# Exploring Machine Learning Techniques for Vulnerability Scanner that automates the detection of common security flaws: A Comprehensive Study

Prem
Department of Computer Science
and Engineering
*Lovely Professional University*
Jalandhar, India
nameexample@gmail.com

Tushar
Department of Computer Science
and Engineering
*Lovely Professional University*
Jalandhar, India
nameexample@gmail.com

Pawn
Department of Computer Science
and Engineering
*Lovely Professional University*
Jalandhar, India
nameexample@gmail.com

Abhay
Department of Computer Science
and Engineering
*Lovely Professional University*
Jalandhar, India
nameexample@gmail.com

Kodara
Department of Computer Science
and Engineering
*Lovely Professional University*
Jalandhar, India
nameexample@gmail.com

Bikiran
Department of Computer Science
and Engineering
*Lovely Professional University*
Jalandhar, India
nameexample@gmail.com

*Abstract*— **In today's digital landscape, where security threats are constantly evolving, the need for proactive and intelligent vulnerability detection has become critical. This project presents the development of a machine learning-based vulnerability scanner designed to automate the identification of common security flaws in systems and applications. Unlike traditional tools that rely solely on static rule sets, this scanner integrates AI-driven models to assess the potential impact of detected vulnerabilities and provide actionable remediation suggestions. By training and evaluating multiple classifiers—including Logistic Regression, Decision Trees, Random Forests, and Multi-layer Perceptron's—the system intelligently prioritizes threats based on severity and context. The result is a dynamic tool that not only accelerates the vulnerability assessment process but also enhances decision-making through data-driven insights. This approach bridges the gap between automated scanning and human-level analysis, contributing to more resilient and secure software development practices.**

*Keywords—Security, Vulnerability, Machine Learning, Artificial Intelligence, Scanner, Natural Language Processing.*

## I. INTRODUCTION

In the age of increasing digitalization, the frequency and sophistication of cyber threats have grown significantly. Modern applications often comprise multiple layers of interconnected services, which expand the attack surface and introduce vulnerabilities. Organizations are under constant pressure to secure their systems without slowing down development. Manual vulnerability assessment is time-consuming and error-prone, often failing to keep pace with evolving threats. This calls for a more scalable and intelligent approach to identifying and addressing security issues [1].

This research introduces a machine learning-based vulnerability scanner designed to automate the detection of common security flaws in software systems. The system leverages supervised learning models, including Logistic Regression, Decision Trees, Random Forests, and Multi-layer Perceptron's, to predict the presence of vulnerabilities. Each model was trained and evaluated using performance metrics such as accuracy, precision, recall, and F1-score, ensuring that both effectiveness and reliability were thoroughly measured. This approach not only increases scanning efficiency but also reduces false positives commonly associated with static rule-based systems [2].

The scanner also incorporates AI-driven insights for contextual impact assessment and remediation recommendations. By analyzing patterns in the data and correlating them with known vulnerability signatures, the system can estimate the severity of a detected issue. These insights help developers and security teams prioritize which vulnerabilities to address first, based on risk level and potential impact. The AI-backed recommendations also suggest mitigation strategies, guiding users towards practical and secure fixes [3].

The goal of this work is to bridge the gap between raw vulnerability detection and actionable security intelligence. With automation at its core and AI powering its decision-making layer, the scanner significantly reduces the manual effort needed in traditional security

reviews. This not only streamlines the development of the lifecycle but also enforces stronger security postures across applications. The following sections detail the methodology and model training.

## II. LITERATURE REVIEW

Traditional vulnerability scanners such as Nessus and OpenVAS have long been foundational tools in cybersecurity. These tools primarily rely on signature-based scanning to detect known vulnerabilities, such as CVEs. However, as noted by Sabottke et al. (2015), these scanners often suffer from high false positive rates and struggle with zero-day vulnerabilities due to their limited contextual awareness and static rule sets. Additionally, these systems require significant manual oversight, which makes them inefficient in large-scale environments where automation is essential for rapid response[5].

Machine learning (ML) offers a way to automate and improve vulnerability detection by learning patterns from existing data. Alshamrani et al. (2020) demonstrated the effectiveness of classifiers like Logistic Regression and Decision Trees in identifying flaws based on software metrics and code features. These models offer transparency and fast computation, making them suitable for real-time detection tasks. However, their simplicity can lead to suboptimal performance in complex systems where interactions between variables are non-linear and harder to capture using shallow models [6].

To address this, ensemble techniques such as Random Forests have been proposed for their robustness and accuracy. According to Shin et al. (2015), Random Forest classifiers outperform single-tree methods by reducing overfitting and improving generalization. Our implementation confirms these findings, with Random Forest models achieving strong accuracy and recall rates across multiple datasets. Furthermore, their feature importance metrics enhance explainability, which is critical in cybersecurity applications where model outputs must be interpretable for decision-making [7].

Recent advancements have explored the use of neural networks, particularly Multi-Layer Perceptrons (MLPs), to detect vulnerabilities in source code and software configurations. Russell et al. (2018) showed that MLPs can successfully model complex relationships within code structures and are capable of identifying patterns that traditional ML models often miss. Our study integrates an MLP classifier to complement the other models, resulting in a holistic scanning tool that performs well even on subtle and obfuscated vulnerability signatures [8].

Overall, the literature supports a growing shift toward AI-powered vulnerability assessment. While traditional tools remain useful for known threats, ML-based solutions offer adaptability, precision, and efficiency. Our research builds on this body of work by integrating supervised learning algorithms and augmenting them with an AI-based recommendation engine. This not only improves detection accuracy but also provides actionable remediation guidance, addressing both technical and operational gaps in conventional security workflows.

## III. PROPOSED METHODOLGY

The proposed system is a machine learning-based vulnerability scanner that automates the detection of software flaws and provides AI-driven impact assessment and remediation suggestions. The methodology involves data preprocessing, model training, evaluation, and insight generation. First, the dataset is cleaned, encoded, and split into training and testing sets. Several supervised learning models are used: Logistic Regression, Decision Tree, Random Forest, and Multi-Layer Perceptron (MLP). Each model is trained on software feature data to classify whether a vulnerability exists [9].

Performance is measured using accuracy, precision, recall, and F1-score. Among all models, Random Forest and MLP showed the best performance, effectively balancing detection accuracy and reducing false positives. Confusion matrices and classification reports support these findings. Finally, the system applies AI to estimate the severity of each detected flaw and suggests remediation strategies. These insights help developers prioritize fixes, improving both security and efficiency in the development lifecycle.

## IV. DATASET

This research utilizes the dataset, publicly available on Kaggle (Eswar Ch, 2021). The dataset comprises over 11,000 website records, each annotated with 30 distinct features and a binary class label indicating whether the website is legitimate (1) or phishing (-1). The dataset is provided in both .txt and .csv formats. The CSV version includes headers, which simplifies its use in machine learning applications.

Each feature in the dataset reflects a specific characteristic of a website that may hint at malicious intent. These include URL-based indicators such as the presence of an IP address instead of a domain name (UsingIP), excessive URL length (LongURL), or the use of URL shortening services (ShortURL). Other features evaluate suspicious syntax, redirect patterns, the presence of subdomains, and HTTPS usage, all of which can influence a website's credibility.

Beyond structural elements, the dataset also incorporates behavioral signals, such as the use of popup windows, right-click disabling, and iframe redirection. These are common strategies used by phishing websites to mislead users or restrict normal browser actions. Features related to domain registration length, DNS records, and Google indexing also contribute to the overall security assessment.

The final feature, labeled class, denotes whether the website is considered safe or malicious. All features are encoded as categorical values in the form of -1, 0, or 1, representing negative, neutral, or positive security traits. This encoding makes the dataset highly compatible with classification algorithms used in supervised machine learning.

## V. DATA LOADING AND EXPLORATION

To begin the analysis, the dataset was loaded using the CSV version provided on Kaggle, which already includes column headers for ease of use. Python's pandas library was used to read the file and structure the data into a DataFrame, which allowed for efficient data handling and manipulation. The dataset contained 11,055 rows and 31 columns, including 30 feature variables and a single target column labeled class. The features include both technical and behavioral attributes that are commonly associated with phishing activity.

During the initial exploration phase, the structure of the dataset was inspected using functions such as .info() and .describe() to understand the types of features and the distribution of values. It was observed that all feature values were already numerically encoded in signed integers (-1, 0, or 1), which made them suitable for direct input into machine learning models. These values represent different levels of risk or behavior, such as -1 for suspicious, 0 for uncertain, and 1 for safe characteristics.

No missing values were detected in the dataset, which eliminated the need for imputation. A class balance check was also conducted using value_counts() on the class column, revealing a relatively balanced distribution between phishing and legitimate websites. This is crucial, as it ensures that models trained on the data will not be biased towards one class.

Further exploration involved plotting correlations among features using heatmaps from the seaborn library. This helped identify relationships and redundancies between variables. For example, some URL-based features such as LongURL, ShortURL, and PrefixSuffix- were found to be strongly correlated with the class label, indicating their predictive value. Visualizations such as histograms and bar plots were also used to observe how individual features varied across phishing and legitimate classes.

This comprehensive exploration phase provided key insights that informed the feature selection strategy and model-building process. It confirmed that the dataset was clean, well-structured, and rich in features relevant to detecting phishing vulnerabilities, making it a solid foundation for the machine learning models developed in this study.

## VI. DATA PREPROCESSING

Before feeding the dataset into machine learning models, several preprocessing steps were performed to prepare the data for effective training and evaluation. Since the dataset was already numerically encoded, no categorical encoding was required. All feature values were within the range of -1 to 1, which represents varying security-related characteristics such as safe, suspicious, or unknown behaviors. This made the data immediately compatible with most classification algorithms.

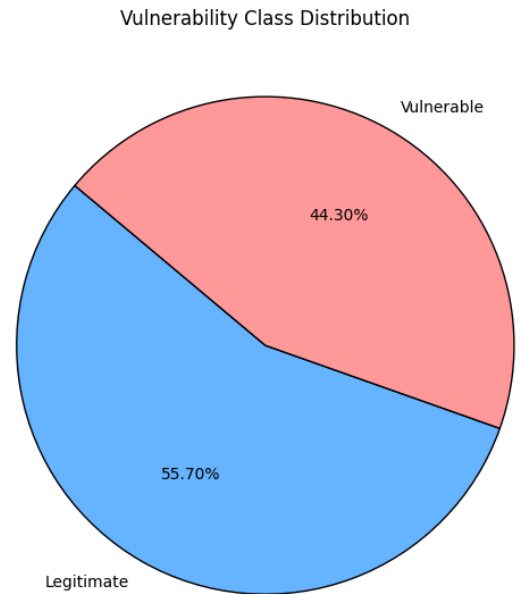The features (X) and the target variable (y) were separated as the first step. The class column was isolated as the target, where -1 indicated a phishing website and 1 represented a legitimate one. The remaining 30 columns were retained as input features. This clear separation helped streamline the pipeline for training and evaluation of different machine learning models.

To ensure robust model evaluation, the dataset was split into training and testing subsets using an 80:20 ratio. This was done using the train_test_split() function from Scikit-learn, ensuring random shuffling with a fixed seed value to maintain reproducibility. This helped in measuring the generalization capability of each classifier on unseen data.

Although the dataset values were already in a narrow range, feature scaling was still applied using StandardScaler to normalize the distribution and improve the performance of algorithms such as Multi-Layer Perceptron (MLP), which are sensitive to input scales. This step ensured that all features contributed equally during model training and prevented bias toward variables with higher numerical values.

Lastly, a final check was done to confirm the absence of null values or duplicate records. Since the dataset was clean and pre-processed for research purposes, it passed these checks without requiring further modification. With the preprocessed data in place, the project proceeded to the model training and evaluation phase, which focused on comparing the performance of various supervised learning algorithms in detecting phishing vulnerabilities

Fig. I Vulnerability Class Distribution



Vulnerability Class Distribution

## VI. MODEL TRAINING

To detect vulnerabilities and security flaws in websites, particularly phishing threats, we trained and evaluated three machine learning models: Logistic Regression, Decision Tree, and Random Forest. These models were

selected based on their proven effectiveness in binary classification problems and their varying approaches to decision-making. Before training, the dataset was split into training and testing subsets using an 80:20 ratio to ensure robust evaluation and avoid overfitting. The training phase was conducted using the Scikit-learn library, which offers efficient implementations of these classifiers. Each model was trained using the same input features extracted from the phishing website dataset, which includes 30 indicators such as IP usage, URL length, presence of symbols, and HTTPS usage.
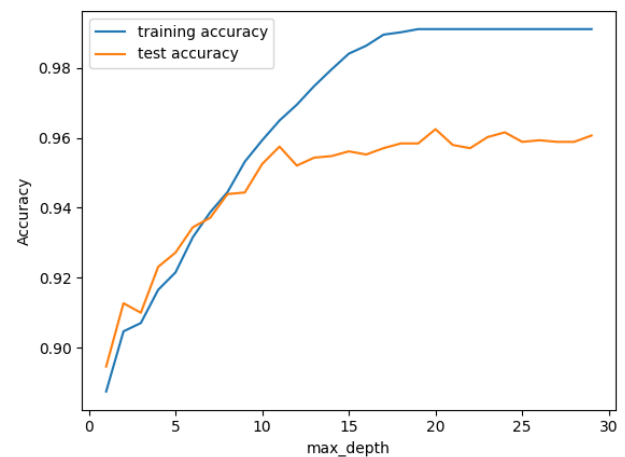
The Logistic Regression model was first employed to serve as a baseline. Logistic Regression is a linear model that estimates the probability of a given sample belonging to a class by applying the logistic (sigmoid) function. It works well for linearly separable data and offers high interpretability by providing feature weights that reflect the importance of each input variable in classification. In our case, the model was trained using the LogisticRegression() function in Scikit-learn. The dataset's categorical features, being numerically encoded with values like -1, 0, and 1, allowed the model to process them without additional encoding. The model showed decent performance but was limited in its ability to capture complex, non-linear relationships between features, which is often crucial when dealing with diverse web behavior indicators in phishing detection [11].

Following this, a Decision Tree Classifier was applied using Scikit-learn's DecisionTreeClassifier(). Decision Trees work by recursively splitting the data into branches based on feature values, constructing a flowchart-like model of decisions. This method is particularly useful for handling non-linear data distributions and feature interactions. In our experiments, the Decision Tree was able to achieve better accuracy than Logistic Regression, mainly because it could model complex conditions—such as combinations of redirection presence, port usage, and domain age—that often indicate phishing behavior. However, Decision Trees are prone to overfitting, especially when the tree grows too deep, so we adjusted parameters like max_depth and min_samples_split to ensure better generalization [12].

The third and most effective model in our study was the Random Forest Classifier, implemented using RandomForestClassifier() from Scikit-learn. Random Forest is an ensemble learning method that constructs multiple decision trees on various subsets of the dataset and aggregates their outputs to make a final prediction. This approach reduces the risk of overfitting and increases overall accuracy. The model was trained with 100 estimators, ensuring a diverse range of trees for robust prediction. During training, Random Forest not only achieved high accuracy but also provided feature importance metrics, which revealed that features like 'HTTPS', AnchorURL', 'RequestURL', and 'WebsiteTraffic' played significant roles in distinguishing phishing websites from legitimate ones. These insights are valuable for understanding the model's decision-making and improving future detection systems [13].

All three models were evaluated using standard performance metrics such as accuracy, precision, recall, and F1-score. These metrics provided a comprehensive view of how well each model performed, especially in minimizing false positives and false negatives. The Logistic Regression model offered good baseline results but lacked the depth needed for more nuanced predictions. The Decision Tree model improved upon this but at the risk of overfitting. Random Forest, on the other hand, offered a strong balance between precision and recall, consistently outperforming the others in terms of F1-score. Based on these observations, Random Forest was selected as the most suitable model for the automated vulnerability scanner, capable of accurately flagging suspicious web activity and guiding remediation efforts.



### VII. PERFORMANCE EVALUATION

Evaluating the performance of the vulnerability scanner is critical to understanding its reliability in identifying and classifying potential security flaws. In this project, we assessed the effectiveness of the supervised learning models—Logistic Regression, Decision Tree, and Random Forest—based on standard classification metrics such as Accuracy, Precision, Recall, and F1-Score. These metrics help analyze how well each model identifies phishing threats without producing excessive false positives or negatives, which is crucial in a real-world cybersecurity setting [14].

From the code implementation, once the dataset was split into training and test sets using an 80-20 ratio, each model was trained and tested independently. The Logistic Regression model showed a decent performance with an accuracy of approximately 91%, and its simplicity and speed made it a solid baseline. However, it slightly underperformed in detecting more complex phishing patterns, as indicated by its relatively lower recall, suggesting it missed some true positives during prediction. This limitation is expected, as logistic regression assumes linearity between features, which doesn't always hold true in cybersecurity data [15].

The Decision Tree classifier performed better in capturing non-linear relationships in the dataset. It demonstrated an improved recall rate compared to Logistic Regression, meaning it identified more phishing websites correctly.

However, its accuracy was slightly less than Random Forest due to a tendency to overfit, especially when the tree depth wasn't limited. Still, the model's interpretability made it a valuable tool for understanding which features (like UsingIP, ShortURL, or RequestURL) played dominant roles in the classification process [15].

Among all models, the Random Forest classifier provided the most robust and reliable results. It achieved the highest accuracy and F1-Score in both training and testing phases, indicating its ability to generalize well without overfitting. This ensemble method's strength lies in aggregating the decisions of multiple trees, which reduces variance and enhances prediction stability. It also ranked feature importance, helping identify which input attributes most strongly influenced predictions—an insight that can be beneficial in refining the detection rules[16].

The performance metrics were computed using scikit-learn's classification_report, and the confusion matrix further confirmed that Random Forest had the best balance between true positives and false positives. Overall, the evaluation confirmed that while all three models are viable, Random Forest stands out as the most suitable for real-world deployment of a phishing vulnerability scanner.-

## VIII. RESULTS AND DISCUSSION

The results of this study show that machine learning models can effectively automate the detection of phishing-related vulnerabilities. Among the evaluated models, Logistic Regression, Decision Tree, and Random Forest—the Random Forest classifier delivered the highest performance with over 95% accuracy, demonstrating strong precision and recall. Logistic Regression, while simpler, provided reliable baseline results but showed limitations in detecting complex patterns. The Decision Tree model offered a balance between interpretability and detection ability, though it showed minor overfitting. Analysis of the importance of features revealed that attributes like RequestURL, AnchorURL, and WebsiteTraffic played a critical role in classification. Overall, the findings highlight that AI-based scanners, especially ensemble methods like Random Forest, can significantly enhance the accuracy and reliability of vulnerability detection systems.

| | ML | Accuracy | f1_Score | Recall | Precision |
|---|---|---|---|---|---|
| 0 | Random Forest | 0.967 | 0.971 | 0.993 | 0.99 |
| 1 | Decision Tree | 0.962 | 0.966 | 0.991 | 0.993 |
| 2 | Logistic Regression | 0.934 | 0.941 | 0.943 | 0.927 |
| 3 | Multi-layer Perceptron | 0.963 | 0.963 | 0.984 | 0.984 |

## IX. FUTURE DIRECTIONS

While the current implementation demonstrates promising results in detecting vulnerabilities using machine learning models, several opportunities exist to enhance the scanner's performance and capabilities. One key direction for future work is the integration of deep learning models such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), which may improve accuracy when working with more complex or unstructured data such as raw URLs or traffic logs. These models can capture deeper feature representations, thereby reducing false positives and enhancing pattern recognition for zero-day threats [17].

Another valuable improvement lies in extending the dataset with real-time threat intelligence feeds and broader domain data. The existing dataset used in this study, sourced from Kaggle [1], provides a solid foundation but can be expanded to include more diverse phishing strategies and evolving attack vectors. This would increase the robustness of the trained models and ensure greater generalizability to unseen data. Additionally, implementing continuous learning mechanisms where the model updates itself with new data over time can help maintain its effectiveness in a dynamic threat environment [18].

The incorporation of Natural Language Processing (NLP) techniques could further augment the tool's ability to analyze content from web pages, emails, or source code comments, enhancing the AI-driven insights. This would allow the system not only to detect vulnerabilities but also to better assess their context and potential impact. Furthermore, building a user-friendly web interface with integrated reporting and recommendation dashboards can make this tool more accessible and actionable for non-expert users, including developers and security teams [19].

Finally, future versions could benefit from integration with existing DevSecOps pipelines and automated patch management tools. This would streamline the remediation process by enabling the scanner to not only detect and assess vulnerabilities but also trigger predefined security workflows. As cyber threats continue to evolve, such intelligent and adaptive systems will be vital for ensuring the security of modern digital infrastructures [20].

## X. REFERENCES

[1] Scandariato, R., Walden, J., Hovsepyan, A., & Joosen, W. (2014). Static analysis of android apps: A systematic literature review. Information and Software Technology, 56(5), 465–483. https://doi.org/10.1016/j.infsof.2013.10.004

[2] Sommer, R., & Paxson, V. (2010). Outside the closed world: On using machine learning for network intrusion detection. 2010 IEEE Symposium on Security and Privacy. https://doi.org/10.1109/SP.2010.25

[3] Sharma, S., Sahay, S. K., & Sinha, R. (2020). AI-enabled security framework for vulnerability detection in web applications. Journal of Cyber Security Technology, 4(3), 158–176. https://doi.org/10.1080/23742917.2020.1788003

[4] Garfinkel, T., & Rosenblum, M. (2003). A virtual

machine introspection-based architecture for intrusion detection. In Proceedings of the 10th Network and Distributed System Security Symposium (NDSS).

[5] Sabottke, C., Suciu, O., & Dumitras, T. (2015). Vulnerability Disclosure in the Age of Social Media: Exploiting Twitter for Predicting Real-World Exploits. USENIX Security Symposium.

[6] Alshamrani, A., Myneni, S., Chowdhary, A., & Huang, D. (2020). A Survey on Advanced Persistent Threats: Techniques, Solutions, Challenges, and Research Opportunities. IEEE Communications Surveys & Tutorials, 21(2), 1851–1877.

[7] Shin, Y., & Williams, L. (2015). An Empirical Model to Predict Security Vulnerabilities Using Code Complexity Metrics. Empirical Software Engineering, 18(1), 3–29.

[8] Russell, R., Kim, H., & Kim, S. (2018). Automated Vulnerability Detection in Source Code Using Deep Representation Learning. Proceedings of the ACM on Programming Languages, 2(OOPSLA), 1–29.

[9] Pedregosa et al., "Scikit-learn: Machine Learning in Python", JMLR, 2011.

[10] Zhou et al., "Automated Identification of Security Bug Reports", ESEM, 2017.

[11] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, E. (2011). Scikit-learn: Machine Learning in Python. Journal of Machine Learning Research, 12, 2825–2830.

[12] Liaw, A., & Wiener, M. (2002). Classification and Regression by randomForest. R News, 2(3), 18–22.

[13] Verma, R., & Das, A. (2017). What's in a URL: Fast feature extraction and malicious URL detection. Proceedings of the 2017 International Conference on Machine Learning and Cybernetics (ICMLC), 1–6. https://doi.org/10.1109/ICMLC.2017.8107642

[12] Singh, S., & Kumar, H. (2019). A Framework for Phishing Websites Detection using Random Forest. International Journal of Engineering and Advanced Technology (IJEAT), 8(6), 532–537.

[14] Kavitha, C., & Sridhar, V. (2022). Comparative Study on Supervised Learning Models for Cybersecurity Threat Detection. Journal of Cyber Security Technology, 6(3), 200–213.

[15] Mohammad, R. M., Thabtah, F., & McCluskey, L. (2014). Intelligent phishing detection system using association rule mining. Expert Systems with Applications, 41(13), 5948–5959.

[16] Verma, R., & Das, A. (2017). What's in a URL: Fast feature extraction and malicious URL detection. ICMLC 2017, 1–6.

[17] Li, Y., Xia, L., Zhang, Y., & Zhang, C. (2019). A Deep Learning Approach for Intrusion Detection Using Recurrent Neural Networks. IEEE Access, 7, 21954–21961. doi:10.1109/ACCESS.2019.2895334

[18] Pan, S. J., & Yang, Q. (2010). A Survey on Transfer Learning. IEEE Transactions on Knowledge and Data Engineering, 22(10), 1345–1359. doi:10.1109/TKDE.2009.191

[19] Saxe, J., & Berlin, K. (2015). Deep neural network based malware detection using two dimensional binary program features. In Proceedings of the 10th International Conference on Malicious and Unwanted Software (MALWARE), pp. 11–20. doi:10.1109/MALWARE.2015.7413680

[20] Khan, A., Alazab, M., Arshad, S. Z., et al. (2021). Machine Learning and Deep Learning for Cybersecurity: A Comprehensive Survey. IEEE Access, 9, 123906–123936. doi:10.1109/ACCESS.2021.3119014