# HW 6 - Analyzing Disinformation Domains
Swathi Venkatesh
11/21/2021

# Q1

Datasets D1 and D2 include the number of tweets that each domain was shared in (found in the last column/field of the dataset).

*Q*: For each of D1 and D2, what were the top 10 domains in terms of tweets?

For each of the top 10 domains from the previous Q:

*Q*: Which ones are no longer live? *Q*: How would you classify the domain? Show this information in a table like the one below, sorted by number of tweets. You should have 2 tables, one for the top 10 in D1 and one for the top 10 in D2.

## Answer

```python
#!/usr/local/bin/python3
import pandas as pd
import numpy as np
#read the files in pandas dataframe
file1= pd.read_csv("D1.csv")
file2 = pd.read_csv("D2.csv")

#sort them in order
file1 = file1.sort_values(by='# Citations in our Alternative Narrative
    Tweets',ascending=False)
file2 = file2.sort_values(by='Tweet count', ascending=False)

search = file1.copy()

#get only 10 items
file1 = file1.head(10)
file2 = file2.head(10)


#drop unwanted columns
file1.drop(['Primary Orientation (Determined through Content Analysis)'
    , 'How Cited in Alternative Narrative of Shooting Events'],axis= 1,
    inplace=True)
file2.drop(['URL count'],axis =1, inplace=True)
```

```python
22  #rename columns
23  file1.rename(columns={"# Citations in our Alternative Narrative Tweets"
       :"Tweets","Media Type (Determined through Content Analysis)":"
       Website Type"},inplace=True)
24  file2.rename(columns={"Tweet count":"Tweet"},inplace=True)
25
26  #swap order for first d1
27  columns_swap = ["Domain","Tweets","Website Type"]
28  file1 = file1.reindex(columns=columns_swap)
29
30  #add new columns to the data with NAN values
31  file1['status']= np.nan
32  file2['Website Type']= np.nan
33  file2['status']= np.nan
34
35  #change column types to string
36  file2['Website Type'] = file2['Website Type'].astype(str)
37  file2['status'] = file2['status'].astype(str)
38
39  numCount = 0
40  temp =""
41   #Match domains in top 10 D2 dataframe with D1 to obtain Website Media
       Type
42  for index, row in file2.iterrows():
43      #find a match(es) and store as a dataframe
44      temp = search[search['Domain'].str.contains(row['Domain'])]
45      #check if data frame is empty
46      if(len(temp) == 0):
47          #assign NaN value
48          final = np.nan
49      else:
50          #assigne Media Type to final value
51          final = temp['Media Type (Determined through Content Analysis)'
       ].iloc[0]
52      #insert into file2 dataframe
53      file2.at[index,"Website Type"] = final
54
55  file1.to_csv("D1_new.csv", index = False, header=True)
56  file2.to_csv("D2_new.csv", index = False, header=True)
```

**Listing 1:** one.py

**Table 1:** Top 10 High Number of Tweets Domains (d1_new.csv)

| Domain | Tweets | Media | Status |
|---|---|---|---|
| therealstrategy.com | 7113 | Alternative Media | not live |
| infowars.com | 1741 | Alternative Media | not live |
| newsbusters.org | 1217 | Alternative Media | live |
| washingtonpost.com | 1108 | MSM | live |
| nodisinfo.com | 774 | Alternative Media | not live |
| nytimes.com | 759 | MSM | live |
| veteranstoday.com | 586 | Alternative Media | live |
| beforeitsnews.com | 580 | Alternative Media | live |
| rawstory.com | 308 | Alternative Media | live |
| hoax.trendolizer.com | 299 | fact checker | live |

**Table 2:** Top 10 High Number of Tweets Domains(d2_new.csv)

| Domain | Tweets | Media | Status |
|---|---|---|---|
| 21stcenturywire.com | 3088 | Alternative Media | live |
| clarityofsignal.com | 2352 | Not found(Alternative Media) | live |
| rt.com | 1598 | Foreign Government Media | live |
| newsweek.com | 1249 | Not found(MSM) | live |
| alternet.org | 1221 | Not found(Alternative Media) | Live |
| sputniknews.com | 1076 | Foreign Government Media | live |
| mintpressnews.com | 919 | Not found(Alternative Media) | live |
| cnn.com | 756 | MSM | live |
| globalresearch.ca | 724 | Alternative Media | live |
| theantimedia.org | 682 | Alternative Media | live |

## Discussion

For D1_new.csv I read the D1.csv file in a pandas dataFrame and was able to easily to sort the data according to the number of Tweets in Highest to lowest. As for D2_new.csv I did the same thing but I read in D2.csv and also compared the data with D1.csv to obtain the Media Values. Those values that are not found shows not found and what i manually classified them. I also manual checked the website on the browser to check if the where active website.

The Table 1 shows Top 10 High Number of Tweets Domains for D1 and Table 2 shows Top 10 High Number of Tweets Domains for D2. The Domains therealstrategy.com, infowars.com, nodisinfo.com are not live and all the other domains are live in d1_new.csv. In d2_new.csv all the top 10 domains are live. In Table 1 the highest number of tweets is 7113 and the top domain is therealstrategy.com. In Table 2 the highest number of tweets is 3088 and the top domain is 21stcenturywire.com. Most of the websites that have the Top 10 tweets belong to either Alternative

Media or Main Stream Media. But in Table 2 two of the domains belong to Foreign Government Media.

*Q: For each of D1 and D2, what were the top 10 domains in terms of tweets?*

Ans. The top 10 domains for D1 and D2 are as shown in Table 1 and Table 2.

*Q: Which ones are no longer live?*

Ans. The Domains therealstrategy.com, infowars.com, nodisinfo.com are not live and all the other domains are live in d1_new.csv. In d2_new.csv all the top 10 domains are live.

*Q: How would you classify the domain?* Ans. The domain classification is as shown in Table 1 and Table 2.

*Q: What can you say about the domains that have been frequently shared on Twitter?*

Comparing both the tables for D1 and D2 files, the top domains either belong to Alternate Media or MSM. For D1 top 10 domains, the ones that are not live belong to Alternate Media. The highest number of tweets in both the tables belong to website type Alternative Media. And the lowest number of tweets in Table 1 belong to Media type fast checker and to Alternate Media in Table 2.

# Q2

Compare the amount of overlap between the three datasets. Remember that:

- D1 - domains shared in tweets related to mass shootings

- D2 - domains shared in tweets related to the White Helmets in Syria

- D3 - domains found to be publishing false Coronavirus information

- a. domains that are present in both D1 and D2

- b. domains that are present in both D2 and D3

- c. domains that are present in both D1 and D3

- d. domains that are present in all three datasets

Create a table showing the number of domains in each of the datasets above. List out the domains in each of the datasets in your report. Upload each of the datasets to your GitHub repo.

## Answer

```python
1  #!/usr/local/bin/python3
2  #import necessary libraries
3  import pandas as pd
4
5  #read the files in pandas dataframe
6  file1= pd.read_csv("D1.csv")
7  file2 = pd.read_csv("D2.csv")
8  file3 = pd.read_csv("D3.csv")
9
10 def compareThisB(lowerCase,upperCase):
11     #create an empty final dataframe
12     number = 0
13     column_name = ["Domain"]
14     final = pd.DataFrame(columns = column_name)
15     for index, row in upperCase.iterrows():
16         #find a match(es) and store as a dataframe
17         #set Uppercases domain to lowercase so that it can propermatch
18         temp = lowerCase[lowerCase['Domain'] == row['Domain'].lower()]
19         #check if data frame is empty
20         if(len(temp) == 0):
21             pass
22         else:
23             final.at[number,"Domain"] = temp['Domain'].iloc[0]
24             number +=1
25     return final
26
27 """
28    Q2 Part a compare D1 and D2
29 """
30 a_final= compareThisB(file1,file2)
31
32 """
33    Q2 Part b compare D2 and D3
34 """
35 b_final = compareThisB(file2,file3)
36 """
37    Q2 Part c compare D1 and D3
38 """
39 c_final = compareThisB(file1,file3)
40
41 """
42    Q2 Part d compare a_final  and D3
43 """
44 d_final = compareThisB(a_final,file3)
45
46 #convert to csv files
```

```
47 a_final.to_csv("a_final.csv", index = False, header=True)
48 b_final.to_csv("b_final.csv", index = False, header=True)
49 c_final.to_csv("c_final.csv", index = False, header=True)
50 d_final.to_csv("d_final.csv", index = False, header=True)
```

**Listing 2:** two.py

## Discussion

The function called compareThisB(lowerCase,upperCase) to handle all the process for comparisons, for a,b,c,and d. This functions takes the two parameters. The second parameters get converted in all lower case before comparisons. Table 3 shows the domains that are present in both D1 and D2. All the data sets that are common either belong to Alternative Media or MSM. Table 4 shows the domains that are present in both D2 and D3. All the data sets that are common in D2 and D3 belong to some or other Media type. Table 5 shows the domains that are present in both data sets D1 and D3. All the data sets that are common in D1 and D3 belong to either Alternative Media or Foreign Government Media. Table 6 shows the domains that are present in both data sets D1, D2 and D3. s. All the data sets that are common in D1, D2 and D3 belong to either Alternative Media or Foreign Government Media.

*Q: Is there anything interesting about the domains that are present in multiple datasets?*

D1 talks about related to mass shootings, D2 is related to work of White Helmets in Syria and D3 about the domains publishing false Coronavirus information. Since we are considering the data sets that are not related to each other so its interesting to see that they still have the common domains. For mass shootings these domains have been listed to be having the correct information but D3 is all about the domains that have spread the false news about the corona virus. Its interesting to see the domains that have reported the correct information to be reporting all the false information about the corona virus.

**Table 3:** Domains that are present in both D1 and D2 from a_final.csv

| Domain |
| --- |
| rt.com |
| breitbart.com |
| theeventchronicle.com |
| themillenniumreport.com |
| beforeitsnews.com |
| thefreethoughtproject.com |
| veteranstoday.com |
| theintercept.com |
| theguardian.com |
| 21stcenturywire.com |
| infowars.com |
| thedailybeast.com |
| heavy.com |
| blacklistednews.com |
| presstv.com |
| dcclothsline.com |
| theantimedia.org |
| upi.com |
| investmentwatchblog.com |
| dailymail.co.uk |
| nydailynews.com |
| fellowshipoftheminds.com |
| thetruthseeker.co.uk |
| abovetopsecret.com |
| cnn.com |
| worldtruth.tv |
| sputniknews.com |
| lewrockwell.com |
| nytimes.com |
| intellihub.com |
| thedailysheeple.com |
| globalresearch.ca |
| foxnews.com |
| thestar.com |
| activistpost.com |
| nbcnews.com |

**Table 4:** Domains that are present in both D2 and D3 from b_final.csv

| Domain |
| --- |
| activistpost.com |
| beforeitsnews.com |
| breitbart.com |
| collective-evolution.com |
| dcclothesline.com |
| gellerreport.com |
| humansarefree.com |
| infowars.com |
| intellihub.com |
| ronpaulinstitute.com |
| sott.net |
| thewashingtonstandard.com |
| worldtruth.com |
| 21stcenturywire.com |
| davidicke.com |
| off-guardian.com |
| presstv.com |
| ukcolumn.org |
| rubikon.news |
| globalresearch.ca |
| theduran.com |

**Table 5:** Domains that are present in both D1 and D3 from c_final.csv

| Domain |
| --- |
| activistpost.com |
| beforeitsnews.com |
| breitbart.com |
| dcclothesline.com |
| infowars.com |
| intellihub.com |
| wakingtimes.com |
| worldtruth.com |
| zerohedge.com |
| 21stcenturywire.com |
| presstv.com |
| globalresearch.ca |

**Table 6:** Domains that are present in all three datasets from d_final.csv

| Domain |
| --- |
| activistpost.com |
| beforeitsnews.com |
| breitbart.com |
| dcclothesline.com |
| infowars.com |
| intellihub.com |
| worldtruth.com |
| 21stcenturywire.com |
| presstv.com |
| globalresearch.ca |

# Q4

There have been several online games created to educate people about disinformation and how it spreads on social media. Play one of the games at https://fakey.osome.iu.edu/, https://www.getbadnews.com, or https://goviralgame.com. Write a paragraph about your experience and some lessons you learned by playing the game. Take some screenshots as you play to include in your report.
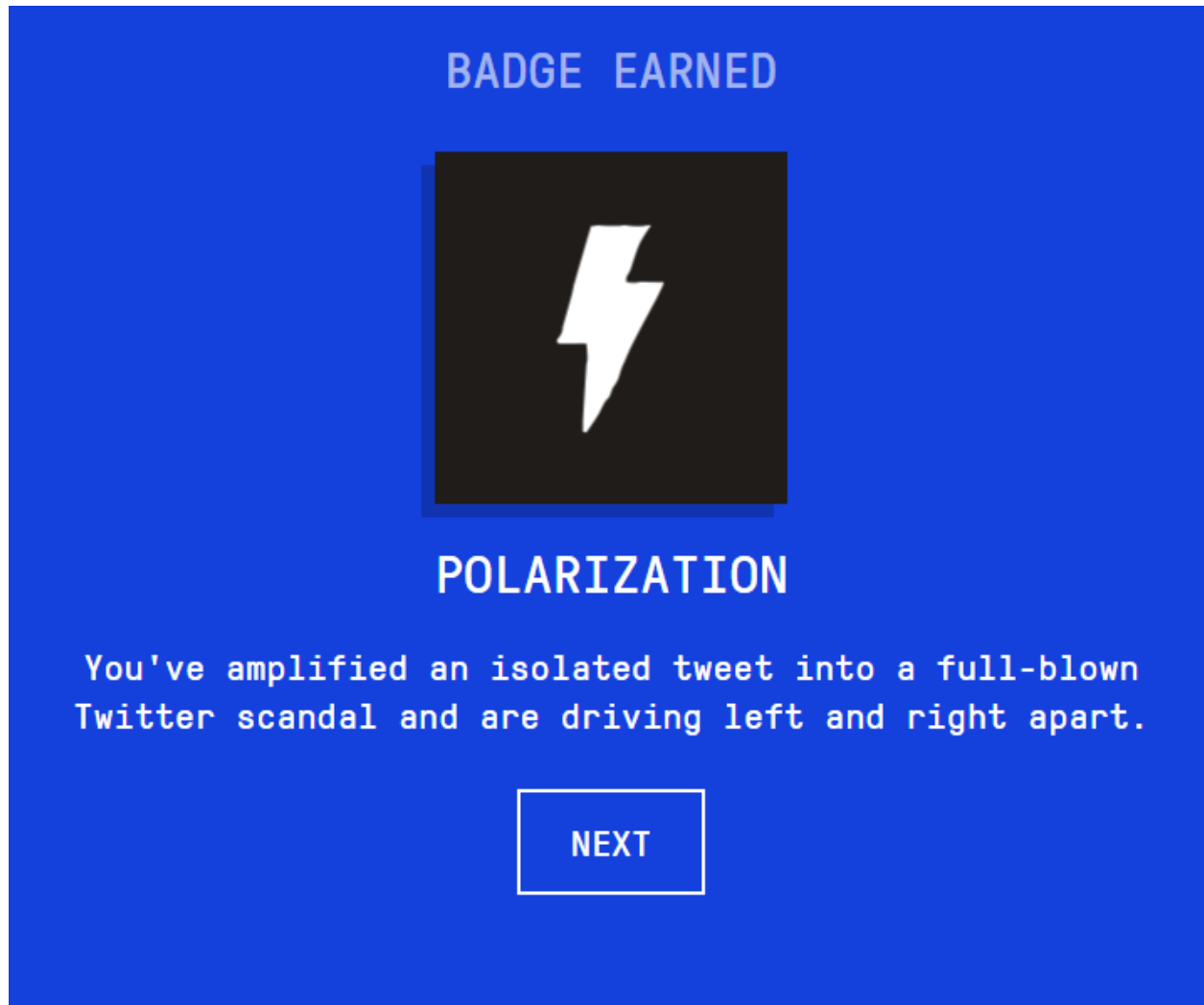
## Answer



**Figure 1:** Badge1
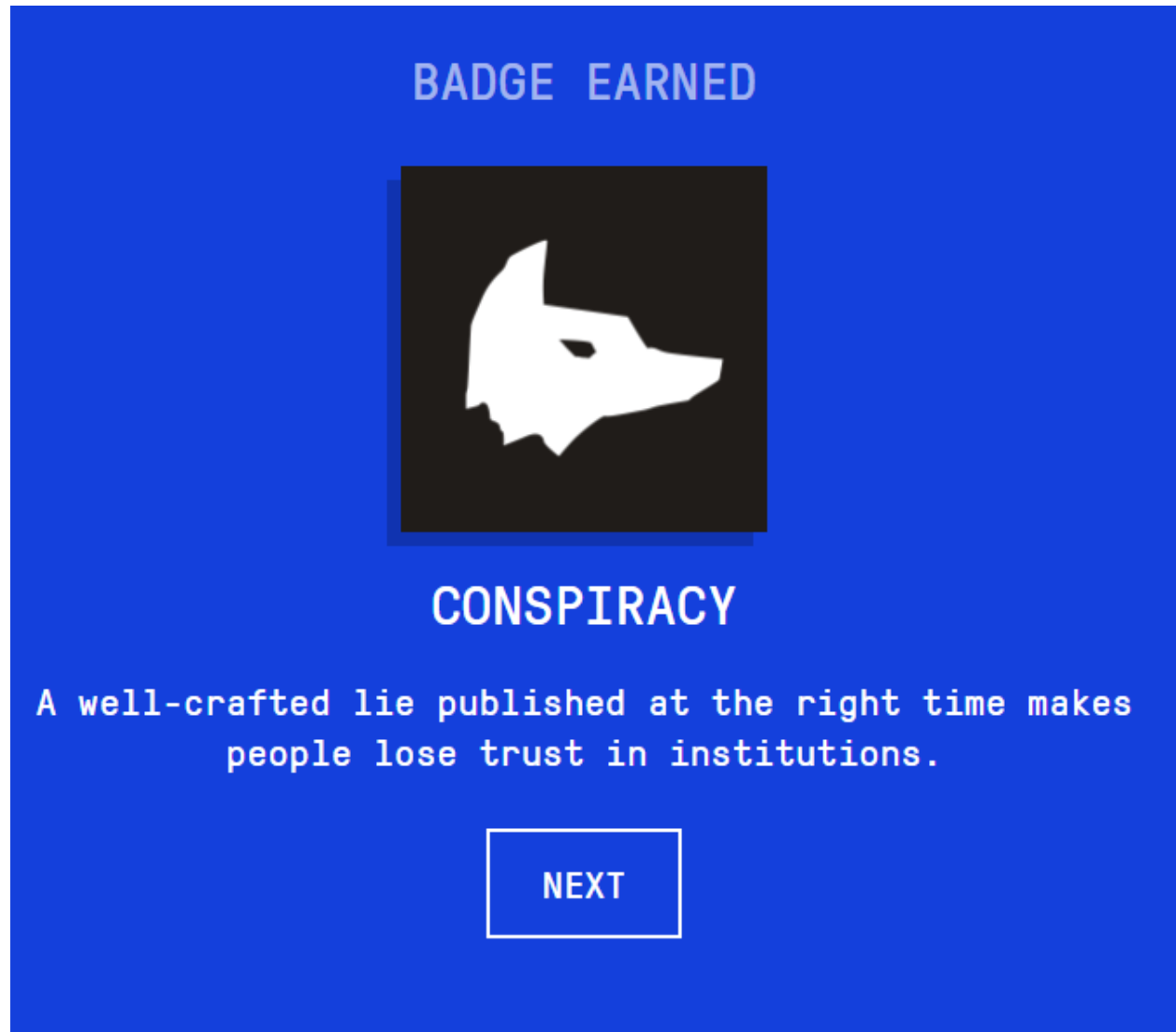
**Figure 2:** Badge2

**Figure 3:** Badge3

**Figure 4:** Badge4

**Figure 5:** Badge5

**Figure 6:** Badge6

**Figure 7:** Final Badge

## Discussion

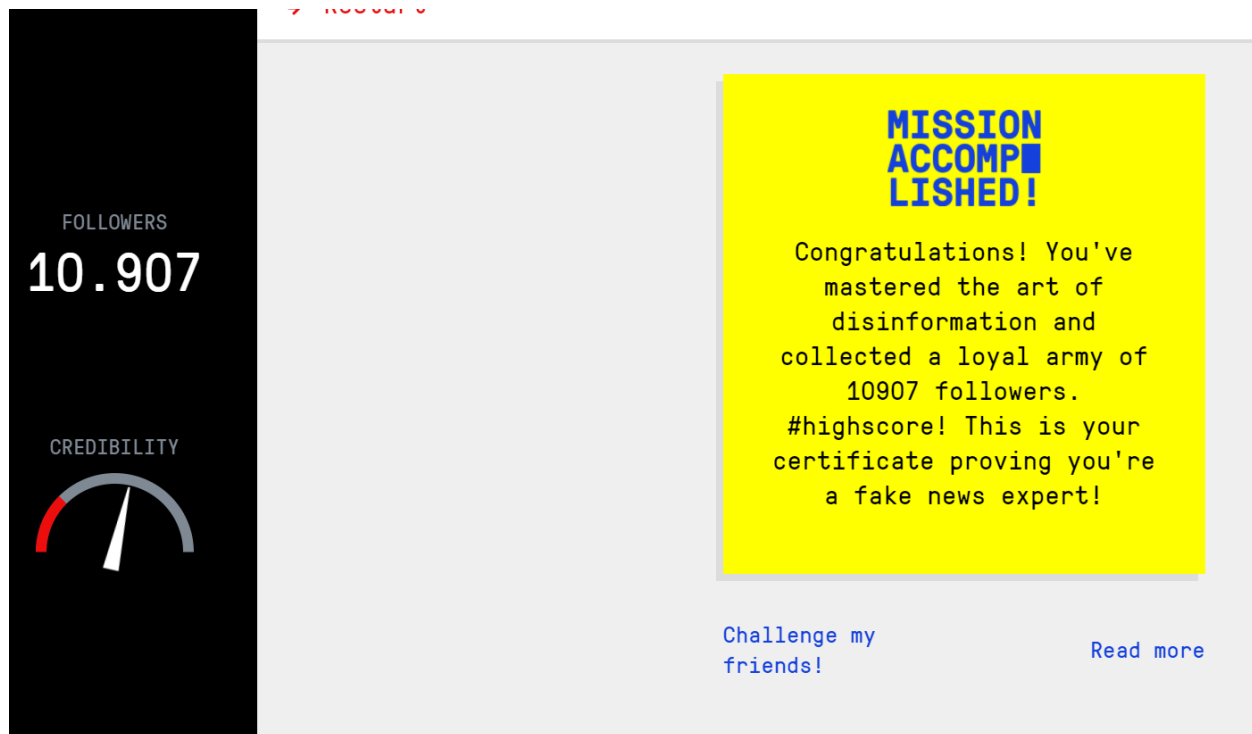*I used https://www.getbadnews.com/ to play the came and was able to earn the above badges below:*

- Impersonation in Figure 1 (Impersonating someone else and disguising myself as a credible news source which was highly effective in increase my followers)

- Emotion in Figure 2 (Playing to people's emotion out of fear, anger or compassing was a great tool for spreading my messages)

- Polarization in Figure 3 (By finding existing grievance and blowing them out of proportion, drove people apart and made think a story is much more important that it really was.)

- Conspiracy in Figure 4(I can use people's desires for the 'truth' as a tool to lure them into my band of followers)

- Discredit in Figure 5(When someone is attacking my credibility i strike back. I do not apologize nor do I play nice and above all I do not retreat)

- Trolling in Figure 6 ( Is a tool that evokes an emotional response such as anger, irritation or sadness.)

- The final score was 10907 followers in Figure 7.

# References

- `https://www.datacamp.com/community/tutorials/pandas-read-csv`

- `https://stackoverflow.com/questions/45164537/filter-pandas-data-frame-based-on-exact-string-match`

- `https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.sort_values.html`

- `https://stackoverflow.com/questions/12021754/how-to-slice-a-pandas-data-frame-by-position`

- `https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.rename.html`

- `https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.drop.html`

- `https://stackoverflow.com/questions/38288372/unable-to-drop-a-column-from-pandas-dataframe`

- `https://stackoverflow.com/questions/45164537/filter-pandas-data-frame-based-on-exact-string-match`

- `https://stackoverflow.com/questions/39092067/pandas-dataframe-convert-column-type-to-string-or-categorical`