# Malicious URL Detection

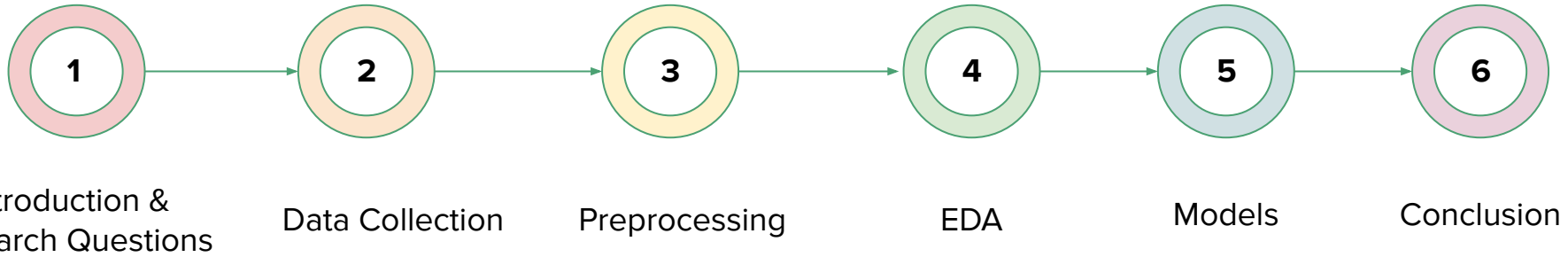## — Using Machine Learning

Data Mining & Predictive Analytics - BUDT758T

Group: Alekya Ghanta, Weian Shi, Xingjian Qin, Yanpu Wang, Jie Gao

# Agenda



1 — Introduction & Research Questions

2 — Data Collection

3 — Preprocessing

4 — EDA

5 — Models

6 — Conclusion

# What is a URL?

Domain name

Path

Scheme

2nd level domain

Directory

Parameter

https://www.example.com/category/webpage.html?id=12345

Subdomain

Top-level domain

Filename

# What is a URL?

**Scheme/Protocol:** Specifies the method used to access the resource. The most common protocols used in URLs are HTTP (Hypertext Transfer Protocol) and HTTPS (HTTP Secure)

**Domain:** Identifies one or more IP addresses on the internet. Domain is the name used to identify a website. Rather than remembering the IP, users can simply type the domain name to access the website.

**Path:** Part of the URL that follows the domain name and identifies the specific file that the resource is located in. The path can include multiple directories, subdirectories, and file names, separated by ("/").

http://www.jenkov.com/books/jquery/index.html

Protocol        Domain          Resource
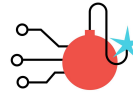                                Path

# Threats associated with URLs

**Types of Cyber Attacks**

Malware
Ransomware
Man in the Middle attack
Zero-day exploit
Phishing
DNS Tunneling
SQL injection
DoS and DDoS attack
Social engineering
XSS attacks
Cryptojacking

There are several threats associated with URLs that users should be aware of. However, we are limiting the scope of this project to **Malware and Phishing** attacks.

**Phishing:** Attackers can **create fake URLs** that appear to be legitimate ones to trick users into giving sensitive information such as **login credentials** or **bank details**.

**Malware:** URLs can be used to deliver viruses, trojans to a user's device for **stealing personal information, damaging computer's hardware**. Malware can be in the URL itself or on the website that the URL leads to.

# Research Area of Focus

Develop an accurate and efficient **machine learning model** that can classify URLs as safe, malware or phishing with high accuracy.

Identify **potentially dangerous URLs** and prevent users from accessing them, thereby reducing the **risk of cybersecurity** threats.

Compare the **performance** of various **algorithms**, as well as the **impact of different features** on the classification accuracy.
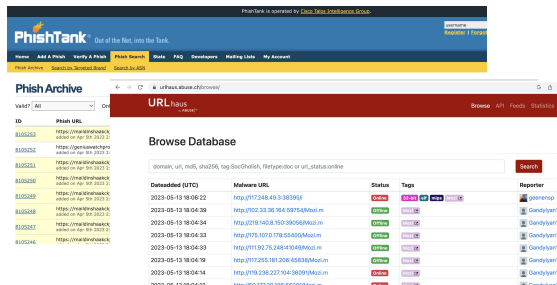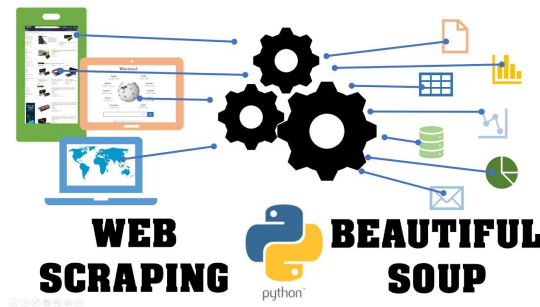
# Data Collection



## Sources

Collected **Phish & safe URLs** from PhishTank, a community-based phish verification system where users submit suspected phishes and other users "vote" if it is a phish or not)

Collected **Malware URLs** from URLhaus a project operated by abuse.ch to collect, track, and share malware URLs, helping network administrators & security analysts to protect their networks from cyber threats.



**WEB SCRAPING** **BEAUTIFUL SOUP**



## Challenges

- **Limited data import speed** – multiple requests may lead to server blocking or flagging the scraper as suspicious activity.
- **Data Size** - Over 1GB Data, too much to handle.
- **Class Imbalance** – Number of Safe URLs in the extracted dataset were just 25 for the time period.

## Solutions

- **Delay** - Used time lapse to avoid too many requests at same time.
- Limited the time frame of data import from **April 1st – May 1st**
- Extracted sufficient Safe URLs instead of making genetic copies to better train the model.
- Ran over multiple systems & combined the data in **csv format.**

| Introduction | Data Collection | Preprocessing | EDA | Models | Conclusion |

# Data Preprocessing - Feature Engineering



**15 Features** have been extracted from the raw URLs to develop a predictive model which can classify URLs into Phish, Malware & Safe URLs

**Features based on Length :**
- URL Length
- HostName Length
- Length of Top Level Domain

**Features based on Abnormality : Presence of**
- Multiple Domains
- Multiple www's
- Multiple Directories

# Data Preprocessing - Feature Engineering



**Features based on Abnormality : (contd)**
- Embedded domains
- Multiple HTTP & Presence of HTTPS
- Use of URL Shortening
- Digit Count
- Suspicious words

**Features based on Special Characters:**
- Count of ?
- Count of =
- Presence of @

# Exploratory Data Analysis

Average length of Safe URLs is highest unlike the usual expectations. Most of the URLs fall in the 30-35 length bin.



URL Length by Label Type



Distribution of URL Lengths

# Exploratory Data Analysis

Malware URLs mostly have multiple WWWs. Phish URLs tend to have dubious words when compared to Malware URLs



Proportion of URLs with multiple www's



URLs with Suspecious Words

# Exploratory Data Analysis

Most of the Malware URLs have no Secure Protocol & presence of Multiple HTTPs in the URL text.

# Predictive Model: Naive Bayes

| Metric | Malware | Phish | Safe |
|---|---|---|---|
| *Sensitivity* | 0.85 | 0.84 | 0.44 |
| *Specificity* | 0.89 | 0.75 | 0.94 |
| *Positive Pred Value* | 0.82 | 0.63 | 0.75 |
| *Negative Pred Value* | 0.91 | 0.90 | 0.79 |
| *Prevalence* | 0.37 | 0.33 | 0.30 |
| *Detection Rate* | 0.31 | 0.28 | 0.13 |
| *Detection Prevalence* | 0.38 | 0.44 | 0.17 |
| *Balanced Accuracy* | 0.87 | 0.80 | 0.69 |

**Accuracy : 72.2%**

# Predictive Model: Decision Tree

| Metric | *Malware* | *Phish* | *Safe* |
|---|---|---|---|
| *Sensitivity* | 0.98 | 0.70 | 0.51 |
| *Specificity* | 0.90 | 0.82 | 0.90 |
| *Positive Pred Value* | 0.86 | 0.65 | 0.68 |
| *Negative Pred Value* | 0.99 | 0.84 | 0.81 |
| *Prevalence* | 0.37 | 0.33 | 0.30 |
| *Detection Rate* | 0.36 | 0.23 | 0.15 |
| *Detection Prevalence* | 0.42 | 0.35 | 0.23 |
| *Balanced Accuracy* | 0.94 | 0.76 | 0.71 |

**Accuracy : 74.6%**

# Predictive Model: Gradient Boosting Model

| Metric | Malware | Phish | Safe |
|---|---|---|---|
| Sensitivity | 0.90 | 0.78 | 0.76 |
| Specificity | 0.97 | 0.87 | 0.89 |
| Positive Pred Value | 0.95 | 0.74 | 0.75 |
| Negative Pred Value | 0.94 | 0.89 | 0.90 |
| Prevalence | 0.37 | 0.33 | 0.30 |
| Detection Rate | 0.33 | 0.26 | 0.23 |
| Detection Prevalence | 0.35 | 0.34 | 0.31 |
| Balanced Accuracy | 0.93 | 0.82 | 0.82 |

Accuracy : 81.7%

# Predictive Model: Random Forest - Bagging

| Metric | Malware | Phish | Safe |
|--------|---------|-------|------|
| *Sensitivity* | 0.97 | 0.85 | 0.84 |
| *Specificity* | 0.98 | 0.93 | 0.93 |
| *Positive Pred Value* | 0.97 | 0.86 | 0.83 |
| *Negative Pred Value* | 0.98 | 0.93 | 0.93 |
| *Prevalence* | 0.37 | 0.33 | 0.30 |
| *Detection Rate* | 0.36 | 0.28 | 0.25 |
| *Detection Prevalence* | 0.37 | 0.33 | 0.30 |
| *Balanced Accuracy* | 0.98 | 0.89 | 0.88 |

**Accuracy : 89.3%**

# Predictive Model: Random Forest

| Metric (m=8) | Malware | Phish | Safe |
|---|---|---|---|
| *Sensitivity* | 0.98 | 0.85 | 0.84 |
| *Specificity* | 0.98 | 0.94 | 0.93 |
| *Positive Pred Value* | 0.96 | 0.87 | 0.84 |
| *Negative Pred Value* | 0.99 | 0.93 | 0.93 |
| *Prevalence* | 0.37 | 0.33 | 0.30 |
| *Detection Rate* | 0.36 | 0.28 | 0.25 |
| *Detection Prevalence* | 0.37 | 0.32 | 0.30 |
| *Balanced Accuracy* | 0.98 | 0.89 | 0.89 |

**Accuracy : 89.7%**

# Model Performance Comparison

| Model | Accuracy | FPR (Classifying Safe URL as Unsafe) | FNR (Classifying Unsafe URL as Safe) |
|---|---|---|---|
| *Naive Bayes* | 72.2% | 56.3% | 6.3% |
| *Decision Tree* | 74.6% | 48.4% | 10.3% |
| *Gradient Boosting* | 81.7% | 23.8% | 11.1% |
| *Random Forest - Bagging* | 89.3% | 16.7% | 7.1% |
| *Random Forest (m=8)* | 89.7% | 15.8% | 6.7% |

**Random Forest** performs well Overall, in terms of **Accuracy, FPR & FNR**. However, **Naive Bayes** model also turns out to be an important model as False Negatives are much costlier than False Positives in this case.

Introduction  Data Collection  Preprocessing  EDA  Models  Conclusion

# Feature Importance

## Feature Importance: by Mean Decrease in Accuracy

| Feature | Mean Decrease in Accuracy |
|---|---|
| Count_WWW | 255.2 |
| Multiple_Dir_Count | 161.3 |
| URL_length | 137.5 |
| Top_Level_Domain_Length | 132.2 |
| dot_count | 119.4 |
| suspicious_words | 76.6 |
| shortening_services | 67.2 |
| URL_HostName_length | 58.4 |
| Count_of_equals | 46.6 |
| Count_of_question_marks | 44.6 |
| Digit_Count | 40.9 |
| Presence_of_Embedded_Domains | 36.2 |
| Count_of_http | 34.9 |
| Presence_of_https | 33.5 |
| @_Presence | 11.4 |

A **high MDA** value suggests that the feature is more **important** in predicting the target variable.

Count of WWWs, Multiple Directories count, URL Length, Top Domain Length, multiple domain count seem to be the **top 5 features** which help in classifying the URLs into safe, phish & malware.
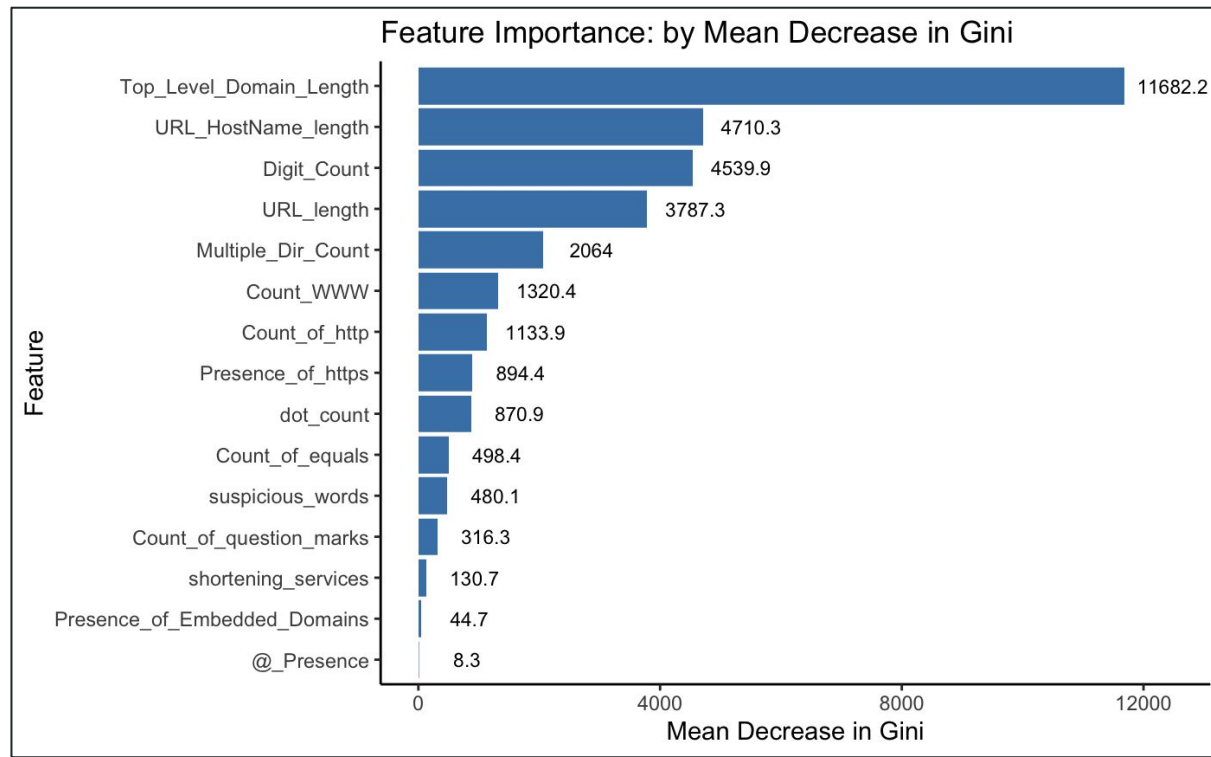
# Feature Importance

## Feature Importance: by Mean Decrease in Gini

| Feature | Mean Decrease in Gini |
|---|---|
| Top_Level_Domain_Length | 11682.2 |
| URL_HostName_length | 4710.3 |
| Digit_Count | 4539.9 |
| URL_length | 3787.3 |
| Multiple_Dir_Count | 2064 |
| Count_WWW | 1320.4 |
| Count_of_http | 1133.9 |
| Presence_of_https | 894.4 |
| dot_count | 870.9 |
| Count_of_equals | 498.4 |
| suspicious_words | 480.1 |
| Count_of_question_marks | 316.3 |
| shortening_services | 130.7 |
| Presence_of_Embedded_Domains | 44.7 |
| @_Presence | 8.3 |

The **mean decrease in Gini importance** of a feature is calculated by summing up total **decrease in Gini index** for each node that **splits on the feature**, then taking the average across all trees in the random forest.

**Higher decrease** in mean Gini index for a feature indicates that the feature is **more important** in splitting the data and creating a decision tree that accurately predicts the target variable.

Top Level Domain Length, Digit Count, Host Name Length, URL Length, Multiple Directory count are the **top 5** important features.

# Conclusion

- Average length of Safe URLs is highest unlike the usual expectations. Most of the URLs fall in the 30-35 length bin.
- Malware URLs mostly have **multiple WWWs**. Phish URLs tend to have **dubious words** when compared to Malware URLs.
- Most of the Malware URLs have **no Secure Protocol** & presence of **Multiple HTTPs** in the URL text.
- **Random Forest** performs well Overall, in terms of **Accuracy, FPR & FNR**
- **Naive Bayes** model also turns out to be an important model as False Negatives are more costlier than False Positives in this case.
- Count of WWWs, Multiple Directories count, URL Length, Top Domain Length, multiple domain count seem to be the **top 5 features** by **Mean Decrease in Accuracy.**
- Top Level Domain Length, Digit Count, Host Name Length, URL Length, Multiple Directory count are the **top 5** important features by **Mean Decrease in Gini.**

# References

- https://wisdomml.in/malicious-url-detection-using-machine-learning-in-python/

- https://www.kaggle.com/datasets/shashwatwork/web-page-phishing-detection-dataset

- https://siddharthakancharla.medium.com/classification-of-urls-f8253dee914

- https://www.kaggle.com/datasets/sid321axn/malicious-urls-dataset

| Introduction | Data Collection | Preprocessing | EDA | Models | Conclusion |