

## Project Development Phase

### Delivery of Sprint – 3

<b>Date</b>	12 November 2022
<b>Team ID</b>	PNT2022TMID18067
<b>Project Name</b>	Web Phishing Detection

## Feature Extraction Source Code:

```
import ipaddress
import re
import urllib.request
from bs4 import BeautifulSoup
import socket
import requests
from googlesearch import search
import whois
from datetime import date, datetime
import time
from dateutil.parser import parse as date_parse
from urllib.parse import urlparse
```

```
class FeatureExtraction:
    features = []
    def __init__(self,url):
        self.features = []
        self.url = url
        self.domain = ""
        self.whois_response = ""
```

```
self.urlparse = ""
self.response = ""
self.soup = ""
self.check = ""
try:
    self.response = requests.get(url)
    self.soup = BeautifulSoup(self.response.text, 'html.parser')
except:
    pass
```

```
try:
    self.urlparse = urlparse(url)
    self.domain = self.urlparse.netloc
    print("1")
except:
    pass
```

```
try:
    print("2")
    self.whois_response = whois.whois(self.domain)
    print("3")
except:
    print("here")
    pass
```

```
self.features.append(self.UsingIp())
self.features.append(self.longUrl())
self.features.append(self.shortUrl())
self.features.append(self.symbol())
self.features.append(self.redirecting())
self.features.append(self.prefixSuffix())
self.features.append(self.SubDomains())
self.features.append(self.Hppts())
self.features.append(self.DomainRegLen())
self.features.append(self.Favicon())
```

```
self.features.append(self.NonStdPort())
self.features.append(self.HTTPSDomainURL())
self.features.append(self.RequestURL())
self.features.append(self.AnchorURL())
self.features.append(self.LinksInScriptTags())
self.features.append(self.ServerFormHandler())
self.features.append(self.InfoEmail())
self.features.append(self.AbnormalURL())
self.features.append(self.WebsiteForwarding())
self.features.append(self.StatusBarCust())
```

```
self.features.append(self.DisableRightClick())
self.features.append(self.UsingPopupWindow())
self.features.append(self.IframeRedirection())
self.features.append(self.AgeofDomain())
self.features.append(self.DNSRecording())
self.features.append(self.WebsiteTraffic())
self.features.append(self.PageRank())
self.features.append(self.GoogleIndex())
self.features.append(self.LinksPointingToPage())
self.features.append(self.StatsReport())
```

# 1.UsingIp

```
def UsingIp(self):
```

```
    try:
```

```
        ipaddress.ip_address(self.url)
```

```
        return -1
```

```
    except:
```

```
        print("11")
```

```
        return 1
```

# 2.longUrl

```
def longUrl(self):
```

```
    if len(self.url) < 54:
```

```
        return 1
```

```
    if len(self.url) >= 54 and len(self.url) <= 75:
```

```
        return 0
```

```
return -1
```

```
# 3.shortUrl
```

```
def shortUrl(self):
```

```
    match = re.search('bit\.ly|goo\.gl|shorte\.st|go2l\.ink|x\.co|ow\.ly|t\.co|tinyurl|tr\.im|is\.gd|cli\.gs|'
        'yfrog\.com|migre\.me|ff\.im|tiny\.cc|url4\.eu|twit\.ac|su\.pr|twurl\.nl|snipurl\.com|'
        'short\.to|BudURL\.com|ping\.fm|post\.ly|Just\.as|bkite\.com|snipr\.com|fic\.kr|loopt\.us|'
        'doiop\.com|short\.ie|kl\.am|wp\.me|rubyurl\.com|om\.ly|to\.ly|bit\.do|t\.co|lnkd\.in|'
        'db\.tt|qr\.ae|adf\.ly|goo\.gl|bitly\.com|cur\.lv|tinyurl\.com|ow\.ly|bit\.ly|ity\.im|'
        'q\.gs|is\.gd|po\.st|bc\.vc|twitthis\.com|u\.to|j\.mp|buzurl\.com|cutt\.us|u\.bb|yourls\.org|'
```

```
'x\.co|prettylinkpro\.com|scrnch\.me|filoops\.info|vzturl\.com|qr\.net|1url\.com|tweez\.me|v\.gd|tr\.im|link\.zip\
.net', self.url)
```

```
    if match:
```

```
        return -1
```

```
    return 1
```

```
# 4.Symbol@
```

```
def symbol(self):
```

```
    if re.findall("@",self.url):
```

```
        return -1
```

```
    return 1
```

```
# 5.Redirecting//
```

```
def redirecting(self):
```

```
    if self.url.rfind('/')>6:
```

```
        return -1
```

```
    return 1
```

```
# 6.prefixSuffix
```

```
def prefixSuffix(self):
```

```
    try:
```

```
        match = re.findall('-', self.domain)
```

```
        if match:
```

```
            return -1
```

```
        return 1
```

```
    except:
```

```
        return -1
```

### # 7.SubDomains

```
def SubDomains(self):
    dot_count = len(re.findall("\.", self.url))
    if dot_count == 1:
        return 1
    elif dot_count == 2:
        return 0
    return -1
```

### # 8.HTTPS

```
def Hppts(self):
    try:
        https = self.urlparse.scheme
        if 'https' in https:
            return 1
        return -1
    except:
        return 1
```

### # 9.DomainRegLen

```
def DomainRegLen(self):
    try:
        expiration_date = self.whois_response.expiration_date
        creation_date = self.whois_response.creation_date
        try:
            if(len(expiration_date)):
                expiration_date = expiration_date[0]
        except:
            pass
        try:
            if(len(creation_date)):
                creation_date = creation_date[0]
        except:
            pass

        age = (expiration_date.year-creation_date.year)*12+ (expiration_date.month-creation_date.month)
        if age >=12:
            return 1
        return -1
```

```
except:
    return -1
```

# 10. Favicon

```
def Favicon(self):
    try:
        for head in self.soup.find_all('head'):
            for head.link in self.soup.find_all('link', href=True):
                dots = [x.start(0) for x in re.finditer('\.', head.link['href'])]
                if self.url in head.link['href'] or len(dots) == 1 or domain in head.link['href']:
                    return 1
        self.check = "rr"
        return -1
    except:
        return -1
```

# 11. NonStdPort

```
def NonStdPort(self):
    try:
        port = self.domain.split(":")
        if len(port)>1:
            return -1
        return 1
    except:
        return -1
```

# 12. HTTPSDomainURL

```
def HTTPSDomainURL(self):
    try:
        if 'https' in self.domain:
            return -1
        return 1
    except:
        return -1
```

# 13. RequestURL

```
def RequestURL(self):
    try:
        for img in self.soup.find_all('img', src=True):
```

```

        dots = [x.start(0) for x in re.finditer('\.', img['src'])]
        if self.url in img['src'] or self.domain in img['src'] or len(dots) == 1:
            success = success + 1
        i = i+1

    for audio in self.soup.find_all('audio', src=True):
        dots = [x.start(0) for x in re.finditer('\.', audio['src'])]
        if self.url in audio['src'] or self.domain in audio['src'] or len(dots) == 1:
            success = success + 1
        i = i+1

    for embed in self.soup.find_all('embed', src=True):
        dots = [x.start(0) for x in re.finditer('\.', embed['src'])]
        if self.url in embed['src'] or self.domain in embed['src'] or len(dots) == 1:
            success = success + 1
        i = i+1

    for iframe in self.soup.find_all('iframe', src=True):
        dots = [x.start(0) for x in re.finditer('\.', iframe['src'])]
        if self.url in iframe['src'] or self.domain in iframe['src'] or len(dots) == 1:
            success = success + 1
        i = i+1

    try:
        percentage = success/float(i) * 100
        if percentage < 22.0:
            return 1
        elif((percentage >= 22.0) and (percentage < 61.0)):
            return 0
        else:
            return -1
    except:
        return 0
except:
    return -1

```

# 14. AnchorURL

```

def AnchorURL(self):
    try:

```

```

i,unsafe = 0,0
for a in self.soup.find_all('a', href=True):
    if "#" in a['href'] or "javascript" in a['href'].lower() or "mailto" in a['href'].lower() or not (url in
a['href'] or self.domain in a['href']):
        unsafe = unsafe + 1
    i = i + 1

```

```

try:
    percentage = unsafe / float(i) * 100
    if percentage < 31.0:
        return 1
    elif ((percentage >= 31.0) and (percentage < 67.0)):
        return 0
    else:
        return -1
except:
    return -1

```

```

except:
    return -1

```

#### # 15. LinksInScriptTags

```
def LinksInScriptTags(self):
```

```

    try:
        i,success = 0,0

        for link in self.soup.find_all('link', href=True):
            dots = [x.start(0) for x in re.finditer('\.', link['href'])]
            if self.url in link['href'] or self.domain in link['href'] or len(dots) == 1:
                success = success + 1
            i = i+1

        for script in self.soup.find_all('script', src=True):
            dots = [x.start(0) for x in re.finditer('\.', script['src'])]
            if self.url in script['src'] or self.domain in script['src'] or len(dots) == 1:
                success = success + 1
            i = i+1

```

```

    try:

```



```

percentage = success / float(i) * 100
if percentage < 17.0:
    return 1
elif((percentage >= 17.0) and (percentage < 81.0)):
    return 0
else:
    return -1
except:
    return 0
except:
    return -1

```

# 16. ServerFormHandler

```
def ServerFormHandler(self):
```

```

    try:
        if len(self.soup.find_all('form', action=True))==0:
            return 1
        else :
            for form in self.soup.find_all('form', action=True):
                if form['action'] == "" or form['action'] == "about:blank":
                    return -1
                elif self.url not in form['action'] and self.domain not in form['action']:
                    return 0
                else:
                    return 1
    except:
        return -1

```

# 17. InfoEmail

```
def InfoEmail(self):
```

```

    try:
        if re.findall(r"[mail\\(\)|mailto:?}", self.soup):
            return -1
        else:
            return 1
    except:
        return -1

```

# 18. AbnormalURL

```
def AbnormalURL(self):
    try:
        if self.response.text == self.whois_response:
            return 1
        else:
            return -1
    except:
        return -1
```

# 19. WebsiteForwarding

```
def WebsiteForwarding(self):
    try:
        if len(self.response.history) <= 1:
            return 1
        elif len(self.response.history) <= 4:
            return 0
        else:
            return -1
    except:
        return -1
```

# 20. StatusBarCust

```
def StatusBarCust(self):
    try:
        if re.findall("<script>.+onmouseover.+</script>", self.response.text):
            return 1
        else:
            return -1
    except:
        return -1
```

# 21. DisableRightClick

```
def DisableRightClick(self):
    try:
        if re.findall(r"event.button ?== ?2", self.response.text):
            return 1
        else:
            return -1
    except:
```

```
return -1
```

```
# 22. UsingPopupWindow
```

```
def UsingPopupWindow(self):
```

```
    try:
```

```
        if re.findall(r"alert\(", self.response.text):
```

```
            return 1
```

```
        else:
```

```
            return -1
```

```
    except:
```

```
        return -1
```

```
# 23. IframeRedirection
```

```
def IframeRedirection(self):
```

```
    try:
```

```
        if re.findall(r"[<iframe>|<frameBorder>]", self.response.text):
```

```
            return 1
```

```
        else:
```

```
            return -1
```

```
    except:
```

```
        return -1
```

```
# 24. AgeofDomain
```

```
def AgeofDomain(self):
```

```
    try:
```

```
        creation_date = self.whois_response.creation_date
```

```
    try:
```

```
        if(len(creation_date)):
```

```
            creation_date = creation_date[0]
```

```
    except:
```

```
        pass
```

```
    today = date.today()
```

```
    age = (today.year-creation_date.year)*12+(today.month-creation_date.month)
```

```
    if age >=6:
```

```
        return 1
```

```
    return -1
```

```
    except:
```

```
        return -1
```

# 25. DNSRecording

```
def DNSRecording(self):
```

```
    try:
```

```
        creation_date = self.whois_response.creation_date
```

```
    try:
```

```
        if(len(creation_date)):
```

```
            creation_date = creation_date[0]
```

```
    except:
```

```
        pass
```

```
    today = date.today()
```

```
    age = (today.year-creation_date.year)*12+(today.month-creation_date.month)
```

```
    if age >=6:
```

```
        return 1
```

```
    return -1
```

```
    except:
```

```
        return -1
```

# 26. WebsiteTraffic

```
def WebsiteTraffic(self):
```

```
    try:
```

```
        rank = BeautifulSoup(urllib.request.urlopen("http://data.alexa.com/data?cli=10&dat=s&url=" +  
url).read(), "xml").find("REACH")["RANK"]
```

```
        if (int(rank) < 100000):
```

```
            return 1
```

```
        return 0
```

```
    except :
```

```
        return -1
```

# 27. PageRank

```
def PageRank(self):
```

```
    try:
```

```
        prank_checker_response = requests.post("https://www.checkpagerank.net/index.php", {"name":  
self.domain})
```

```
        global_rank = int(re.findall(r"Global Rank: ([0-9]+)", prank_checker_response.text)[0])
```

```
        if global_rank > 0 and global_rank < 100000:
```

```
            return 1
```

```
        return -1
    except:
        return -1
```

# 28. GoogleIndex

```
def GoogleIndex(self):
    try:
        site = search(self.url, 5)
        if site:
            return 1
        else:
            return -1
    except:
        return 1
```

# 29. LinksPointingToPage

```
def LinksPointingToPage(self):
    try:
        number_of_links = len(re.findall(r"<a href=", self.response.text))
        if number_of_links == 0:
            return 1
        elif number_of_links <= 2:
            return 0
        else:
            return -1
    except:
        return -1
```

# 30. StatsReport

```
def StatsReport(self):
```

```
    try:
        url_match = re.search(
            'at\.ua|usa\.cc|baltazarpresentes\.com\.br|pe\.hu|esy\.es|hol\.es|sweddy\.com|myjino\.ru|96\.lt|ow\.ly', url)
        ip_address = socket.gethostbyname(self.domain)
        ip_match =
```

```
re.search('146\.112\.61\.108|213\.174\.157\.151|121\.50\.168\.88|192\.185\.217\.116|78\.46\.211\.158|181\.174\.165\.13|46\.242\.145\.103|121\.50\.168\.40|83\.125\.22\.219|46\.242\.145\.98|'
```

'107\151\148\44|107\151\148\107|64\70\19\203|199\184\144\27|107\151\148\108|107\151\148\109|119\28\52\61|54\83\43\69|52\69\166\231|216\58\192\225|'

'118\184\25\86|67\208\74\71|23\253\126\58|104\239\157\210|175\126\123\219|141\8\224\221|10\10\10\10|43\229\108\32|103\232\215\140|69\172\201\153|'

'216\218\185\162|54\225\104\146|103\243\24\98|199\59\243\120|31\170\160\61|213\19\128\77|62\113\226\131|208\100\26\234|195\16\127\102|195\16\127\157|'

'34\196\13\28|103\224\212\222|172\217\4\225|54\72\9\51|192\64\147\141|198\200\56\183|23\253\164\103|52\48\191\26|52\214\197\72|87\98\255\18|209\99\17\27|'

'216\38\62\18|104\130\124\96|47\89\58\141|78\46\211\158|54\86\225\156|54\82\156\19|37\157\192\102|204\11\56\48|110\34\231\42', ip\_address)

if url\_match:

return -1

elif ip\_match:

return -1

return 1

except:

return 1

def getFeaturesList(self):

return self.features

# url = "https://www.google.com/"

# obj = FeatureExtraction(url)

# print(obj.features)

# print(obj.url)

# print(obj.domain)

# print(obj.whois\_response)

# print(obj.urlparse)

# print(obj.response)

# print(obj.check)