# Zach Fechko

## Finding Phishing Websites Using Machine Learning

## Introduction

### What is Phishing

```python
#importing basic packages
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from scipy.io import arff
```

The dataset gives each category a value of -1, 0, or 1

- -1 signifies a phishing website
- 0 signifies a website doesn't contain a given property
- 1 signifies a legitimate website

```python
#reading in the dataset from arff file
data = arff.loadarff('data/Training Dataset.arff')
df = pd.DataFrame(data[0])


#convert values with b'1' to 1 and b'-1' to -1
df.replace(b'1', 1, inplace=True)
df.replace(b'-1', -1, inplace=True)
df.replace(b'0', 0, inplace=True)
df.head()
```

| | having_IP_Address | URL_Length | Shortining_Service | having_At_Symbol | double_slash_redirecting |
|---|---|---|---|---|---|
| 0 | -1 | 1 | 1 | 1 | -1 |
| 1 | 1 | 1 | 1 | 1 | 1 |
| 2 | 1 | 0 | 1 | 1 | 1 |
| 3 | 1 | 0 | 1 | 1 | 1 |
| 4 | 1 | 0 | -1 | 1 | 1 |

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11055 entries, 0 to 11054
Data columns (total 31 columns):
 #   Column                       Non-Null Count  Dtype
---  ------                       --------------  -----
 0   having_IP_Address            11055 non-null  int64
 1   URL_Length                   11055 non-null  int64
 2   Shortining_Service           11055 non-null  int64
 3   having_At_Symbol             11055 non-null  int64
 4   double_slash_redirecting     11055 non-null  int64
 5   Prefix_Suffix                11055 non-null  int64
 6   having_Sub_Domain            11055 non-null  int64
 7   SSLfinal_State               11055 non-null  int64
 8   Domain_registeration_length  11055 non-null  int64
 9   Favicon                      11055 non-null  int64
 10  port                         11055 non-null  int64
 11  HTTPS_token                  11055 non-null  int64
 12  Request_URL                  11055 non-null  int64
 13  URL_of_Anchor                11055 non-null  int64
 14  Links_in_tags                11055 non-null  int64
 15  SFH                          11055 non-null  int64
 16  Submitting_to_email          11055 non-null  int64
 17  Abnormal_URL                 11055 non-null  int64
 18  Redirect                     11055 non-null  int64
 19  on_mouseover                 11055 non-null  int64
 20  RightClick                   11055 non-null  int64
 21  popUpWidnow                  11055 non-null  int64
 22  Iframe                       11055 non-null  int64
 23  age_of_domain                11055 non-null  int64
 24  DNSRecord                    11055 non-null  int64
 25  web_traffic                  11055 non-null  int64
 26  Page_Rank                    11055 non-null  int64
 27  Google_Index                 11055 non-null  int64
 28  Links_pointing_to_page       11055 non-null  int64
 29  Statistical_report           11055 non-null  int64
 30  Result                       11055 non-null  int64
dtypes: int64(31)
memory usage: 2.6 MB
```