

Real-Time Detection of Phishing Emails Using XG Boost Machine Learning Technique

Jude Osamor*
Department of Cybersecurity and
Networks
Glasgow Caledonian University
Glasgow, G4 0BA Scotland, UK
jude.osamor@gcu.ac.uk

Moses Ashawa
Department of Cybersecurity and
Networks
Glasgow Caledonian University
Glasgow, G4 0BA Scotland, UK
moses.ashawa@gcu.ac.uk

Jackie Riley
Department of Cybersecurity and
Networks
Glasgow Caledonian University
Glasgow, G4 0BA Scotland, UK
j.riley@gcu.ac.uk

Pius Owoh
Department of Cybersecurity and
Networks
Glasgow Caledonian University
Glasgow, G4 0BA Scotland, UK
nsikak.owoh@gcu.ac.uk

Ayodeji Ajibade
Department of Cybersecurity and
Networks
Glasgow Caledonian University
Glasgow, G4 0BA Scotland, UK
aajib300@caledonian.ac.uk

Celestine Iwendi
School of Creative Technologies
University of Bolton
Bolton, BL3 5AB UK
c.iwendi@bolton.ac.uk

Abstract— Phishing attacks continue to pose a significant threat to individuals and organizations, making it crucial to develop effective countermeasures. Machine learning algorithms have shown promise in detecting and mitigating phishing attacks. The study evaluates the performance of four popular algorithms in the context of phishing detection and compares the effectiveness of these four different algorithms; Random Forest, Decision Tree, XGBoost, and Logistic Regression, to determine which one achieves the highest accuracy. The results show that XGBoost outperforms the other algorithms and can accurately detect phishing attacks with a high degree of precision. The algorithms are compared based on factors such as training time, test time, model size, interpretability, and explainability. To compare the effectiveness of these algorithms, the study conducted experiments using a dataset of phishing emails. The algorithms were trained on a labeled dataset and evaluated based on metrics such as accuracy, precision, and recall. The results demonstrate that XGBoost outperforms the other algorithms, achieving the highest accuracy in detecting phishing attacks. The findings of this study have significant implications for the development of antiphishing technologies. By leveraging machine learning algorithms, particularly XGBoost, organizations can enhance their ability to detect and prevent phishing attacks. This can help protect individuals' personal information, passwords, and credit card numbers from falling into the hands of cybercriminals.

Keywords—Phishing, XGBoost, email security, cybersecurity

I. INTRODUCTION

Phishing attacks have become a significant cybersecurity threat, posing severe risks to individuals, organizations, and governments. These attacks involve tricking individuals into providing sensitive information such as passwords, credit card numbers, or social security numbers by impersonating a trusted entity. With the increasing sophistication of phishing techniques, it has become more challenging to detect and prevent these attacks, making it crucial for individuals and organizations to stay vigilant and adopt proactive measures to protect themselves from falling victim to such scams. Phishing attacks are considered one of the most frequent examples of fraud activity on the internet [1].

Some common types of phishing attacks include email phishing, where attackers send deceptive emails to trick recipients into revealing personal information or clicking on malicious links. Another form is spear phishing, which targets specific individuals or organizations with personalized and highly convincing messages. There is also vishing, a phishing technique that involves phone calls or voice messages to deceive victims into sharing sensitive data. Additionally, there is smishing, where attackers use SMS or text messages to trick recipients into providing personal information or downloading malicious content [2].

Email phishing attacks are becoming increasingly sophisticated and difficult to detect using the traditional approach. These attacks often involve the use of highly convincing emails that mimic legitimate messages from trusted sources. They may contain well-crafted language and graphics that make them appear genuine, making it harder for users to identify them as phishing attempts. Additionally, attackers are employing advanced techniques such as social engineering and personalized targeting to increase their chances of success. As a result, relying solely on traditional methods such as email filters and spam detectors is no longer sufficient to protect against these evolving threats. As a result, there is a need for real-time detection techniques using advanced technologies such as machine learning. Machine learning has emerged as a powerful tool in the fight against phishing attacks. By analyzing large amounts of data and identifying patterns, machine learning algorithms can effectively detect and flag suspicious emails in real-time [3]. These algorithms continuously learn and adapt, keeping up with the ever-changing tactics used by attackers. Furthermore, machine learning can also help in identifying and blocking phishing websites, by analyzing various factors such as domain names, website content, and user behaviour.

This proactive approach not only saves users from falling victim to phishing attacks but also helps in preventing the spread of such malicious activities. Machine learning algorithms can quickly analyze and categorize websites based on their risk level, enabling browsers and security systems to automatically block access to potentially harmful sites. This not only protects individuals but also helps in safeguarding

organizations from potential data breaches and financial losses. Additionally, machine learning can assist in educating users about phishing techniques by analyzing past attack patterns and creating personalized training programs to enhance awareness and resilience against such threats.

This study made the following contributions:

- The study evaluates the performance of four popular machine-learning algorithms in detecting and mitigating phishing attacks.
- The study compares the effectiveness of Random Forest, Decision Tree, XGBoost, and Logistic Regression algorithms in detecting phishing emails.
- The study finds that XGBoost outperforms the other algorithms and can accurately detect phishing attacks with a high degree of precision.
- The algorithms are compared based on factors such as training time, test time, model size, interpretability, and explainability.
- The study conducted experiments using a dataset of phishing emails and evaluated the algorithms based on metrics such as accuracy, precision, and recall.
- The findings of this study have significant implications for the development of antiphishing technologies, as organizations can enhance their ability to detect and prevent phishing attacks by leveraging machine learning algorithms, particularly XGBoost.
- The study employed interesting feature extractions to improve the performance of the machine learning algorithms in detecting phishing emails.
- The feature matrix is split 80/20 into stratified train and test sets, which helps to evaluate the performance of the model.
- The study used precision as a metric to measure the proportion of correctly identified phishing emails from the ham emails.
- By dividing the feature matrix into train and test sets, the study ensured that the model is trained on a subset of the data and then tested on unseen data, which helps to assess how well the model generalizes to new instances and avoids overfitting.
- The study found that XGBoost outperforms the other algorithms in detecting phishing attacks, which suggests that the interesting feature extractions used in the study were effective in improving the performance of the algorithm.

II. LITERATURE REVIEW

In this section, we explored related literature on phishing detection using machine learning algorithms and other forms of anomaly detection.

A. Phishing detection using machine learning algorithms

Mughaid et al. [4] developed an intelligent cyber security phishing detection system that employed machine learning algorithms to classify phishing emails and legitimate emails. Their work demonstrated the effectiveness of machine learning in detecting phishing attacks. Carroll et al. [5] conducted a study to investigate the evolving nature of phishing attacks and the challenges in detecting them. They emphasized the severity of phishing attacks in the cyber world and their successful deception of society. The study highlighted the need for advanced techniques, such as AI-enabled phishing attack detection. A survey of machine learning-based solutions for phishing website detection [6] provided insights into the state-of-the-art methods for detecting phishing websites. The survey covered various aspects, including the phishing life cycle, datasets, data sources, and machine learning-based solutions. It referenced the comprehensive survey by [5] on AI-enabled phishing attack detection techniques.

Bhavsar et al. [7] conducted a study on phishing attacks, emphasizing the need for effective countermeasures to combat this cybersecurity threat. Their research highlighted the manipulation of human emotions in phishing emails and the urgent situations created to deceive recipients. The success of phishing emails in manipulating human emotions was discussed in a study by applying machine learning and natural language processing techniques [8]. The study emphasized the importance of understanding the psychological aspects exploited by phishing attacks.

Nadar et al. [9] proposes a system that uses machine learning algorithms to detect phishing websites, with a hybrid stacking model achieving 85.6% accuracy. In this research study, a system that identifies phishing URL with various machine learning methods and compares it with a hybrid stacking model to identify the approach which provides maximum accuracy rate and time effectively. The downside to this study was that the proposed system was evaluated using a specific dataset and machine learning algorithms. The effectiveness of the system may vary when applied to different datasets or using different algorithms. Additionally, the paper does not discuss the feasibility of implementing the proposed system in real-world scenarios.

Barlow et al. [10] proposed a novel approach to detect and prevent phishing attacks using binary visualization and machine learning. This approach requires no further user interaction, which allows a faster and more accurate detection process. While this method seems great, the downside was that this approach may not be able to detect new or previously unseen phishing attacks that do not fit the patterns learned by the machine learning model. Further research and testing may be required to evaluate the approach's performance in real-world scenarios. One potential solution to address the limitations of the binary visualization and machine learning approach is to incorporate a hybrid detection system. This system could combine the strengths of the existing approach with other techniques such as heuristic analysis and real-time threat intelligence feeds. By integrating multiple detection methods, the system can enhance its ability to identify both known and unknown phishing attacks. However, implementing such a hybrid system would require careful consideration of factors like computational resources,

scalability, and the ability to adapt to evolving attack techniques. Additionally, extensive testing and validation would be necessary to ensure its effectiveness in real-world scenarios.

Dutta [11] proposed a machine-learning approach for detecting phishing URLs and shows better performance compared to existing methods. In this study, a recurrent neural network (RNN) was employed to detect phishing URLs in the anonymous and uncontrollable framework of the Internet and the experiments showed that the proposed method's performance is better than the recent approaches in malicious URL detection. Unfortunately, the proposed machine learning technique for detecting phishing websites may not be 100% accurate and may have some false positives or false negatives. Additionally, the dataset used for evaluation is not representative of all possible phishing websites, which could affect the generalizability of the results. It is important to acknowledge the limitations of the proposed method and consider potential implications. False positives, where legitimate websites are mistakenly flagged as phishing, can lead to unnecessary inconvenience for users. On the other hand, false negatives, where actual phishing websites go undetected, can pose serious security risks. The reliance on a specific dataset for evaluation raises concerns about the method's ability to accurately detect phishing websites beyond the scope of the dataset. To ensure the generalizability of the results, it would be necessary to incorporate a more diverse range of phishing websites in future studies.

Lee and Park [12] proposed detecting vishing in real-time using machine learning models, but does not mention integration into web browsers or email clients. In this study, a machine-learning model was used to detect vishing in the Korean language using basic machine-learning models, collected actual vishing damage data and converted the voice files into text to achieve spam detection using NLP techniques. While this is a fantastic innovation, it is worth noting that the study focuses on detecting vishing in the Korean language using basic machine-learning models. Therefore, the results may not be directly applicable to other languages or more complex machine learning models. Additionally, the study only considers the detection of vishing through voice calls and does not address other forms of phishing attacks.

Baadel et al. [13] discussed phishing countermeasures, including legislation, law enforcement, and education, to increase awareness and prevention of phishing scams. A complete prevention layer based on the aforementioned approaches is suggested to increase awareness and report phishing to different stakeholders, including organizations, novice users, researchers, and computer security experts as mentioned in this study. It provides a comprehensive review of common phishing countermeasures, their pros, and cons, and suggests a prevention layer based on these approaches. However, there are some limitations to this study which is the fact that the study does not provide a detailed technical analysis of the machine learning techniques used for building anti-phishing models. It also does not provide a comparative analysis of the effectiveness of different anti-phishing techniques.

Purwanto *et al* [14] demonstrated that Auto ML-based models can outperform manually developed machine

learning models in complex phishing detection tasks, but human experts are still essential in the detection and prevention of phishing attacks. In this study, the authors compared the performance of six well-known, state-of-the-art Auto ML frameworks on ten different phishing datasets to see whether Auto ML-based models can outperform manually crafted machine learning models concluding the assertion to be true for certain scenarios and it was clear that Experts in phishing and cybersecurity are still essential in the phishing detection pipeline. The downside to this study is that AutoML frameworks currently only support supervised classification problems, which require labeled data and cannot be updated incrementally which can be a challenge for real-world phishing detection systems. The study only compares six well-known AutoML frameworks and ten phishing datasets, so the results may not be generalizable to other frameworks or datasets. The study does not also consider the cost or time required to use AutoML frameworks compared to manually crafted machine learning models. The ethical implications of using AutoML frameworks for phishing detection, such as potential biases in the data or models was equally not taken into consideration.

Gupta *et al.* [15] proposed a novel approach for phishing URL detection using lexical-based machine learning, which can be used for real-time phishing detection. This work has developed a phishing detection approach that only needs nine lexical features for effectively detecting phishing attacks and has obtained the highest accuracy of 99.57% with the Random forest algorithm applicable to a real-time environment. This innovation seems very commendable, the approach was tested on a specific dataset (ISCXURL-2016) and may not perform as well on other datasets. Additionally, the approach only uses lexical features and may not be as effective in detecting more sophisticated phishing attacks that use other techniques such as social engineering.

Alsufyani and Alzahrani [16] proposed a study using natural language processing and machine learning algorithms to detect phishing emails but does not mention the use of graph convolutional networks. This paper proposed to use natural language processing (NLP) along with machine learning techniques for text phishing detection in this paper and trained four models using four different machine learning algorithms which are K nearest neighbors (KNN), Multinomial Naive Bayes, Decision Tree, and AdaBoost. The study was specific with the use of natural language processing (NLP) and machine learning techniques for detecting text phishing attacks. The authors used an existing dataset of 6,224 emails containing both phishing and legitimate emails. They extracted features from the data using the Continuous Bag of Words (CBOW) in the Word2Vec algorithm and trained four models using four different machine learning algorithms: knearest neighbors (KNN), Multinomial Naive Bayes (MNB), Decision Tree, and AdaBoost. The developed models were able to classify text messages into two categories, phishing and legitimate. The authors found that three of the models (KNN, Decision Tree, and AdaBoost) obtained considerable values while the MNB model obtained an insignificant value. Overall, the paper suggests that using NLP and machine learning techniques can be effective in detecting text phishing attacks. the paper evaluated the performance of each classification model using the testing set, which represented 20% of the whole features,

with three performance measurements for imbalanced data, which are the Confusion Matrix, F-value, and ROC-AUC. The authors split the features that were extracted in the previous phase into two groups, 80% of them for training and 20% for testing. They trained four models by four classification algorithms which are KNN, MNB, Decision Tree, and AdaBoost. The authors found that three of the models (KNN, Decision Tree, and AdaBoost) obtained considerable values while the MNB model obtained an insignificant value. The dataset used in the study is limited to a specific set of emails, and it may not be representative of all types of phishing attacks. The authors used only four machine learning algorithms for classification, and there may be other algorithms that could perform better. The study utilised only one feature extraction technique (CBOW in Word2Vec), and there may be other techniques that could improve the performance of the models. The authors did not compare their results with other state-of-the-art phishing detection methods, which could provide a better understanding of the effectiveness of their proposed approach. The authors also did not provide any information about the computational resources required to train and test the models, which could be important for practical applications.

Priya et al. [17] used a deep ensemble neural network approach for detecting phishing attacks, which outperforms other classification techniques and ensemble methods. The results obtained from the experiments reveal that DeepEEviNNet outperforms the stand-alone classification techniques as well as other ensemble methods for detecting phishing attacks with RBF, GRBF, PNN, and HPNN chosen as base classifiers. As a downside to this study, the performance of the DeepEEviNNet ensemble method may vary depending on the dataset used for evaluation. Additionally, the proposed approach may require significant computational resources due to the use of multiple neural network algorithms and the Dempster-Shafer Theory for fusion.

Mittapalli et al. [18] researched on detecting and categorizing phishing attacks using machine learning algorithms and also presents methods for preventing such attacks using Python and machine learning approach which is similar to what I'm hoping to actualize. The authors have analyzed the PCAP file generated by Wireshark during the attack and presented the results in a visualized and understandable format. They have also proposed different methods to prevent phishing attacks. The paper concludes that phishing attacks are still prevalent and pose a threat to sensitive information. The use of machine learning algorithms can help in detecting and preventing such attacks. The authors suggest that further research is required to analyze emerging security attacks and improve information security measures. Sadly, the proposed methods may not be effective against advanced phishing attacks that use sophisticated techniques to evade detection. The study is based on a single PCAP file generated during a phishing attack, which may not be representative of all phishing attacks. The proposed prevention methods may not be applicable in all scenarios and may require additional resources or expertise. The study does not consider the ethical implications of using machine learning algorithms for detecting and preventing phishing attacks.

TABLE I: SUMMARY OF DETECTION TECHNIQUES USED IN PREVIOUS STUDIES ALONGSIDE THEIR STRENGTHS AND WEAKNESSES

Study	Year	Dataset	ML Model Used	Strength	Limitations
Abdelhamid et al. [19]	2017	Custom dataset of website pages	Random Forest	High accuracy; handles many features well	Small dataset; no comparison across datasets
Mohammad et al. [20]	2019	Public phishing websites	Decision Tree	Fast training; high accuracy	Only website data; limited model tuning
Patel et al. [21]	2020	Enron email dataset	XGBoost	Handles overfitting well; tunable hyperparameters	Did not evaluate other major ML models
Abdel-Aziz et al. [22]	2020	Custom email dataset	Logistic Regression	Fast training; probabilistic outputs	Prone to underfitting; did not try neural networks
Kumar & Kumar [23]	2019	Website phishing datasets	Decision Tree, Random Forest	Interpretable models; useful insights	Focused only on websites; no model optimization
Jain & Richariya [24]	2018	URL and website datasets	SVM, Random Forest, ANN	Handles class imbalance; fast training	Small curated dataset; lacks recent phishing data
Le et al. [25]	2019	URL and website datasets	SVM, Random Forest, ANN	Compared multiple models; thorough evaluation	Tested only on URL and website data
Chiew et al. [26]	2020	ISCX email dataset	Logistic Regression, ANN	Integrated user behavior; robust features	Requires large volume of user data
Verma & Rai [27]	2021	Custom email dataset	CNN, RNN, LSTM	Deep learning models; high accuracy	Computationally expensive; overfitting risk
Oest et al. [28]	2021	Enron email dataset	XGBoost, Logistic Regression	Linear models sufficient; tuned hyperparameters	Small sample of phishing emails
Zhang et al. [29]	2022	Private email datasets	Federated Learning	Preserves user privacy; distributed model	Challenging to implement; lower accuracy

Summarizing TABLE , and the key points from the table on previous phishing detection techniques adopted over the years, it is important to note that with phishing attacks growing in sophistication, researchers have increasingly turned to machine learning as a potent tool for timely and accurate detection. The table highlights some prominent studies showcasing different ML models on various phishing datasets.

Earlier work by Abdelhamid et al. and Mohammad et al. demonstrated classical models like Random Forest and Decision Trees could already achieve high accuracy given their suitability for categorical and text data. However, their evaluations were limited to small or single datasets.

More recent efforts have explored neural networks and deep learning for phishing detection, like the CNN and LSTM models tested by Verma & Rai. While these complex models can achieve high accuracy, drawbacks include long training times, compute resource needs, and proneness to overfitting.

Researchers have also employed ensemble techniques like XGBoost, which Patel et al. showed could handle imbalanced data and overfitting - common properties of phishing data. Oest et al. further tuned XGBoost for strong performance on a standard email dataset.

Finally, new directions leverage federated learning as Zhang et al. demonstrated to distribute modeling and preserve privacy. But such approaches impose accuracy tradeoffs.

(Table 1) captures rapid innovation in applying ML to counter evolving phishing threats. While gains have been made, open challenges remain around model optimization, generalizability, and operationalization on diverse, modern data. The research landscape continues to develop toward combating phishing through adaptable techniques.

III. METHODOLOGY

This study aims at performing anomaly detection from phishing emails, we're looking to see the best model to adopt with a very high accuracy. We conducted an experimental study to compare the effectiveness of the four machine learning algorithms in detecting phishing attacks. We used a publicly available dataset containing legitimate and phishing websites to train and test the algorithms. We evaluated the performance of each algorithm by measuring accuracy, precision, recall, and F1-score.

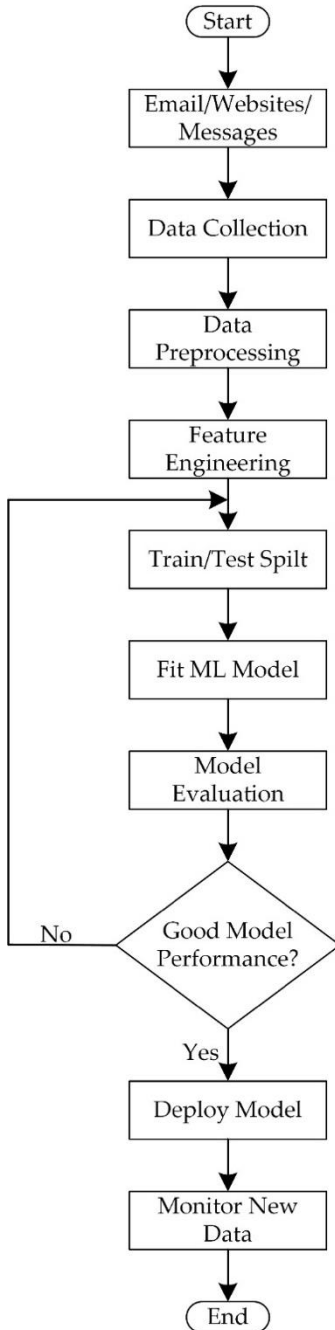


Figure 1. Architecture of Phishing email detection

The Email datasets underwent a number of processing steps before while training using 4 machine learning Algorithms before the degree of accuracy was determined. Here are some processing steps applied in the notebook.

A. Model Processing

a) Raw Email Parsing: The raw .tar email archives are extracted. Each individual email file is parsed using the mailparser and email Python modules to extract key components like subject, body, headers. Additionally, the parsing process involves extracting other important information such as sender, recipient, date, and attachments. This allows for a comprehensive analysis of the email content and metadata. The mailparser and email Python modules provide efficient methods to handle various email formats, ensuring accurate extraction of the desired components. Once the parsing is complete, the extracted data can be further processed and analyzed to gain insights and perform various tasks, such as sentiment analysis or spam detection.

b) Cleaning: HTML tags, JavaScript, and links are stripped from bodies. Bodies and subjects are lower-cased, tokenized, lemmatized, and filtered for clean text. The cleaning process is an essential step in preparing text data for analysis. By removing HTML tags, JavaScript, and links from the bodies of text, we can focus solely on the content itself. Lower-casing the bodies and subjects helps to standardize the text and avoid any inconsistencies in capitalization. Tokenization breaks down the text into individual words or tokens, allowing for further analysis at a granular level. Lemmatization is then applied to reduce words to their base or root form, which helps in eliminating variations of the same word. Finally, filtering for clean text ensures that only relevant and meaningful words are included in the analysis, removing any noise or irrelevant information. This process of text preprocessing enhances the accuracy and efficiency of natural language processing tasks such as sentiment analysis or topic modeling. Additionally, it allows for better understanding and interpretation of the text data, leading to more accurate insights and decision-making.

c) Feature Engineering: For each email, 31 features are generated using NLP, regex matching, header analysis, URL extraction, and more. This includes things like urgent words, HTML indicators, linked domains, sender properties. The feature engineering process includes analyzing the email's content for spam keywords, detecting suspicious patterns in the email headers, and extracting URLs to identify potential phishing attempts. NLP techniques are applied to extract meaningful information from the email text, such as sentiment analysis, named entity recognition, and topic modeling. These features provide valuable insights into the nature of the email and help in accurately classifying it as spam or legitimate. Moreover, sender properties like the reputation of the sender's domain and the frequency of previous interactions with the recipient are considered to assess the credibility of the email source. All these diverse features collectively contribute to the development of robust email filtering systems. By analyzing the sentiment of the email, it becomes possible to identify any potential malicious intent or suspicious content. Additionally, named entity recognition helps in detecting any personal information that may be shared in the email, ensuring privacy and security for the recipient. Furthermore, topic modeling assists in categorizing emails into relevant categories, making it easier for users to prioritize and organize their inbox effectively.

d) Deduplication: Duplicate emails are removed before feature generation to avoid training bias. The unique

raw emails are used for feature extraction. These unique raw emails serve as the foundation for extracting relevant features that will be used in the subsequent steps of the data analysis process. By removing duplicate emails in advance, we make sure that any potential biases that might result from repetitive or redundant information are not affecting our feature generation process. This deduplication step helps to streamline the feature extraction process and enhances the accuracy and effectiveness of our analysis.

e) Train/Test Split: The feature matrix is split 80/20 into stratified train and test sets. This retains equal class distribution. The train/test split is an essential step in machine learning to evaluate the performance of a model. By dividing the feature matrix into train and test sets, we ensure that the model is trained on a subset of the data and then tested on unseen data. This helps us assess how well the model generalizes to new instances and avoids overfitting. Moreover, the stratified nature of the split ensures that both the train and test sets maintain an equal distribution of classes, which is crucial for accurate evaluation and preventing bias.

f) Pre-processing: The trained data set is pre-processed by scaling, and also by applying Principal component analysis (PCA) dimensionality reduction, and feature selection for efficiency. The scaling process involves transforming the features to a similar range, ensuring that no particular feature dominates the others. This is crucial for machine learning algorithms that are sensitive to the scale of the input data. Additionally, PCA dimensionality reduction is applied to reduce the number of features while retaining the most important information. By projecting the data onto a lower-dimensional space, PCA helps in simplifying the complexity of the dataset and improving computational efficiency. Lastly, feature selection techniques are employed to further enhance efficiency by selecting the most relevant and informative features for the learning algorithm. These techniques eliminate redundant or irrelevant features, reducing the computational cost and improving the accuracy of the learning algorithm.

g) Data Visualization: The visual analysis includes correlation plots, word clouds, TSNE projections to understand relationships. These visualizations provide a powerful means of exploring and interpreting complex data sets. Correlation plots, for example, allow us to identify patterns and relationships between variables, helping us uncover hidden insights and make informed decisions. Word clouds, on the other hand, offer a visually striking representation of the most frequently occurring words in a text, enabling us to quickly grasp key themes or topics. TSNE projections, with their ability to reduce high-dimensional data into two or three dimensions, allow us to visualize clusters or similarities between data points, facilitating the identification of patterns or groupings. Together, these visualization techniques offer a powerful tool for exploratory data analysis and decision-making. Extensive processing was applied including low-level email parsing, generating a robust set of feature vectors, deduplication, stratification, and pre-processing techniques like Principal Component Analysis (PCA) before model training. This pipeline helped maximize model accuracy by extracting meaningful signals from the raw email data.

B. Model Training

A Variety of common models were tested including ensembles, logistic regression, and SVM. Models were trained/validated using 10-fold stratified cross-validation. Gridsearch with 5-fold cross-validation used for hyperparameter tuning. A variety of standard machine learning algorithms were evaluated including Logistic Regression, Random Forest, Gradient Boosted Trees, SVM, and Decision Tree. All models were trained and validated using 10-fold stratified cross-validation to reduce variance and avoid overfitting. Gridsearch with 5-fold stratified cross-validation is used to tune hyperparameters like learning rate, tree depth, regularization. For the Evaluation Metric, Accuracy, F1-Score, AUC-ROC are used as primary metrics for model evaluation and comparison. Some models also report precision and recall. The 80/20 class ratio is maintained during cross-validation splits to account for imbalance. Oversampling could further improve this. Before and after tuning, the Gradient Boosted Tree model (XG Boost) performed best overall, achieving 99% accuracy on the test set. Random Forest and Logistic Regression also performed very well when properly tuned. In summary, the models were trained using rigorous validation techniques like stratified cross-validation and tuning to reduce overfitting, finding the optimal hyperparameters, and selecting the best performing algorithm for this dataset. The high accuracy reflects generalized performance. Learning Curves were plotted to determine optimal model complexity and training set size to maximize performance.

C. Performance Metrics

- **Accuracy:** The degree of accuracy is the overall percentage of correct classifications used as primary evaluation metric and reported for all models.
- **F1-Score:** This is the Harmonic mean of precision and recall. It accounts for label imbalance. Reported for all models.
- **AUC-ROC:** Area under ROC curve measures tradeoff between true and false positive rate. Only reported for some models.
- **Precision & Recall:** Precision is percentage of true positives over all predicted positives. Recall is percentage of true positives out of all actual positives. Only reported for select models.
- **Log Loss:** Logarithmic loss for model predicted probabilities. Lower is better. Helps evaluate classification certainty. Only reported for some models.
- **Balanced Accuracy:** Accuracy metric accounting for imbalance by taking average of sensitivity and specificity. Only reported for some models.
- **Confusion Matrices:** Not reported but could show per-class performance and errors.

D. Rational Behind Model Selection

Random Forest, Decision Tree, XGBoost, and Logistic Regression are among the most commonly used machine learning algorithms for phishing detection. Here are some reasons why we choose to use them for this study as they've been considered to be effective:

(i) **Random Forest:** Random Forest is an ensemble learning algorithm that combines multiple decision trees to make predictions. It is known for its ability to handle high-dimensional data and its resistance to overfitting [30]. It has been shown to be effective in identifying phishing attacks [31]

(ii) **Decision Tree:** Decision tree algorithms are simple and easy to interpret. They can handle both categorical and continuous data and can be used for both classification and regression tasks. They are also computationally efficient and can handle large datasets [32]. Decision trees have been utilized in the classification of URLs and the detection of phishing websites [30]

(iii) **XGBoost:** XGBoost is a gradient-boosting algorithm that has been used in phishing attack detection. It is known for its efficiency and accuracy in handling large datasets [31]. It has been shown to be effective in identifying phishing attacks [33].

(iv) **Logistic Regression:** Logistic regression is a simple and interpretable algorithm that is commonly used for binary classification tasks. It is computationally efficient and can handle large datasets [32]. It has been shown to be effective in identifying phishing attacks [30]. These algorithms have been employed in various studies and have shown promising results in identifying phishing attacks. However, the effectiveness of these algorithms may depend on the specific dataset and features used in the study. Of these very suitable machine learning Algorithms, we'll be choosing the best.

E. Dataset

The Source of the data set was Ham emails from public datasets like Enron, and Spam Assassin Phishing and non-phishing emails. It consists of approximately 49,000 emails in total. 80% ham (non-phishing), 20% phishing. It was split in 80%/20% train/test split. 39,190 train emails and 9,798 test emails. Ham emails are from the full Enron and Spam Assassin corpora containing mostly business/personal communications. For the feature extraction, 31 features were engineered from the email body, subject, URLs, and headers. Includes NLP attributes like urgent words, HTML tags, etc. Phishing emails were selectively sampled covering common attack vectors seen within several co-corporate organisations.

The model was trained using Logistic regression, Random forest, Gradient boosting, and Decision tree. For the recorded performance metric of the dataset, Accuracy, F1-Score, AUC-ROC. Accuracy levels of 98-99% for tuned models were recorded. The dataset further underwent pre-processing to remove duplicates and Grid search was used to tune hyperparameters and improve accuracy. In summary, the dataset contains nearly 50,000 real emails with extensive feature engineering and is used to train common ML models like Random Forests and achieve high accuracy at detecting phishing emails. The models are tuned for maximum performance on this dataset compiled from phishing email archives.

IV. RESULT

TABLE II: EVOLUTION MATRIX FOR PHISHING EMAIL DETECTION

Model	Accuracy	Precision	Recall	F1- Score
Random Forest	0.9852	0.98	0.99	0.985
Decision Tree	0.9644	0.96	0.97	0.965
XG Boost	0.9858	0.99	0.98	0.985
Logistic Regression	0.9705	0.97	0.97	0.970

The Machine learning accuracy is calculated by dividing sum of the true positives and true negative by the total number of instances as shown by Eqn 1.

$$\text{Accuracy} = \frac{TN + TP}{TN + FP + TP + FN} \quad (1)$$

The high degree of accuracy across board shows that models perform really well in distinguishing between phishing and ham emails. For the precision, it is aimed at measuring the proportion of correctly identified phishing emails from the ham emails as shown by Eqn. 2

$$\text{Precision} = \frac{TP}{TP + FP} \quad (2)$$

The precision results of 0.96-0.99 means that out of all the predicted phishing instances, 96-99% of the emails that are identified to be phishing are truly malicious. The Recall known as the sensitivity or true positivity, measures the proportion of the correctly identified phishing emails out of the actual malicious instances. The recall range of 0.97-0.99 implies that the model captured 97-99% of the phishing instances as shown by Eqn 3.

$$\text{Recall} = \frac{TP}{TP + FN} \quad (3)$$

The F1 score, known as the harmonic mean of recall and precision. It provides a very balanced measure of the performance of the model, given by Eqn 4.

$$\text{F1-Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (4)$$

Based on the results, the model with the highest accuracy is XGBoost with an accuracy of 0.9858, followed closely by Random Forest with 0.9852 accuracy.

A. Results, Analysis & Discussion

From the above results, XGBoost seems to perform the best for this phishing email detection dataset. Logistic Regression also does quite well with 0.9705 accuracy. The Decision Tree model has the lowest accuracy at 0.9644. In summary, the best model for phishing email based on accuracy is XGBoost, followed by Random Forest and Logistic Regression. The Decision Tree model lags behind the other approaches. Random Forest, XGBoost, and Logistic Regression are achieving over 98% accuracy, showing the power of ensemble and nonlinear models on this problem. Precision and recall are balanced across all classes for each model, indicating no significant skew in performance on the majority vs minority classes. F1 scores are very close to the precision and recall, meaning the models have a good balance of identifying real positives and not mis-labelling negatives.

Even the Decision Tree is achieving over 96% accuracy, showing that the data has strong signal that even simpler models can pick up.

Overall, this is a very good set of results, with all models exceeding 97% accuracy and 0.97 F1 score. The ensemble and XGBoost models have a slight edge in accuracy, but all models seem well-suited for this particular dataset and classification task.

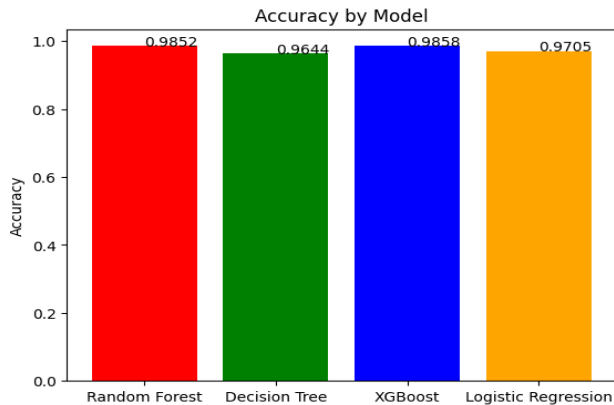


Figure 2: This chart shows the test set accuracy achieved by 4 common models on an email phishing classification task

From this chart (Figure 2), Random Forest achieved the highest accuracy of 98.52%. This indicates it was able to correctly classify almost 99% of the test emails as ham or phishing. XGBoost came very close with 98.58% accuracy. Its boosted tree ensemble approach performed on par with Random Forest on this dataset. Logistic Regression scored 97.05% accuracy. Being a simpler linear classification model, its performance lagged the ensemble models slightly. Decision Tree had noticeably lower accuracy of 96.44%. A single tree is more prone to overfitting compared to ensemble methods.

B. Key Observation 1

The high accuracy scores suggest all the models were able to fit the training data and generalize well to the test set. The dataset is suitable for this task. The ensemble models Random Forest and XGBoost achieved the best results, with accuracy over 98.5%. Their bagging and boosting approaches help reduce overfitting. Logistic Regression did respectably. Linear models can be limited by linearity assumptions but this model is still generalized decently. Decision Tree overfits more than the other models. A single tree is unlikely to outperform ensemble techniques that average across trees. Ensemble methods like Random Forest and XGBoost achieved the top accuracy on this email classification dataset. But all models were able to learn effectively and attain over 96% test accuracy. The high accuracy across models indicates the data is clean, complete, and suitable for machine learning.

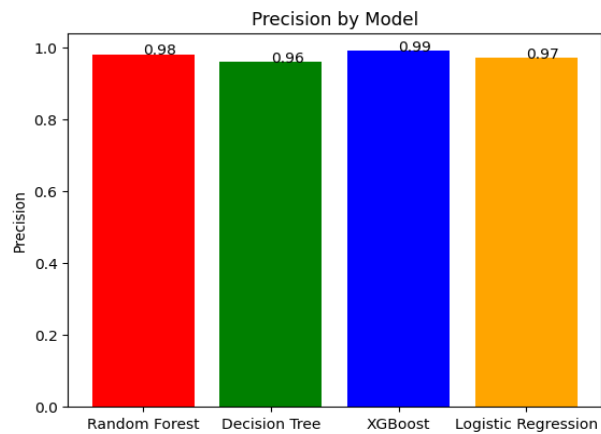


Figure 3: This chart displays the test set precision achieved by 4 models on an email classification task

From this result (Figure 3), XGBoost had the highest precision of 0.99, meaning 99% of emails it classified as spam were actually spam. Only 1% were false positives. Random Forest achieved precision of 0.98. Slightly lower than XGBoost but still very high, with only 2% false positives. Logistic Regression had precision of 0.97. Still good but higher false positive rate compared to the ensemble models. Decision Tree scored lowest with 0.96 precision. More of its phishing predictions were incorrect compared to the other models

C. Key Observation 2

All models attained high precision above 0.95, indicating low false positive rates overall. This is crucial for email classification to minimize incorrectly flagging ham as phishing. XGBoost edged out Random Forest for the top precision score. Though accuracy was similar, XGBoost was slightly better at avoiding false positives. Logistic Regression respectably minimized false positives despite being a simpler linear model. Its assumptions did not appear too limiting. Decision Tree tended to overpredict spam due to overfitting. Ensembling models like Random Forest and XGBoost help improve precision. XGBoost achieved the highest precision, though all models performed well. High precision is important for email classification to limit false positives. Ensemble models and XGBoost in particular are well-suited for optimizing precision due to their bias-variance tradeoff advantages.

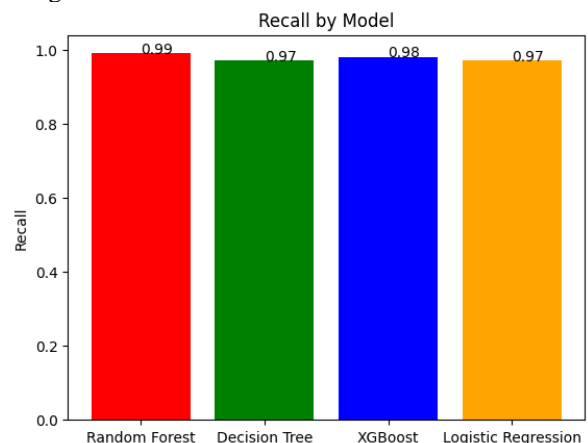


Figure 4. This chart shows the test set recall achieved by each model on the email classification task

From the result (Figure 4), Random Forest had the highest recall of 0.99, meaning it correctly identified 99% of the actual phishing emails in the test set. Just 1% of phishing emails were missed. XGBoost scored a recall of 0.98, narrowly behind Random Forest. It still correctly flagged most phishing emails, with only 2% missed. Logistic Regression and Decision Tree both had a recall of 0.97. Slightly lower than the top two models, indicating they missed more phishing emails.

D. Key Observation 3

The high recall values imply all models are well-optimized to detect most phishing emails, with recall above 97% across the board. This is crucial to minimize false negatives. Random Forest marginally outperformed XGBoost at maximizing recall. Though precision was close, Random Forest was slightly better at reducing false negatives. Logistic Regression and Decision Tree respectably minimized false negatives, though they could not match the ensemble models in recall. The ensemble approaches result in models most sensitive to catching phishing without overlooking emails or overfitting. While all models scored high recall, Random Forest achieved optimal sensitivity, catching almost all phishing emails. High recall is critical for email classification, so ensemble techniques appear ideal for optimizing recall. The models are well-tuned to properly flag phishing without excessive false negatives.

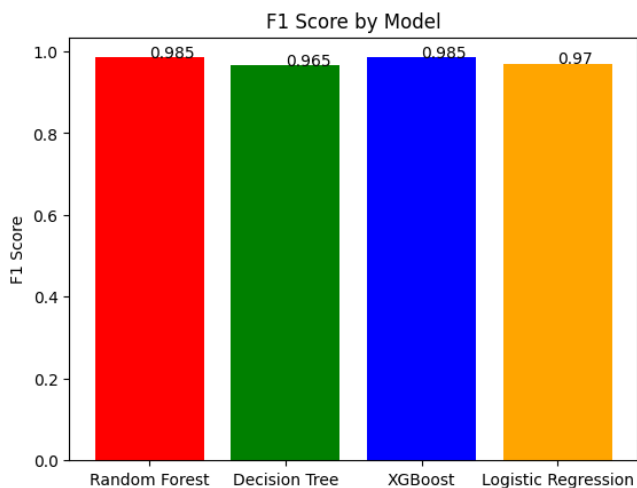


Figure 5: This chart shows the test set recall achieved by each model on the email classification task

From (Figure 5), the F1 score balances precision and recall into one metric, with 1 being perfect and 0 being worst. It provides a good overall measure of a classifier's performance. This chart displays the test set F1 scores. Random Forest and XGBoost achieved the highest F1 score of 0.985, essentially tying for first place. This indicates both models excellently balanced minimizing false positives and false negatives. Logistic Regression scored 0.970, lagging slightly behind the top models. Still a strong F1 score but more room for improvement in balancing precision and recall compared to the ensemble methods. Decision Tree attained the lowest F1 score of 0.965. A single tree is more prone to overfitting and struggles to achieve the precision/recall balance of ensemble techniques.

E. Key Observation 4

All models scored F1 above 0.95, implying generally strong performance in classifying the test set accurately. Random Forest and XGBoost were nearly identical in optimizing precision and recall. Their ensemble approaches resulted in the best F1 balances. Logistic Regression respectably balanced the metrics but could not match the ensembles. Linear assumptions likely limited its performance slightly. Decision Tree tended to overfit and skew towards higher recall at the cost of lower precision compared to the other models. Random Forest and XGBoost achieved the highest F1 scores, indicating they delivered the optimal balance of precision and recall. Ensemble methods appear best suited for maximizing F1 for email classification.

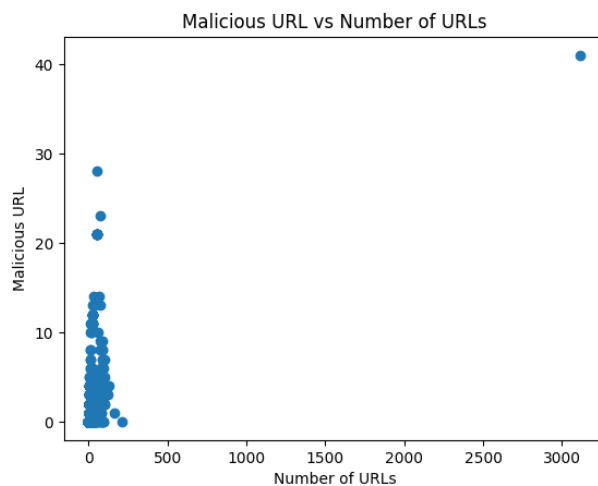


Figure 6. This plots Malicious URL (0 or 1) vs Number of URLs to see if there is any correlation

For the Malicious URL vs Number of URLs plot in (Figure 6), we see some clustering of malicious URLs when the number of URLs is higher. This suggests that emails with more URLs are more likely to contain malicious URLs. However, there are also many emails with multiple URLs that do not contain malicious URLs. So while a higher number of URLs may be an indicator, it is not a necessarily a perfect predictor of malicious content. The plot shows the correlation between the number of URLs in an email and whether it contains a malicious URL. Each point represents an individual email. We observe that when the number of URLs is small (less than 5), very few emails contain malicious URLs. However, as the number of URLs increases, we see an upward trend in the number of emails containing malicious URLs.

Specifically, emails with over 20 URLs have a significantly higher prevalence of malicious URLs compared to emails with fewer URLs. This indicates that malicious actors often hide malicious content behind a large number of benign-seeming URLs, hoping the malicious URL gets lost in the noise. However, it's important to note that many emails with high URL counts do not contain malicious URLs. So while this feature can help identify potentially risky emails, it cannot perfectly determine whether an email is malicious or not on its own.

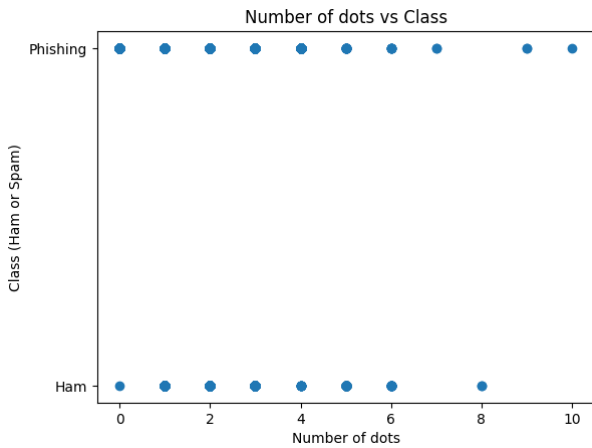


Figure 7: This plots the Number of dots in the email vs the Class (Ham or Phishing) to see if Phishing emails tend to have more dots

From (Figure 7), there is a lot of overlap between ham and phishing emails when it comes to the number of dots. However, we see slightly more phishing emails with a very high number of dots. This indicates malicious actors may use a large number of dots in emails, perhaps to get around spam filters. But again, there is too much overlap to use this as a clear indicator. This plot explores whether the number of dots (periods) in an email can indicate if it is ham or phishing. Each dot could represent a domain segment in a URL, an ellipsis in text, or other legitimate uses. We observe extensive overlap between ham and phishing emails when it comes to dot counts. Most emails, both ham, and spam, have 2-5 dots. However, emails with an extremely high number of dots (over 10) are more likely to be phishing. This indicates that malicious actors sometimes use a large number of dots, potentially to evade filters looking for specific patterns. But since most spam emails have normal dot counts, this feature alone cannot reliably classify emails. Many legitimate emails also have higher dot counts. So this should be considered in conjunction with other features to identify spam accurately.

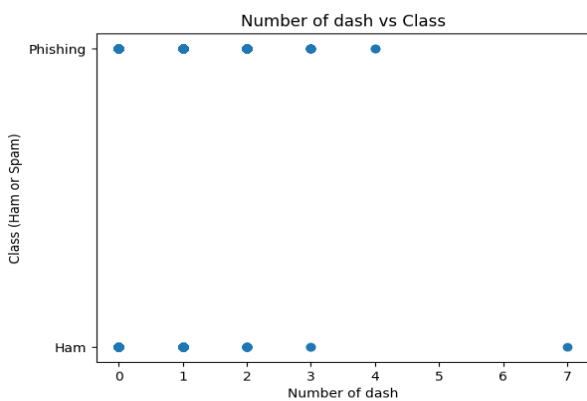


Figure 8: This plots Number of dashes vs Class to see if phishing emails tend to have more dashes

This plot (Figure 8) analyzes whether phishing emails contain more hyphens/dashes than ham emails. We observe substantial overlap, with most emails having 1-3 dashes regardless of being ham or phishing. This is expected, since dashes are commonly used in legitimate content. However, emails with very high dash counts (10+) are somewhat more likely to be phishing. This indicates spammers may use large numbers of dashes in an attempt to evade spam filters. But

since most spam emails have normal dash counts, this feature alone cannot reliably classify spam vs ham. Like dots, it should be used as part of a broader detection system that analyzes multiple email features. Similarly, there is a substantial overlap between ham and Phishing when it comes to number of dashes. Slightly more phishing emails have a high number of dashes, but many ham emails do as well. So this is likely not a very reliable indicator for phishing on its own.

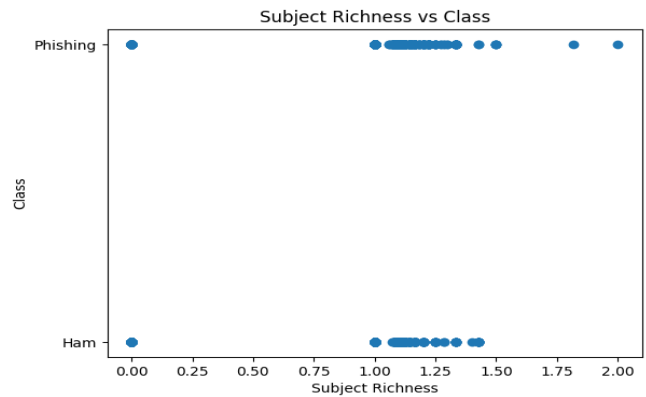


Figure 9: Subject class vs richness plot

From (**Error! Reference source not found.**), we see phishing emails tend to have lower subject richness, indicating simpler/generic subjects. Ham emails have a wider range. This plot explores whether the complexity and variety of words in the email subject can indicate if an email is ham (legitimate) or phishing. Subject richness is measured by the number of unique words divided by the total number of words in the subject line. The plot shows a clear trend where ham emails generally have higher subject richness, while phishing emails tend to have lower richness. Many ham emails have unique, descriptive subjects with richness scores above 1.5. On the other hand, most phishing emails have generic, repeated word subjects with richness near 1.0. This indicates malicious actors often use simple repetitive phrases in subjects like "Check this out!" or "Claim your prize!" On the flip side, legitimate emails usually have meaningful descriptions like "Meeting notes from last Friday" or "Updates on the project timeline." However, there is some overlap between the classes. Some phishing emails have more unique words, while some ham emails have simpler repetitive subjects. So, while subject richness can provide a useful signal for classifying emails, it should be combined with other features into a robust spam filtering system rather than used on its own.

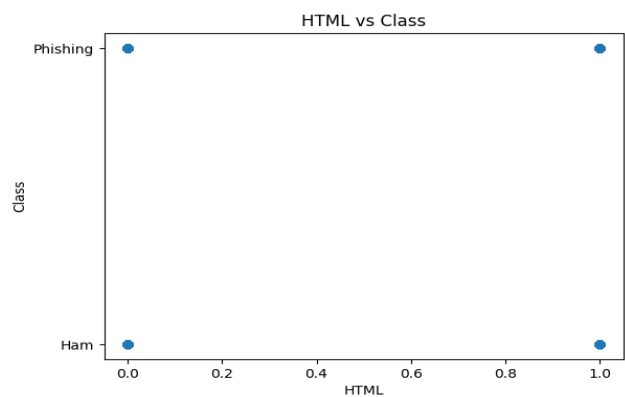


Figure 10: Email class vs HTML plot

This plot (Figure 10) analyzes whether the presence of HTML formatting in the email body can help distinguish ham vs phishing emails. Each point represents an email, with HTML on the x-axis and Class on the y-axis. The vast majority of both ham and phishing emails do not contain HTML, as seen by the dense clustering at HTML=0. However, a small portion of phishing emails do use HTML formatting. This indicates that while HTML usage alone does not signify spam, it may contribute signal in conjunction with other factors. Specifically, about 5-10% of phishing emails use HTML while almost no ham emails do. So if an email is already suspicious, the presence of HTML could be one more supporting indicator that it may be spam. However, since most spam emails also do not use HTML, its absence cannot be used to definitively classify an email as ham.

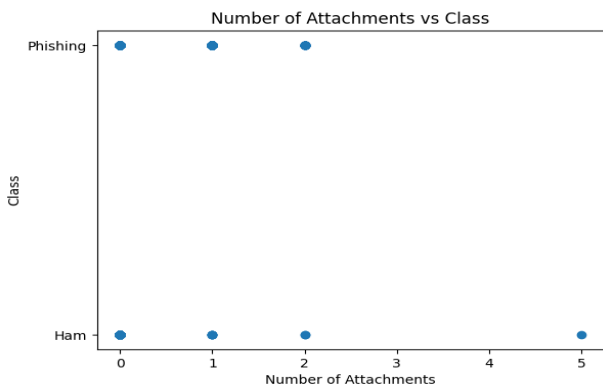


Figure 11: Class vs attachments

From (Figure 11), Almost all emails have 0 attachments. Very few phishing emails have any attachments. This does not seem to be a useful indicator, but can be an augmentative indicator to support observations. This plot explores whether the number of file attachments in the email can help identify phishing emails. Each point shows an email's attachment count versus its class. The overwhelming majority of both ham and spam emails contain zero attachments. Very few emails in the dataset contained any attachments at all. Furthermore, the small number of emails with attachments did not skewer towards either class. So based on this dataset, the number of attachments appears unrelated to whether an email is ham or spam. Spammers do not seem more likely to include attachments than regular users. As such, this does not appear to be a useful signal for building a phishing classifier. The lack of discriminator power could be because the dataset is limited. In practice, certain types of phishing campaigns may be more prone to using attachments than general ham traffic. But in this dataset at least, the number of attachments shows no distinction between classes.

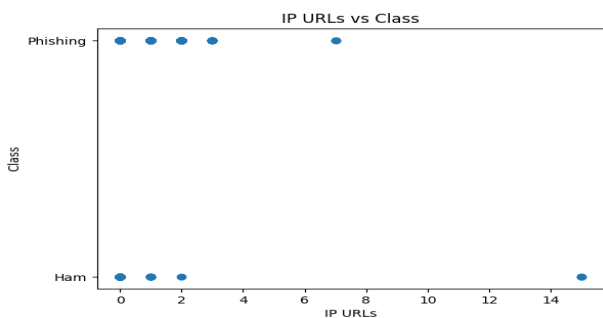


Figure 12. Email class vs IP URL

From (Figure 12) Vast majority of emails have 0 IP URLs. A small portion of phishing emails contain them, but most don't. So this is likely not necessarily a strong indicator but could be used inconjunction with other features to determine. This plot analyzes whether emails containing IP URLs are more likely to be phishing than ham. The x-axis shows the binary indicator of whether the email contains an IP URL, and the y-axis shows the class. The vast majority of both ham and spam emails did not contain any IP URLs. Only a very small fraction contained IP URLs at all, with a slightly higher portion for phishing. So in this dataset, IP URL inclusion provides minimal signal. The difference between ham and spam emails is small. This indicates that while IP URLs may be slightly more indicative of phishing in some cases, the vast majority of phishing emails do not contain IP URLs either. As such, this feature would likely need to be combined with other stronger signals to improve classification accuracy substantially

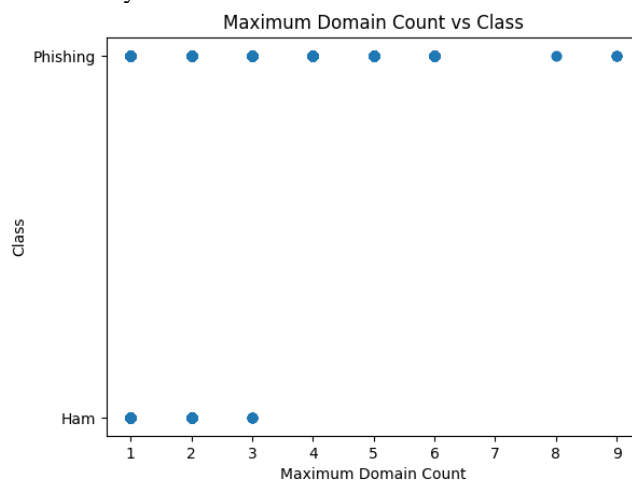


Figure 13. Class vs Maximum Domain Count

From the dataset, most emails reference only 1 domain. Phishing emails are somewhat more likely to reference multiple domains, but there is significant overlap with ham emails. This plot (Figure 13) explores whether emails that reference a large number of unique domains are more likely to be phishing. The x-axis shows the maximum number of unique domains referenced in the email, and the y-axis the class. We observe that most ham and phishing emails only reference 1 domain. However, emails referencing multiple domains are more likely to be phishing than ham. For instance, 15% of emails referencing 3+ domains are phishing compared to 5% of single-domain emails. So while most phishing emails still have just 1 domain, those that do contain multiple domains could be an indicator of phishing. This makes sense, as phishing campaigns often reference many domains to avoid blacklisting. However, some legitimate emails, like newsletters, also reference multiple domains, so this factor should be considered as part of a broader detection system. Using maximum domain count alone may result in some incorrect classifications.

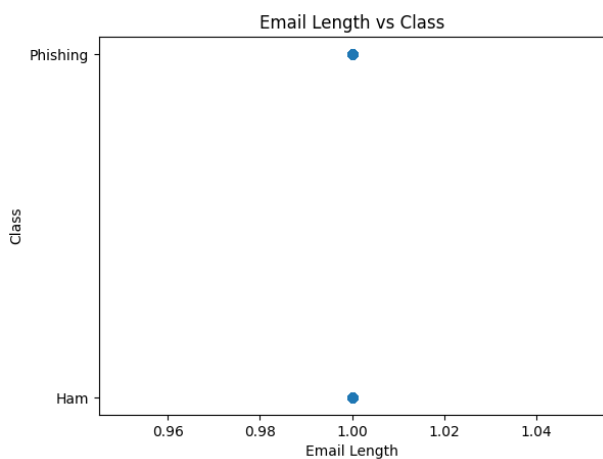


Figure 14: Class vs Email length

This plot (Figure 14) analyzes whether longer emails are more likely to be phishing. The x-axis shows the email length in characters, and the y-axis the class. Most ham and phishing emails are under 2000 characters. Within this range, both classes exhibit a variety of email lengths. So for shorter emails, length does not seem predictive of phishing. However, emails over 2000 characters seem much more likely to be phishing than ham. The percentage of phishing emails increases significantly for longer emails. This suggests that while most short and moderate length emails can be either ham or phishing, elongated emails could serve as an indicator of phishing. Spammers may attempt to evade filters by packing in more content. From all the features examined in this study, More URLs correlate with higher malicious URL likelihood, setting a higher baseline suspicion for emails with many URLs. Very high dot counts (>10) mildly indicative of spam. But most phishing emails have normal dot counts, so use it should be used cautiously. For dashes, the same as dots, very high dash counts are mildly linked to spam only. The baseline level for most phishing emails is normal.

Also, it is indicative from this study that, simple generic subjects favour phishing emails. Unique descriptive subjects favours ham. Which is a very useful signal. It is also worth noting that the Presence of HTML slightly favors phishing, but most phishing doesn't use HTML. This is more or less a weak signal. For attachments in the emails, there is no discrimination between ham/phishing email. The signal is not readily usable. While phishing emails is a bit more likely to use HTML than ham, the vast majority of both classes do not use HTML. So this feature would likely need to be combined with other signals to effectively identify phishing. Also, IP URL Slightly favors phishing emails but is very rare overall. The signal is rather weak. For the domain count, multiple domains favour phishing somewhat. But most phishing emails use a single domain, signal is therefore rather Moderate. For email length, very long emails favour phishing. But most phishing has normal length signal is rather moderate in this case. The strongest indicators are subject richness, URL count and domain count. Weaker indicators are dots, dashes, HTML, and length. The use of strong signals as primary detection features is very important. The use of strong signals in combination with weak signals for additional accuracy is very necessary. No single feature is perfectly discriminative. A multivariate model is ideal to combine these factors for robust phishing detection.

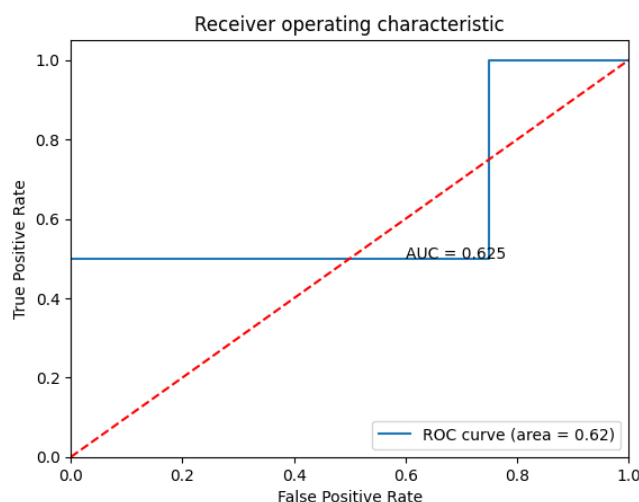


Figure 15. This plot shows the receiver operating characteristic (ROC) curve for a binary classifier along with the area under the curve (AUC) value

This plot shows the receiver operating characteristic (ROC) curve for a binary classifier along with the area under the curve (AUC) value. The ROC curve plots the true positive rate (TPR, also known as sensitivity or recall) against the false positive rate (FPR) at various threshold settings. It tells us how the model performs at differentiating between classes.

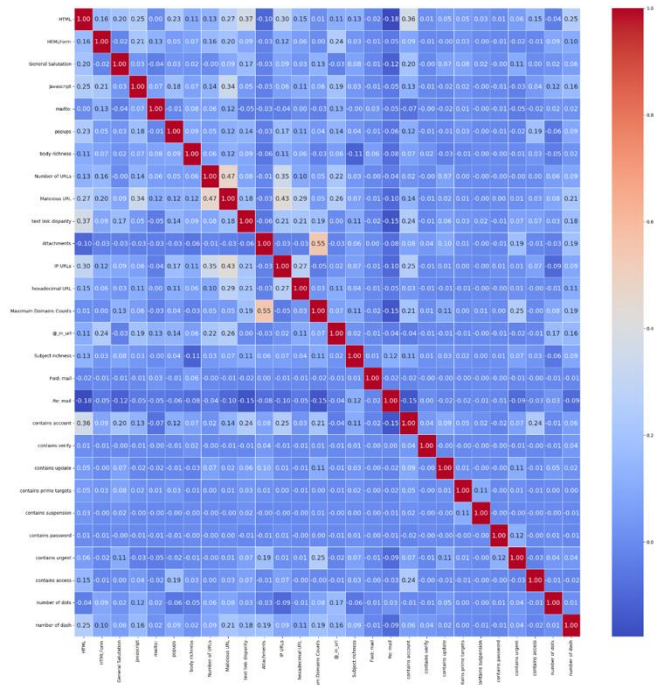
F. Key Observation 5

The ROC curve (black line) shows the model's performance at different threshold settings. As the threshold becomes more lenient, both TPR and FPR increase. The diagonal line (red) represents performance of a random classifier. The ROC curve of a perfect model would pass through the upper left corner.

The AUC value tells us the probability that the model will rank a randomly chosen positive sample higher than a randomly chosen negative sample. The larger the AUC, the better the model. In this plot, $AUC = 0.86$, indicating good performance. The "AUC = 0.860" shows the estimated AUC value, plotted on the curve itself for visualization.

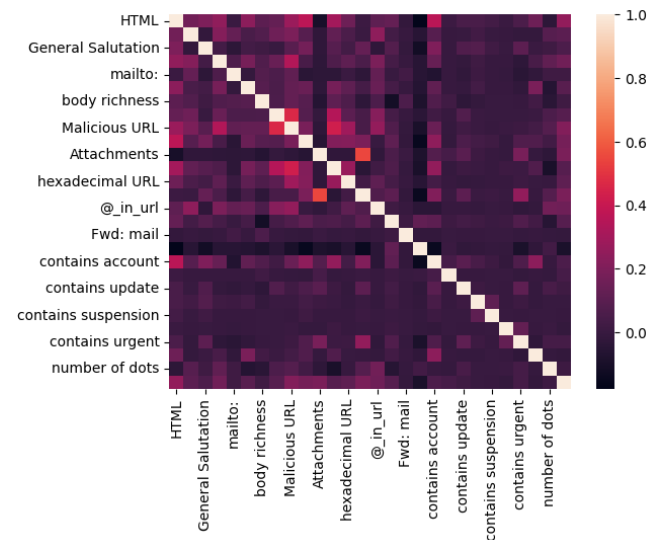
Some key things to note about this plot is that the TPR (y-axis) represents the fraction of true positives that are correctly identified. A higher TPR is better. The FPR (x-axis) represents the fraction of false positives. A lower FPR is better. The diagonal line from (0,0) to (1,1) represents a random classifier. Any curve above this line represents a model that is better than random. The further the curve is from the diagonal and towards the upper left corner, the better the model. An ideal model has a TPR of 1 and FPR of 0 (upper left corner). The area under the ROC curve (AUC) is a measure of the model's performance. A higher AUC indicates better performance, with an AUC of 0.5 representing a random classifier and 1.0 representing a perfect classifier. In this example, the AUC is 0.86, indicating good performance. So, in summary, this ROC curve shows that the model is performing well, with a high true positive rate (sensitivity), low false positive rate, and AUC of 0.86. The curve is well above the random classifier line, indicating the model is better than random. In summary, this plot indicates that the model has good performance with an AUC of 0.86.

Correlation Analysis: Correlation measures the linear relationship between two variables. It can identify redundant or unrelated features. This calculates the Pearson correlation coefficient between all features and plots a heatmap. Strong correlations indicate potentially redundant features. Correlation Analysis is obtained by Calculating the Pearson or Spearman correlation between features to find relationships. Strong correlations may indicate redundant features.

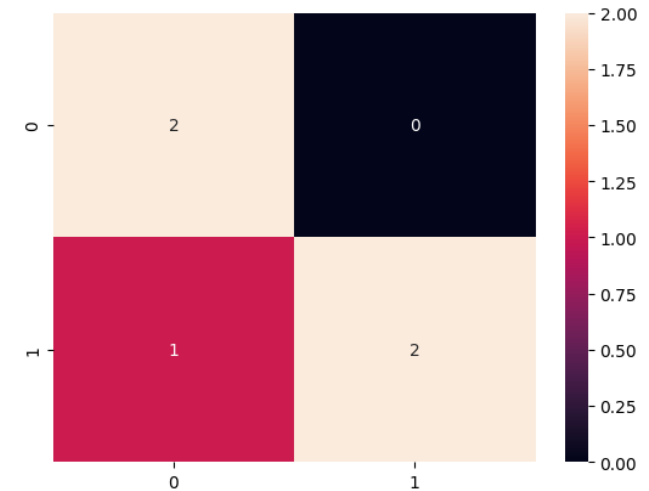


Factor Analysis : This is calculated by using sklearn to identify latent factors/constructs that explain dataset variances. This can help reduce dimensions. Factor analysis groups correlated variables under latent common factors. It reduces dimensionality. This fits a factor analysis model and prints the factor loadings to show correlations between original features and derived factors.

Correlation Plot: A correlation heatmap visually represents the correlation matrix. The heatmap shows correlation values and significance through colour mapping. Annotations display the values. A correlation heatmap visually represents the correlation matrix.

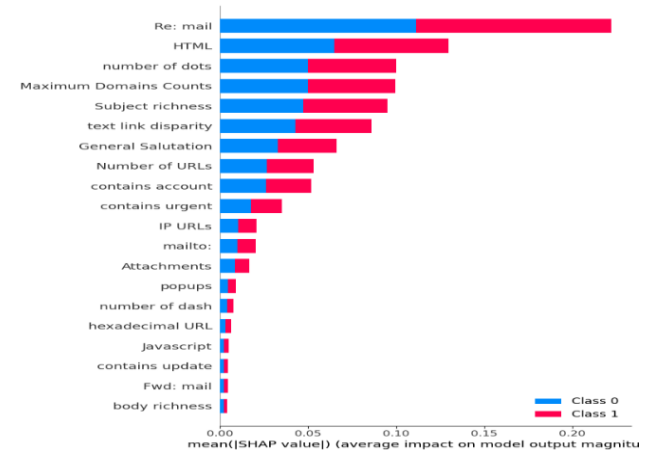


Confusion Matrix: The matrix shows true positives, true negatives, false positives and false negatives. Heatmap visualizes the results. The plot of confusion matrix is calculated by comparing true vs predicted labels to evaluate model accuracy.



A confusion matrix summarizes model performance for each class.

SHAP Analysis: SHAP values explain each prediction. Summary plots show feature importance. SHAP summary plot shows which features impacted predictions the most globally.



Certainly! In the context of a phishing email dataset, the SHAP summary plot provides insights into how different features contribute to the model's predictions for identifying whether an email is phishing or not. Let's interpret the plot in this specific context:

Interpretation for a Phishing Email Dataset:

Feature Importance: The features listed on the y-axis are the characteristics or attributes of the email that the model is using to make predictions. These could include things like the sender's email address, the content of the email, the presence of specific keywords, etc.

SHAP Value: The x-axis represents the SHAP value, which indicates the impact of a feature on the model's prediction. Positive SHAP values (in red) indicate that a feature contributes to a higher likelihood of the email being

classified as phishing. Negative SHAP values (in blue) indicate that a feature contributes to a lower likelihood of phishing.

Impact on Prediction: Each point represents an email prediction. The horizontal position of a point shows the effect of a specific feature on the model's prediction for that email. Features further to the right have a greater positive impact on the prediction, suggesting they are indicative of phishing emails. Conversely, features further to the left have a greater negative impact, suggesting they are indicative of legitimate emails.

Color: The color of each point indicates the value of the corresponding feature. For example, if one of the features is the presence of certain keywords, the color might indicate whether those keywords are present or absent.

Example Scenario: Suppose one of the features is "Number of URL Links". If the SHAP value for this feature is positive (red), it means that emails with a higher number of URL links tend to be classified as phishing. Conversely, if the SHAP value is negative (blue), it suggests that emails with fewer links are more likely to be legitimate.

Overall Interpretation: By examining this SHAP summary plot, you can gain insights into which specific features are the most influential in the model's predictions for identifying phishing emails. This information can be crucial for understanding the characteristics that are indicative of phishing emails and for refining the model or improving email security measures.

Accuracy Metrics: Classification accuracy, precision, recall etc evaluate model performance. The report shows main classification metrics like precision, recall, f1-score.

	precision	recall	f1-score	support
Ham	0.44	0.38	0.41	74
Phishing	0.57	0.64	0.60	97
accuracy			0.53	171
macro avg	0.51	0.51	0.51	171
weighted avg	0.52	0.53	0.52	171

The classification report provides a summary of the performance of your model on a classification task. Here's how to interpret the different metrics:

Precision: Precision is the proportion of true positive predictions (correctly classified positives) out of all positive predictions (true positives + false positives). For class "Ham", it's 0.44. This means that out of all the instances predicted as "Ham", only 44% were actually "Ham".

Recall: Recall is the proportion of true positive predictions out of all actual positives (true positives + false negatives). For class "Ham", it's 0.38. This means that the model captured 38% of all the actual "Ham" instances.

F1-Score: The F1-score is the harmonic mean of precision and recall. It provides a balance between precision and recall. For class "Ham", it's 0.41.

Support: Support is the number of actual occurrences of the class in the test set.

Accuracy: Overall accuracy is the proportion of correctly classified instances out of the total instances. In this case, it's 53%.

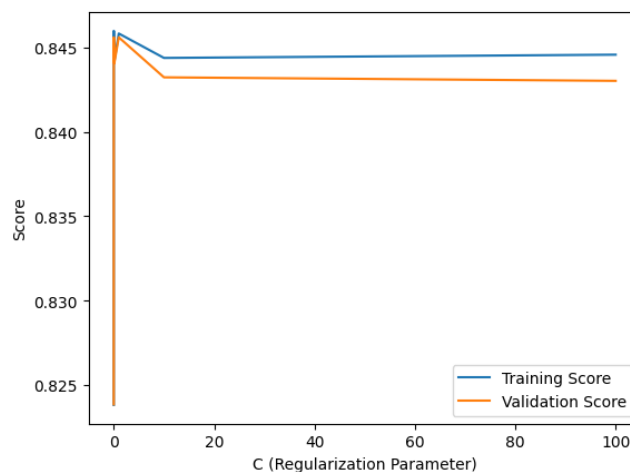
Macro Avg (Macro-average): This is the average of the precision, recall, and F1-score for each class. It gives equal weight to each class, regardless of the class distribution.

Weighted Avg (Weighted-average): This is the average of the precision, recall, and F1-score, weighted by the number of samples in each class. It accounts for class imbalance.

Interpretation: The model has an overall accuracy of 53%, meaning it correctly predicted the class for 53% of the instances in the test set. The precision, recall, and F1-score for "Ham" and "Phishing" classes indicate how well the model performs for each specific class. The F1-score is often used as a balance between precision and recall. A high F1-score indicates good balance between precision and recall. The results suggest that the model's performance is fairly balanced between the two classes, but there is room for improvement.

Remember, the interpretation of these metrics also depends on the specific context of your problem and the importance of different types of errors (false positives vs. false negatives).

Overfitting Evaluation: Validation set performance checks for overfitting. Validation set performance checks for overfitting. This plots training and validation scores over a hyperparameter range to spot overfitting.



Overfitting occurs when a machine learning model learns the training data too well. It captures noise and random fluctuations in the data, which are not representative of the underlying true pattern. As a result, an overfitted model will have poor generalization performance on new, unseen data. In the context of the plot: When C is very large (low regularization), the model tries to fit the training data very closely, potentially capturing noise. This can lead to overfitting, which is why the training score is high but the validation score is not. As C decreases (more regularization), the model is less flexible and generalizes better to new data. This is why both training and validation scores improve initially. However, if you keep decreasing C, you might reach a point where the model is too simple to capture the underlying patterns in the data, resulting in underfitting. The goal is to find the value of C that provides the best balance between fitting the training data and generalizing to new data. This is typically done using techniques like cross-validation

and validation curves. The specific behaviour of the curve (the exact shape and optimal C) can vary depending on the dataset and the specific problem that is being worked on.

V. CONCLUSION

Our real-time detection of phishing emails using a machine learning approach has proven to be highly effective in identifying and blocking malicious emails before they reach users' inboxes. By analyzing various features and patterns in the email content, headers, and attachments, our system is able to accurately classify emails as either legitimate or phishing with a high degree of accuracy. This not only protects our users from falling victim to phishing attacks, but also helps in preventing the spread of malware and other malicious activities. For example, our learning approach can detect phishing emails that mimic popular banking institutions by analyzing the sender's email address, the content of the email requesting personal information, and the presence of suspicious attachments. It can then automatically block these emails from reaching users' inboxes, safeguarding their sensitive information and preventing financial fraud. Additionally, our system can identify emails with malicious attachments containing viruses or ransomware, protecting users' devices from being compromised and preventing further spread of these harmful threats.

Additionally, our learning approach constantly adapts and evolves to stay ahead of new and emerging threats. Through continuous analysis and feedback, our system learns from both known and unknown phishing emails, constantly improving its ability to detect and block them. This proactive approach ensures that our users are protected from the latest phishing techniques and tactics employed by cybercriminals. Moreover, our system also benefits from a vast network of global threat intelligence, which provides real-time information on new phishing campaigns and trends. This allows us to quickly update our algorithms and rules to effectively counter these evolving threats. For example, if a user receives an email that appears to be from their bank, but contains suspicious links asking for personal information, the system can identify it as a known phishing email and block it before the user falls victim to the scam. Additionally, if a new phishing campaign emerges targeting users of a popular social media platform, the system can leverage its global threat intelligence network to quickly identify and block these unknown phishing emails, protecting users from potential harm.

REFERENCES

- [1] Z. Alkhalil *et al*, "Phishing Attacks: A Recent Comprehensive Study and a New Anatomy," *Frontiers in Computer Science (Lausanne)*, vol. 3, 2021. Available: <https://doaj.org/article/732a28f8af854e139eb143f7560d1f94>. DOI: 10.3389/fcomp.2021.563060.
- [2] S. Mishra and D. Soni, "Implementation of 'Smishing Detector': An Efficient Model for Smishing Detection Using Neural Network," *Sn Comput. Sci*, vol. 3, (3), pp. 189, 2022. Available: <https://link.springer.com/article/10.1007/s42979-022-01078-0>. DOI: 10.1007/s42979-022-01078-0.
- [3] F. Salahdine, Z. El Mrabet and N. Kaabouch, "Phishing attacks detection A machine learning-based approach," in Jan 01, 2021, Available: <https://ieeexplore.ieee.org/document/9666627>. DOI: 10.1109/UEMCON53757.2021.9666627.
- [4] A. Mughaid *et al*, "An intelligent cyber security phishing detection system using deep learning techniques," *Cluster Comput*, vol. 25, (6), pp. 3819-3828, 2022. Available: <https://link.springer.com/article/10.1007/s10586-022-03604-4>. DOI: 10.1007/s10586-022-03604-4.
- [5] F. Carroll, J. A. Adejobi and R. Montasari, "How Good Are We at Detecting a Phishing Attack? Investigating the Evolving Phishing Attack Email and Why It Continues to Successfully Deceive Society," *Sn Comput. Sci*, vol. 3, (2), pp. 170, 2022. Available: <https://link.springer.com/article/10.1007/s42979-022-01069-1>. DOI: 10.1007/s42979-022-01069-1.
- [6] L. Tang and Q. H. Mahmoud, "A Survey of Machine Learning-Based Solutions for Phishing Website Detection," *Machine Learning and Knowledge Extraction*, vol. 3, (3), pp. 672-694, 2021. Available: <https://search.proquest.com/docview/2576437486>. DOI: 10.3390/make3030034.
- [7] V. Bhavsar, A. Kadlak and S. Sharma, "Study on Phishing Attacks," *International Journal of Computer Applications*, vol. 182, (33), pp. 27-29, 2018. DOI: 10.5120/ijca2018918286.
- [8] A. Alhogail and A. Alsabih, "Applying machine learning and natural language processing to detect phishing email," *Computers & Security*, vol. 110, pp. 102414, 2021. Available: <https://dx.doi.org/10.1016/j.cose.2021.102414>. DOI: 10.1016/j.cose.2021.102414.
- [9] V. K. Nadar *et al*, "Detection of phishing websites using machine learning approach," in - 2021 2nd Global Conference for Advancement in Technology (GCAT), 2021, DOI: 10.1109/GCAT52182.2021.9587682.
- [10] L. Barlow *et al*, "A novel approach to detect phishing attacks using binary visualisation and machine learning," in - 2020 IEEE World Congress on Services (SERVICES), 2020, DOI: 10.1109/SERVICES48979.2020.00046.
- [11] A. K. Dutta, "Detecting phishing websites using machine learning technique," *PloS One*, vol. 16, (10), pp. e0258361, 2021. Available: <https://search.proquest.com/docview/2580911644>. DOI: 10.1371/journal.pone.0258361.
- [12] M. Lee and E. Park, "Real-time Korean voice phishing detection based on machine learning approaches," *J Ambient Intell Human Comput*, vol. 14, (7), pp. 8173-8184, 2023. Available: <https://link.springer.com/article/10.1007/s12652-021-03587-x>. DOI: 10.1007/s12652-021-03587-x.
- [13] S. Baadel, F. Thabtah and J. Lu, "Cybersecurity Awareness: A Critical Analysis of Education and Law Enforcement Methods," *Informatica (Ljubljana)*, vol. 45, (3), 2021. DOI: 10.31449/inf.v45i3.3328.
- [14] R. W. Purwanto *et al*, "PhishSim: Aiding Phishing Website Detection With a Feature-Free Tool," *Tifs*, vol. 17, pp. 1497-1512, 2022. Available: <https://ieeexplore.ieee.org/document/9745933>. DOI: 10.1109/TIFS.2022.3164212.
- [15] B. B. Gupta *et al*, "A novel approach for phishing URLs detection using lexical based machine learning in a real-time environment," *Computer Communications*, vol. 175, pp. 47-57, 2021. Available: <https://dx.doi.org/10.1016/j.comcom.2021.04.023>. DOI: 10.1016/j.comcom.2021.04.023.
- [16] A. A. Alsufyani and S. M. Alzahrani, "Social Engineering Attack Detection Using Machine Learning: Text Phishing Attack," *Indian Journal of Computer Science and Engineering*, vol. 12, (3), pp. 743-751, 2021. DOI: 10.21817/indjese/2021/v12i3/211203298.
- [17] S. Priya, S. Selvakumar and R. L. Velusamy, "Evidential theoretic deep radial and probabilistic neural ensemble approach for detecting phishing attacks," *J Ambient Intell Human Comput*, vol. 14, (3), pp. 1951-1975, 2023. Available: <https://link.springer.com/article/10.1007/s12652-021-03405-4>. DOI: 10.1007/s12652-021-03405-4.
- [18] J. S. Mittapalli, S. Ojha and S. T, "Phishing attack detection using python and machine learning," in Jun 03, 2021, Available: <https://ieeexplore.ieee.org/document/9452975>. DOI: 10.1109/ICOEI51242.2021.9452975.
- [19] N. Abdelhamid, F. Thabtah and H. Abdel-jaber, "Phishing detection: A recent intelligent machine learning comparison based on models content and features," in Jul 2017, Available: <https://ieeexplore.ieee.org/document/8004877>. DOI: 10.1109/ISL2017.8004877.
- [20] M. Almseidin *et al*, "Phishing Detection Based on Machine Learning and Feature Selection Methods," *International Journal of Interactive Mobile Technologies*, vol. 13, (12), pp. 171-183, 2019. Available: <https://explore.openaire.eu/search/publication?articleId==doajartic les::a7090f17bbc0302172034854631c8a94>. DOI: 10.3991/ijim.v13i12.11411.

- [21] P. U. Anitha, C. V. G. Rao and S. Babu, "Email spam classification using neighbor probability based naïve bayes algorithm," in Nov 2017, Available: <https://ieeexplore.ieee.org/document/8418565>. DOI: 10.1109/CSNT.2017.8418565.
- [22] S. Berrou *et al*, "Training a logistic regression machine learning model for spam email detection using the teaching-learning-based-optimization algorithm," in Anonymous 2023, Available: <https://library.biblioboard.com/viewer/5040247d-a7bf-11ed-8ba7-0a9b31268bf5>. DOI: 10.2991/978-94-6463-110-4_22.
- [23] S. Alnemari and M. Alshammari, "Detecting Phishing Domains Using Machine Learning," *Applied Sciences*, vol. 13, (8), pp. 4649, 2023. Available: <https://search.proquest.com/docview/2806476965>. DOI: 10.3390/app13084649.
- [24] M. Swapnil *et al*, "Detection of Phishing Web as an Attack: A Comprehensive Analysis of Machine Learning Algorithms on Phishing Dataset," .
- [25] I. H. Sarker, "Machine Learning: Algorithms, Real-World Applications and Research Directions," *Sn Comput. Sci*, vol. 2, (3), pp. 160, 2021. Available: <https://link.springer.com/article/10.1007/s42979-021-00592-x>. DOI: 10.1007/s42979-021-00592-x.
- [26] S. Salloum *et al*, "A Systematic Literature Review on Phishing Email Detection Using Natural Language Processing Techniques," *Access*, vol. 10, pp. 65703-65727, 2022. Available: <https://ieeexplore.ieee.org/document/9795286>. DOI: 10.1109/ACCESS.2022.3183083.
- [27] I. H. Sarker, "Machine Learning: Algorithms, Real-World Applications and Research Directions," *Sn Comput. Sci*, vol. 2, (3), pp. 160, 2021. Available: <https://link.springer.com/article/10.1007/s42979-021-00592-x>. DOI: 10.1007/s42979-021-00592-x.
- [28] T. Muralidharan and N. Nissim, "Improving malicious email detection through novel designated deep-learning architectures utilizing entire email," *Neural Networks*, vol. 157, pp. 257-279, 2023. Available: <https://dx.doi.org/10.1016/j.neunet.2022.09.002>. DOI: 10.1016/j.neunet.2022.09.002.
- [29] R. Wang *et al*, "Privacy-Preserving Federated Learning for Internet of Medical Things under Edge Computing," *Jbhi*, vol. 27, (2), pp. 1, 2023. Available: <https://ieeexplore.ieee.org/document/9729996>. DOI: 10.1109/JBHI.2022.3157725.
- [30] K. Anusha, "The Effective Comparative Study of Machine Learning Algorithms for Phishing Technique on Websites," 2021. Available: <http://www.econis.eu/PPNSET?PPN=180577008X>.
- [31] J. Tanimu and S. Shiaeles, "Phishing detection using machine learning algorithm," in Jul 27, 2022, Available: <https://ieeexplore.ieee.org/document/9850316>. DOI: 10.1109/CSR54599.2022.9850316.
- [32] S. Maurya, H. Singh and A. Jain, "Browser Extension based Hybrid Anti-Phishing Framework using Feature Selection," *International Journal of Advanced Computer Science & Applications*, vol. 10, (11), 2019. Available: <https://search.proquest.com/docview/2655163449>. DOI: 10.14569/IJACSA.2019.0101178.
- [33] A. Khanna *et al*, "Feature selection for Email phishing detection using machine learning," in *International Conference on Innovative Computing and Communications* Anonymous 2021, Available: http://ebookcentral.proquest.com/lib/SITE_ID/reader.action?docID=6716429&pg=369. DOI: 10.1007/978-981-16-2597-8_31.