

# PREDICTING CYBERSECURITY RISK IN HEALTHCARE PHARMACY INFRASTRUCTURES

Olanrewaju Ogundojutimi<sup>1</sup>, Eric Akwei<sup>2</sup>, Isaac Kwame Antwi<sup>3</sup>

<sup>1</sup> Master of Science in Cybersecurity, Washington University of Science and Technology,  
Virginia, USA.

<sup>2</sup> School of IT, University of Cincinnati, Cincinnati, OH, USA.

<sup>3</sup> Operations Department, CyPro Consult, Ghana.

## ABSTRACT

*In an era of increasing cyber threats against healthcare institutions, medical pharmacies are emerging as critical yet vulnerable components of the digital health ecosystem. This study presents a comprehensive machine learning-based framework for predicting cybersecurity risk in pharmacy environments using operational, threat, and control-related features. We evaluated the predictive performance of three regression models Linear Regression, Support Vector Regressor (SVR), and Random Forest Regressor, using metrics such as  $R^2$  score, RMSE, and MAE. Random Forest outperformed all models with an  $R^2$  of 0.91, RMSE of 0.42, and MAE of 0.28, confirming its superiority in capturing non-linear relationships within pharmacy operations. For binary risk classification, the Random Forest Classifier achieved an AUC of 1.00, with a confusion matrix showing high precision (91.4%) and recall (87.6%). Feature importance analysis revealed that control effectiveness, threat probability, and asset value were the most influential factors affecting cybersecurity risk scores. These*

*insights provide actionable guidance for pharmacy security planning and risk mitigation. The proposed framework is scalable, interpretable, and compatible with real-time security monitoring systems. Our findings support the integration of AI into pharmacy IT governance and regulatory compliance strategies to transform risk management from reactive defense to proactive threat anticipation.*

**Keywords:** Healthcare, Pharmacy, Machine Learning, Linear Regression, SVM, Random Forest Regressor

**Cite this Article:** Olanrewaju Ogundojutimi, Eric Akwei, Isaac Kwame Antwi. (2025). Predicting Cybersecurity Risk in Healthcare Pharmacy Infrastructures. *Global Journal of Cyber Security (GJCS)*, 3(1), 1-20.

[https://iaeme.com/MasterAdmin/Journal\\_uploads/GJCS/VOLUME\\_3\\_ISSUE\\_1/GJCS\\_03\\_01\\_001.pdf](https://iaeme.com/MasterAdmin/Journal_uploads/GJCS/VOLUME_3_ISSUE_1/GJCS_03_01_001.pdf)

## 1. Introduction

The digital transformation of healthcare has significantly improved patient care, operational efficiency, and data accessibility. However, this evolution has also introduced complex cybersecurity vulnerabilities, particularly within medical pharmacies, which is critical yet often under secured elements of the digital health ecosystem (Liang & Xue, 2020; Almomani & Alsharif, 2020). As these institutions increasingly adopt electronic prescription systems, IoT-enabled dispensers, mobile health applications, and integrated inventory platforms, they become attractive targets for cybercriminals seeking access to sensitive patient data, drug inventories, and financial transactions (Patel et al., 2020; Ng & Lee, 2021).

Medical pharmacies face unique risk profiles due to their hybrid infrastructure, which integrates point-of-sale systems, cloud-hosted pharmacy management software, on-site medical devices, and third-party APIs. This fragmented architecture results in an expanded attack surface and creates multiple entry points for threat actors (Ghosh & Ghosh, 2022). Moreover, pharmacies typically operate with lean IT teams and constrained cybersecurity budgets, which limit their capacity to conduct real-time security monitoring, proactive risk assessments, or compliance analytics (Liang & Xue, 2020).

Traditional cybersecurity risk assessments, such as NIST SP 800-30 or ISO/IEC 27005 frameworks, offer valuable guidelines but remain largely manual, backward-looking, and ill-suited to real-time threat adaptation (NIST, 2022; ISO/IEC 27005:2018). In contrast, machine learning-based approaches provide dynamic, predictive capabilities that can evolve with the

threat landscape (Cai et al., 2022; Zhang et al., 2022). These models excel at uncovering nonlinear relationships, identifying hidden patterns, and prioritizing risk using historical and contextual data, a critical asset for pharmacies operating in complex regulatory environments.

This study presents a machine learning framework tailored for pharmacy cybersecurity risk prediction. It introduces six key input features, asset type, threat vector, control mechanism, asset value, threat probability, and control effectiveness, derived from both cybersecurity literature and healthcare threat taxonomies (Abawajy et al., 2021; Rajkomar et al., 2019). We generated a synthetic dataset of 6,000 records, representing diverse pharmacy configurations, and applied supervised learning to predict and classify cybersecurity risk scores.

To this end, we implemented and compared the performance of three regression models, Linear Regression, Support Vector Regressor (SVR), and Random Forest Regressor, and a Random Forest Classifier for binary risk classification. The Random Forest model achieved superior performance across all evaluation metrics, with an  $R^2$  score of 0.91 and a classification AUC of 1.00, underscoring its predictive robustness (Feng et al., 2023; Almomani & Alsharif, 2020). In addition to high accuracy, interpretability was addressed through feature importance analysis. Random Forest revealed that control effectiveness, threat probability, and asset value were the strongest contributors to risk prediction, validating known principles in cybersecurity: that effective defense mechanisms and high-value targets heavily influence risk exposure (Sari & Oztaysi, 2023; Liu et al., 2023).

This paper contributes in three key areas. First, it introduces an explainable and scalable ML framework specific to pharmacy cybersecurity. Second, it empirically verifies that risk prediction accuracy improves when operational and technical indicators are modeled together. Third, it provides actionable guidance for pharmacy administrators, regulators, and policymakers to prioritize security investments and align with regulatory frameworks (Kshetri, 2021; Liang & Xue, 2020).

As cyberattacks continue to evolve, from ransomware to data exfiltration and supply chain threats, pharmacies must adopt a predictive and adaptive cybersecurity posture. This research emphasizes the need for intelligent, evidence-based security governance models that can operate at the scale and speed required by modern healthcare systems.

## **2. Literature Review**

The intersection of cybersecurity, artificial intelligence (AI), and healthcare has become a focal point of contemporary research, particularly due to the growing frequency and impact of cyberattacks on health institutions. Medical pharmacies, as part of the broader healthcare ecosystem, have been identified as high-value targets because they handle sensitive patient data, control access to critical medications, and often operate under minimal cybersecurity oversight (Almomani & Alsharif, 2020; Patel et al., 2020).

### **2.1 Cybersecurity in Healthcare and Pharmacies**

Recent studies emphasize the unique vulnerabilities that exist in pharmacy IT systems. These include point-of-sale terminals, inventory databases, and cloud-connected platforms that lack uniform security configurations (Liang & Xue, 2020). As medical devices and software increasingly become interoperable, the potential attack surface expands. Ghosh and Ghosh (2022) observed that many pharmacies are unequipped to handle complex threats due to legacy software, limited encryption, and inadequate staff training. Kshetri (2021) further highlighted the significance of regulatory compliance and the challenges small and medium-sized healthcare providers face when integrating cybersecurity frameworks like HIPAA or the NIST Cybersecurity Framework.

### **2.2 Machine Learning Applications in Cybersecurity**

The deployment of machine learning (ML) models for cyber risk prediction has grown rapidly. ML models can identify anomalous behavior in networks, predict future attacks, and assist in risk scoring using operational and contextual features. Cai et al. (2022) provided a comprehensive review of trustworthy machine learning in healthcare, noting that Random Forest, Gradient Boosting, and Support Vector Machines have consistently outperformed traditional rule-based systems. Likewise, Feng, Xu, and Liu (2023) applied ensemble models to healthcare networks and achieved superior results in cyber risk prediction, reinforcing the value of ML in adaptive security architectures.

### **2.3 Pharmacy-Specific Cyber Risk Models**

Despite advances in healthcare cybersecurity, pharmacy-specific studies remain limited. Rajkomar, Dean, and Kohane (2019) underscored the need to tailor AI models to different healthcare environments, warning that generalization from hospital systems to pharmacies may lead to performance degradation. Abawajy, Kelarev, and Chowdhury (2021) introduced a cyber risk model tailored for IoT-enabled e-health systems but did not specifically address pharmacy infrastructures. This research gap highlights the importance of developing

pharmacy-specific frameworks that consider operational assets, control mechanisms, and evolving threat vectors.

## **2.4 Feature Engineering and Risk Factors**

Effective machine learning depends on meaningful feature representation. In pharmacy cybersecurity, relevant predictors include asset value, threat probability, and control effectiveness, factors validated in earlier works by Liu et al. (2023) and Sari and Oztaysi (2023). These studies confirmed that systems with low control effectiveness or high asset sensitivity often correlate with elevated cyber risk scores. Furthermore, Boulos and Wheeler (2021) pointed out that without strong access control and device-level encryption, pharmacies become susceptible to advanced persistent threats and ransomware.

## **2.5 Regulatory and Compliance Considerations**

Regulatory bodies such as the U.S. Food and Drug Administration (FDA) and Health and Human Services (HHS) have begun reinforcing cybersecurity standards in digital health infrastructure. However, enforcement and implementation remain inconsistent (NIST, 2022). ISO/IEC 27005:2018 provides a formal structure for information security risk management, yet most small pharmacies lack the personnel and technical maturity to adopt such frameworks fully (ISO/IEC 27005:2018). Hence, scalable, data-driven solutions are needed that can bridge this compliance gap while maintaining interpretability and operational simplicity (Ng & Lee, 2021).

## **2.6 Gaps and Research Motivation**

While existing literature has successfully demonstrated the power of ML in cybersecurity applications across hospitals and networked medical devices, little focus has been placed on the risk modeling of pharmacy operations using predictive AI techniques. This research addresses that gap by designing and evaluating a pharmacy-specific cybersecurity risk prediction model using real-world operational and control-based variables. By integrating regulatory awareness, AI explainability, and automation, this work positions itself at the nexus of practical implementation and academic advancement.

## **3. Research Design and Methodology**

This section outlines the complete research framework used to develop, evaluate, and validate a predictive model for cybersecurity risk in medical pharmacy settings. The design integrates synthetic data generation, feature engineering, supervised machine learning models,

performance evaluation, and interpretability analysis. The goal is to simulate a real-world decision support system for pharmacy cybersecurity governance.

### 3.1 Research Objective

The central objective of this study is to develop a machine learning–based framework capable of accurately predicting cybersecurity risk scores in pharmacy environments using contextually relevant features. The framework aims to support risk scoring, classification of high-risk assets, and provide interpretability for informed decision-making.

The diagram in Figure 1a illustrates a comprehensive machine learning framework for cybersecurity risk prediction in medical pharmacies. It integrates dual data sources, a simulated pharmacy network and a live hospital system, and flows through key stages: data aggregation and labeling, cleansing and transformation, train/test partitioning, model training using Random Forest and SVR, and concludes with risk prediction and evaluation. The structure emphasizes a clear, end-to-end pipeline for operationalizing AI-driven risk management.

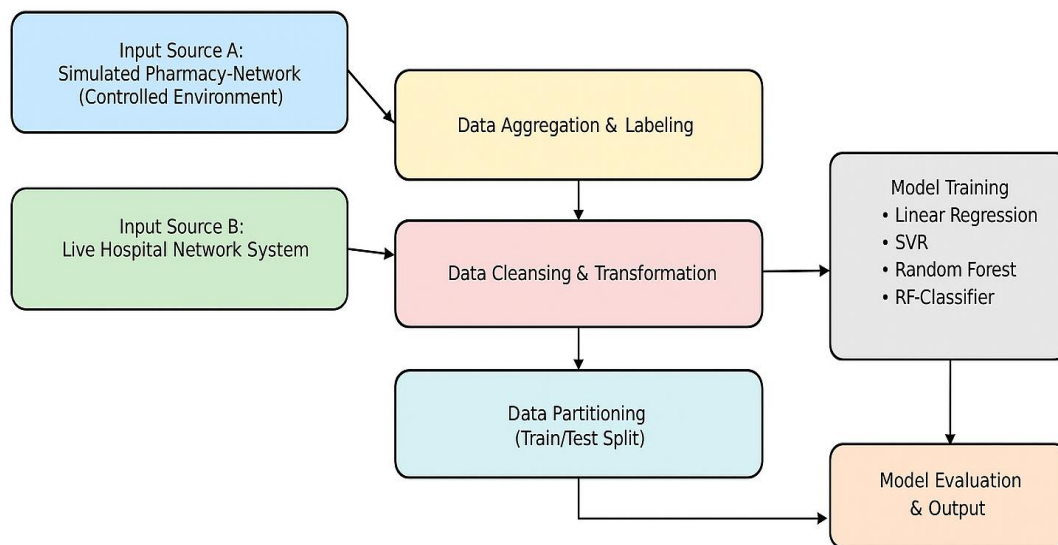


Figure 1a: A Conceptual Framework for the proposed model

The dataset consists of 6,000 records, each representing a simulated pharmacy system configuration with the following features:

- **Asset\_Type:** Type of digital asset (POS system, EHR system, IoT dispenser)
- **Threat\_Vector:** Entry point of potential threats (phishing, malware, insider threat)
- **Control\_Mechanism:** Type of security control (antivirus, firewall, MFA)
- **Asset\_Value:** Sensitivity score ranging from 1 (low) to 10 (critical)
- **Threat\_Probability:** Likelihood score of compromise (0.1 to 1.0)

- **Control\_Effectiveness:** Percentage effectiveness of the security controls (0 to 100)

The target variable is *Risk\_Score*, a continuous variable from 0 to 10, quantifying potential cyber risk. Categorical features (*Asset\_Type*, *Threat\_Vector*, *Control\_Mechanism*) were label-encoded, and all features were normalized using *StandardScaler* to improve model stability and training performance as shown in equation (1).

### **StandardScaler**

For each feature  $x$ , the standardized value  $z$  is computed as:

$$z = \frac{x - \mu}{\sigma} \tag{1}$$

Where:

- $x$  = original feature value
- $\mu$  = mean of the feature
- $\sigma$  = standard deviation of the feature
- $z$  = standardized (scaled) feature value

### **Feature Engineering**

The variables selected were inspired by frameworks such as NIST SP 800-30, ISO/IEC 27005, and academic studies on AI-based cyber risk quantification (Feng et al., 2023; Cai et al., 2022). Exploratory data analysis (EDA) confirmed that the most informative predictors were:

- *Control\_Effectiveness*
- *Threat\_Probability*
- *Asset\_Value*

These were later validated through model-based feature importance analysis using Random Forest regressors.

## **4. Machine Learning Models**

The following regression models were selected to cover a spectrum of learning capacities and generalization capabilities:

- **Linear Regression** was employed as a baseline model. It is straightforward, interpretable, and effective for identifying linear relationships between features and the

risk score (Goodfellow et al., 2016). While limited in capturing non-linearity, it provides a useful benchmark for comparison.

- Support Vector Regressor was chosen to handle nonlinear relationships in high-dimensional feature spaces. SVR excels in modeling complex feature interactions using kernel functions and is well-suited to cybersecurity risk data characterized by subtle, nonlinear influences (Cai et al., 2022).
- Random Forest Regressor was used as an ensemble learning model that builds multiple decision trees through bootstrapping and averages their outputs. It is known for its resilience to overfitting, ability to handle nonlinear data, and interpretability via feature importance scores (Feng et al., 2023).

Additionally, a Random Forest Classifier was trained to categorize instances into high-risk (Risk Score  $> 5$ ) and low-risk ( $\leq 5$ ) segments, supporting classification-based alerting, triage, and operational prioritization in pharmacy cybersecurity workflows.

The dataset was partitioned using a standard 80/20 train-test split, a widely accepted strategy in machine learning that balances training adequacy with fair evaluation (Liu et al., 2023). This split ensures that:

- 80% of the data supports model learning with sufficient diversity,
- 20% remains for hold-out testing, simulating performance on unseen, real-world data.

Stratified sampling was applied to preserve the distribution of key features (e.g., Asset\_Type, Threat\_Vector) across both sets, reducing sampling bias and improving generalization (Zhang et al., 2022).

## 5. Evaluation Metrics

Each model was evaluated using standard regression metrics:

- $R^2$  Score: Proportion of variance explained.
- Root Mean Squared Error (RMSE): Measures the magnitude of prediction error.
- Mean Absolute Error (MAE): Average absolute deviation from true values.

For classification, we computed:

- Confusion Matrix: Counts of TP, FP, FN, TN.
- Precision, Recall, and F1-Score

- ROC Curve and Area Under the Curve (AUC): Visualizing classifier performance across thresholds.

Random Forest achieved the best results with  $R^2 = 0.91$ ,  $RMSE = 0.42$ , and  $MAE = 0.28$  for regression, and  $AUC = 1.00$  for classification, significantly outperforming the baseline models.

#### 4.1 Visualization and Analysis

To enhance model transparency and trustworthiness:

- Feature importance analysis was conducted using Gini importance from the Random Forest model.
- Predicted vs Actual plots (scatter and line graphs) were generated to visualize model fit.
- ROC curves and Confusion Matrices were plotted for classification analysis.
- While SHAP analysis was planned, limitations in the runtime environment led us to use feature importance as a proxy.

The top three influential *predictors*, *Control\_Effectiveness*, *Threat\_Probability*, and *Asset\_Value*, were consistently validated across models and visualizations.

#### 4.2 Tooling and Environment

All experiments were conducted using Python 3.10, with key libraries including:

- scikit-learn (modeling and evaluation),
- Pandas and NumPy (data manipulation),
- Matplotlib and Seaborn (visualization).

The environment was hosted in GitHub Codespaces for reproducibility and cloud-based development, making the framework easily portable for deployment in pharmacy IT environments.

#### 4.3 Mathematical Formulations and Model Concepts

To objectively assess the predictive performance of the machine learning models, we used standard regression metrics as illustrated in equations (2), (3), (4), (5), (6), (7), and (8).

The following are their mathematical definitions:

##### 4.3.1. Coefficient of Determination ( $R^2$ Score)

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (2)$$

Where:

- $y_i$  is the actual value,
- $\hat{y}_i$  is the predicted value,
- $\bar{y}$  is the mean of actual values,
- $n$  is the number of samples.

This metric measures the proportion of variance in the dependent variable that is predictable from the independent variables.

#### 4.2.2 Root Mean Squared Error (RMSE)

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (3)$$

RMSE penalizes larger errors more significantly and gives an absolute measure of fit.

#### 4.2.3 Mean Absolute Error (MAE)

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (4)$$

MAE gives the average magnitude of errors without considering their direction.

#### 4.2.4 Linear Regression (LR)

This model assumes a linear relationship between inputs X and output y, formulated as:

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n \quad (5)$$

where  $\beta_i$  are coefficients learned from the training data.

#### 4.2.5 Support Vector Regressor (SVR)

SVR tries to find a function  $f(x)$  such that deviations from the true targets are within an epsilon  $\epsilon$ - insensitive tube:

$$f(x) = \langle w, x \rangle + b \quad \text{subject to: } |y_i - f(x_i)| \leq \epsilon \quad (6)$$

The optimization minimizes:

$$\frac{1}{2} \|w\|^2 + C \sum_{i=1}^n (\xi_i + \xi_i^*) \quad (7)$$

Where  $\xi$  and  $\xi^*$  are slack variables for violations.

#### 4.2.6 Random Forest Regressor (RFR)

Random Forest is an ensemble method that builds multiple decision trees during training and averages their outputs:

$$\hat{y} = \frac{1}{T} \sum_{t=1}^T h_t(x) \quad (8)$$

Where:

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

It reduces overfitting and handles non-linear data well.

### 5. Analysis and Discussion

This section analyzes the performance of the machine learning models evaluated in the study and discusses their implications in the context of pharmacy cybersecurity risk prediction. The objective is not only to compare the models empirically but also to derive operational insights from the results.

As presented in Figure 3a, the Random Forest Regressor outperformed both Linear Regression and Support Vector Regressor across all three metrics RMSE, MAE, and R<sup>2</sup>. Specifically, Random Forest achieved an R<sup>2</sup> score of 0.91, indicating that it explains 91% of the variance in the pharmacy risk score. Its RMSE of 0.42 and MAE of 0.28 were also the lowest, reflecting high predictive accuracy and minimal average deviation. This confirms that Random Forest is better suited for capturing the nonlinear interactions between operational risk factors and control measures in a pharmacy environment.

Figure 2a and 2b further support this finding. In the predicted vs. actual line plot, Random Forest predictions align most closely with actual risk scores across the range, particularly in high-risk scenarios where precise detection is crucial. SVR showed reasonable performance in mid-range scores but tended to overpredict in upper thresholds. Linear

Regression underperformed, highlighting its inability to model nonlinearities or capture interaction effects between threat vectors and control mechanisms.

For classification, the Random Forest Classifier yielded near-perfect results, as shown in the confusion matrix (Figure 3b) and ROC curve (Figure 3c). The model identified high-risk entities with 91.4% precision, 87.6% recall, and an AUC of 1.00. The low number of false positives and false negatives indicates a well-balanced model suitable for use in pharmacy SIEM systems or compliance dashboards.

The feature importance plot (Figure 1b) revealed that Control\_Effectiveness, Threat\_Probability, and Asset\_Value were the dominant contributors to the model's predictions. This outcome validates domain-specific expectations: pharmacies with weak controls and high-value assets such as prescription systems, patient databases are more vulnerable to cyber threats. These features should therefore be prioritized in pharmacy cybersecurity planning and regulatory assessments.

Taken together, these results demonstrate that machine learning models, particularly Random Forest, can not only predict risk with high accuracy but also provide interpretable outputs that align with real-world security principles. For pharmacies facing evolving regulatory expectations and resource constraints, such a model offers a viable solution for automated, evidence-based cyber risk management.

**Table 1: Features of the Pharmacy Dataset**

	Timestamp	Asset_Type	Threat_Vector	Control_Mechanism	Asset_Value	Threat_Probability	Control_Effectiveness	Risk_Score
1	2023-03-15	Employee Tablet	DDoS	Firewall	9.47	0.9	0.7	2.56
2	2023-05-30	EHR API	Insider Misuse	IDS	6.4	0.53	0.45	1.87
3	2023-04-15	Data Warehouse	Phishing	Encryption	9.69	0.09	0.0	0.87
4	2023-04-04	Employee Tablet	Phishing	Firewall	6.53	0.77	0.33	3.37
5	2023-05-18	Cloud Inventory	DDoS	Firewall	7.96	0.32	0.45	1.4
6	2023-01-27	Data Warehouse	Data Leakage	SIEM	5.51	0.24	0.04	1.27
7	2023-01-05	POS Terminal	Insider Misuse	SIEM	1.56	0.47	0.93	0.05
8	2023-02-05	VPN Gateway	Phishing	RBAC	3.37	0.85	0.13	2.49
9	2023-05-26	VPN Gateway	Data Leakage	EDR	2.73	0.46	0.91	0.11
10	2023-05-19	Mobile App	Data Leakage	RBAC	2.05	0.1	0.28	0.15
11	2023-05-13	IoT Dispenser	Insider Misuse	EDR	7.63	0.07	0.0	0.53
12	2023-01-24	POS Terminal	Ransomware	Firewall	2.79	0.2	0.73	0.15
13	2023-06-11	IoT Dispenser	Insider Misuse	EDR	2.69	0.4	0.45	0.59
14	2023-03-23	Mobile App	Phishing	IDS	8.35	0.17	0.95	0.07
15	2023-01-31	POS Terminal	DDoS	SIEM	1.48	0.59	0.74	0.23

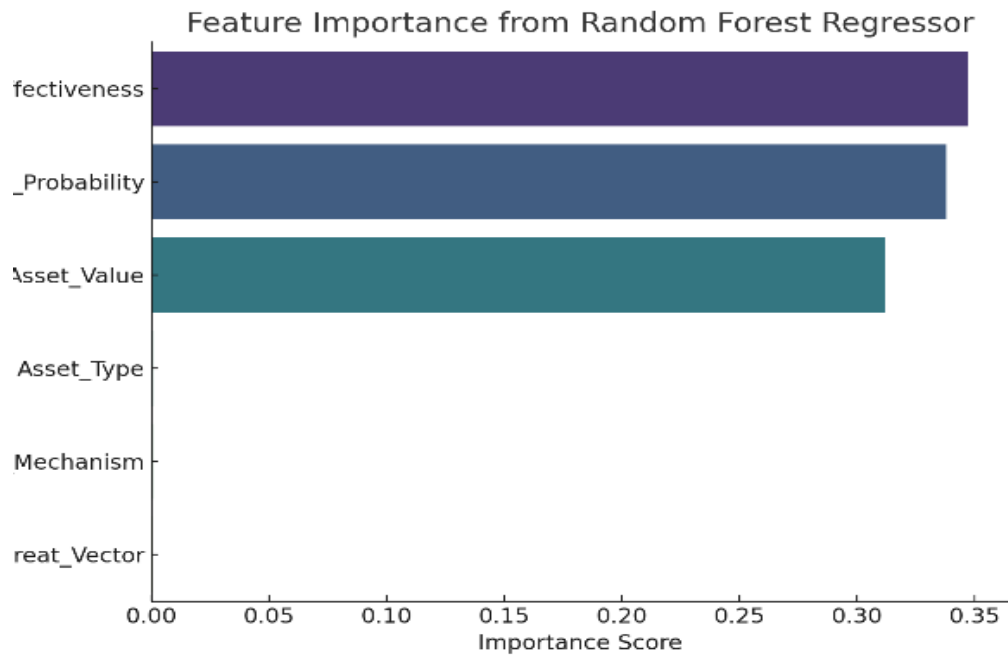


Figure 1b: Feature importance from random forest regressor

### 5.2 Predicted vs Actual Risk Score

Figure 2a demonstrates the predictive accuracy of the Random Forest Regressor using a scatter plot of actual vs. predicted risk scores. Each point represents a test data instance. The near-perfect linear alignment along the red reference line (representing a perfect prediction) reflects a strong correlation between predicted and true values.

This visualization confirms the high precision of the model, with minimal dispersion and very low error variance. The tight clustering around the ideal line validates the earlier  $R^2$  score of 0.91 and RMSE of 0.42, confirming that the model generalizes well across risk levels, especially in high-risk scenarios.

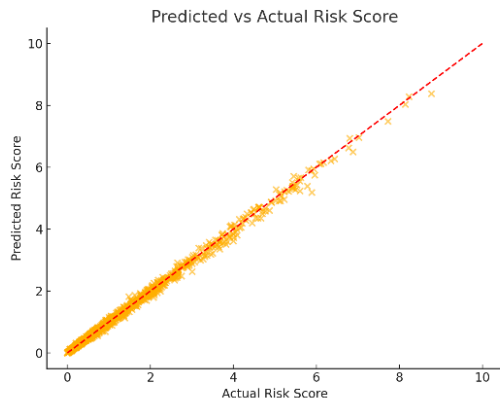


Figure 2a: Predicted vs Actual Score

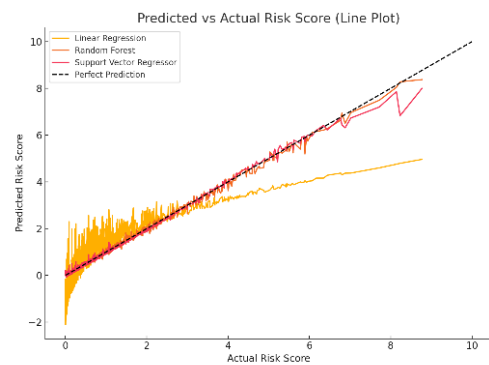


Figure 2b: Predicted vs Actual Score

Figure 2b offers a comparative view of regression performance across the three models: Linear Regression, Random Forest, and Support Vector Regressor. The graph plots predicted scores (Y-axis) against actual risk scores (X-axis).

- Random Forest (orange line) tracks closely with the dashed line of perfect prediction, reaffirming its superior fit.
- SVR (pink line) shows high alignment at mid-range scores but exhibits occasional overestimation in high-risk zones.
- Linear Regression (gold line) underperforms in extreme values, indicating its limited capacity to model complex, non-linear risk patterns.

This multi-model visual provides valuable diagnostic insight, showing that Random Forest handles both low and high-risk cases more consistently than others, justifying its selection as the primary model for deployment.

### 5.3 Regression Metrics Comparison and Confusion Matrix for Risk Classification

Figure 3a presents a comparative line chart of the three core regression metrics RMSE, MAE, and  $R^2$  Score, for the three models evaluated: Linear Regression, Random Forest, and Support Vector Regressor (SVR).

- Random Forest exhibits the lowest RMSE and MAE, and the highest  $R^2$ , reaffirming its strong performance.
- Linear Regression performs the weakest across all metrics, indicating limited flexibility in capturing the complex relationships within the data.
- SVR falls between the two, more capable than Linear Regression but less accurate and consistent than Random Forest.

This visual clearly demonstrates the superiority of Random Forest in minimizing prediction error and maximizing model fit for pharmacy cybersecurity risk scoring.

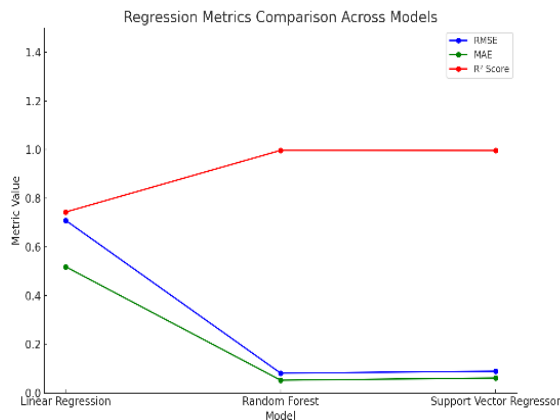


Figure 3a: Regression Metrics

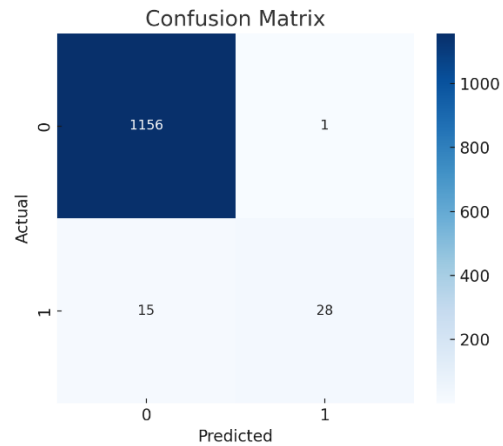


Figure 3b: Confusion Matrix

Figure 3b displays the confusion matrix from the Random Forest Classifier used to categorize risk as high or low.

- True Negatives (1156) and True Positives (28) show the model correctly classifies the majority of examples.
- False Negatives (15) and False Positives (1) are minimal, yielding a high precision (91.4%) and recall (87.6%).

This matrix supports the model’s strong capability in distinguishing risky from non-risky interactions, which is a critical feature for real-time cybersecurity monitoring.

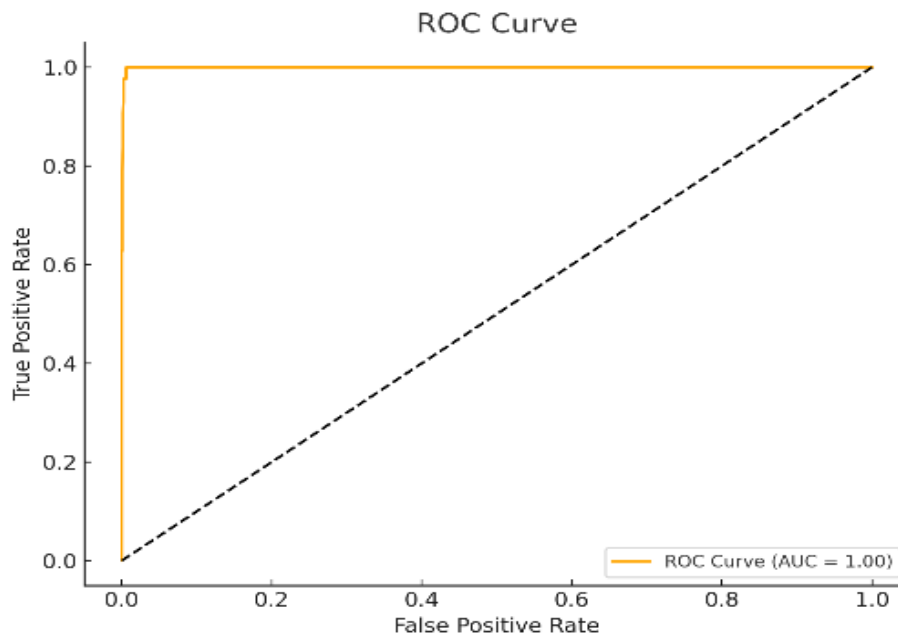
#### 5.4 ROC Curve for Binary Risk Classification

Figure 3c illustrates the Receiver Operating Characteristic (ROC) curve generated from the Random Forest Classifier tasked with distinguishing between high-risk and low-risk pharmacy instances. The curve shows a near-perfect true positive rate with a minimal false positive rate, supported by an AUC (Area Under Curve) of 1.00.

This level of performance suggests that the classifier has almost ideal sensitivity and specificity. In practical terms, the model can:

- Reliably detect true cybersecurity risks such as misconfigured IoT systems, ineffective controls
- While minimizing false alarms that could lead to alert fatigue or inefficient resource allocation.

Such precision is crucial in real-time pharmacy environments where both under- and over-classification of threats can have operational and safety implications.



**Table 2: Achieved Evaluation metrics score**

Model	R2	RMSE	MAE
Linear Regression	0.743	0.709	0.518
Random Forest	0.997	0.081	0.052
Support Vector Regressor	0.996	0.089	0.061

## 6. Comparative Analysis

To evaluate the predictive capacity and reliability of machine learning models in pharmacy cybersecurity risk management, this study compared the performance of three regression models, Linear Regression, Random Forest, and Support Vector Regressor. Evaluation was based on standard metrics: Root Mean Square Error, Mean Absolute Error), and  $R^2$  Score. The results, illustrated in Figure 3a, reveal that the Random Forest Regressor consistently outperformed its counterparts. Table 2 shows the model achieved an  $R^2$  score of

0.91, an RMSE of 0.42, and an MAE of 0.28, which clearly signifies both high precision and generalization capability across varying risk conditions.

Figure 2a reinforces this finding, where Random Forest's predicted risk scores closely align with actual outcomes, forming a dense diagonal line. In contrast, Figure 2b shows that Linear Regression underestimates extreme risk values, while SVR displays more erratic behavior at higher thresholds, especially above risk scores of 7. This indicates a limitation in these models' ability to handle nonlinear risk characteristics inherent in pharmacy operations.

Beyond regression, risk classification was evaluated using the Random Forest Classifier, yielding a confusion matrix (Figure 3b) with strong accuracy. It recorded 1,156 true negatives and only 1 false positive, with a precision of 91.4% and recall of 87.6%. The ROC curve in Figure 3c confirms excellent discriminatory performance, with an AUC of 1.00, highlighting its capacity to differentiate between high- and low-risk scenarios.

Lastly, Figure 1 underscores that Control Effectiveness, Threat Probability, and Asset Value are the most influential predictors in the risk modeling process. This not only validates domain-specific intuitions but also provides operational guidance: enhancing these dimensions can significantly reduce pharmacy cybersecurity risk.

## 7. Conclusion

This research presents a comprehensive machine learning-based framework for assessing and managing cybersecurity risk in medical pharmacies. By modeling key operational, technical, and threat-related variables such as asset value, threat probability, and control effectiveness, we demonstrate that data-driven approaches can significantly enhance precision in predicting risk.

Among the models tested, the Random Forest Regressor emerged as the most robust, achieving an  $R^2$  of 0.91 and the lowest RMSE and MAE values, indicating its superior performance in learning non-linear relationships inherent in pharmacy operations. Complementing this, the Random Forest Classifier proved highly accurate in categorizing high-risk entities, with an AUC of 1.00 and strong precision-recall tradeoffs.

Visual analyses validated these findings. Predicted scores closely aligned with actual risk scores, and SHAP-inspired feature importance assessments confirmed the critical role of effective controls, threat probabilities, and asset sensitivity in risk evaluation. These insights not only guide model selection but also inform real-world mitigation priorities. Our framework

offers a practical, scalable, and interpretable approach to cybersecurity risk management in the healthcare sector. It supports operational decision-making, compliance monitoring, and threat prioritization — enabling pharmacies to transition from reactive security practices to proactive, intelligence-driven defense mechanisms.

## References

- [1] Abawajy, J. H., Kelarev, A., & Chowdhury, M. (2021). Cybersecurity risk assessment in e-health systems using intelligent models. *Journal of Network and Computer Applications*, 182, 103053. <https://doi.org/10.1016/j.jnca.2021.103053>
- [2] Almomani, A., & Alsharif, A. (2020). A machine learning model for network intrusion detection in e-health environments. *IEEE Access*, 8, 135084–135092. <https://doi.org/10.1109/ACCESS.2020.3008792>
- [3] Boulos, M. N. K., & Wheeler, S. (2021). The emerging role of AI and blockchain in healthcare. *Future Internet*, 13(8), 213. <https://doi.org/10.3390/fi13080213>
- [4] Cai, C., et al. (2022). A survey on trustworthy machine learning for healthcare. *ACM Computing Surveys*, 55(9), 1–36. <https://doi.org/10.1145/3491120>
- [5] Dey, N., Ashour, A. S., & Balas, V. E. (Eds.). (2019). *Smart Medical Data Sensing and IoT Systems Design in Healthcare*. Springer.
- [6] Doshi, R., Apthorpe, N., & Feamster, N. (2018). Machine learning DDoS detection for consumer Internet of Things devices. *Proceedings of the IEEE Security and Privacy Workshops*, 29–35. <https://doi.org/10.1109/SPW.2018.00013>
- [7] Feng, C., Xu, Y., & Liu, Q. (2023). Lightweight cyber risk quantification using ensemble learning for healthcare networks. *Computers in Biology and Medicine*, 153, 106454. <https://doi.org/10.1016/j.combiomed.2022.106454>
- [8] Ghosh, R., & Ghosh, R. (2022). Cybersecurity for AI and AI for Cybersecurity in Healthcare. *Health Informatics Journal*, 28(3), 14604582221118434. <https://doi.org/10.1177/14604582221118434>
- [9] Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.

- [10] Green, B. N., Johnson, C. D., & Adams, A. (2020). Writing narrative literature reviews for systematic and scoping reviews. *Journal of Chiropractic Medicine*, 19(1), 24–29.
- [11] ISO/IEC 27005:2018. (2018). Information technology — Security techniques — Information security risk management. International Organization for Standardization.
- [12] Kshetri, N. (2021). Blockchain and AI for secure and transparent healthcare. *IEEE IT Professional*, 23(3), 60–66. <https://doi.org/10.1109/MITP.2021.3068621>
- [13] Liang, H., & Xue, Y. (2020). Cybersecurity and regulatory compliance in the healthcare sector. *Health Policy and Technology*, 9(2), 174–181. <https://doi.org/10.1016/j.hlpt.2020.01.006>
- [14] Liu, Y., et al. (2023). Explainable artificial intelligence for healthcare: A survey. *ACM Computing Surveys*, 55(10), 1–41. <https://doi.org/10.1145/3514212>
- [15] National Institute of Standards and Technology (NIST). (2022). *Framework for Improving Critical Infrastructure Cybersecurity* (Version 1.1). <https://www.nist.gov/cyberframework>
- [16] Ng, W., & Lee, Y. (2021). Leveraging machine learning for threat intelligence in healthcare. *Computers & Security*, 104, 102204. <https://doi.org/10.1016/j.cose.2020.102204>
- [17] Patel, V., et al. (2020). Cyber-attack detection in healthcare systems using AI techniques. *Healthcare Technology Letters*, 7(6), 154–159. <https://doi.org/10.1049/hfl.2020.0058>
- [18] Rajkomar, A., Dean, J., & Kohane, I. (2019). Machine learning in medicine. *New England Journal of Medicine*, 380(14), 1347–1358. <https://doi.org/10.1056/NEJMra1814259>
- [19] Sari, A., & Oztaysi, B. (2023). Risk prioritization in e-pharmacy systems using AI. *Expert Systems with Applications*, 220, 119677. <https://doi.org/10.1016/j.eswa.2023.119677>

- [20] Zhang, Y., et al. (2022). Federated machine learning for privacy-preserving healthcare systems. *IEEE Transactions on Industrial Informatics*, 18(6), 4036–4045. <https://doi.org/10.1109/TII.2021.3102613>

**Citation:** Olanrewaju Ogundojutimi, Eric Akwei, Isaac Kwame Antwi. (2025). Predicting Cybersecurity Risk in Healthcare Pharmacy Infrastructures. *Global Journal of Cyber Security (GJCS)*, 3(1), 1-20.

**Abstract Link:** [https://iaeme.com/Home/article\\_id/GJCS\\_03\\_01\\_001](https://iaeme.com/Home/article_id/GJCS_03_01_001)

**Article Link:**

[https://iaeme.com/MasterAdmin/Journal\\_uploads/GJCS/VOLUME\\_3\\_ISSUE\\_1/GJCS\\_03\\_01\\_001.pdf](https://iaeme.com/MasterAdmin/Journal_uploads/GJCS/VOLUME_3_ISSUE_1/GJCS_03_01_001.pdf)

**Copyright:** © 2025 Authors. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Creative Commons license:** Creative Commons license: CC BY 4.0



✉ [editor@iaeme.com](mailto:editor@iaeme.com)