



A REVOLUTIONARY MACHINE LEARNING BASED CYBER THREAT AND PHISHING DETECTOR

Syed Suhaila S¹

Assistant Professor, Department of Computer Science Engineering, Alagappa Chettiar Government College of
Engineering and Technology, Karaikudi, India¹

Abstract: For security researchers, phishing assaults' increasing frequency continues to be a major worry. Conventional signature-based methods for phishing website detection frequently miss freshly created phishing sites. By utilizing large and varied datasets, researchers are creating machine learning-based systems that can accurately identify and categorize phishing websites in order to address this issue. Logistic Regression (LR), K-Nearest Neighbors (KNN), Support Vector Machine (SVM), Decision Tree (DT), Random Forest (RF), and Artificial Neural Network (ANN) are among the models that are trained after a sequence of procedures to suitably prepare the dataset. In order to wrap up this study, the top-performing model as determined by performance metrics is integrated into a Flask web application. The integration of machine learning models for end-user accessibility is frequently overlooked in existing research, which frequently concentrates only on model performance. This study makes a contribution by outlining the thorough procedures necessary to include the chosen model into an intuitive web application utilizing the Flask framework, in addition to comparing the suggested model's performance with earlier research.

Keywords: Security, Phishing, Conventional based methods, LR, KNN, SVM, DT, RF, ANN, Flask.

1) INTRODUCTION

The digital revolution has altered how people use technology and how businesses operate, leading to advancements in banking, marketing, communication, and service delivery. Unprecedented increases in internet usage have made it easier to acquire information in real time and connect globally, satisfying customers' constantly changing needs. However, this rapid digital revolution has resulted in significant security concerns. Cybercriminals are using increasingly complex tactics, such as malware, phishing, and other exploitative approaches, to compromise the integrity of online interactions. One of the most common cyberthreats, phishing, uses dishonest methods to pose as trustworthy websites or communication platforms in an attempt to trick individuals into divulging personal information. Hackers are increasingly exploiting security flaws to make phony websites appear genuine. According to a 2024 report by the Anti-Phishing Working Group (APWG), there were more than five million reported instances of phishing attacks, setting a new high. Of them, social media outlets accounted for 42.8% of the efforts. These trends demonstrate the pressing need for robust cybersecurity defenses, increased user awareness, and ongoing research to lower the dangers associated with these threats and provide a safer online environment.

Using Uniform Resource Locators, one can locate the website depicted. The seven parts of it are as follows protocol, domain at the top level, Malevolent Query, child domain, parameter, path, and domain name. A protocol controls the exchange of information between a web server and a browser.

Web Transfer Protocol (WTP), Post Office Protocol (POP), and Simple Mail Transfer Protocol (SMTP), HTTPS, and Internet Message Access Protocol (IMAP) are a few of the most widely used protocols. Additionally, a domain name acts as a unique internet reference to help identify a website. A web server's path, like /home/address/image.jpeg, indicates the precise location where a given directory or file resides. A branch domain is contained within the main domain name. A domain name always contains the top-level domain (TLD) in the case of stanford.edu, the TLD is still edu. Webpages that a redynamic include questions. An inquiry is always followed by a question mark. When a client asks a server for a



page, it uses a query string to run program. The URL <https://example.com/completed/track/there?name=alexa> is one example. Name = alexa is a query in this URL.

To improve the model's ability to find pertinent patterns for precise phishing website identification, feature selection algorithms such as Principal Component Analysis (PCA), Chi-Square, and Recursive Feature Elimination (RFE) are used. Logistic regression, Random Forest (RF), Support Vector Machine (SVM), Decision Tree (DT), K-Nearest Neighbors (KNN) and Artificial Neural Network are among the algorithms that are being compared across various feature sets. Increasing classification accuracy with the use of very effective algorithms from various models. Using performance evaluation criteria including accuracy, precision, recall, and F1-score, all models' efficacy is evaluated.

II) RELATED WORK

Phishing, a type of social engineering attack, remains one of the most significant and evolving cybersecurity threats, targeting sensitive user information such as login credentials, credit card details, and personal data. Despite advancements in cybersecurity, the dynamic nature of phishing attacks makes them challenging to detect using traditional methods. Attackers often create fraudulent webpages that mimic legitimate websites, such as banking or social media platforms, deceiving users into divulging sensitive information. Traditional detection mechanisms, including blacklist-based, rule-based, and anomaly-based systems, have proven inadequate against increasingly sophisticated phishing techniques. These challenges have led researchers to explore more adaptive and robust solutions, particularly machine learning (ML)-based approaches, which excel in identifying patterns and adapting to new attack methods.

Recent studies have demonstrated the effectiveness of ML in addressing the challenges posed by phishing. For instance, M. M. Bala Krishna's research highlights the growing vulnerability of websites as prime targets for phishing attacks due to advancements in internet technologies. Traditional security measures, such as firewalls and antivirus software, primarily focus on technical vulnerabilities but fail to address phishing's exploitation of human weaknesses. Bala Krishna's work emphasizes that ML-based anomaly detection systems are particularly well-suited for combating phishing because of their dynamic capabilities to identify "zero-day" attacks by analyzing patterns and distinguishing between legitimate and malicious activities. The study proposed a fusion classifier that combines the outcomes of two priority-based ML algorithms, achieving an impressive accuracy of 97% on a dataset from the UCI Machine Learning Repository.

Further research in phishing detection explores innovative techniques, such as image visualization and feature extraction from malicious URLs. Jibrilla Tanimu's work introduces an ML-based model that leverages the visual characteristics of website code to identify fraudulent sites. This approach, combined with URL feature extraction, provides a multi-layered detection mechanism capable of addressing phishing attempts that evade traditional systems. Similarly, Rascha Zieni, Luisa Massari, and Maria Carla Calzarossa propose a dual approach that integrates advanced image visualization and URL feature analysis, highlighting the effectiveness of combining multiple methodologies for robust phishing detection.

The rapid increase in phishing attacks has resulted in significant financial losses and raised concerns over data confidentiality and privacy. Attackers continually adapt their strategies to bypass existing safeguards, underscoring the need for innovative and adaptive solutions. Machine learning, with its ability to analyze complex patterns and adapt to evolving threats, has emerged as a cornerstone of modern phishing detection. This paper explores the application of ML in phishing detection, utilizing eight different algorithms and three datasets for a comparative analysis. The results demonstrate exceptional performance, showcasing the effectiveness of ML models in combating phishing with greater accuracy and adaptability. By integrating techniques such as fusion classifiers, image visualization, and URL feature extraction, this research advances the state of the art in phishing detection and provides a promising direction for future cybersecurity efforts.

III) METHODOLOGY

The suggested framework for phishing website identification follows a systematic and structured approach to ensure accuracy and reliability in detection. The process begins with selecting an appropriate phishing dataset, which serves as the foundation for training and evaluating the machine learning models. These datasets typically contain a mix of phishing



and legitimate website data, providing a balanced representation of real-world scenarios. The quality and relevance of the dataset are crucial, as they significantly impact the performance of the classifiers.

Once the dataset is chosen, feature selection methods are applied to identify the most significant attributes that contribute to distinguishing phishing websites from legitimate ones. Techniques such as **Principal Component Analysis (PCA)**, **Chi-Square**, and **Recursive Feature Elimination (RFE)** are employed for this purpose. PCA reduces dimensionality by identifying the most influential principal components, simplifying the data without losing critical information. The Chi-Square test assesses the independence between features and the target variable, highlighting attributes with the strongest correlation. RFE, on the other hand, iteratively removes less significant features while retaining the most impactful ones, ensuring that the model is trained on only the most relevant data. These feature selection methods not only enhance the performance of classifiers but also reduce computational overhead by minimizing unnecessary data.

After determining the significant features, the data is split into **training and testing sets**, typically in an **80:20 ratio**. The training set, comprising 80% of the data, is used to train the machine learning models, while the remaining 20% serves as the testing set to evaluate the model's performance. This split ensures that the models are trained on a large portion of the data while leaving sufficient samples for unbiased testing.

To assess the effectiveness of the classifiers, **Performance Evaluation Measures** are implemented. Accuracy provides a general measure of the classifier's correctness, while precision and recall evaluate its performance in identifying phishing websites specifically.

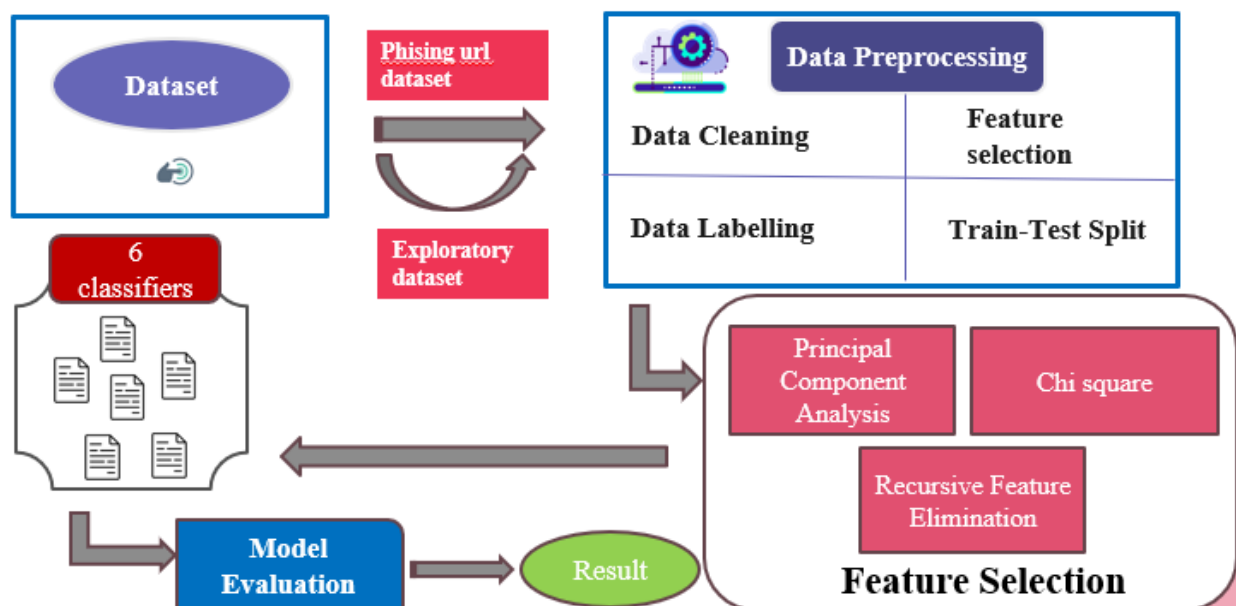


FIGURE 3 WORKFLOW OF PROPOSED METHODOLOGY

3.1 DATA PREPROCESSING

Data preprocessing is a crucial step in any machine learning project, as it ensures the dataset is clean, structured, and ready for analysis. For this study, the dataset was sourced from Kaggle, containing 11,054 rows and 32 columns. The dataset is well-structured, with no null values, eliminating the need to address missing data challenges. This allows for a seamless transition to exploratory data analysis (EDA) and feature selection processes.

To identify the most effective independent variables, statistical techniques such as Principal Component Analysis (PCA), Recursive Feature Elimination (RFE), and Chi-Square tests were employed. These methods help pinpoint features that have a significant impact on phishing website detection while removing redundant or irrelevant attributes, optimizing the dataset for model training and improving overall accuracy.



TABLE 3.1 DATASET

| SI NO | FEATURES | DESCRIPTION |
|-------|--------------|---|
| 1. | UsingIP | If the website's address is an IP address instead of a domain name, it could be a phishing website. |
| 2. | LongURL | A URL that is excessively long, often containing unnecessary characters or parameters. Phishing websites often use long URLs to obfuscate their true purpose. |
| 3. | ShortURL | A shortened URL, typically generated using URL shortening services like Bitly or TinyURL. While convenient, short URLs can sometimes mask malicious websites. |
| 4. | Symbol@ | The presence of the "@" symbol in the URL, which is typically used for email addresses. Its presence in a URL can be a red flag. |
| 5. | Redirect | A process where a website automatically redirects the user to another URL. Phishing websites often use redirects to hide their true location. |
| 6. | PrefixSuffix | The presence of unusual prefixes or suffixes in the URL, which might indicate a phishing website. |
| 7. | SubDomain: | A subdomain is a part of a domain name that comes before the main domain (e.g., "mail" in "mail.google.com"). Phishing websites might use subdomains to appear legitimate. |
| 8. | HTTPS | Hypertext Transfer Protocol Secure. It encrypts communication between the user's browser and the website, making it more secure. However, phishing websites can sometimes use HTTPS to appear legitimate. |
| 9. | DomainReg | The registration date of the domain. Newer domains might be more likely to be associated with phishing websites. |
| 10. | StatusBar | This likely refers to the browser's status bar, which can display information about the website. Phishing websites might try to manipulate the status bar to appear legitimate. |

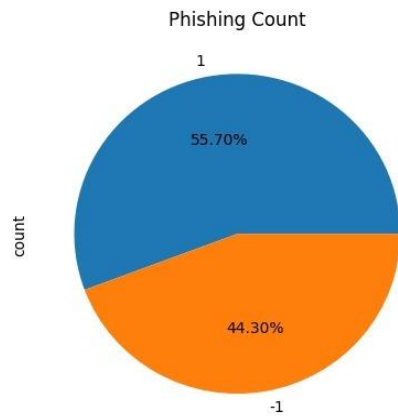


FIGURE 3.1 A) BALANCED DATASET

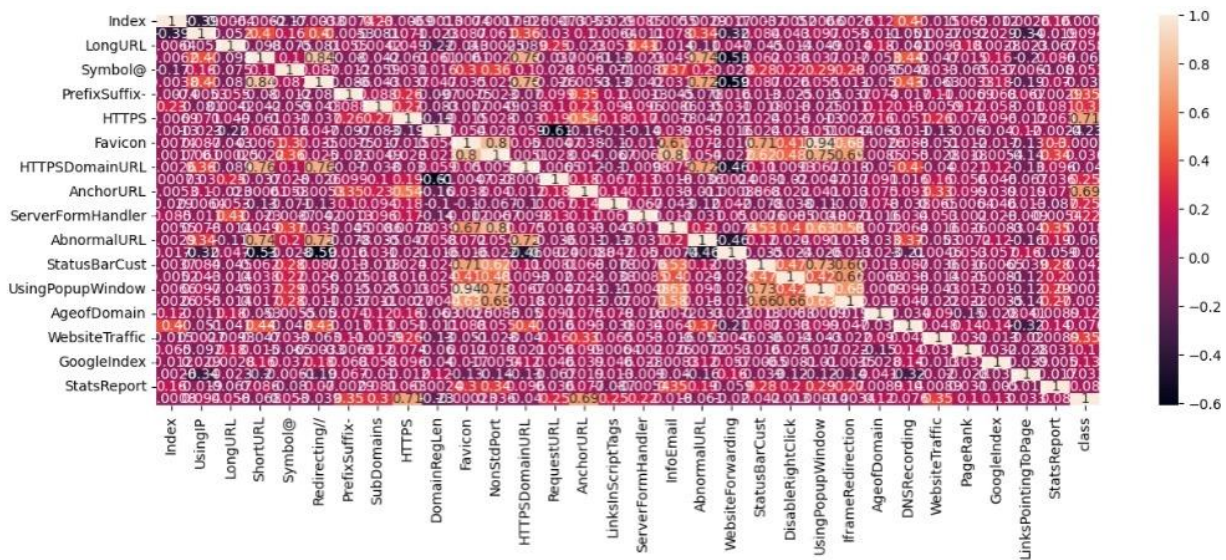


FIGURE 3.1 B) HEAT MAP

3.2 FEATURE SELECTION METHODS

During the feature selection procedure, the characteristics are whittled down to those that pertain to the dependent variable, either manually or automatically. This step is performed because it has a significant impact on the performance of the model in terms of the time required to construct it and its level of accuracy. Due to the fact that they force the model to learn on irrelevant data, irrelevant dataset characteristics may have a negative influence on training. Due to the fact that irrelevant data acts as noise, poor feature selection will result in unreliable accuracy. When features are picked, there is a reduction in overfitting, an increase in accuracy, and a decrease in the training time required. Through the use of feature selection, it is feasible to reduce the dimensions of a dataset.

3.2.1 MANUAL FEATURE SELECTION

Id is one of the features of the dataset. However, it simply represents numbering of each cell of the entries. This does not have any effect on the outcome of the model prediction. This particular feature is quickly removed.

3.2.2 AUTOMATIC FEATURE SELECTION

During this research, Principal Component Analysis (PCA), Recursive Feature Elimination (RFE), and Chi-Square



method was implemented . By employing this method, one can gain a deeper knowledge of the relationship between the dependent and independent variables.

3.2.2.1 PRINCIPAL COMPONENT ANALYSIS (PCA)

Principal Component Analysis (PCA) is a dimensionality reduction technique widely used in machine learning and data analysis. It transforms a high-dimensional dataset into a smaller set of uncorrelated variables called principal components, while preserving as much variability (information) in the data as possible. PCA works by identifying the directions (or axes) in which the variance in the dataset is maximized and projecting the data onto these new axes. This technique is particularly useful for reducing the complexity of large datasets, minimizing redundancy, and improving computational efficiency.

3.2.2.2 RECURSIVE FEATURE ELIMINATION (RFE)

Recursive Feature Elimination (RFE) is a feature selection technique that works by recursively removing the least important features from a dataset to identify the most relevant ones for model training. It is particularly effective in reducing the complexity of a dataset and improving the performance of machine learning models.

RFE operates by fitting a model (commonly linear or tree-based) and assigning weights or importance scores to features. Features with the lowest importance are eliminated in each iteration, and the model is refitted with the remaining features. This process continues until the desired number of features is reached.

3.2.2.3 CHI SQUARE

The Chi-Square test is a statistical method used to determine if there is a significant association between categorical features and a target variable. It is suitable for categorical features and a categorical target.

3.3 FEATURE ENGINEERING

Since the dataset is balanced, strategies for feature engineering that address the issue of an unbalanced dataset are not covered. In contrast, normalization approaches allow the dataset to be contained inside a specific, closed range, such as $[0,1]$, $[-1,1]$, etc., so that it is equitable and has a greater chance of being successfully predicted. In this study, the interval $[-1,1]$ is used .

3.4 MODEL DEVELOPMENT

Since this research utilize Logistic Regression (LR), K-Nearest Neighbors (KNN), Support Vector Machine (SVM), Decision Tree (DT), Random Forest (RF), and Artificial Neural Network (ANN). This section documents the important facts about these models. Similarly, the steps required to build Flask web application is also documented. Feature selection is an important step in machine learning, where irrelevant or redundant features are removed to improve model performance, reduce overfitting, and decrease computational costs. Let's dive into some popular feature selection methods, including Principal Component Analysis (PCA), Recursive Feature Elimination (RFE), and Chi-Square method.

3.4.1 LOGISTIC REGRESSION

When the dependent variable is dichotomous (binary), the supervised learning technique known as logistic regression can be used. Unlike the previously described linear regression, real-world issues usually call for non-linear models. Real-world issues can be logistic, quadratic, or exponential. Potential results from logistic regression are categorical as opposed to quantitative. We may forecast categorical results using logistic regression, such as "yes or no will be a phishing website" or "0 or 1 will not be a phishing website ". Decisions in the context of this study usually boil down to a straightforward yes/no choice. We can generate far more basic predictions with logistic regression, like will this URL at all be a phishing website? Regarding the gradient descent in this investigation.

3.4.2 K-NEAREST NEIGHBORS (KNN)



K-Nearest Neighbors (KNN) is one of the simplest Machine Learning algorithms for regression and classification challenges. KNN algorithms utilize existing data to classify new data points based on similarity measures (e.g. distance function). Classification is chosen by a vote of the majority of adjacent communities. The data is assigned to the class that is the closest neighbor. Increasing the number of nearest neighbors, or the value of k , can improve accuracy.

3.4.3 ARTIFICIAL NEURAL NETWORK (ANN)

Each algorithm in the hierarchy nonlinearly alters its input and produces a statistical model as a result. Iterations will keep going until the results are precise and usable While feature extraction is a time-consuming process in machine learning, ANN only employs weights to produce the best accurate prediction. The learning rate of this model was accelerated using the Adam optimization algorithms over the course of 400 iterations and an 80-element hidden layer size.

3.4.4 SUPPORT VECTOR MACHINE (SVM)

Support Vector Machine (SVM) is a supervised machine learning algorithm primarily used for classification and regression tasks. It works by finding the optimal hyperplane that best separates data points belonging to different classes in a high-dimensional space. SVM is particularly effective for handling linearly separable data and can also manage non-linear data by using kernel functions, such as the polynomial, radial basis function (RBF), or sigmoid kernel. SVM excels in high-dimensional spaces and is robust against overfitting, especially when the number of features exceeds the number of samples.

3.4.5 RANDOM FOREST(RF)

Random Forest is an ensemble learning method that builds multiple decision trees during training and aggregates their predictions to improve accuracy and reduce overfitting. Each tree in the forest is trained on a randomly sampled subset of the data (with replacement), and at each split, a random subset of features is considered. This randomness ensures diversity among the trees and helps mitigate the limitations of individual decision trees, such as overfitting. Random Forest is highly versatile and performs well on both classification and regression tasks while being resistant to noise and capable of handling missing data effectively.

3.4.6 DECISION TREE(DT)

A Decision Tree is a supervised learning algorithm used for classification and regression tasks. It works by splitting the dataset into subsets based on feature values, creating a tree-like structure where each internal node represents a feature, each branch represents a decision rule, and each leaf node represents an outcome or class. Decision Trees are intuitive, easy to interpret, and can handle both numerical and categorical data. However, they are prone to overfitting, especially with deep trees, which is why techniques like pruning or ensemble methods (e.g., Random Forest) are often applied to improve their performance.

IV)EVALUATION METRICS

Well-known measuring metrics such as recall, accuracy, f1-score, roc-curve, etc. will be used.

4.1) Confusion Matrix

Confusion matrix is used to understand what the model is getting correctly and what it is getting wrongly. Figure 4.1 shows how confusion matrix

True Positive (TP) and True Negative (TN): Real value at the diagonal axis

False Negative (FN): Other values along the horizontal

False Positive (FP): Other value present in the vertical column

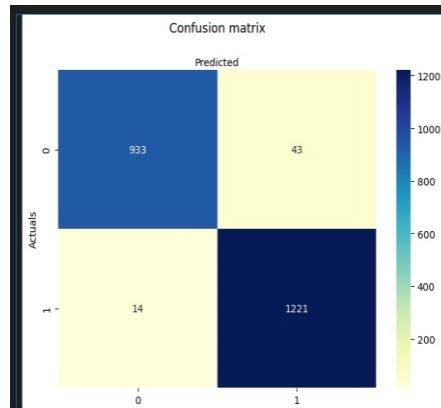


FIGURE 4.1) CONFUSION MATRIX

4.1.2 Classification Report

This is the recall, precision, accuracy, f1-score, and support as explained below.

a) Recall: The TP divided by how many times the classifier predicted that class.

FORMULA:

$$\frac{TP}{TP + FP}$$

b) Precision: Number of correct predictions divided by how many occurrences of that class were in the test data.

FORMULA:

$$\frac{TP}{TP + FP}$$

c) F1-score: The weighted harmonic means of the precision and recall values for the test is the F1-score. A high f1-score indicates that the precision is more balanced

FORMULA:

$$\frac{2 * Precision * recall}{Precision + recall}$$

d) Support: The total number of true response samples in the class.

FORMULA:

$$TP + FN$$

e) Accuracy: Combination of TP and TN divided TP, TN, FP and FN. The higher the better

FORMULA:

$$\frac{TP + TN}{TP + TN + FP + FN}$$

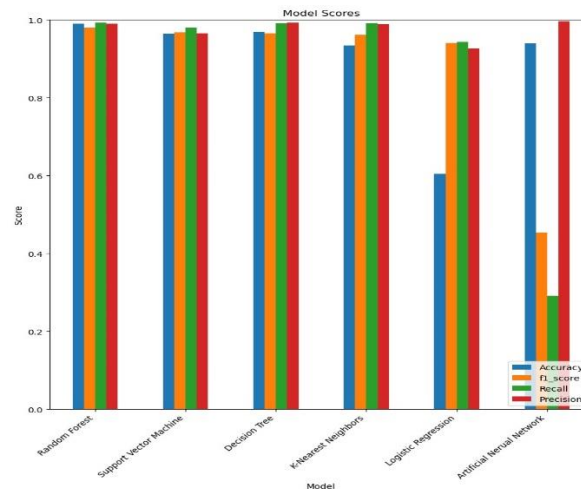


FIGURE 4.1.2 RESULT

V) DEPLOYMENT AND USER INTERFACE

For real-time accessibility, the project's user interface was developed using Flask, a lightweight framework for building web applications. Flask provides an intuitive environment for deploying machine learning models, allowing users to easily interact with the heart disease prediction tool. The interface is designed to be simple and accessible, ensuring that users regardless of their technical background can input their health metrics and receive immediate predictions related to heart disease.

The deployment structure of the application involves a server-side processing model, which takes care of handling user inputs, running the prediction algorithms, and delivering results back to the user in real time. This ensures the system can be accessed online, providing a seamless experience for users to get immediate feedback on their heart disease risk. The use of Flask enables smooth interaction between the user interface and the underlying machine learning model, maintaining efficient performance and scalability for multiple users.

This deployment setup makes the prediction tool not only accessible but also practical, as users can quickly assess their health status without needing to install or configure complex software.

VI) CONCLUSION

The increasing frequency of phishing attacks remains a major concern for cybersecurity researchers, as traditional signature-based methods often fail to detect newly created phishing sites. Machine learning-based systems, leveraging large and diverse datasets, offer a promising solution to accurately identify and categorize phishing websites. This study evaluates several models, including Logistic Regression (LR), K-Nearest Neighbors (KNN), Support Vector Machine (SVM), Decision Tree (DT), Random Forest (RF), and Artificial Neural Network (ANN), optimizing the dataset for each. The study concludes by integrating the top-performing model, as determined by performance metrics, into a user-friendly Flask web application. This integration provides end-users with real-time access to the model, a valuable contribution often overlooked in previous research that focuses solely on model performance. The study not only highlights the importance of integrating machine learning models into practical applications but also paves the way for future research to refine these models and address emerging phishing tactics.

REFERENCES

- [1] R. Alabdian, "future internet Phishing Attacks Survey: Types, Vectors, and Technical Approaches", doi: 10.3390/fi12100168.
- [2] L. Barlow, G. Bendiab, S. Shiaeles, and N. Savage, "A Novel Approach to Detect Phishing Attacks using Binary Visualisation and Machine Learning," in Proceedings - 2020 IEEE World Congress on Services, SERVICES 2020, Oct. 2020, pp. 177–182. doi: 10.1109/SERVICES48979.2020.00046.



- [3] M. A. Adebowale, K. T. Lwin, E. Sánchez, and M. A. Hossain, "Intelligent web-phishing detection and protection scheme using integrated features of Images, frames and text," *Expert Systems with Applications*, vol. 115. Elsevier Ltd, pp. 300–313, Jan. 01, 2019. doi: 10.1016/j.eswa.2018.07.067.
- [4] C. Iuga, J. R. C. Nurse, and A. Erola, "Baiting the hook: factors impacting susceptibility to phishing attacks," *Human-centric Computing and Information Sciences*, vol. 6, no. 1, Dec. 2016, doi: 10.1186/s13673-016-0065-2.
- [5] M. Almseidin, A. M. Abu Zuraiq, M. Al-kasassbeh, and N. Alnidami, "Phishing detection based on machine learning and feature selection methods," *International Journal of Interactive Mobile Technologies*, vol. 13, no. 12, pp. 71–183, 2019, doi: 10.3991/ijim.v13i12.11411.
- [6] E. S. Gualberto, R. T. de Sousa, T. P. B. de Vieira, J. P. C. L. da Costa, and C. G. Duque, "From Feature Engineering and Topics Models to Enhanced Prediction Rates in Phishing Detection," *IEEE Access*, vol. 8, pp. 76368–76385, 2020, doi: 10.1109/ACCESS.2020.2989126.
- [7] R. S. Rao, • Alwyn, and R. Pais, "Detection of phishing websites using an efficient feature-based machine learning framework," *Neural Computing and Applications*, vol. 31, doi: 10.1007/s00521-017-3305-0.
- [8] A. Kumar Jain and B. B. Gupta, "Towards detection of phishing websites on client-side using machine learning based approach," vol. 68, pp. 687–700, 2018, doi: 10.1007/s11235-017-0414-0.
- [9] A. Basit, M. Zafar, A. R. Javed, and Z. Jalil, "A Novel Ensemble Machine Learning Method to Detect Phishing Attack," Nov. 2020. doi: 10.1109/INMIC50486.2020.9318210.
- [10] M. Volkamer, K. Renaud, B. Reinheimer, and A. Kunz, "User experiences of TORPEDO: TOoltip-poweRed Phishing Email DetectiOn," *Computers and Security*, vol. 71, pp. 100–113, Nov. 2017, doi: 10.1016/j.cose.2017.02.004.