



# PHISHING DOMAIN DETECTION (Machine Learning)

## ARCHITECT DOCUMENT

### Project Members:

1. Rishabh
2. Shivansh Srivastava
3. Ashish Diwakar

## **INTRODUCTION:**

Phishing Domain Detection is a technique by which we should be able to predict whether the domain is real or fake.

## **PROBLEM STATEMENT:**

Phishing is a type of fraud in which an attacker impersonates a reputable company or person in order to get sensitive information such as login credentials or account information via email or other communication channels. Phishing is popular among attackers because it is easier to persuade someone to click a malicious link that appears to be authentic than it is to break through a computer's protection measures.

The main goal is to predict whether the domains are real or malicious.

## **APPROACH:**

The classical machine learning tasks like Data Exploration, Data Cleaning, Feature Engineering, Model Building and Model Testing as been done on the project. Tried with different machine learning algorithms such as Logistic Regression, SVM, Gradient Boosting, KNN, Random Forest and found out best fit model in the project as Random Forest.

## **DATASET:**

These data consist of a collection of legitimate as well as phishing website instances. Each website is represented by the set of features which denote, whether website is legitimate or not. Data can serve as an input for machine learning process.

The dataset had two variants of the Phishing Dataset are presented.

Full variant - dataset\_full.csv

- Short description of the full variant dataset:
- Total number of instances: 88,647
- Number of legitimate website instances (labeled as 0): 58,000
- Number of phishing website instances (labeled as 1): 30,647
- Total number of features: 111

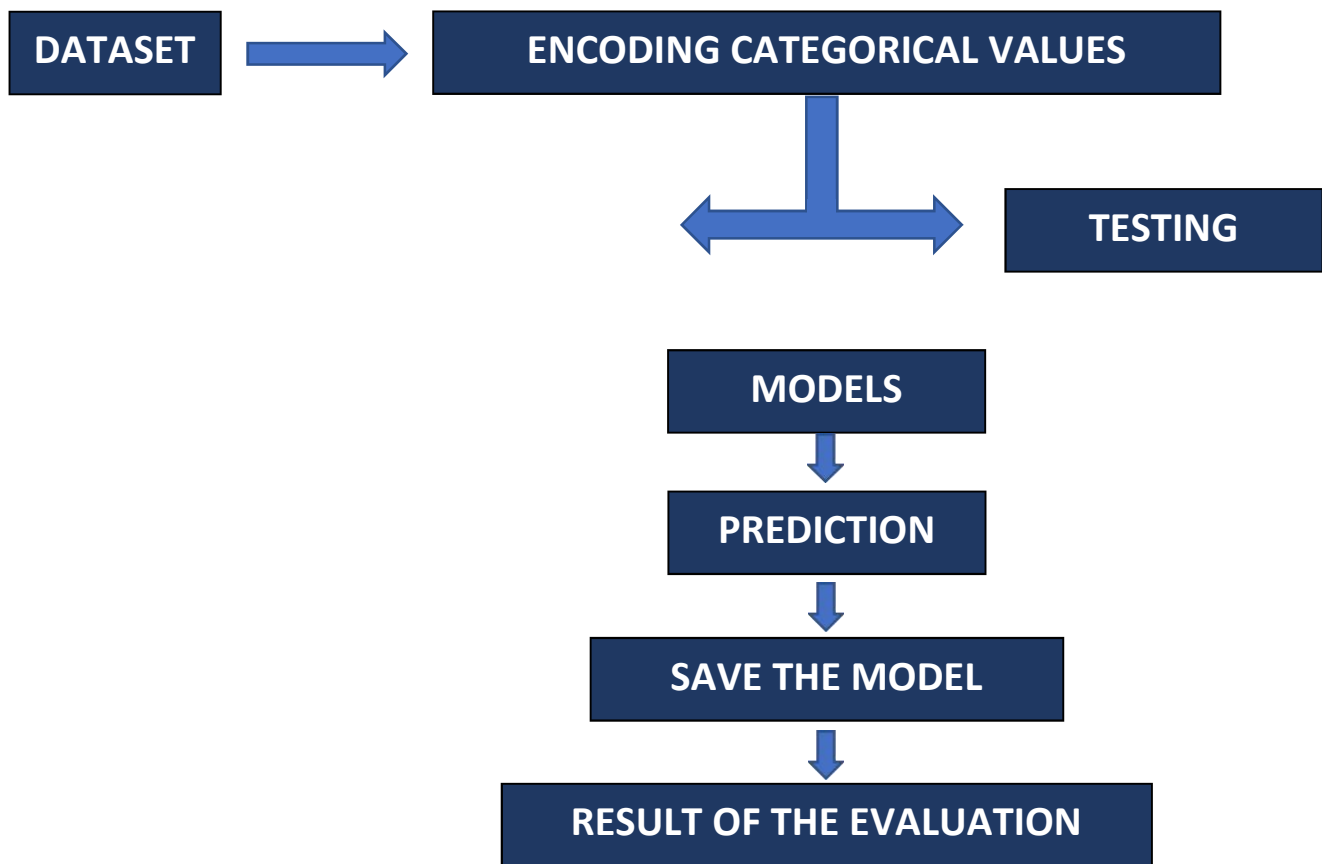
Small variant - dataset\_small.csv

- Short description of the small variant dataset:
- Total number of instances: 58,645
- Number of legitimate website instances (labeled as 0): 27,998
- Number of phishing website instances (labeled as 1): 30,647
- Total number of features: 111

## USER INPUT/OUTPUT FLOW:



## DESIGN FLOW:



## CONCLUSION:

In conclusion, the results obtained from the phishing domain detection using the random forest model have been quite promising. With an accuracy of 98%, precision of 98%, and recall of 97.6%, the model has demonstrated excellent performance in correctly identifying phishing domains. These metrics indicate that the model has a high degree of accuracy in detecting malicious domains and is reliable in predicting whether a domain is phishing or not. Overall, this model can be a valuable tool in protecting users against phishing attacks, and its high accuracy and precision make it a strong candidate for use in real-world scenarios.