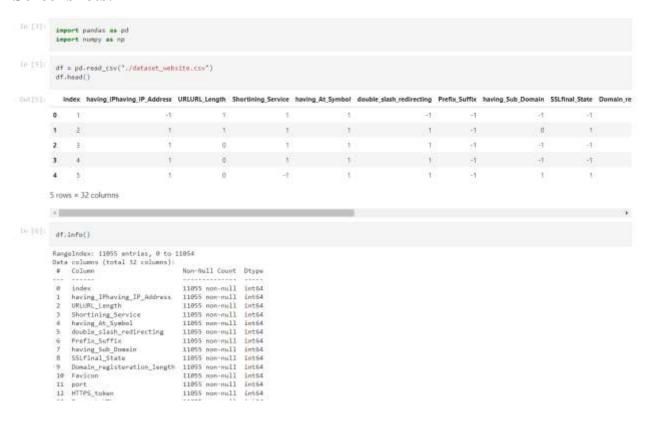
## **SPRINT 1**

| Date         | 29 October 2022        |
|--------------|------------------------|
| Team ID      | PNT2022TMID15890       |
| Project Name | Web Phishing Detection |

## **Sprint Objective:**

In Sprint 1 the team members download the given dataset from the given link. Then the dataset should be pre-processed by checking all the information in the given datasets like the features given for the URL, what the values mean for each Feature (URL's feature) in the dataset; then, they should check for the null values, if there is any null value they should have them.

## **Screenshots:**



```
l= [7]: off.iseull().any()
                         index
having iPhaving iP Address
having iPhaving iP Address
having Service
having At Symbol
double slash, addrecting
Prefix Suffia
having Sob Jomain
SSifinal State
Domain registeration length
Favicos
                                                                                                                              False
                                                                                                                               False
                                                                                                                              False
False
                                                                                                                              False
False
                                                                                                                              False
False
                           Favicon
                                                                                                                               False
                          port
HTTPS_token
                         HTTPS_token
Request_UML
URL_of_Sechor
Links_in_tags
SFN
Submitting_to_email
Abovemal_URL
Hadimart
on_mouseover
HightClick
popUsphidnom
Ifrace
age_of_domain
DMSRacord
web_traffic
Page_Rank
Google_Index
Links_pointing_to_page
Statistical_report
Result
                                                                                                                               False
                                                                                                                              False
False
                                                                                                                              False
False
                                                                                                                              False
False
                                                                                                                              False
False
                                                                                                                              Folse
False
False
                                                                                                                               False
                                                                                                                              False
False
                                                                                                                              False
False
                                                                                                                              False
False
                           Result
dtype: bool
```

## **Sprint Result:**

Team members downloaded the given data. The dataset contains 11,055 data samples. Each data samples have 30 different features and these features has either of these 3 different values (-1, 0, 1). The dataset is checked for null values and there is none.