

# URLGuard: A Holistic Hybrid Machine Learning Approach for Phishing Detection

Mr. Siddheshwar Patil

Mr. Rohit Mehetre

Mr. Abhishek Kubde

Mr. Shubham Gaikwad

under supervision of

Dr. Pradip Paithane

**Department of Artificial Intelligence and Data Science ,  
Vidya Pratishthans Kamalnayan Bajaj Institute of Engineering and Technology,  
Vidyanagari, Baramati-413133**

**Academic Year: 2023-24**

# Contents

- Introduction
- Problem Statement
- Objectives
- Architecture
- Mathematical Modeling
- Algorithm Development
- Results
- Performance Evaluation Parameters
- GUI
- Conclusion
- References

# Introduction

- The phishing attacks poses a significant threat to individuals and organizations alike. To combat the escalating threat of phishing attacks on the internet effectively, the comprehensive URL-based phishing detection system utilizing a variety of machine learning algorithms.
- As we explore the intricacies of such a system, we uncover the pivotal role it plays in ensuring a secure digital environment for users, ultimately contributing to the ongoing battle against cyber threats.
- It introduces the concept of hybrid machine learning, emphasizing its potential to enhance detection accuracy by combining the strengths of multiple algorithms.

- Phishing detection using hybrid machine learning models outlines the growing threat of phishing attacks in the digital era and highlights the need for advanced, accurate detection methods.
- a hybrid machine learning model, which includes Decision Tree, Random Forest, and XGB is employed to safeguard against phishing URLs.
- The proposed approach undergoes evaluation with key metrics such as precision(0.96), accuracy(0.96), recall(0.97),and F1-score(0.96). Results demonstrate that the proposed method surpasses the state of art.

# Problem statement

To combat the escalating threat of phishing attacks on the internet effectively, the comprehensive URL-based phishing detection system utilizing a variety of machine learning algorithms.

# Objectives

- To create a ongoing detection system to spot phishing URLs.
- To introduce a hybrid model (DT+RF+XGB) for enhanced phishing detection.
- Discuss and compare evaluation parameters to demonstrate the superiority of the proposed approach.

# Architecture

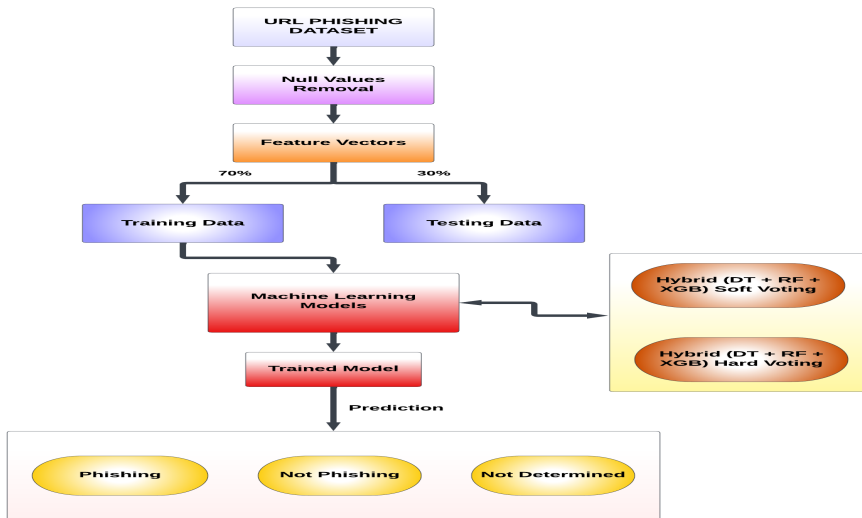


Figure: Architecture

- **Random Forest:** Random Forest is an ensemble method using decision trees on random data subsets, with randomness in sampling and feature selection, preventing overfitting and improving robustness. It handles high-dimensional data and offers feature importance scores.
- Ensemble Method

$$(F(x_t) = \frac{1}{B} \sum_{i=0}^B F_i(x_t) \quad (1)$$

- $F(x_t)$  is the output of the random forest for the input  $x_t$ ,
- $B$  is the number of trees in the random forest,
- $F_i(x_t)$  is the output of the  $i$ -th tree in the forest for the input  $x_t$ ,
- $\sum$  denotes the summation over the specified range.



**XGB:** Extreme Gradient Boosting (XGBoost) enhances gradient boosting for supervised learning, known for scalability and efficiency. It builds an ensemble of weak models, typically decision trees, optimizing an objective function to correct previous models errors.

- Objective Function

$$\text{Objective} = \sum_{i=1}^n \text{Loss}(\hat{y}_i, y_i) + \sum_{k=1}^K \Omega(f_k) \quad (2)$$

- $n$  is the number of samples.
- $\hat{y}_i$  is the predicted value for the  $i$ -th sample.
- $y_i$  is the true label for the  $i$ -th sample.
- $K$  is a parameter representing some set of values.

- Gradient Boosting

$$\hat{y}(t) = \sum_{k=1}^t f_k(x) \quad (3)$$

- $t$  is the iteration or boosting round.
- $f_k(x)$  is the prediction of the  $k$ -th model at input  $x$ .

- Regularization

$$\Omega(f_k) = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2 + \alpha \sum_{j=1}^T |w_j| \quad (4)$$

- $T$  is the number of leaves in the tree.

- $w_j$  is the weight assigned to the  $j$ -th leaf.
- $\gamma$ ,  $\lambda$ , and  $\alpha$  are regularization parameters.

**Decision tree:** A decision tree is a tree-like model in data analysis and machine learning. It splits data based on key attributes, creating an interpretable flowchart. It handles categorical and numerical data, effective for classification and regression tasks.

- **Entropy ( $E(S)$ ):**  $E(S) = -\sum_{i=1}^c p_i \log_2(p_i)$
- **Conditional Entropy ( $E(T, X)$ ):**

$$E(T, X) = - \sum_{c \in X} p(c) \cdot E(c) \quad (5)$$

- **Information Gain (IG(T, X)):**

$$IG(T, X) = E(T) - E(T, X) \quad (6)$$

- $S$  is a set,
- $c$  represents the classes or values in the set,
- $p_i$  is the probability of occurrence of element  $i$  in set  $S$ ,
- $T$  is a set, and  $X$  is an attribute or feature,
- $p(c)$  is the probability of occurrence of value  $c$  in set  $X$ ,
- $E(c)$  represents the entropy of the subset of  $T$  associated with value  $c$ ,  
and
- $\sum$  denotes the summation over the specified range.

# Algorithm Development

---

**Algorithm 1:** Hybrid Machine Learning Model (RF + DT + XGB)

---

**Input:** Training data  $X_{train}$ ,  $y_{train}$ , Test data  $X_{test}$ ,  $y_{test}$

**Output:** Predicted class label

```
train_and_test_hybrid_model( $X_{train}$ ,  $y_{train}$ ,  $X_{test}$ ,  $y_{test}$ , num_of_epochs,
param_grid):
    initialize rf_model  $\leftarrow$  RandomForestClassifier()
    initialize xgb_model  $\leftarrow$  XGBClassifier()
    initialize dt_model  $\leftarrow$  DecisionTreeClassifier()
    initialize estimators  $\leftarrow$  [('rf', rf_model), ('xgb', xgb_model),
('dt', dt_model)]
    initialize best_model  $\leftarrow$  None
    for epoch  $\leftarrow$  1 to num_of_epochs do
        initialize grid_search  $\leftarrow$ 
GridSearchCV(estimator=VotingClassifier(estimators),
param_grid=param_grid, cv=5)
        grid_search.fit( $X_{train}$ ,  $y_{train}$ )
        best_model  $\leftarrow$  grid_search.best_estimator_
        best_model.fit( $X_{train}$ ,  $y_{train}$ )
        // Evaluate on the test set and print results
        validation_results  $\leftarrow$  best_model.evaluate( $X_{test}$ ,  $y_{test}$ )
        print(Accuracy: validation_results)
    end
    return best_model.predict( $X_{test}$ )
```

---

Figure: Proposed Algorithm

- The dataset was gathered and saved as a CSV file from the well-known Kaggle dataset repository, which offers benchmark datasets for academic use.
- The collection included 33 attributes and 11054 items that were taken from over 11,000 websites.
- Datasetlink: <https://www.kaggle.com/code/eswarchandt/website-phishing/input?select=phishing.csv>

# Results

Algorithm	Accuracy	Precision	Recall	F1 score
Decision Tree (DT)	0.94	0.94	0.94	0.94
Random Forest (RF)	0.96	0.92	0.93	0.92
Extreme Gradient Boosting (XGB)	0.96	0.96	0.96	0.96
DT + RF + XGB (Soft Voting)	0.96	0.96	0.96	0.96
DT + RF + XGB (Hard Voting)	0.96	0.97	0.96	0.96

Figure: Result Comparison

# Results

Sr.no	Author	Algorithm	Evaluation Parameters
1	Proposed algorithm	Hybrid Machine Learning Model(DT+RF+XGB)	<ul style="list-style-type: none"> <li>• Accuracy= 0.96</li> <li>• precision=0.96</li> <li>• Recall= 0.97</li> <li>• F1-Score=0.96</li> </ul>
2	ABDULKARIM1, MOBEEN SHAHROZ 2, KHABIB MUSTOFA 1, SAMIR BRAHIM BEL-HAOUARI 3, ANDS. RAMANA KUMAR JOGA(2023)	Hybrid Machine Learning Model(SVC + DT +LR)r	<ul style="list-style-type: none"> <li>• Accuracy= 0.95</li> <li>• Precision=0.95</li> <li>• Recall=0.96</li> <li>• F1-Score=0.95</li> </ul>
3	Sonowal, G., Kuppusamy, K(2020)	five-layer phishing detection model called PhiDMA	<ul style="list-style-type: none"> <li>• Accuracy= 0.92</li> <li>• Precision=0.91</li> <li>• Recall=0.90</li> <li>• F1-Score=0.90</li> </ul>



# Results

4	KangLengChiew a, ChoonLinTan a,, KokSheik Wong b, Kelvin S.C.Yongc,Wei KingTionga(2019)	Hybrid Ensem- ble Feature Selec- tion(HEFS),Random ForestClassifier	<ul style="list-style-type: none"> <li>• Accuracy= 0.94</li> </ul>
5	YONG FANG, CHENG ZHANG, CHENG HUANG, LIANG LIU,ANDYUE YANG (2019)	Deep learning model named THEMIS(advanced version of RCNN)	<ul style="list-style-type: none"> <li>• Accuracy= 0.99</li> <li>• Precision=0.99</li> <li>• Recall=0.99</li> <li>• F1-Score=0.99</li> </ul>
6	Ozgun Koray Sahin- goz, Ebubekir Bu- ber b,OnderDemirb, BanuDiri(2017)	Using Random Forest with NLP-based fea- tures	<ul style="list-style-type: none"> <li>• Accuracy= 0.97</li> <li>• Precision=0.97</li> <li>• Sensitivity=0.99</li> <li>• F1-Score=0.98</li> </ul>

Figure: Result Comparison

# Results

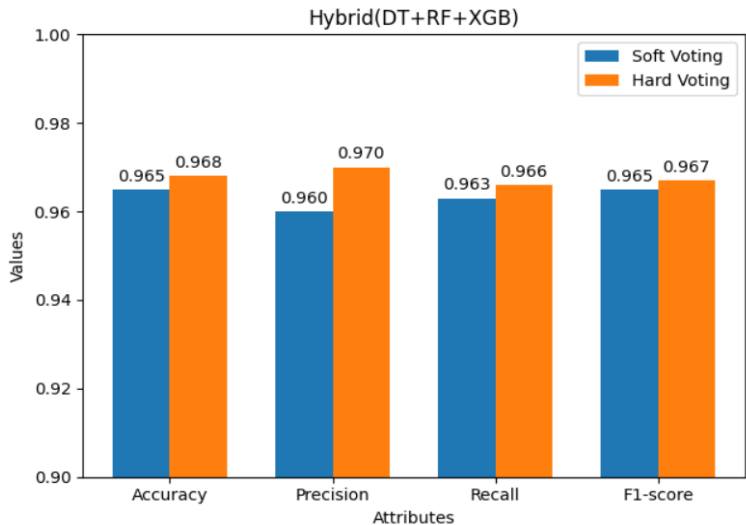


Figure: Results

UsingIP -1 ▾	Subdomains -1 ▾	RequestURL -1 ▾	WebsiteForwarding -1 ▾	DNSRecording -1 ▾
LongURL -1 ▾	HTTPS -1 ▾	AnchorURL -1 ▾	StatusBarCust -1 ▾	WebsiteTraffic -1 ▾
ShortURL -1 ▾	DomainRegLen -1 ▾	LinksInScriptTags -1 ▾	DisableRightClick -1 ▾	PageRank -1 ▾
Symbol@ -1 ▾	Favicon -1 ▾	ServerFormHandle -1 ▾	UsingPopupWindo w -1 ▾	GoogleIndex -1 ▾
Redirecting -1 ▾	NonStdPort -1 ▾	InfoEmail -1 ▾	IframeRedirection -1 ▾	LinksPointingToPag e -1 ▾
Prefixsuffix -1 ▾	HTTPSDomainURL -1 ▾	AbnormalURL -1 ▾	AgeofDomain -1 ▾	StatsReport -1 ▾

Predict

Figure: GUI

# Conclusion

- <https://github.com/Rohit9860/Phishing-detection>
- The hybrid machine learning models ( DT+RF+XGB ) boost phishing detection accuracy, paving the way for future cybersecurity research.

- A. Karim, M. Shahroz, K. Mustofa, S. B. Belhaouari and S. R. K. Joga, "Phishing Detection System Through Hybrid Machine Learning Based on URL," in IEEE Access, vol. 11, pp. 36805-36822, 2023, doi: 10.1109/ACCESS.2023.3252366.
- V. Shahrivari, M. M. Darabi, and M. Izadi, "Phishing detection using machine learning techniques," 2020, arXiv:2009.11116.
- O. K. Sahingoz, E. Buber, O. Demir, and B. Diri, "Machine learning based phishing detection from URLs," Expert Syst. Appl., vol. 117, pp. 345–357, Mar. 2019

- Y. Fang, C. Zhang, C. Huang, L. Liu, and Y. Yang, “Phishing email detection using improved RCNN model with multilevel vectors and attention mechanism,” IEEE Access, vol. 7, pp. 56329–56340, 2019.
- K. L. Chiew, C. L. Tan, K. Wong, K. S. C. Yong, and W. K. Tiong, “A new hybrid ensemble feature selection framework for machine learning-based phishing detection system,” Inf. Sci., vol. 484, pp. 153–166, May 2019.



THANK YOU