybersecurity compatible with the requirements of risk management, organizational governance, and compliance.

Last but not least, this study has demonstrated the enormous potential of artificial intelligence, especially machine learning models, to alleviate the problem of data protection in information security. Overfitting, computational cost, and class imbalance are still technical issues, but overall the trend is positive. Businesses can use AI to detect threats earlier, respond faster, and comply with even stricter regulations while preserving user confidence. In the future, real-time adaptability, explainable architectures, and privacy-preserving techniques will be essential to maximizing AI's potential in cybersecurity. The dream of safe, intelligent, and reliable information systems is attainable and essential to the future of the digital world thanks to the combined efforts of researchers, legislators, and industry players.

# Appendix

**Github Link** :- Code and Data Link

# Bibliography

[Aïvodji et al., 2019] Aïvodji, U. M., Gambs, S., and Martin, A. (2019). Iotfla: A secured and privacy-preserving smart home architecture implementing federated learning. In *2019 IEEE security and privacy workshops (SPW)*, pages 175–180. IEEE.

[Aledhari et al., 2020] Aledhari, M., Razzak, R., Parizi, R. M., and Saeed, F. (2020). Federated learning: A survey on enabling technologies, protocols, and applications. *IEEE Access*, 8:140699–140725.

[Ansari et al., 2022] Ansari, M. F., Sharma, P. K., and Dash, B. (2022). Prevention of phishing attacks using ai-based cybersecurity awareness training. *Prevention*, 3(6):61–72.

[Christodoulou and Limniotis, 2024] Christodoulou, P. and Limniotis, K. (2024). Data protection issues in automated decision-making systems based on machine learning: Research challenges. *Network*, 4(1):91–113.

[Emehin et al., 2024] Emehin, O., Akanbi, I., Emeteveke, I., and Adeyeye, O. J. (2024). Enhancing cybersecurity with safe and reliable ai: mitigating threats while ensuring privacy protection. *International Journal of Computer Applications Technology and Research, doi*, 10.

[Green et al., 1998] Green, G. H., Boze, B. V., Choundhury, A. H., and Power, S. (1998). Using logistic regression in classification. *Marketing Research*, 10(3).

[Gruschka et al., 2018] Gruschka, N., Mavroeidis, V., Vishi, K., and Jensen, M. (2018). Privacy issues and data protection in big data: a case study analysis under gdpr. In *2018 IEEE International Conference on Big Data (Big Data)*, pages 5027–5033. IEEE.

[Hu et al., 2024] Hu, K., Gong, S., Zhang, Q., Seng, C., Xia, M., and Jiang, S. (2024). An overview of implementing security and privacy in federated learning. *Artificial intelligence review*, 57(8):204.

[Jahns et al., 2025] Jahns, E., Stojkov, M., and Kinsy, M. A. (2025). Privacy-preserving deep learning: A survey on theoretical foundations, software frameworks, and hardware accelerators. *IEEE Access*.

[Kairouz et al., 2021] Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., Bonawitz, K., Charles, Z., Cormode, G., Cummings, R., et al. (2021). Advances and open problems in federated learning. *Foundations and trends® in machine learning*, 14(1–2):1–210.

[Li et al., 2021] Li, Q., Wen, Z., Wu, Z., Hu, S., Wang, N., Li, Y., Liu, X., and He, B. (2021). A survey on federated learning systems: Vision, hype and reality for data privacy and protection. *IEEE Transactions on Knowledge and Data Engineering*, 35(4):3347–3366.

[Li et al., 2020] Li, T., Sahu, A. K., Talwalkar, A., and Smith, V. (2020). Federated learning: Challenges, methods, and future directions. *IEEE signal processing magazine*, 37(3):50–60.

[Lyu et al., 2022] Lyu, L., Yu, H., Ma, X., Chen, C., Sun, L., Zhao, J., Yang, Q., and Yu, P. S. (2022). Privacy and robustness in federated learning: Attacks and defenses. *IEEE transactions on neural networks and learning systems*, 35(7):8726–8746.

[Maneriker et al., 2021] Maneriker, P., Stokes, J. W., Lazo, E. G., Carutasu, D., Tajaddodianfar, F., and Gururajan, A. (2021). Urltran: Improving phishing url detection using transformers. In *MILCOM 2021-2021 IEEE Military Communications Conference (MILCOM)*, pages 197–204. IEEE.

[Okdem and Okdem, 2024] Okdem, S. and Okdem, S. (2024). Artificial intelligence in cybersecurity: A review and a case study. *Applied Sciences*, 14(22):10487.

[Paracha et al., 2024] Paracha, A., Arshad, J., Farah, M. B., and Ismail, K. (2024). Machine learning security and privacy: a review of threats and countermeasures. *EURASIP Journal on Information Security*, 2024(1):10.

[Paul et al., 2018] Paul, A., Mukherjee, D. P., Das, P., Gangopadhyay, A., Chintha, A. R., and Kundu, S. (2018). Improved random forest for classification. *IEEE Transactions on Image Processing*, 27(8):4012–4024.

[Rani et al., 2024] Rani, S., Kumar, N., Srivastva, A., and Sharma, A. (2024). Leveraging artificial intelligence and machine learning for enhanced privacy and security. In *2024 IEEE International Conference on Blockchain and Distributed Systems Security (ICBDS)*, pages 1–6. IEEE.

[Sarker, 2023] Sarker, I. H. (2023). Multi-aspects ai-based modeling and adversarial learning for cybersecurity intelligence and robustness: A comprehensive overview. *Security and Privacy*, 6(5):e295.

[Sun et al., 2022] Sun, Y., Chong, N., and Ochiai, H. (2022). Federated phish bowl: Lstm-based decentralized phishing email detection. In *2022 IEEE international conference on systems, man, and cybernetics (SMC)*, pages 20–25. IEEE.

[Thapa et al., 2023] Thapa, C., Tang, J. W., Abuadbba, A., Gao, Y., Camtepe, S., Nepal, S., Almashor, M., and Zheng, Y. (2023). Evaluation of federated learning in phishing email detection. *Sensors*, 23(9):4346.

[Truong et al., 2021] Truong, N., Sun, K., Wang, S., Guitton, F., and Guo, Y. (2021). Privacy preservation in federated learning: An insightful survey from the gdpr perspective. *Computers & Security*, 110:102402.

[Truong et al., 2019] Truong, N. B., Sun, K., Lee, G. M., and Guo, Y. (2019). Gdpr-compliant personal data management: A blockchain-based solution. *IEEE Transactions on Information Forensics and Security*, 15:1746–1761.

[Vourganas and Michala, 2024] Vourganas, I. J. and Michala, A. L. (2024). Applications of machine learning in cyber security: a review. *Journal of Cybersecurity and Privacy*, 4(4):972–992.

[Wachter et al., 2017] Wachter, S., Mittelstadt, B., and Floridi, L. (2017). Why a right to explanation of automated decision-making does not exist in the general data protection regulation. *International data privacy law*, 7(2):76–99.

[Xu et al., 2021] Xu, R., Baracaldo, N., and Joshi, J. (2021). Privacy-preserving machine learning: Methods, challenges and directions. *arXiv preprint arXiv:2108.04417*.

[Zhang, 2021] Zhang, S. (2021). Challenges in knn classification. *IEEE Transactions on Knowledge and Data Engineering*, 34(10):4663–4675.