## HW#6 - Analyzing Disinformation Domains
PRASHANT TOMAR
04/09/2023

# Q1 - Data Analysis

## Answer

To conduct this analysis, I created a python script that extracts the top 10 websites based on the number of tweets they have. Once the data is collected, I utilized the request library to assess the status of the website and check whether it is still operational.

**Table 1:** D1 analysis

| Domains | Type | Tweets | Status |
|---|---|---|---|
| therealstrategy.com | Alternative Media | 7113.0 | Live |
| infowars.com | Alternative Media | 1741.0 | Live |
| newsbusters.org | Alternative Media | 1217.0 | Live |
| washingtonpost.com | MSM | 1108.0 | Live |
| nodisinfo.com | Alternative Media | 774.0 | Not Live |
| nytimes.com | MSM | 759.0 | Live |
| veteranstoday.com | Alternative Media | 586.0 | Live |
| beforeitsnews.com | Alternative Media | 580.0 | Live |
| rawstory.com | Alternative Media | 308.0 | Live |
| hoax.trendolizer.com | fact checker | 299.0 | Live |

**Table 2:** D2 analysis

| Domains | Type | Tweets | Status |
|---|---|---|---|
| 21stcenturywire.com | Blogging Site | 3088 | Not Live |
| clarityofsignal.com | Geopolitical Info | 2352 | Live |
| rt.com | News Channel | 1598 | Live |
| newsweek.com | Magazine site | 1249.0 | Not Live |
| alternet.org | Magazine | 1221 | Live |
| sputniknews.com | Newsfeed site | 1076.0 | Live |
| mintpressnews.com | Journalism | 919.0 | Not Live |
| cnn.com | News channel | 756.0 | Live |
| globalresearch.ca | research site | 724.0 | Live |
| theantimedia.org | News Channel | 682.0 | Live |

```python
 1 import pandas as pd
 2 from pandas import DataFrame
 3 import requests
 4
 5 def check_website_status(domain):
 6     website_status = 'Live'
 7     try:
 8         response = requests.get('http://' + domain, timeout = 10)
 9         if response.status_code != 200:
10             website_status = 'Not Live'
11     except:
12         website_status = 'Not Live'
13     return website_status
14
15 # Implementation for D1
16 D1 = pd.read_csv('D1.csv')
17 D1 = D1.sort_values('# Citations in our Alternative Narrative Tweets',
   ascending=False)
18 D1 = D1.iloc[0:10, [0,1,4]]
19 D1_list = []
20
21 for index, row in D1.iterrows():
22     website_status = check_website_status(row['Domain'])
23     website_object = {
24         'domain': row['Domain'],
25         'type': row['Media Type (Determined through Content Analysis)'
   ],
26         'tweets': row['# Citations in our Alternative Narrative Tweets'
   ],
27         'status': website_status
28         }
29
30     D1_list.append(website_object)
31
32 D1_updated_dataframe = DataFrame(D1_list)
33 D1_updated_dataframe.to_csv('D1_updated_dataset.csv')
34
35 # Implementation for D2
36 D2 = pd.read_csv('D2.csv')
37 D2 = D2.sort_values('Tweet count', ascending=False)
38 D2 = D2.iloc[0:10, [0,2]]
39 D2_list = []
40
41 for index, row in D2.iterrows():
42     website_status = check_website_status(row['Domain'])
43     website_object = {
44         'domain': row['Domain'],
```

```
45        'type': '',
46        'tweets': row['Tweet count'],
47        'status': website_status
48        }
49
50    D2_list.append(website_object)
51
52 D2_updated_dataframe = DataFrame(D2_list)
53 D2_updated_dataframe.to_csv('D2_updated_dataset.csv')
```

**Listing 1:** Generate a graph for various cluster sizes.

## Analysis

To conduct the analysis, I opened each website in a browser and navigated to its "about us" page to determine its media type. The top shared tweets were mainly from large organizations such as cnn.com and rt.com. Upon further analysis, it was revealed that most of the domains were affiliated with news channels and independent blogging industries that commonly share news and critiques regarding governments and geopolitical information.

# Q2 - Dataset Overlap

## Answer

To obtain the unique domains from each dataset DataFrame and identify any overlapping domains between different datasets, I utilized the set feature in Python and employed the sets intersection method.

Result of these overlaps are shared below,

**Table 3:** Dataset Combination

| Dataset Combination | Count |
|---|---|
| D1&D2 | 36 |
| D2&D3 | 21 |
| D1&D3 | 12 |
| D1&D2&D3 | 10 |

The table below displays the domains associated with each dataset. To identify the unique links that belong to each dataset, I utilized a set and have listed them in the table.

**Table 4:** Domains belongs to above dataset

| | | |
|---|---|---|
| ukcolumn.org | intellihub.com | ronpaulinstitute.org |
| beforeitsnews.com | lewrockwell.com | theintercept.com |
| off-guardian.org | blacklistednews.com | presstv.com |
| foxnews.com | thefreethoughtproject.com | rt.com |
| investmentwatchblog.com | theeventchronicle.com | themillenniumreport.com |
| gellerreport.com | thewashingtonstandard.com | theguardian.com |
| infowars.com | sott.net | 21stcenturywire.com |
| breitbart.com | nbcnews.com | collective-evolution.com |
| fellowshipoftheminds.com | globalresearch.ca | cnn.com |
| thestar.com | theantimedia.org | dailymail.co.uk |
| nydailynews.com | thetruthseeker.co.uk | thedailysheeple.com |
| abovetopsecret.com | theduran.com | nytimes.com |
| rubikon.news | humansarefree.com | |
| thedailybeast.com | davidicke.com | |
| upi.com | dcclothesline.com | |
| wakingtimes.com | heavy.com | |
| activistpost.com | worldtruth.tv | |

The following section contains the shared implementation of the dataset mentioned above.

```python
1  import pandas as pd
2  from pandas import DataFrame
3
4  def lowercase_domain_names(dataset):
5      dataset["Domain"] = dataset["Domain"].str.lower()
6      return dataset
7
8  def get_mutual_links(dataset1, dataset2):
9      links_set1 = set(dataset1.iloc[:,0])
10     links_set2 = set(dataset2.iloc[:,0])
11     mutual_links = set.intersection(links_set1, links_set2)
12     return mutual_links
13
14 def save_mutual_dataset(mutual_links, file_name):
15     mutual_dataset = DataFrame(mutual_links, columns=['Domain'])
16     mutual_dataset.to_csv(file_name, index=False)
17
18 def count_domains(file_name):
19     mutual_dataset = pd.read_csv(file_name)
20     count = len(mutual_dataset.index)
21     return count
22
23 def get_unique_domains(*args):
24     list_domains = []
25     for file_name in args:
26         mutual_dataset = pd.read_csv(file_name)
27         for index, row in mutual_dataset.iterrows():
28             list_domains.append(str(row['Domain']))
29     domain_set = set(list_domains)
30     unique_domains = list(domain_set)
31     return unique_domains
32
33 # Lowercase domain names in all datasets
34 D1_dataset = pd.read_csv('D1.csv')
35 D1_dataset = lowercase_domain_names(D1_dataset)
36 D2_dataset = pd.read_csv('D2.csv')
37 D2_dataset = lowercase_domain_names(D2_dataset)
38 D3_dataset = pd.read_csv('D3.csv')
39 D3_dataset = lowercase_domain_names(D3_dataset)
40
41 # Find mutual links between datasets
42 D1_D2_mutual_links = get_mutual_links(D1_dataset, D2_dataset)
43 D2_D3_mutual_links = get_mutual_links(D2_dataset, D3_dataset)
44 D1_D3_mutual_links = get_mutual_links(D1_dataset, D3_dataset)
45 D1_D2_D3_mutual_links = set.intersection(D1_D2_mutual_links, D3_dataset
       .iloc[:,0])
```

```
46
47  # Save mutual datasets to CSV files
48  save_mutual_dataset(D1_D2_mutual_links, 'D1_D2_mutual_links.csv')
49  save_mutual_dataset(D2_D3_mutual_links, 'D2_D3_mutual_links.csv')
50  save_mutual_dataset(D1_D3_mutual_links, 'D1_D3_mutual_links.csv')
51  save_mutual_dataset(D1_D2_D3_mutual_links, 'D1_D2_D3_mutual_links.csv')
52
53  # Count number of domains in each mutual dataset
54  D1_D2_mutual_count = count_domains('D1_D2_mutual_links.csv')
55  D2_D3_mutual_count = count_domains('D2_D3_mutual_links.csv')
56  D1_D3_mutual_count = count_domains('D1_D3_mutual_links.csv')
57  D1_D2_D3_mutual_count = count_domains('D1_D2_D3_mutual_links.csv')
58
59  # Print unique domains
60  unique_domains = get_unique_domains('D1_D2_mutual_links.csv', '
        D2_D3_mutual_links.csv', 'D1_D3_mutual_links.csv', '
        D1_D2_D3_mutual_links.csv')
61  for domain in unique_domains:
62      print(domain)
```

**Listing 2:** Generate a graph for various cluster sizes.

## Analysis

An interesting observation about the list of mutual domains is that all of the website domains included are highly popular.
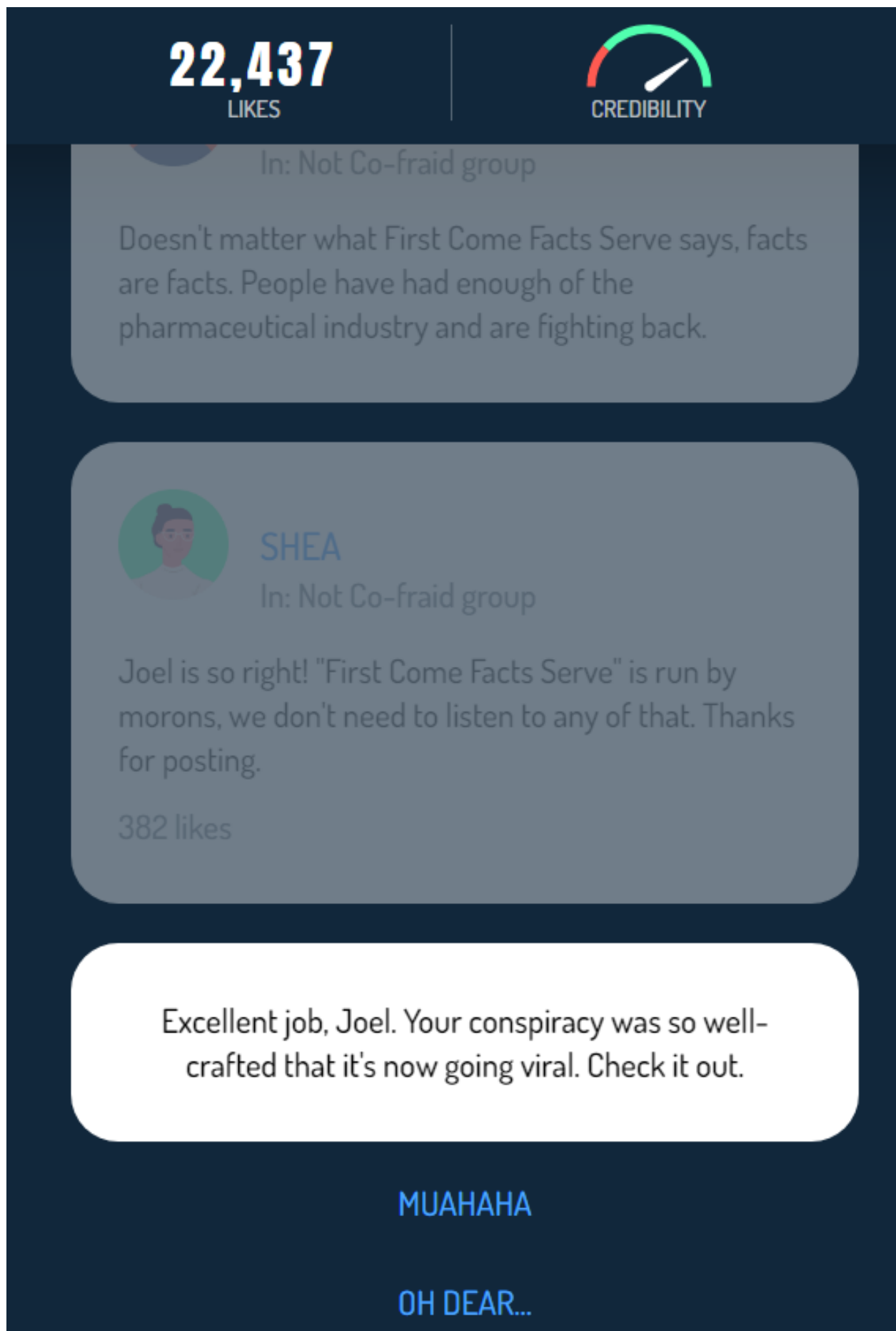
# Q3 - Disinformation Games

## Answer

I played the game on https://goviralgame.com and found it to be an engaging and informative experience. The game is designed to educate players about how disinformation spreads on social media platforms, and it does this by putting the player in the shoes of someone creating and spreading fake news. Through a series of choices and actions, the player can choose to create fake news articles, purchase bots to spread the articles, and use other tactics to increase their reach.

As I played the game, I found myself becoming more and more aware of the tactics used by people who create and spread fake news. I learned about the importance of verifying sources, fact-checking information before sharing it, and being critical of sensational headlines. The game also emphasized the importance of responsible social media use and how our actions online can have real-world consequences.

Overall, I found the game to be a thought-provoking and engaging way to learn about the dangers of disinformation. The game's design and user interface were easy to navigate, and the graphics were well-done. I appreciated the opportunity to experience the creation and spread of disinformation in a safe and controlled environment and would recommend this game to anyone looking to learn more about the topic. Below are some screenshots from my playthrough:

**Figure 1:** Final result of Game

# References

- Bar charts seaborn library, `https://seaborn.pydata.org/generated/seaborn.barplot.html`

- DataFrame operations, `https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.groupby.html`

- Tweepy, `https://docs.tweepy.org/en/latest/`

- Numpy, `https://numpy.org/doc/stable/reference/generated/numpy.log2.html`

- Goviralgame, `https://www.goviralgame.com/en/play`

- Github Datasets, `https://github.com/odu-cs432-websci/spring23-hw6-Badjedi04`