# PHISHING WEBSITE DETECTION USING MACHINE LEARNING ALGORITHMS

BY; MANJARI JAYAN

INTMCA S9

**SEMINAR**

GUIDE - Ms. JETTY BENJAMIN
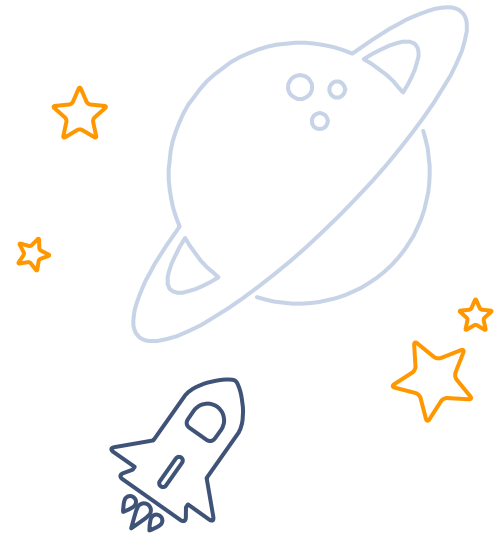
# CONTENTS

# TITLE

## PHISHING WEBSITE DETECTION USING MACHINE LEARNING ALGORITHMS

# ABSTRACT

- Phishing attacks are now more complex, making it harder to differentiate between real and fake email links and websites.

- This paper focuses on using supervised learning, particularly the Gradient Boosting method, to mitigate security risks through intelligent solutions.

- The aim is to create a robust classifier for phishing websites by analyzing their characteristics and determining the best feature combination for model training.

# ABSTRACT cont.

- The proposed system achieves a remarkable accuracy of 97.4% using 32 key features.

- This research strengthens phishing attack detection, providing a valuable resource to enhance user safety and inspire stronger anti-phishing defenses.

# INTRODUCTION

❖ Phishing is the most commonly used social engineering and cyber attack.

❖ Phishers deceive unsuspecting online users into disclosing sensitive information for fraudulent purposes.

❖ In order to avoid getting phished,

  ❖ users should have awareness of phishing websites

  ❖ maintain a blacklist of phishing websites which requires the knowledge of website being detected as phishing.

  ❖ detect them in their early appearance, using machine learning algorithms.

# OBJECTIVE

❖ A phishing website is a common social engineering method that mimics trustful uniform resource locators (URLs) and webpages.

❖ The objective of this project is to train machine learning models on the dataset created to predict phishing websites.

❖ Both phishing and legitimate URLs of websites are gathered to form a dataset and from them required URL and website content-based features are extracted.

❖ The performance level of each model is measures and compared.

# LITERATURE REVIEW

**[ 1 ] Alyssa Anne Ubing, "Phishing Website Detection: An Improved Accuracy through Feature Selection and Ensemble Learning", October 2019**

- Feature selection algorithm is employed and integrated with an ensemble learning methodology,

- Based on majority voting

- Compared with different classification models including Random forest, Logistic Regression, Prediction model etc.

- An accuracy rate between 70% and 92.52%

# LITERATURE REVIEW cont.

**[ 2 ] Mrs. Vaneeta M, Pratik N N, Prajwal D, Pradeep K S, and Suhas Kakade K, "Machine Learning Techniques for Phishing Website Detection", June 2020**

- ❖ Developed a phishing websites detection technique

- ❖ Based on machine learning classifiers with a features selection method

- ❖ Classification algorithms used are Artificial Neural Network, Random Forest and Support Vector Machine.

- ❖ An accuracy rate above 93%.

**[ 3 ]. Ammar Odeh, Ismail Keshta, and Eman Abdelfattah, "An Adaptive Boosting Approach for Phishing Detection", March 2021**

- ❖ Proposed model uses an adaptive boosting approach

- ❖ Consists of multiple classifiers to increase the model's accuracy.

- ❖ Produced an accuracy of approximately 95%.

[ 4 ].  **Amani Alswailem, Bashayr Alabdulla, Nora Alrumayh, Aram Alsedrani, "Detecting Phishing Websites Using Machine Learning", July 2019**

- The system is based on a machine learning method, particularly supervised learning.

- Selected the Random Forest technique .

- Choose the better combination of them to train the classifier.

- An accuracy of 97.7%

[ 5 ]. **Lizhen Tang, Qusay H. Mahmoud**, **" A Deep Learning-Based Framework for Phishing Website Detection", December 2021**

- compared multiple machine learning models using several datasets.

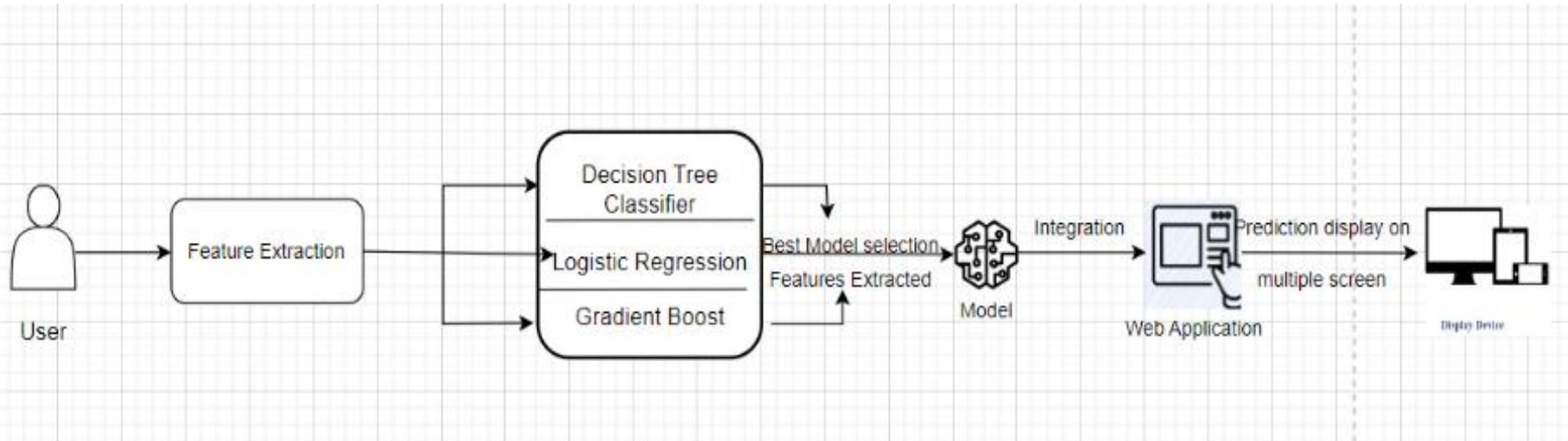- RNN-GRU model obtained the highest accuracy of 97.8%.

# METHODOLOGY

❖ **An extensive review was done on existing works of literature and machine learning models on detecting phishing websites to best decide the classification models to solve the problem of detecting phishing websites.**
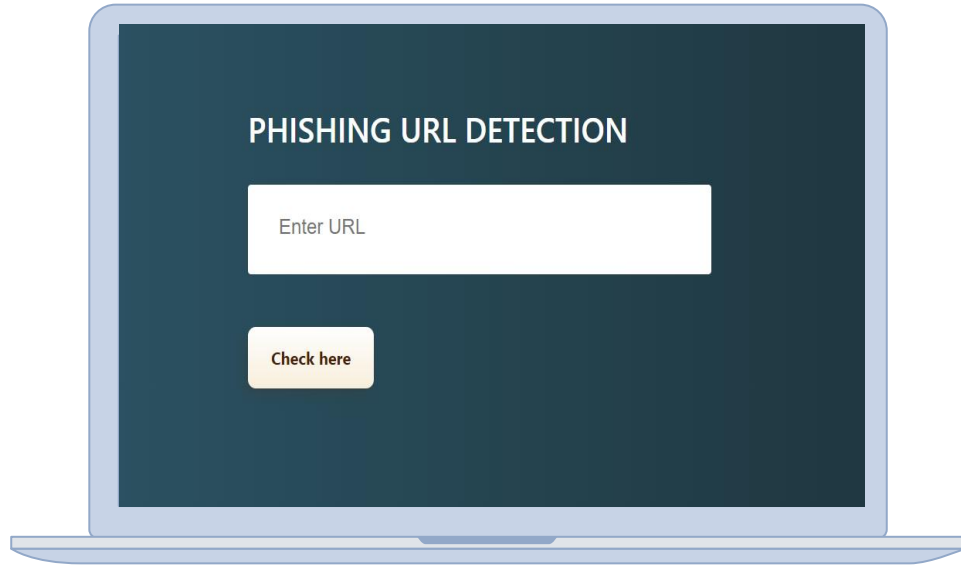


Data Collection → Feeding Data → Extracting Features → Training the Algorithm → Creating a Model

**Architectural design of proposed system**

# IMPLEMENTATION

PHISHING URL DETECTION

Enter URL

Check here

# 1. Loading Data:

❖ The dataset is borrowed from Kaggle, https://www.kaggle.com/eswarchandt/phishing-website-detector .

❖ A collection of website URLs for 11000+ websites. Each sample has 30 website parameters and a class label identifying it as a phishing website or not (1 or -1).

❖ The overview of this dataset is, it has 11054 samples with 31 features.

# 2. Familiarizing with Data & EDA:

❖ In this step, few dataframe methods are used to look into the data and its features.

3. Visualizing the data:

❖ Few plots and graphs are displayed to find how the data is distributed and the how features are related to each other.
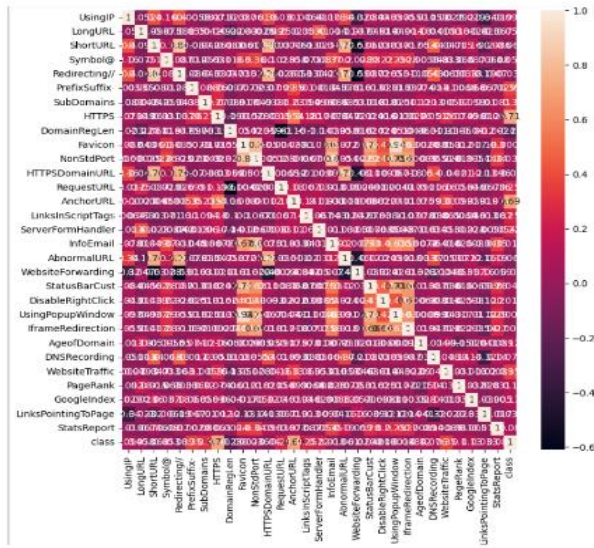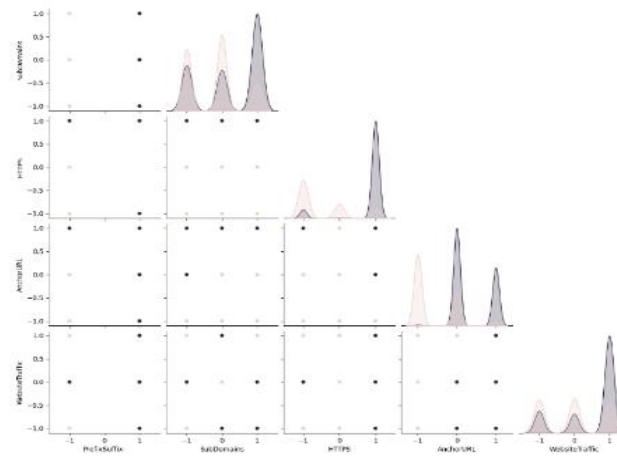


**FIG. 1. CORRELATION HEATMAP**
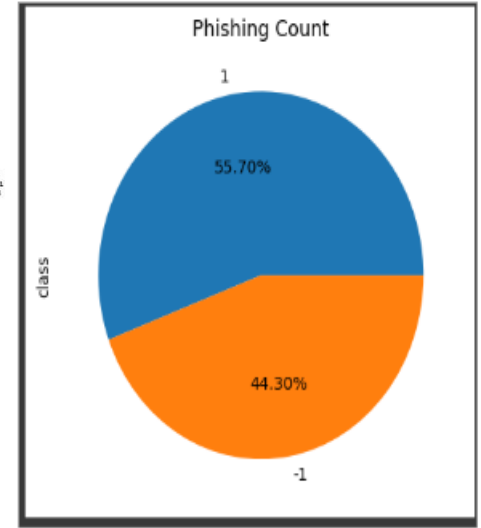
**FIG. 2. PAIRPLOT FOR PARTICULAR FEATURES**

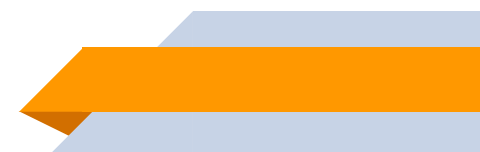**FIG.3. PHISHING COUNT IN PIE CHART**

## 4. Splitting the Data:

The data is split into train & test sets, 80-20 split.

## 5. Model Building & Training:

Supervised machine learning is one of the most commonly used and successful types of machine learning.

The supervised machine learning models (regression) considered to train the dataset in this notebook are:

- Logistic Regression
- Decision Tree
- Gradient Boosting

# 6. Comparision of Models

To compare the models performance, a dataframe is created. The columns of this dataframe are the lists created to store the results of the model.

```
[39]  #creating dataframe
      result = pd.DataFrame({ 'ML Model' : ML_Model,
                              'Accuracy' : accuracy,
                              'f1_score' : f1_score,
                              'Recall'   : recall,
                              'Precision': precision,
      })
```
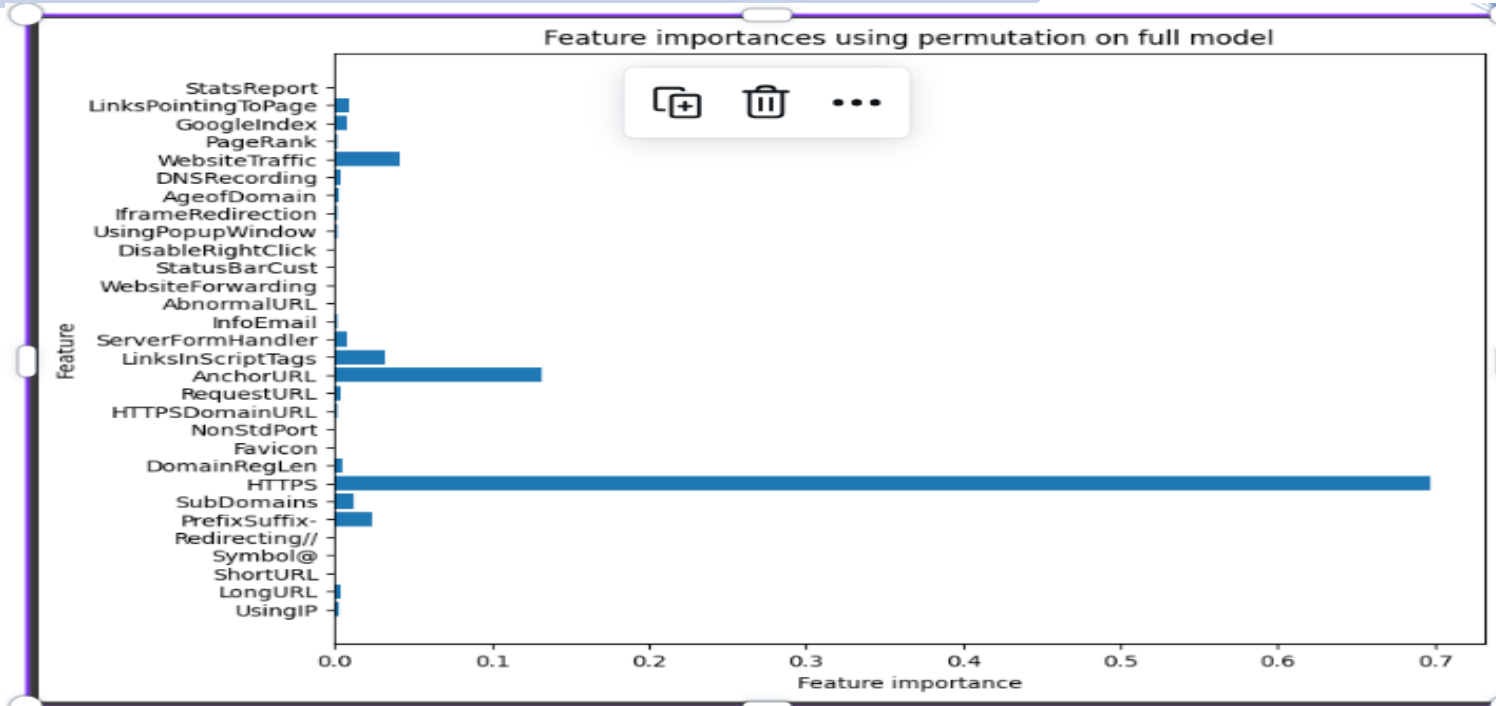
# RESULT-MODEL EVALUATION

❖ The models are evaluated, and the considered metrics are accuracy, f1_score, recall, and precision.

❖ The below figure shows the training and test dataset accuracy, f1_score, recall, and precision for each respective model.

| | ML Model | Accuracy | f1_score | Recall | Precision |
|---|---|---|---|---|---|
| 0 | Gradient Boosting Classifier | 0.974 | 0.977 | 0.994 | 0.986 |
| 1 | Decision Tree | 0.958 | 0.963 | 0.991 | 0.993 |
| 2 | Logistic Regression | 0.934 | 0.941 | 0.943 | 0.927 |

❖ Based on the evaluation, the Gradient Boosting Classifier outperforms the other models in terms of accuracy, with an accuracy score of 0.974.

❖ Therefore, the Gradient Boosting Classifier model has been saved for further usage.

# Feature Importance in the Model



Feature importances using permutation on full model

# INTEGRATION OF MODEL TO WEB APPLICATION

# INTEGRATION OF MODEL TO WEB APPLICATION

# CONCLUSION

❖ The final conclusion on the Phishing dataset is that the some feature like "HTTPS", "AnchorURL", "WebsiteTraffic" have more importance to classify URL is phishing URL or not.

❖ Gradient Boosting Classifier currectly classify URL upto 97.4% respective classes and hence reduces the chance of malicious attachments.

# REFERENCES

**[ 1 ]** Alyssa Anne Ubing, "Phishing WebsiteDetection: An Improved Accuracy throughFeature Selection and Ensemble Learning", October 2019

**[ 2 ]** Mrs. Vaneeta M, Pratik N N, Prajwal D, Pradeep K S, and Suhas Kakade K, "Machine Learning Techniques for Phishing Website Detection", June 2020

**[ 3 ]** Ammar Odeh, Ismail Keshta, and Eman Abdelfattah, "An Adaptive Boosting Approach for Phishing Detection", March 2021

**[ 4 ]** Amani Alswailem, Bashayr Alabdulla, Nora Alrumayh, Aram Alsedrani, "Detecting Phishing Websites Using Machine Learning", July 2019

**[ 5 ]** Lizhen Tang, Qusay H. Mahmoud, " A Deep Learning-Based Framework for   Phishing Website Detection", December 2021

# THANK YOU