

**Team -DataVerse**

# PHISHIUNITO (PHISHING) URL DATASET

DATA SCIENCE MID SEM PROJECT SUBMISSION

[Colab Link](#)

# Dataset Explanation

## DATASET COMPOSITION:

- TOTAL INSTANCES: 235,795
- LEGITIMATE URLs: 134,850
- PHISHING URLs: 100,945
- FEATURES: 54, DERIVED FROM WEBPAGE SOURCE CODE AND URL.
- ASSOCIATED TASK: CLASSIFICATION (PHISHING VS. LEGITIMATE URLs)
- FEATURE TYPES: REAL, CATEGORICAL, INTEGER
- NO MISSING VALUES IN THE DATASET

## KEY FEATURES:

- URL FEATURES: URL LENGTH, DOMAIN, TLD, CHAR CONTINUATION RATE
  - WEBPAGE CONTENT FEATURES: TITLE, HASSOCIALNET, NO POPUP
- LABEL: LABEL (1 = PHISHING, 0 = LEGITIMATE)

# Problem Statement and Importance

- PROBLEM STATEMENT: TO DEVELOP A MODEL THAT CAN ACCURATELY CLASSIFY URLs AS PHISHING OR LEGITIMATE, HELPING TO DETECT AND PREVENT PHISHING ATTACKS.
- IMPORTANCE:
  - PHISHING IS A MAJOR CYBERSECURITY THREAT, OFTEN USED TO STEAL SENSITIVE DATA LIKE PASSWORDS AND FINANCIAL INFORMATION.
  - A ROBUST URL-BASED DETECTION MODEL CAN ACT AS A FRONTLINE DEFENSE, PROVIDING QUICK IDENTIFICATION OF MALICIOUS SITES WITHOUT NEEDING TO VISIT THEM.
- APPLICATION: SUCH A MODEL COULD BE IMPLEMENTED IN BROWSERS, EMAIL FILTERS, AND SECURITY SYSTEMS TO PROTECT USERS FROM PHISHING THREATS.



# Exploratory Data Analyses

## 01 Percentage of unique/total instances in each feature

	count	unique	Distinct %	is_categorical
FILENAME	235795	235795	100.0	NA
URL	235795	235370	99.819759	NA
URLLength	235795.0	482	0.204415	NA
Domain	235795	220086	93.337857	NA
DomainLength	235795.0	101	0.042834	NA
IsDomainIP	235795.0	2	0.000848	Categorical
TLD	235795	695	0.294748	NA
URLSimilarityIndex	235795.0	36360	15.420174	NA
CharContinuationRate	235795.0	898	0.380839	NA
TLDLegitimateProb	235795.0	465	0.197205	NA
URLCharProb	235795.0	227421	96.44861	NA
TLDLength	235795.0	12	0.005089	NA
NoOfSubDomain	235795.0	10	0.004241	NA
HasObfuscation	235795.0	2	0.000848	Categorical
NoOfObfuscatedChar	235795.0	20	0.008482	NA
ObfuscationRatio	235795.0	146	0.061918	NA
NoOfLettersInURL	235795.0	421	0.178545	NA
LetterRatioInURL	235795.0	709	0.300685	NA

We get 5 features that have more than 82% distinct features indicating that they do not contribute much to the target.

FileName :100%  
URL :99.819%  
URLCharProb:96.4%  
Domain : 93.33 %  
Title : 83.91 %

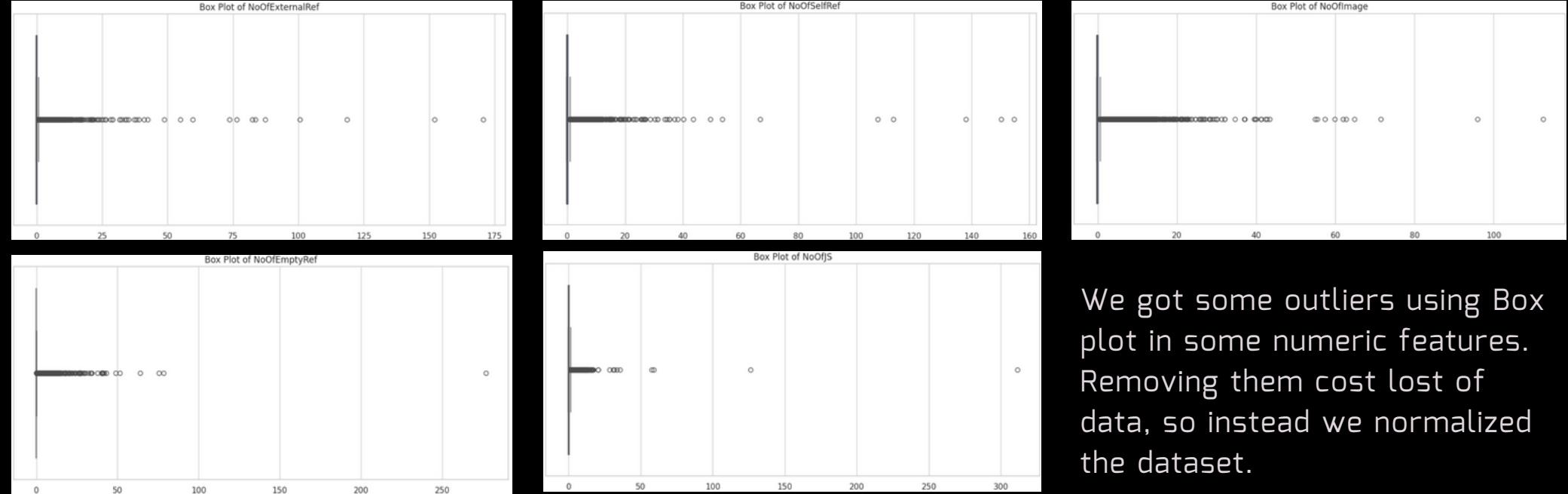
## 02 Categorical features = 21 Unique instances <=5

'IsDomainIP', 'HasObfuscation', 'NoOfQMarkInURL', 'IsHTTPS', 'HasTitle', 'HasFavicon', 'Robots', 'IsResponsive', 'NoOfURLRedirect', 'NoOfSelfRedirect', 'HasDescription', 'HasExternalFormSubmit', 'HasSocialNet', 'HassubmitButton', 'HasHiddenFields', 'HasPasswordField', 'Bank', 'Pay', 'Crypto', 'HasCopyrightInfo', 'label'

## 03 Outlier Detection

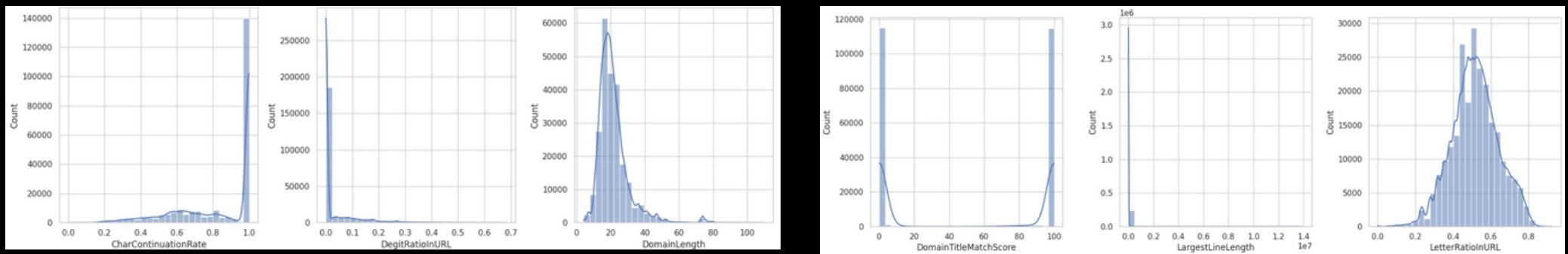
- Total Features: 54
- Missing Values: NA
- Data Type: Object (5) / int64 (41) / float64(10)

Box plot on all the numerical features



We got some outliers using Box plot in some numeric features. Removing them cost lost of data, so instead we normalized the dataset.

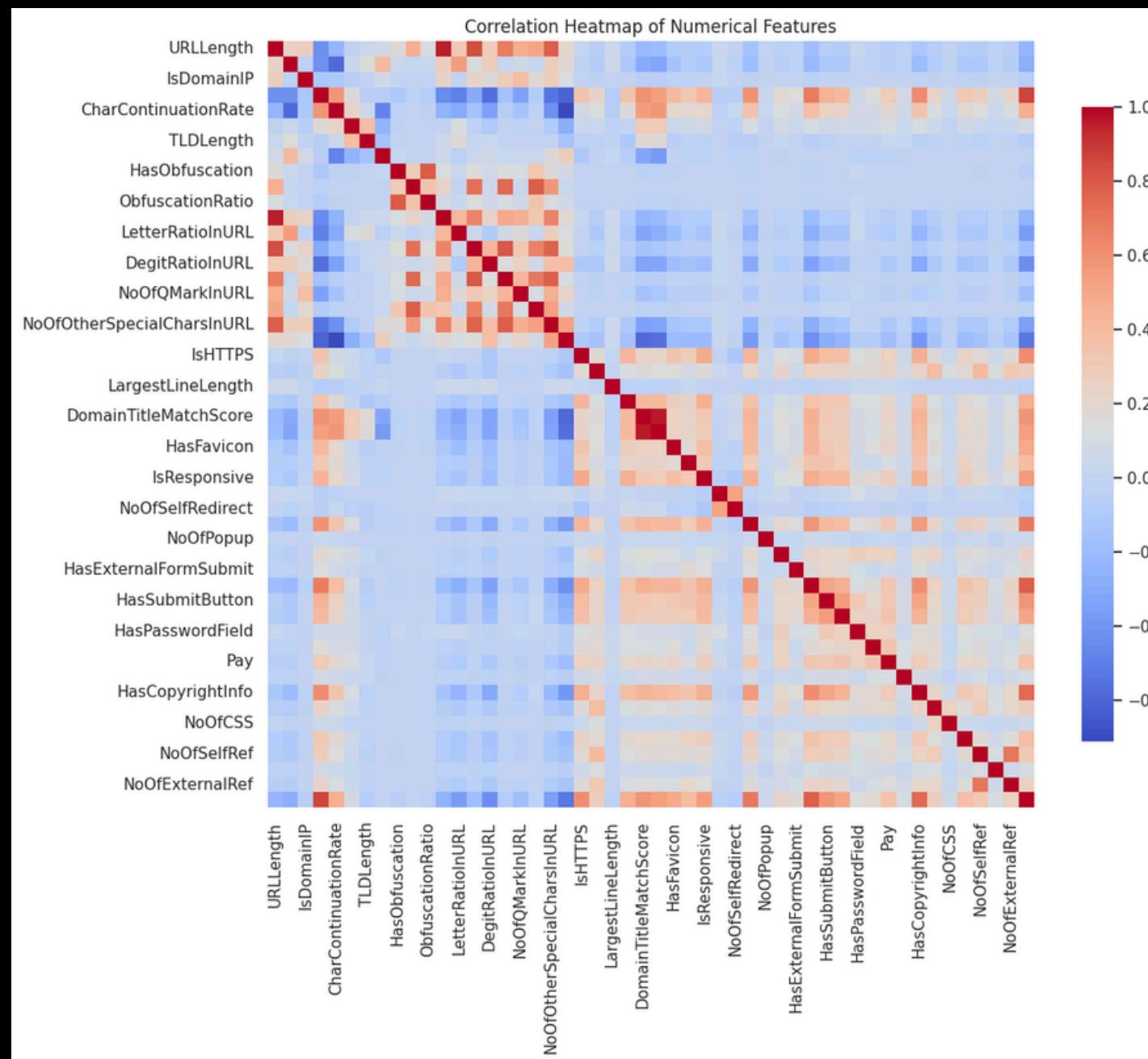
## 04 Data Distribution using histogram



This helped in understanding data distribution and its skewness . Moreover we get categorical features too after visualizing this .

# 05

# Heatmap



## Strong Positive Correlations:

- **URLLength** with **LargestLineLength** and **NoOfQMarkInURL** suggesting that longer URLs may have more special characters and a loner line length.
- **NoOfSelfRedirect** with **NoOfPopup**, indicating that sites with self-redirects may also use pop-ups, possibly a characteristic of phishing sites.

## Strong Negative Correlations:

- **IsDomainIP** with **IsHTTPS** and **HasPasswordField**, indicating that URLs using IP addresses instead of domain names are less likely to use HTTPS or display favicons.
- **HasObfuscation** with **IsHTTPS** and **HasPasswordField** which could suggest that obfuscated URLs are less likely to employ HTTPS and include password fields, perhaps indicative of phishing tendencies.

## Other notable Correlations:

- **DomainTitleMatchScore** and **IsResponsive** have a positive correlation, which may imply that legitimate URLs with a matching title are often responsive, contrasting with some phishing sites that may ignore responsiveness.
- **Pay** and **HasPasswordField** are positively correlated, as payment fields are often found alongside password fields.

# Data Preprocessing

## 1. Encoding categorical features

There are binary and nominal categorical features for example IsDomainIP, TLD. Since binary categorical features are already encoded, we shall apply Label Encoding only on TLD column

## 2. Feature Scaling

Firstly we find all the numerical continuous features excluding categorical data. To those features we apply Standard Scaler function which standardizes those continuous data.

## 3. Handling Imbalanced Data

As the target feature, ie, 'label' might be imbalanced , we use Synthetic Minority Over-sampling Technique (SMOTE) to balance the classes.

# Hypothesis Testing on URL Characteristics

## 01 URL Length Hypothesis

- NULL HYPOTHESIS ( $H_0$ ): THERE IS NO SIGNIFICANT DIFFERENCE IN THE AVERAGE URL LENGTH BETWEEN PHISHING URLs AND LEGITIMATE URLs.
- ALTERNATIVE HYPOTHESIS ( $H_1$ ): THE AVERAGE URL LENGTH DIFFERS SIGNIFICANTLY BETWEEN PHISHING AND LEGITIMATE URLs.

## 02 HTTPS Usage Hypothesis

- NULL HYPOTHESIS ( $H_0$ ): THE PROPORTION OF PHISHING URLs THAT USE HTTPS IS EQUAL TO THE PROPORTION OF LEGITIMATE URLs THAT USE HTTPS.
- ALTERNATIVE HYPOTHESIS ( $H_1$ ): THE PROPORTION OF PHISHING URLs USING HTTPS DIFFERS FROM THAT OF LEGITIMATE URLs

## 03 Number of Subdomains Hypothesis

- NULL HYPOTHESIS ( $H_0$ ): THERE IS NO SIGNIFICANT DIFFERENCE IN THE AVERAGE NUMBER OF SUBDOMAINS BETWEEN PHISHING URLs AND LEGITIMATE URLs.
- ALTERNATIVE HYPOTHESIS ( $H_1$ ): THE AVERAGE NUMBER OF SUBDOMAINS IS SIGNIFICANTLY DIFFERENT BETWEEN PHISHING AND LEGITIMATE URLs.

# Experiments Conducted to Validate Hypothesis Tests

01

## URL LENGTH HYPOTHESIS

- TEST USED: INDEPENDENT T-TEST
- RESULT: SIGNIFICANT DIFFERENCE (T-STATISTIC: -101.05, P-VALUE: 0.0)
- CONCLUSION: SIGNIFICANT DIFFERENCE IN URL LENGTH BETWEEN PHISHING AND LEGITIMATE URLs.

02

## HTTPS USAGE HYPOTHESIS

- TEST USED: CHI-SQUARE TEST
- RESULT: SIGNIFICANT DIFFERENCE ( $\chi^2$ -STATISTIC: 87,486.79, P-VALUE: 0.0)
- CONCLUSION: SIGNIFICANT DIFFERENCE IN HTTPS USAGE BETWEEN PHISHING AND LEGITIMATE URLs.

03

## NUMBER OF SUBDOMAINS HYPOTHESIS

- TEST USED: INDEPENDENT T-TEST
- RESULT: SIGNIFICANT DIFFERENCE (T-STATISTIC: -2.66, P-VALUE: 0.0079)
- CONCLUSION: SIGNIFICANT DIFFERENCE IN THE AVERAGE NUMBER OF SUBDOMAINS BETWEEN PHISHING AND LEGITIMATE URLs.

# Summary of Hypothesis Testing Insights

- SIGNIFICANT FEATURES IDENTIFIED:
  - URL LENGTH: PHISHING URLs HAVE DIFFERENT LENGTHS COMPARED TO LEGITIMATE URLs.
  - HTTPS USAGE: HTTPS IS USED DIFFERENTLY BETWEEN PHISHING AND LEGITIMATE URLs.
  - NUMBER OF SUBDOMAINS: PHISHING URLs TEND TO HAVE FEWER SUBDOMAINS COMPARED TO LEGITIMATE URLs.

OVERALL CONCLUSION: THESE FEATURES PROVIDE STRONG DISCRIMINATORY POWER FOR DISTINGUISHING PHISHING URLs FROM LEGITIMATE ONES, SUPPORTING THE USE OF THESE FEATURES IN A CLASSIFICATION MODEL.

*Thank  
You*