# Malicious
# URL Detector

**Marco Fresco & Fabio Fadda**

Artificial Intelligence for Cybersecurity
Master's Degree in Cybersecurity
Academic Year 2023/24

# Goal Definition

**Drive-By Downloads**
Execute malicious code on the victim's system

**Malware Distribution**
Host or redirect to sites that distribute malware

**Phishing Attacks**
Lead users to fake websites that mimic legitimate ones

**Real-Time Protection**
Identifying and blocking access to harmful URLs

**Early Threat Detection**
Preventing users from interacting with dangerous content

**Reducing Attack Surface**
Prevent potential entry points for cyberattacks

03
04
02
05
01
06

# Common Characteristics of Malicious URLs

**Misspelled Domain Names**

Cybercriminals register domains that are intentionally similar to well-known websites

**IP Addresses**

Raw IP Addresses instead of domain names to bypass domain registration requirements

**Long, Random Strings**

Attempt to obfuscate the true purpose of the URL

**Unusual Characters**

Used to confuse users or evade detection

**Lack of HTTPS**

Phishing sites may not have valid SSL certificates (HTTPS)

**Overuse of Subdirectories**

A technique to obscure the final destination

# Data Gathering

**PhishTank**
Data and information about phishing on the Internet.

**Kaggle**
Platform for data science competitions

Feature Extraction

Length  Special Chars  HTTPS  IPv4 Address  Depth  CCTLD

# Heuristics Confrontation

## Special Chars

'@' Symbol: 93.827%
'-' Symbol: 90.126%
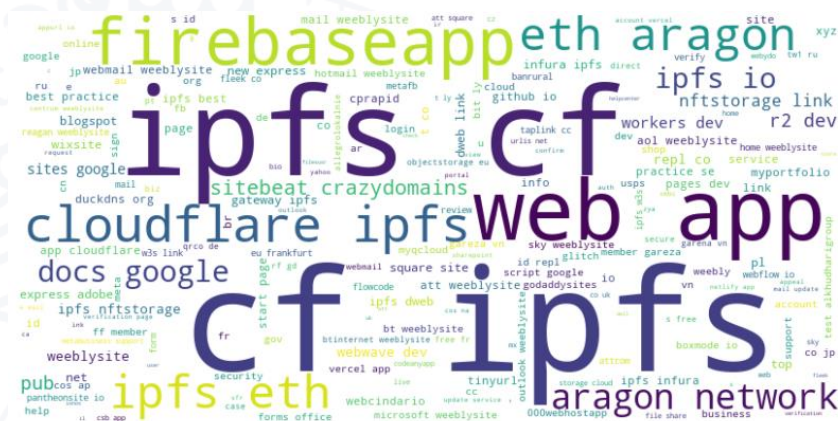'_' Symbol: 53.846%

## IPv4 Address

100% Malicious URLs

Reference: Youness Mourtaji, Mohammed Bouhorma, Daniyal Alghazzawi, Ghadah Aldabbagh, Abdullah Alghamdi, "Hybrid Rule-Based Solution for Phishing URL Detection Using Convolutional Neural Network".
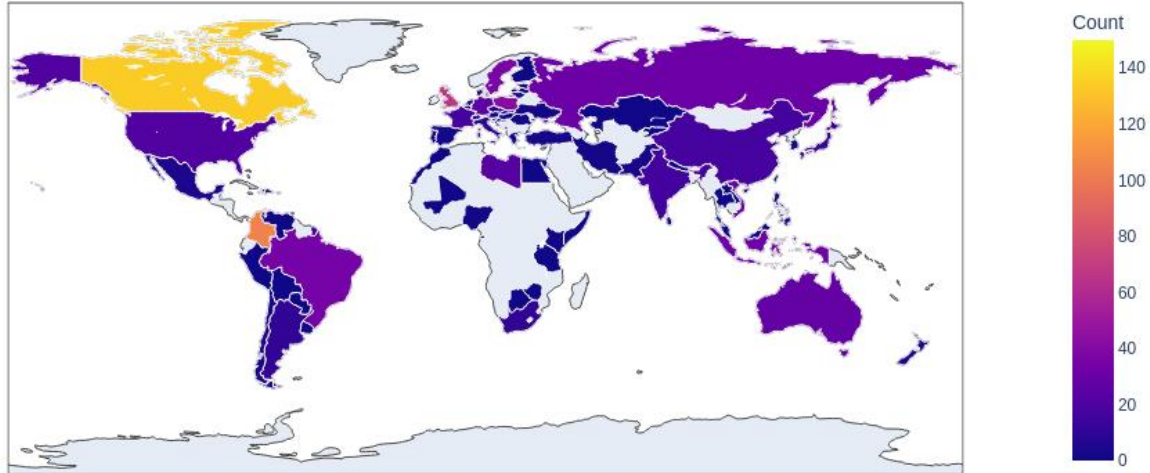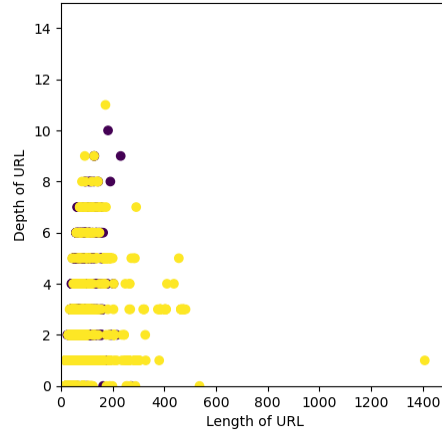
# Preliminary Data Exploration



**Benign**

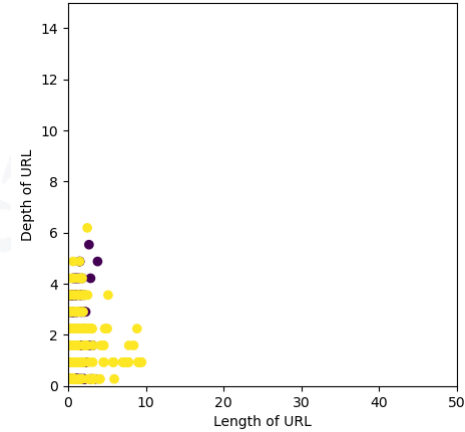**Malicious**

# Preliminary Data Exploration



Distribution of Country-Code Top Level Domains

# Normalization process



Z-Score

# Data Mining

**Decision Tree Classifier**

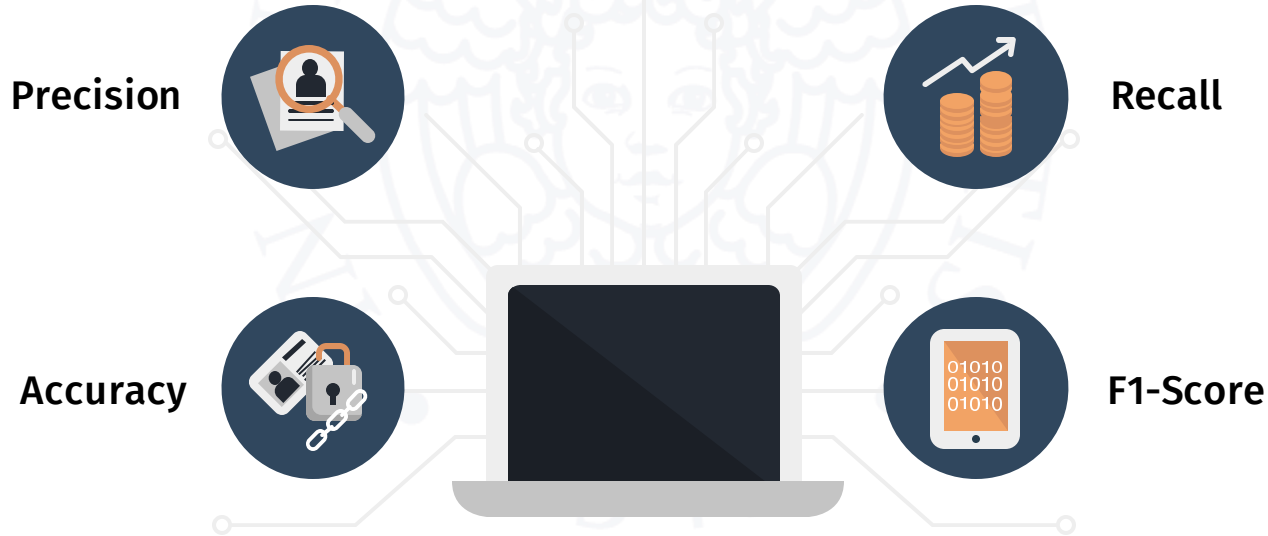**Random Forest Classifier**
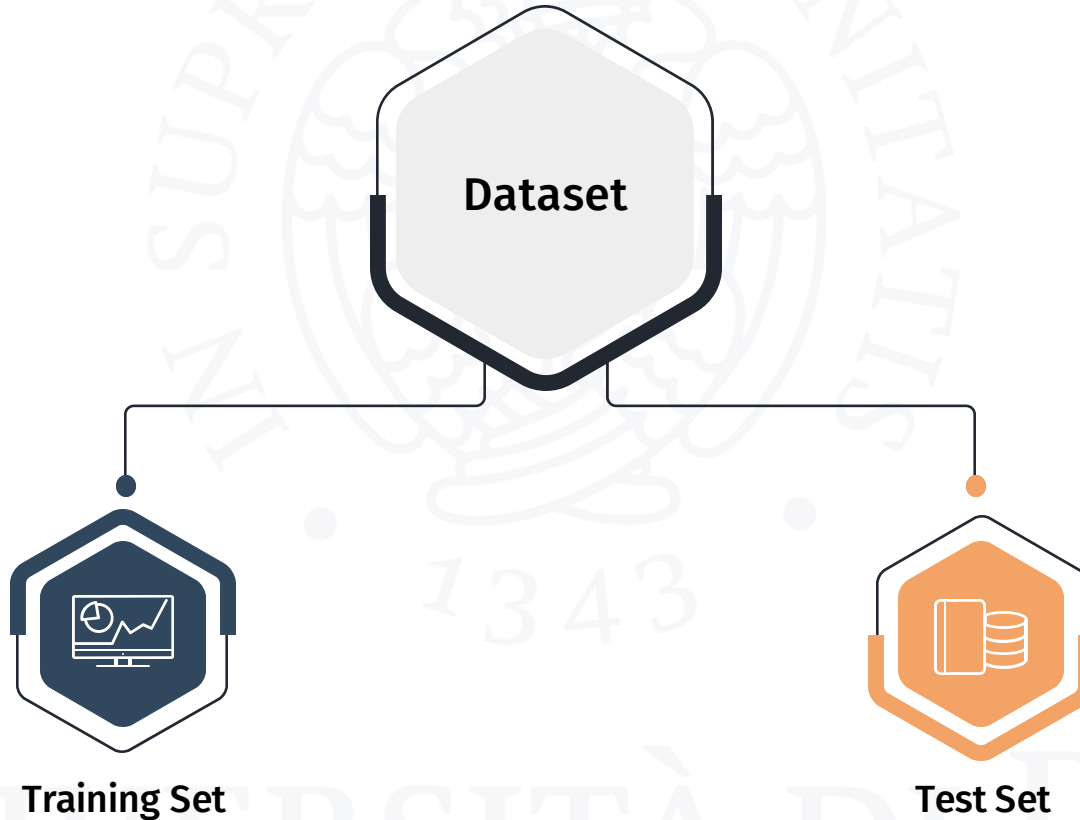
**Logistic Regression**
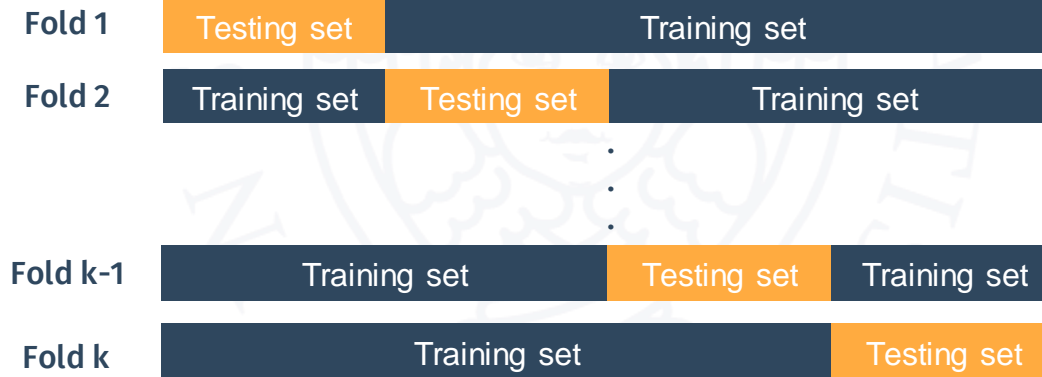
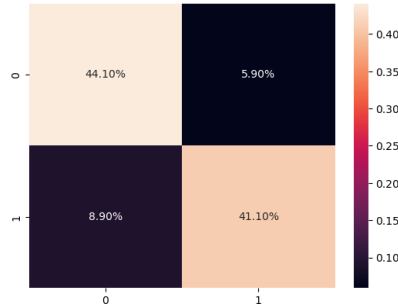**Gaussian Naive Bayes**

**Linear SVC**

# Performance Evaluation

Precision

Recall

Accuracy

F1-Score

# Holdout Method



Dataset

Training Set

Test Set

# Stratified K-Fold

# Results – Hold out

### Decision Tree Classifier



|       | 0      | 1      |
|-------|--------|--------|
| 0     | 44.10% | 5.90%  |
| 1     | 8.90%  | 41.10% |

**85.20%**

### Random Forest Classifier



|       | 0      | 1      |
|-------|--------|--------|
| 0     | 44.10% | 5.90%  |
| 1     | 8.05%  | 41.95% |

**86.05%**

### Logistic Regression



|       | 0      | 1      |
|-------|--------|--------|
| 0     | 43.45% | 6.55%  |
| 1     | 17.10% | 32.90% |

**78.35%**

### Linear SVC



|       | 0      | 1      |
|-------|--------|--------|
| 0     | 44.95% | 5.05%  |
| 1     | 20.00% | 30.00% |

**74.95%**

### Gaussian Naive Bayes



|       | 0      | 1      |
|-------|--------|--------|
| 0     | 49.95% | 0.05%  |
| 1     | 44.15% | 5.85%  |

**55.80%**

# Results – Accuracy

| | Model | Accuracy | Accuracy with SKF | Accuracy with Feature Selection |
|---|---|---|---|---|
| 0 | DecisionTreeClassifier | 0.8520 | 0.842 | 0.7770 |
| 1 | RandomForestClassifier | 0.8605 | 0.852 | 0.7785 |
| 2 | LogisticRegression | 0.7635 | 0.736 | 0.6655 |
| 3 | LinearSVC | 0.7495 | 0.740 | 0.6700 |
| 4 | GaussianNB | 0.5580 | 0.561 | 0.5520 |

# Decision Tree Representation

# System Improvement



HTML encoding

JavaScript code

Request parallelization

Updated/Tested Dataset

# Thanks For The Attention!