



L OVELY
P ROFESSIONAL
U NIVERSITY

DATA MANAGEMENT PROJECT REPORT

ON

Analysis of Cyber Attacks Happened All Over the World

Submitted by:

Manas Bhardwaj

11704118

Under the guidance of:

Sandeep Kaur - 23614

**School of Computer Science and Engineering
Lovely Professional University, Phagwara**

CERTIFICATE

This is to certify that **Manas Bhardwaj** bearing Registration number **11704118** has completed **Data Management (INT217)** project titled, “**Analysis of Cyber Attacks Happened All Over the World**” under my guidance and supervision. To the best of my knowledge, the present work is the result of his/her original development, effort and study.

Signature and Name of the Supervisor
School of Computer Science and Engineering
Lovely Professional University
Date:

DECLARATION

I, **Manas Bhardwaj**, student of **Computer Science and Engineering** under CSE/IT Discipline at, Lovely Professional University, Punjab, hereby declare that all the information furnished in this project report is based on my own intensive work and is genuine.

Date:

Signature:

Registration No. 11704118

Manas Bhardwaj

ACKNOWLEDGEMENT

I take this opportunity to present our votes of thanks to all those guideposts who really acted as lightening pillars to enlighten my way throughout this project that has led to the successful and satisfactory completion of this Project. I am grateful to **Lovely Professional University** for providing us with an opportunity to undertake this project and providing us with all the facilities. I am highly thankful to All for their active support, valuable time and advice, whole-hearted guidance, sincere cooperation and painstaking involvement during the project and in completing the assignment of preparing the said project within the time stipulated. Lastly, I am thankful to all those, particularly the various friends, who have been instrumental in creating a proper, healthy and conducive environment and including new and fresh innovative ideas for me during the project, without their help, it would have been extremely difficult for me to complete the project in a time-bound framework.

INDEX

1. Introduction
 - a. About Data Analysis
 - b. Background of the project
 - c. Data Definition
2. Scope of Analysis
 - a. Objectives of the project
 - b. Scope
3. ETL Process
 - a. What is ETL?
 - b. ETL on the dataset
4. Analysis of objectives
5. Analyzing Results
6. References and Bibliography

INTRODUCTION

1. Data:

It is the raw and isolated facts about an entity or records.

2. Information:

When the data are used to extract some insights then those data which are being used is called informative data and the results that come out of them is called information.

3. Data Analysis:

The process of analysing the data which later turns into information is called data analysis. In this process we observe the data to make some inference out of it, we use several tools and resources to read the jargon datasets into the common and usable form.

Or, It is a process of inspecting, cleansing, transforming, and modelling data with the goal of discovering useful information, informing conclusions, and supporting decision-making. Data analysis has multiple facets and approaches, encompassing diverse techniques under a variety of names while being used in different business, science, and social science domains.

4. Data Science:

The science of studying the behaviours of the data like what they stand for, what can be extracted from them, their nature and how they can be used for future references and used to solve some real-world problems. It is one of the emerging fields of science and technologies where gazillions of data are fetched into a machine and analysed for the beautiful insights.

5. Excel 2016:

Microsoft office tools which are used widely for data analysis and exploration. Excel 2016 is the best tool for data analysis when it comes to handy data sets. Excels contains all the features for converting small uncleaned data into a readable form and also contains different functions which provide an eminent role in data exploration. It's easy to use and cosmopolitan in the distribution.

6. Tableau Prep:

Tableau Prep is another exciting tool which is widely used by data analysts around the world. It is used when we have to deals with a lot of inbuilt errors, like removing wrong spellings, removing the redundancy of the data, removing data using different aggregate functions and visual guidance. Sometimes we have to deal with hidden errors and there this tool comes in handy experiences. Tableau is not only used to explore the errors but also used to explore and analyse the data sets for better insights.

SCOPE OF ANALYSIS

Data visualization is useful in a plethora of areas, whether it be technical sector or not technical sector. Every micro-business or macro large scale industries require this analysis to meet their different standards and goals. Data visualization is not only required in one field but in the today world, it has become quintessential in all most all the sectors.

Data visualization is used in business, managerial skills, for training and testing machine models, to solve real-world problems with the aid of machine learning and artificial intelligence. Moreover, it has been able to solve some of the complex problems of the world.

Some directions to take when exploring the data:

1. Top 5 Cyber-attacks happened on 5th October 2019.
2. Top 5 attacking countries.
3. Attack type analysis.
4. Countries Attacking Too & Fro.
5. Top 5 victim countries.
6. The number of attacks happened on the different hour of the day.
7. Type of Attacks happened on the different hour of the day.

Aim of this project is to answer the above objectives in the form of visualization by creating a dashboard to convey the answers effectively.

ETL PROCESS

ETL is defined as a process that extracts the data from different RDBMS source systems, then transforms the data (like applying calculations, concatenations, etc.) and finally loads the data into the Data Warehouse system. ETL full-form is Extract, Transform and Load.

Extraction:

Extraction is the first step in the ETL process, which basically deals with the data mining and data selection phase. In this step, we have to use different sources through which we are going to extract our files for the analysis. Extraction involves the files like .csv and others which need cleaning. This process plays eminent roles because once we are done with the extraction, everything turns out to be simple for us. Different functions are used to maintain the normality and atomicity of the data sets.

Transformation:

Transformation is the second phase of data cleaning. Once the data are extracted and cleaned it needs to be transformed into a set of cleaned and well-defined data. The transformation takes place in different models. In our coursework, we did transformation with an in-built Pivot Table and manually too. However, Pivot table helped a lot in the data transformation. Different charts models are used to make the visual appearance of the data sets.

Load/Visualization:

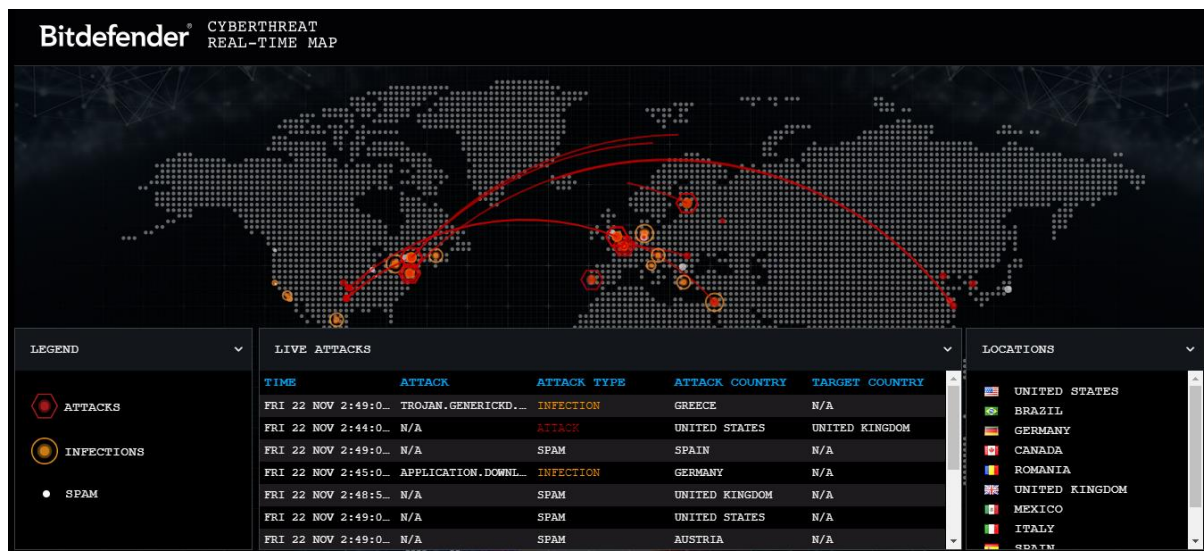
Data once get transformed needs to be analysed. This process is the ultimate process in which results are extracted from the graphs and visual models. Different models represented dynamically can be used to anticipate results, make careful observations and create a fascinating story out of it.

DATASET DESCRIPTION:

Source :	https://threatmap.bitdefender.com/		
Date of Collection :	05-Oct-19		
Time of Collection :	00 : 00 - 23 : 00		
Collected By :	Manas Bhardwaj		
Number of Columns :	6		
Number of Rows :	22,447		
Method of Collection :	Web Scraping using python, refer Scraping-script.py in root directory		
Tools used for cleaning :	Tableau prep, Python, Microsoft Excel 2016		
Method Used to remove NULLS :	Exclusion of rows, Forward Filling in pandas(Python).		
Code Book			
Columns Name	Description	Type	Total NULLS
Time	It contains time of the attack in 24-hour format.	Time(HH) 24-hour	0
Date	It contains date of attack.	Date(DD-mm-YYYY)	0
Attack Name	It contains the name of attack which was performed on the victim country by the attacking country.	Text	0
Attacking Country	It contains the name of the attacking country.	Text	0
Target Country	It contains the name of the target country on which attack was performed.	Text	0
Attack Type	It contains the type of attack, i.e INFECTION, SPAM, ATTACK.	Text	0

DATA SOURCE:

Link: <https://threatmap.bitdefender.com>



Data Collection Script:

This script is written to scrape the data from a dynamic website.
A dynamic website is a site that contains dynamic pages such as templates,
contents, scripts etc. In a nutshell, the dynamic website displays various
content types every time it is browsed. The web page can be changed with the
reader that opens the page, character of consumer interplay, or day time.
Or, Dynamic website is a website in which the data changes very frequently.

```
from bs4 import BeautifulSoup
from selenium import webdriver
from time import sleep
import csv
```

Below code is to automate Chrome only if you want to automate Mozilla you have
to download geckodriver according to the version you are using of Mozilla.
Just download the geckodriver and place it in the directory
and replace 'chromedriver.exe' with the geckodriver in below line.
Or, you can simply uncomment the line 12 and comment the line 11.

```
driver = webdriver.Chrome(executable path='chromedriver.exe')
driver.get('https://threatmap.bitdefender.com/') # Opens the link in the browser.
sleep(60) # It gives the website sufficient time to load all its content.
while(1):
```

Executes the javascript to get the HTML every time it is executed without
 # refreshing the website since we are scraping the dynamic website it is
 # very important to prevent refreshing of the website since the data will not
 # be available as we wanted.

```
res = driver.execute_script("return document.documentElement.outerHTML")
# To convert the response of the above line to readable html.
soup = BeautifulSoup(res,'lxml')
# It find the table-body('tbody') tag since we want the data from the table
# which is updating regularly and the id of the table is 'attacks_data'.
# In futue it can change since website is maintained regularly so you can
# change the below body tag and id tag accordingly.
tags = soup.find_all('tbody',{ 'id':'attacks_data'})
tags = tags[0] # We just want first table.

for tr in tags.find_all('tr'): # Iterate over the all the rows in the table.
    print(tr.contents,end='\n') # Prints the content of the row in form of list.

# To Collect the data we use 'csv' module to write in csv file.
with open("attacks.csv",'a+',newline='') as f:
    writer = csv.writer(f,delimiter=',',quotechar='|', quoting=csv.QUOTE_MINIMAL)
    if len(list(tr.contents)) == 5: # we want that data only which has all the 5 columns data.
        writer.writerow(list(tr.contents)) # Writes the row to the csv file.

sleep(20)
driver.quit() # When you're done just press ctrl+c to quit and the browser will get close.
```

Uncleaned Data Collected in CSV format,

```
[<td>Fri 22 Nov 2:34:14 AM</td>, <td>N/A</td>, <td class="type_attack">attack</td>, <td>Romania</td>, <td>Finland</td>
<td>Fri 22 Nov 2:34:16 AM</td>, <td>#esif_assist_64.exe:0000F001,00008006,00008050,00009903</td>, <td class="type_in
<td>Thu 21 Nov 7:44:31 PM</td>, <td>Trojan.Downloader.VBS.HY</td>, <td class="type_infection">infection</td>, <td>Br
<td>Fri 22 Nov 2:34:41 AM</td>, <td>N/A</td>, <td class="type_spam">spam</td>, <td>Costa Rica</td>, <td>n/a</td>]
<td>Fri 22 Nov 2:33:45 AM</td>, <td>N/A</td>, <td class="type_attack">attack</td>, <td>Vietnam</td>, <td>United Stat
<td>Fri 22 Nov 2:33:45 AM</td>, <td>N/A</td>, <td class="type_attack">attack</td>, <td>France</td>, <td>Germany</td>
<td>Fri 22 Nov 2:34:39 AM</td>, <td>Application.BitCoinMiner.SX</td>, <td class="type_infection">infection</td>, <td>
<td>Fri 22 Nov 2:34:41 AM</td>, <td>N/A</td>, <td class="type_spam">spam</td>, <td>United States</td>, <td>n/a</td>]
```

DATA CLEANING:

For cleaning the messy data tableau prep was used,

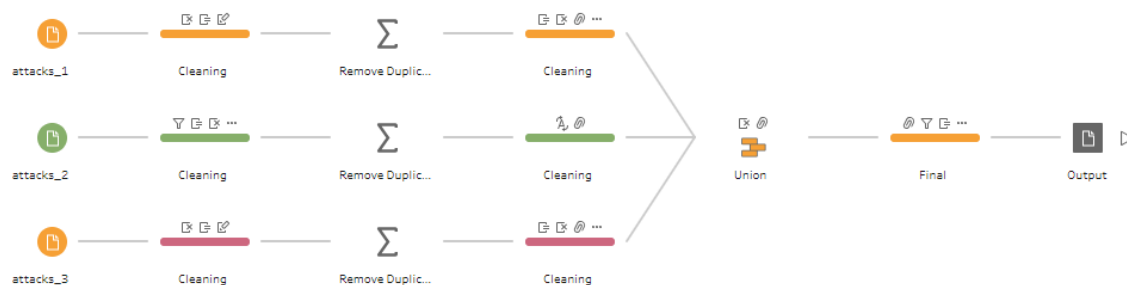


Tableau Flow File

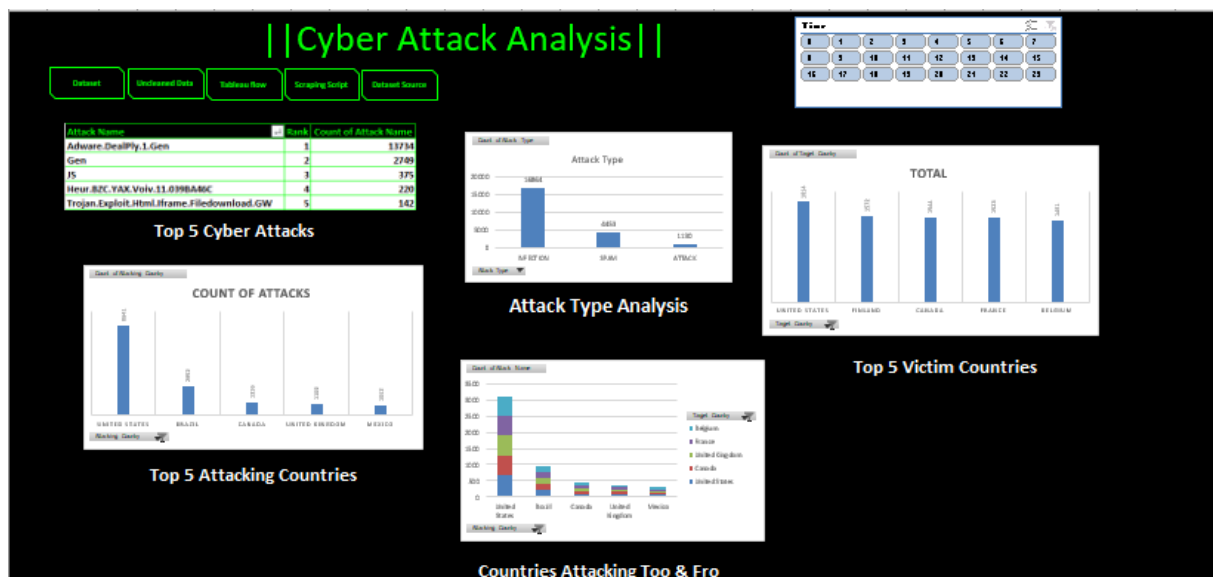
Messy Data:

Cleaning 12 Fields 120K Rows					
Filter Values... Automatic Split Custom Split... Rename Field ... 8 Recommendations Search					
Changes (59)	Abc F1	Abc F2	Abc F3	Abc F4	Abc F5
	<td>Sat 5 Oct 1:19:21 PM</td>	<td>N/A</td>	<td class="type_attack">attack</td>	<td>Vietnam</td>	<td>United Kingdom</td>
	<td>Sat 5 Oct 1:20:13 PM</td>	<td>Trojan.SalityStub.F</td>	<td class="type_infection">infection</td>	<td>Brazil</td>	<td>n/a</td>
	<td>Sat 5 Oct 1:20:16 PM</td>	<td>Adware.DealPly.1.Gen</td>	<td class="type_infection">infection</td>	<td>United States</td>	<td>n/a</td>
	<td>Sat 5 Oct 1:20:14 PM</td>	<td>Application.Miner.AG</td>	<td class="type_infection">infection</td>	<td>Venezuela</td>	<td>n/a</td>
	<td>Sat 5 Oct 1:20:18 PM</td>	<td>N/A</td>	<td class="type_attack">attack</td>	<td>Romania</td>	<td>United States</td>
	<td>Sat 5 Oct 1:20:17 PM</td>	<td>N/A</td>	<td class="type_attack">attack</td>	<td>Romania</td>	<td>France</td>
	<td>Sat 5 Oct 1:20:15 PM</td>	<td>N/A</td>	<td class="type_spam">spam</td>	<td>India</td>	<td>n/a</td>
	<td>Sat 5 Oct 1:19:14 PM</td>	<td>N/A</td>	<td class="type_attack">attack</td>	<td>United States</td>	<td>United Kingdom</td>
	<td>Sat 5 Oct 1:19:23 PM</td>	<td>N/A</td>	<td class="type_attack">attack</td>	<td>Romania</td>	<td>Belgium</td>
<td>Sat 5 Oct 1:19:19 PM</td>	<td>N/A</td>	<td class="type_attack">attack</td>	<td>United States</td>	<td>Poland</td>	
<td>Sat 5 Oct 1:20:11 PM</td>	<td>Win32.Generic.494661</td>	<td class="type_infection">infection</td>	<td>United States</td>	<td>n/a</td>	
<td>Sat 5 Oct 1:20:14 PM</td>	<td>#ipadchg2.exe:0000F001	00008050	0000D007</td>	<td class="type_infectio	
<td>Sat 5 Oct 1:19:16 PM</td>	<td>N/A</td>	<td class="type_attack">attack</td>	<td>United States</td>	<td>Finland</td>	
<td>Sat 5 Oct 1:20:10 PM</td>	<td>Application.JS.Miner.C</td>	<td class="type_infection">infection</td>	<td>United States</td>	<td>n/a</td>	
<td>Sat 5 Oct 1:20:13 PM</td>	<td>Gen.Variant.Zusy.193507</td>	<td class="type_infection">infection</td>	<td>United States</td>	<td>n/a</td>	

Cleaned Data:

Abc Time 11K	Abc Date 1	Abc Attack Name 63	Abc Country/Region Attacking Coun... 145	Abc Country/Region Target Country 16	Abc Attack Type 3
1:00:12 PM 1:00:13 PM 1:00:15 PM 1:00:17 PM 1:00:28 PM 1:00:36 PM 1:00:41 PM 1:00:46 PM 1:00:49 PM 1:00:50 PM 1:00:51 PM 1:00:52 PM	5 Oct 2019	Cloud.Malware.0133.a... CS_GO_Arx_Applet.exe CS1.tmp Datacrime.1480 DeepScan DeleteFolderTask.exe digipass-nativebridge... e0 d f g f f f g f, e f<... Exploit.Iframe.Vulnera... Fiserv.Cleartouch.Desk... formatfactory_29222... FortniteClient-Win64...	Afghanistan Algeria Angola Argentina Armenia Aruba Australia Austria Azerbaijan Bahamas Bahrain Bangladesh	Finland France Germany Ireland Italy Netherlands Poland Portugal Russia Spain United Kingdom United States	ATTACK INFECTION SPAM

DASHBOARD



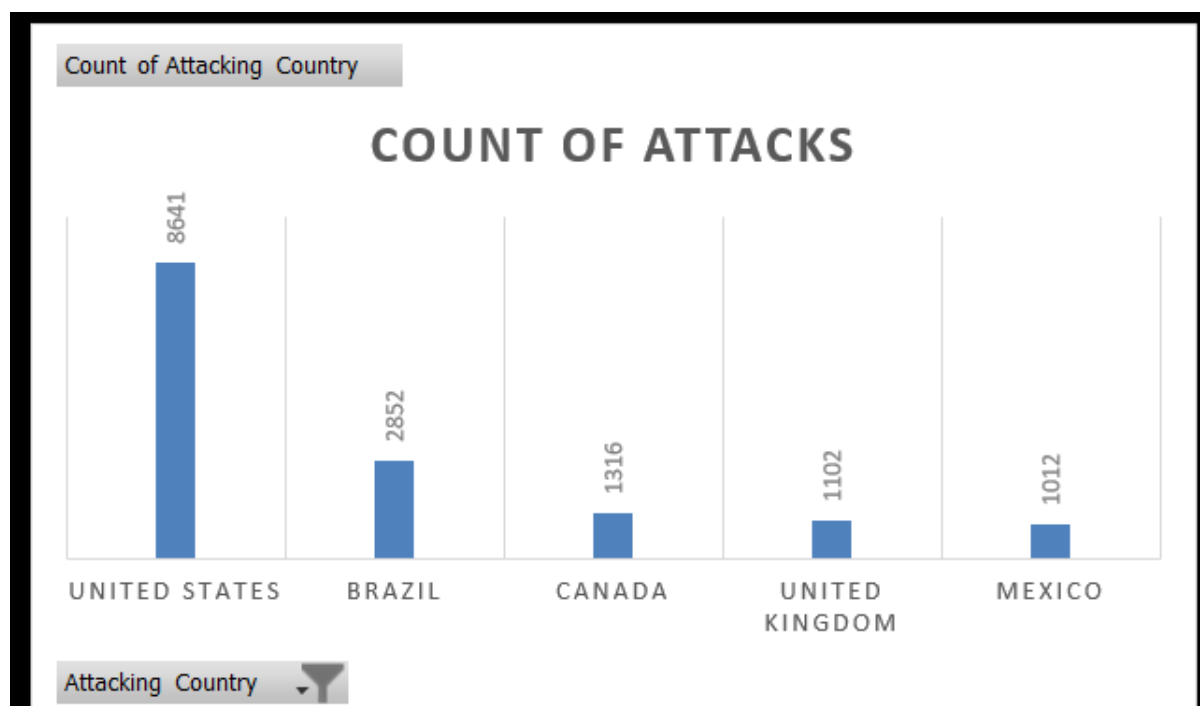
ANALYSIS OF DATASET

1. Top 5 Cyber-attacks happened on 5th October 2019.

Attack Name	Rank	Count of Attack Name
Adware.DealPly.1.Gen	1	13734
Gen	2	2749
JS	3	375
Heur.BZC.YAX.Voiv.11.039BA46C	4	220
Trojan.Exploit.Html.Iframe.Filedownload.GW	5	142

2. Top 5 attacking countries.

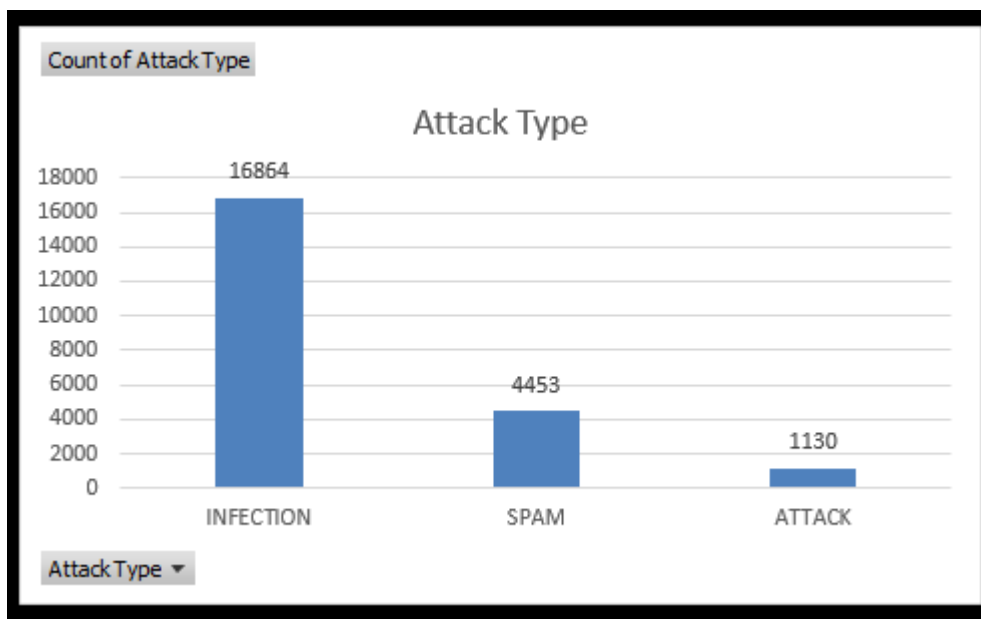
Attacking Countries	Count of Attacking Country
United States	8641
Brazil	2852
Canada	1316
United Kingdom	1102
Mexico	1012
Grand Total	14923



3. Attack type analysis.

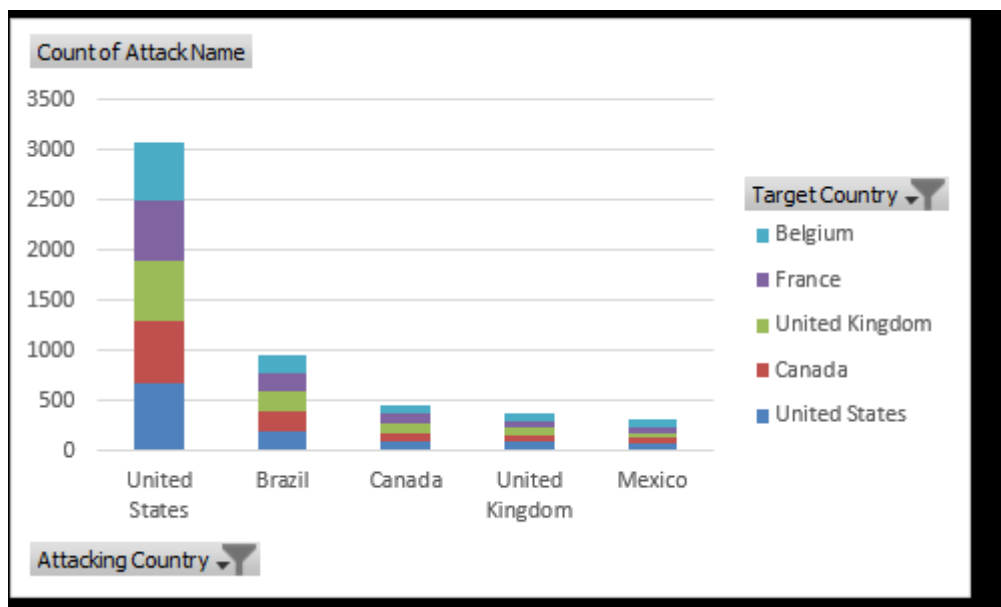
Attack Type	Count of Attack Type
INFECTION	16864
SPAM	4453
ATTACK	1130
Grand Total	22447

Time							
0	1	2	3	4	5	6	7
8	9	10	11	12	13	14	15
16	17	18	19	20	21	22	23



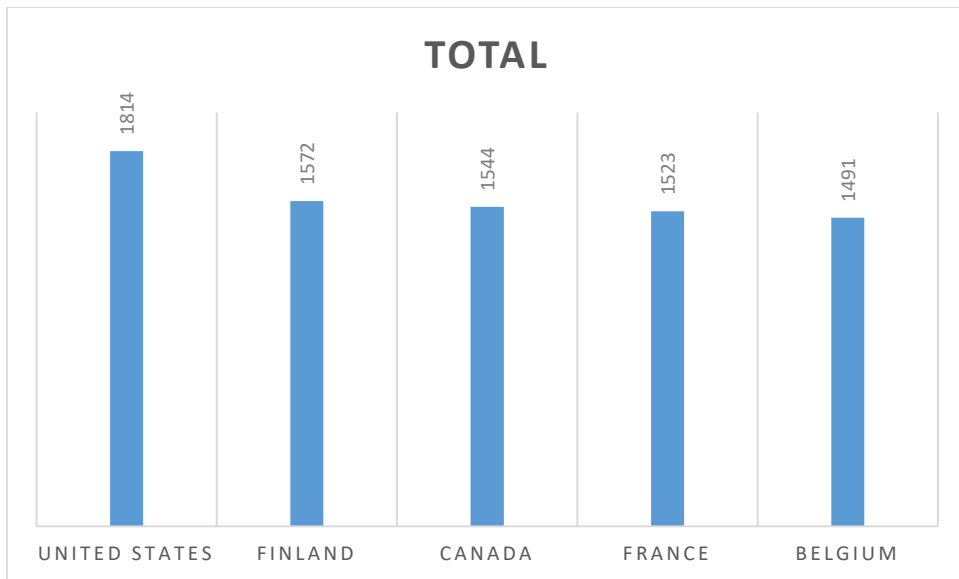
4. Countries Attacking Too & Fro.

Count of Attack Name	Target Countries					
Attacking Countries	United States	Canada	United Kingdom	France	Belgium	Grand Total
United States	670	616	608	604	582	3080
Brazil	198	199	192	180	189	958
Canada	86	94	95	92	77	444
United Kingdom	91	69	66	75	70	371
Mexico	66	57	57	55	84	319
Grand Total	1111	1035	1018	1006	1002	5172



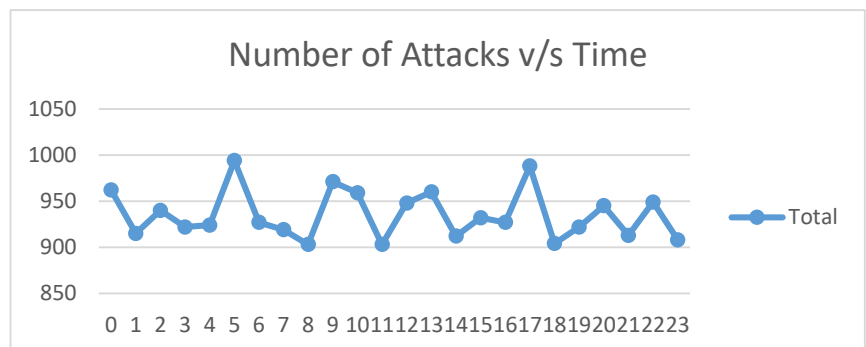
5. Top 5 victim countries.

Row Labels	Count of Target Country
United States	1814
Finland	1572
Canada	1544
France	1523
Belgium	1491
Grand Total	7944



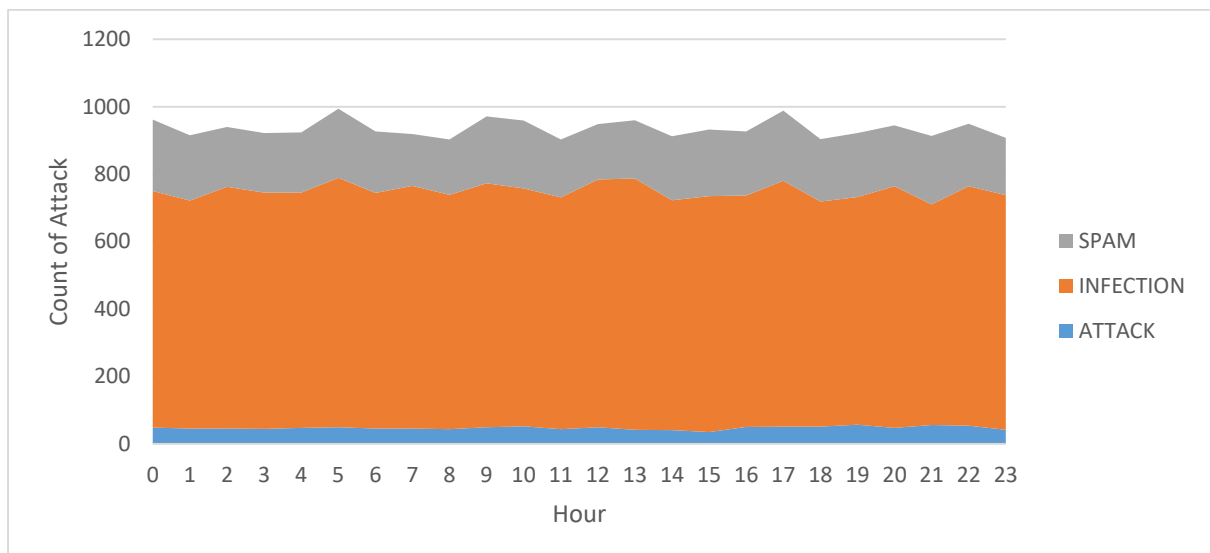
6. The number of attacks happened on the different hour of the day.

Hour	Count of Attacks
0	962
1	915
2	940
3	922
4	924
5	994
6	927
7	919
8	903
9	971
10	959
11	903
12	948
13	960
14	912
15	932
16	927
17	988
18	904
19	922
20	945
21	913
22	949
23	908
Grand Total	22447



7. Type of Attacks happened on the different hour of the day.

Count of Attack Name	Attack Types			
Hours	ATTACK	INFECTION	SPAM	Grand Total
0	48	702	212	962
1	45	676	194	915
2	45	717	178	940
3	44	701	177	922
4	47	698	179	924
5	49	740	205	994
6	45	699	183	927
7	45	720	154	919
8	43	696	164	903
9	49	724	198	971
10	52	706	201	959
11	43	688	172	903
12	49	735	164	948
13	42	745	173	960
14	41	681	190	912
15	35	700	197	932
16	50	687	190	927
17	51	729	208	988
18	51	668	185	904
19	57	675	190	922
20	47	717	181	945
21	56	654	203	913
22	54	710	185	949
23	42	696	170	908
Grand Total	1130	16864	4453	22447



References and Bibliography

1. <https://selenium-python.readthedocs.io/>
2. https://help.tableau.com/current/prep/en-us/prep_about.htm
3. <https://www.tableau.com/learn/tutorials/on-demand/getting-started-tableau-prep>
4. https://www.youtube.com/watch?v=e_o2S3oJnYQ
5. <https://www.tableau.com/learn/tutorials/on-demand/tableau-prep-interface>
6. <https://www.youtube.com/watch?v=YDvpgdy2ox8>
7. <https://www.youtube.com/watch?v=rwbho0CgEAE>
8. https://www.youtube.com/watch?v=RdTozKPY_OQ
9. <https://www.excel-easy.com/>
10. <https://www.tutorialspoint.com/excel/index.htm>
11. <https://www.guru99.com/excel-tutorials.html>