

Silver and Reddit's Wallstreetbets

Overview

Research Question:

How are bots on r/WallStreetBets influencing posts about silver stock?

Context:

On February 1, 2021, silver had an 11% price increase, the largest surge in over a decade. Redditors on the subreddit r/WallStreetBets claim to be uninvolved, but it remains unclear if bot activity is influencing these posts. The subreddit r/WallStreetBets has over 9 million members.

Motivation:

We aim to identify the extent to which misinformation and bots affect r/WallStreetBets, as well as identifying other subreddits that have the potential to generate a GME-level stock upset.

Intial Data Load

In the first phase of our project, we scraped r/WallStreetBets for posts containing the keyword "silver". The following code was used to pull data from reddit using the API.

We began by loading data using pushshift.

```
In [23]: 1 #from google.colab import files
          2 #uploaded = files.upload()
          3 import pandas as pd
          4 final_df = pd.read_csv('updated_final_df.csv')
```

```
In [24]: 1 from urllib.request import urretrieve
          2
          3
          4 urretrieve('https://files.pushshift.io/reddit/comments/RC_2006-02.bz2', fil
          5
          6 import bz2, json
```

Print out some sample information from the reddit file

```

In [25]: 1  #We processed the data by using a push shift. <br>
          2  #This allowed us to gather the exact data that we wanted which included the
          3  #'url','author','title','subreddit','id','num_comments', 'score','upvote_rat
          4  #All these features were collected from submissions from the subreddit with
          5
          6
          7  counter =0
          8  with bz2.BZ2File('./example_reddit_comments.bz2', "r") as fp:
          9      for line in fp:
         10      counter = counter +1
         11      if counter < 5:
         12          job = json.loads(line)
         13          print( "subreddit", job['subreddit'])
         14          print( "author", job['author'])
         15          print( "upvotes", job['ups'])
         16          print( "score", job['score'])
         17          print
         18      else:
         19          break

```

```

In [26]: 1  #!/pip install psaw

```

```

In [27]: 1  from psaw import PushshiftAPI
          2
          3  api = PushshiftAPI()
          4
          5  # Define some keys for submission attribtues you care about -- these are sim
          6  filter_keys = ['url','author','title','subreddit','id','num_comments',
          7                  'score','upvote_ratio','domain','selftext', 'gilded']

```

Below we create the search and filter by the subreddit and the key term which we want to look at, in this case it is silver.

```

In [29]: 1  #https://pushshift.io/api-parameters/
          2
          3  start = int(datetime(2021, 1, 20).timestamp())
          4
          5  #can add a before if we want to
          6  #before = int(datetime(2021, 1, 20).timestamp())
          7
          8
          9  search = api.search_submissions(after=start,
         10                                subreddit='wallstreetbets',
         11                                filter=filter_keys,
         12                                q = 'silver',
         13                                sort='asc',
         14                                #score > 100,
         15                                limit=50000)

```

Create a list of the necessary information from pushshift

```
In [30]: 1 # Storage for the results
2 all_subs = []
3
4 # Loop through the search results to actually get data
5 for i,sub in enumerate(search):
6
7     # Add each result's dictionary (the .d_ attribute) to the all_subs
8     all_subs.append(sub.d_)
9
10    # Print out status updates every 10,000 submissions
11    if i % 10000 == 0:
12
13        # The current time so you know how long in between updates
14        time_now = datetime.now().time().replace(microsecond=0)
15
16        # The date of the submission to give you an idea of how far along you
17        record_date = datetime.utcfromtimestamp(sub.d_['created']).date()
18
19        # Print it out
20        print("{0:},} for {1} received at {2}".format(i,record_date,time_now))
```

```
In [31]: 1 import pandas as pd
2 #all_subs = pd.read_csv('final_df.csv')
3 subs_df = pd.read_csv('updated_final_df.csv')
4 #final_df = subs_df
5 # len(subs_df)
```

The JSON of the pushshift response is a dictionary of every post

```
In [32]: 1 all_subs[0]
2
3 OUTPUT
4
5 {'author': 'slackrooster',
6  'created': 1611130062.0,
7  'created_utc': 1611130062,
8  'domain': 'self.wallstreetbets',
9  'id': 'l14tq9',
10  'num_comments': 0,
11  'score': 1,
12  'selftext': '[removed]',
13  'subreddit': 'wallstreetbets',
14  'title': "Silver lining of Shitron's $GME BS for tomorrow's trading day",
15  'upvote_ratio': 1.0,
16  'url': 'https://www.reddit.com/r/wallstreetbets/comments/l14tq9/silver_lini
```

What we want to achieve here is create a date column that is easier to read instead of a timestamp.

We are then saving our DF as a CSV. We have 8,586 total data entries from r/WallStreetBets. This way, we won't have to go through the time consuming process of pulling from Pushshift, we can just load the CSV that we already have saved.

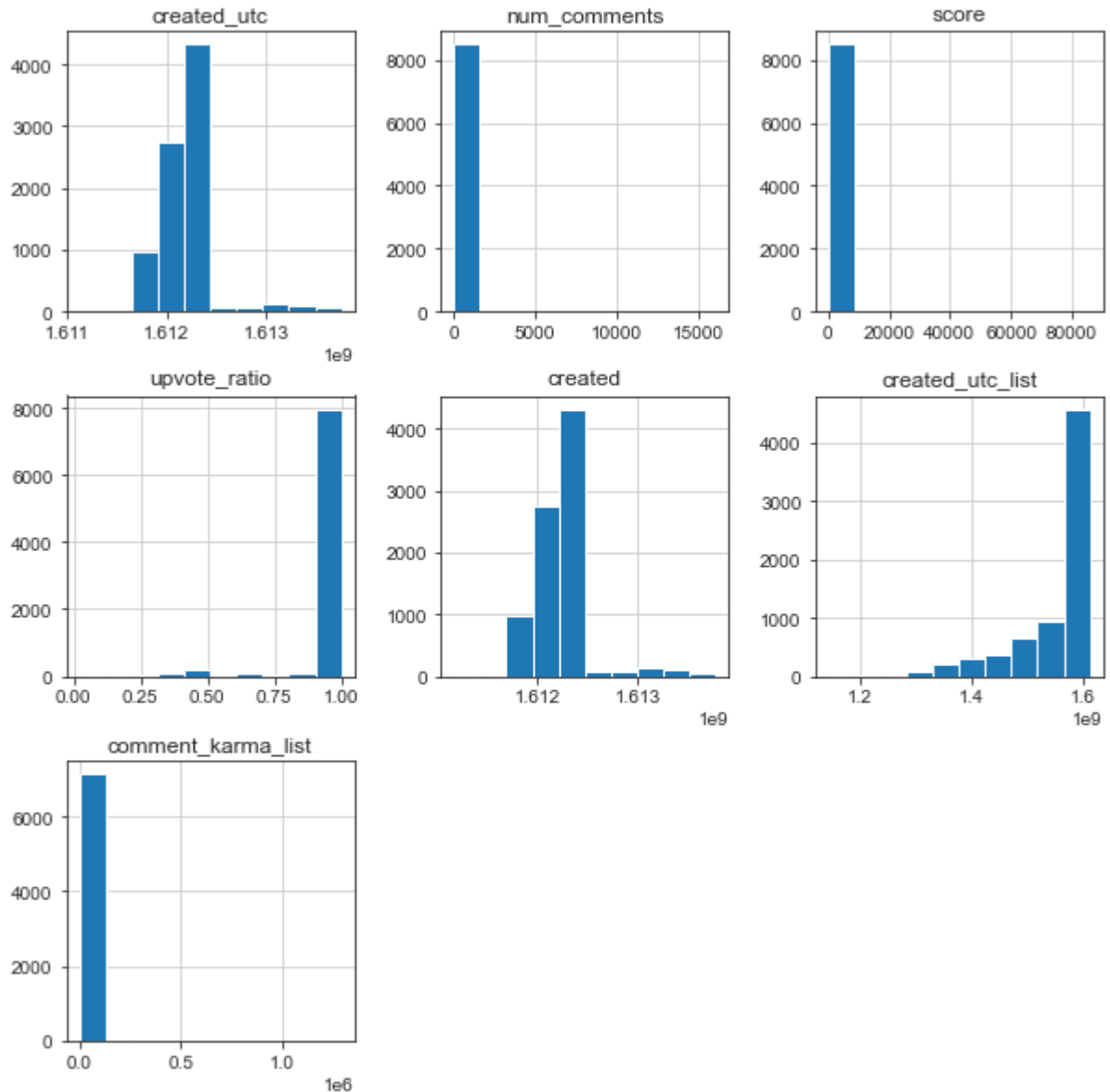
```
In [33]: 1 subs_df = pd.DataFrame(all_subs)
2 print('{:,}'.format(len(all_subs)))
3
4 subs_df['timestamp'] = subs_df['created'].apply(datetime.utcfromtimestamp)
5 subs_df['date'] = subs_df['timestamp'].apply(lambda x:x.date())
6 subs_df.to_csv('wallstreet_push_shift.csv',encoding='utf8',index=False)
7
8 subs_df = pd.read_csv('wallstreet_push_shift.csv')
9 from datetime import datetime
10 subs_df['timestamp'] = subs_df['created'].apply(datetime.utcfromtimestamp)
11 subs_df['date'] = subs_df['timestamp'].apply(lambda x:x.date())
12
13 subs_df.head()
```

Exploratory Analysis

Get a general idea of the shape of the data

```
In [34]: 1 final_df = pd.read_csv('final_df.csv')
          2 final_df.hist(figsize=(10,10))
          3
```

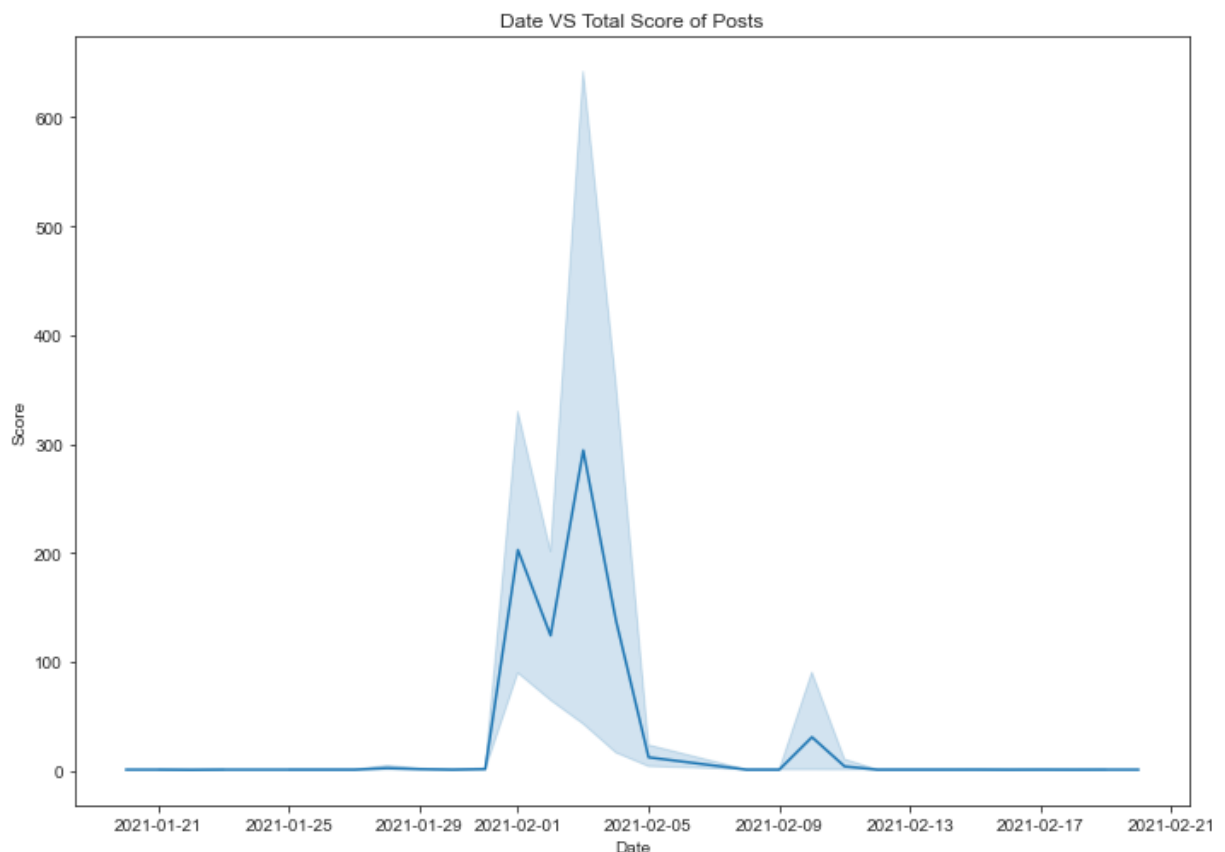
```
Out[34]: array([[<AxesSubplot:title={'center':'created_utc'}>,
                  <AxesSubplot:title={'center':'num_comments'}>,
                  <AxesSubplot:title={'center':'score'}>],
                [<AxesSubplot:title={'center':'upvote_ratio'}>,
                  <AxesSubplot:title={'center':'created'}>,
                  <AxesSubplot:title={'center':'created_utc_list'}>],
                [<AxesSubplot:title={'center':'comment_karma_list'}>,
                  <AxesSubplot:>, <AxesSubplot:>]], dtype=object)
```



We created 8 graphs to begin with to initially visualize our data. This gave us a good starting point to illustrate activity on r/WallStreetBets during the last week of January 2021.

Next, we used the csv to create a visualization. Figure 1 shows us the date and scores of subreddit data over time which helps identify the uptick in activity.

```
In [36]: 1 import seaborn as sns
2
3 from datetime import datetime
4 import matplotlib.pyplot as plt
5
6 # csv_subs_df = pd.read_csv('wallstreet_push_shift.csv')
7 final_df['timestamp'] = final_df['created'].apply(datetime.utcfromtimestamp)
8 final_df['date'] = final_df['timestamp'].apply(lambda x:x.date())
9
10 # #for stretching out the plot
11 # #https://stackoverflow.com/questions/31594549/how-do-i-change-the-figure-s
12
13 sns.set_style('ticks')
14 fig, ax = plt.subplots()
15 # the size of A4 paper
16 fig.set_size_inches(11.7, 8.27)
17
18 sns.lineplot(x='date',y='score',data=final_df)
19 plt.xlabel("Date")
20 plt.ylabel("Score")
21 plt.title("Date VS Total Score of Posts")
22 plt.show()
```

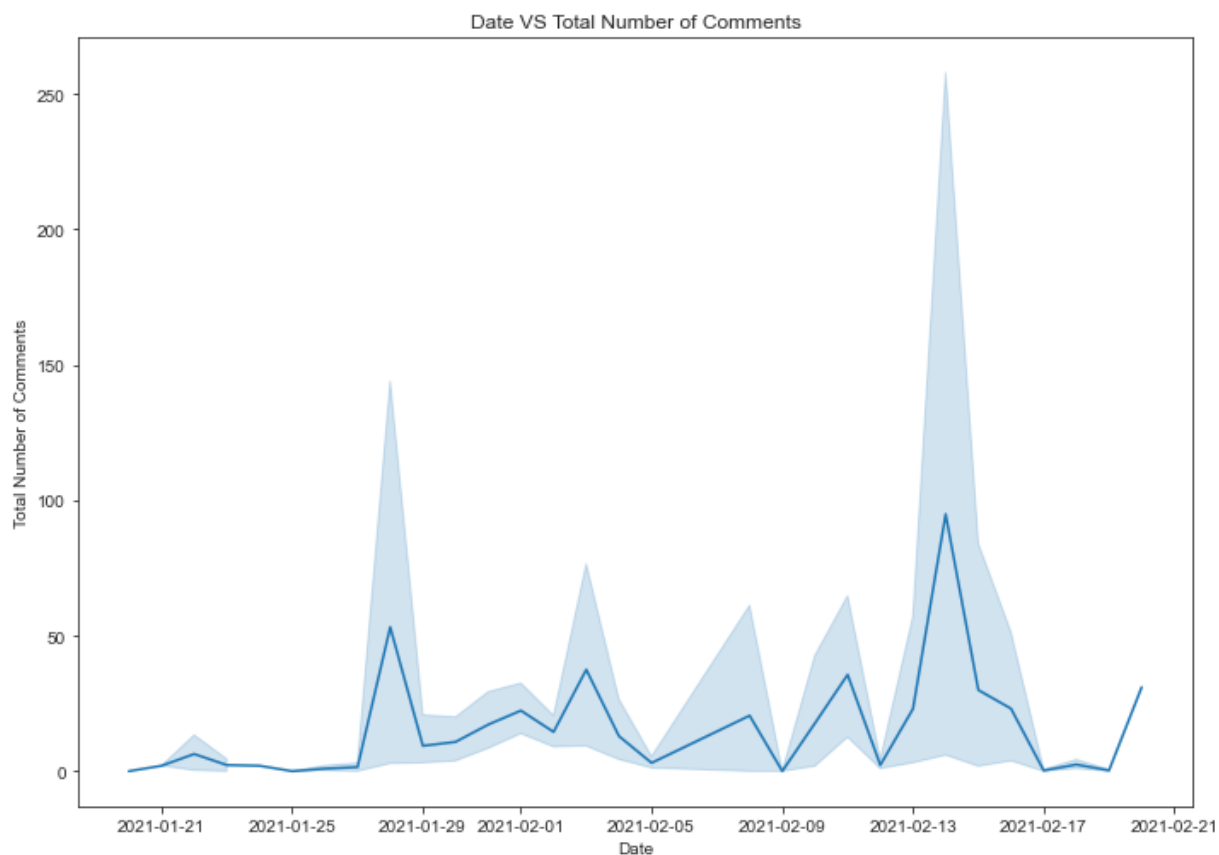


Based on the visualization, February 4th had the highest scoring posts, with another small peak on February 1st. This gives us some information, but score only reflects the ratio of upvoted and downvoted posts. We also want to learn about which accounts are doing the majority of upvotes.

Next, we chose to illustrate the number of comments on r/WallStreetBets over time.

In another visualization what we did was look at the number of comments based on dates over time and we chose to look at the 3 week span from about the January 21st to February 17th. This showed us where there were lots of comments.

```
In [37]: 1 sns.set_style('ticks')
2 fig, ax = plt.subplots()
3 # the size of A4 paper
4 fig.set_size_inches(11.7, 8.27)
5
6 sns.lineplot(x='date',y='num_comments',data=final_df)
7 plt.xlabel("Date")
8 plt.ylabel("Total Number of Comments")
9 plt.title("Date VS Total Number of Comments")
10 plt.show()
```



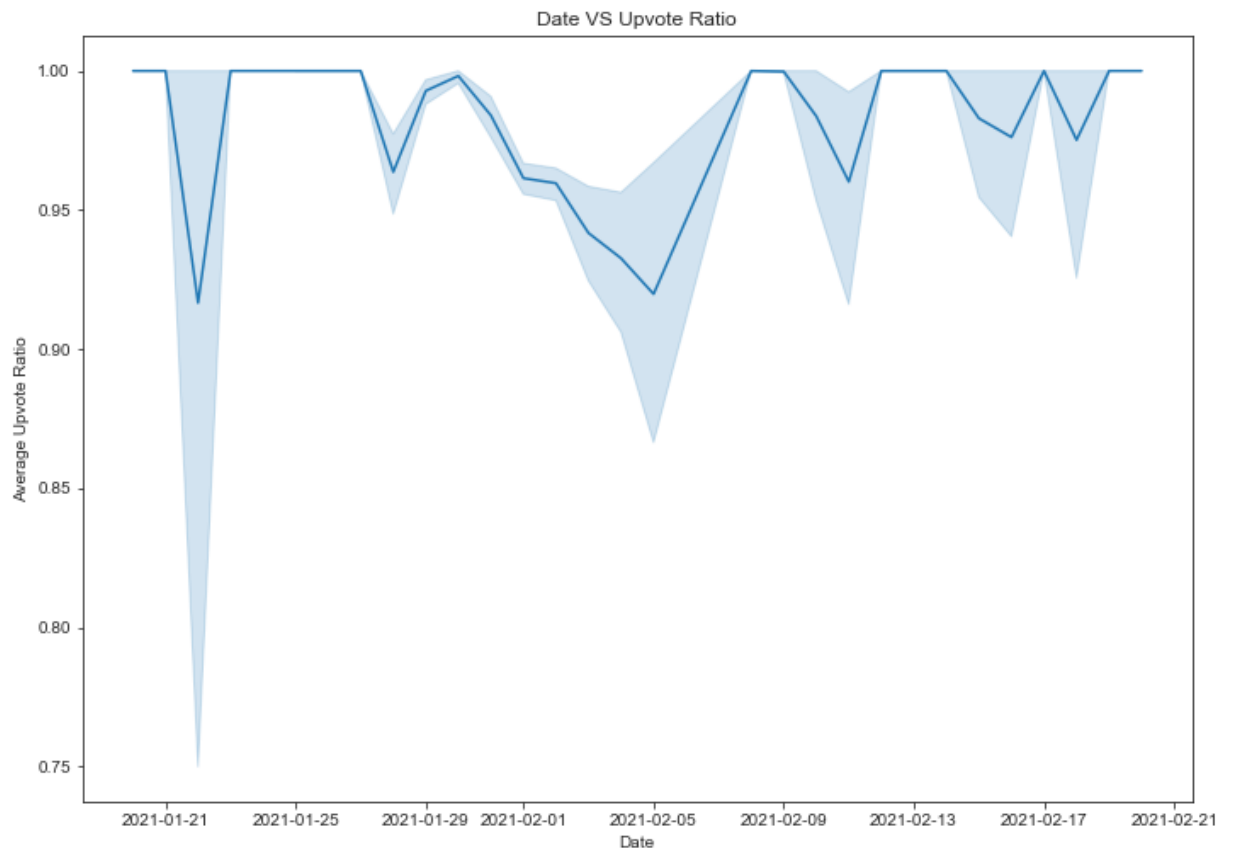
The peaks on January 27 and February 14 were most notable in this visualization. Silver stock reached a relative low point on January 27, prompting many users to buy stock to force a short squeeze (Silver 2021).

Another visualization we have is one which shows us the upvote to downvote ratio over time.


```

In [38]: 1 import seaborn as sns
          2
          3 from datetime import datetime
          4 import matplotlib.pyplot as plt
          5
          6 #for stretching out the plot
          7 #https://stackoverflow.com/questions/31594549/how-do-i-change-the-figure-size
          8
          9 sns.set_style('ticks')
         10 fig, ax = plt.subplots()
         11 # the size of A4 paper
         12 fig.set_size_inches(11.7, 8.27)
         13
         14 sns.lineplot(x='date',y='upvote_ratio',data=final_df)
         15
         16 plt.xlabel("Date")
         17 plt.ylabel("Average Upvote Ratio")
         18 plt.title("Date VS Upvote Ratio")
         19 plt.show()

```

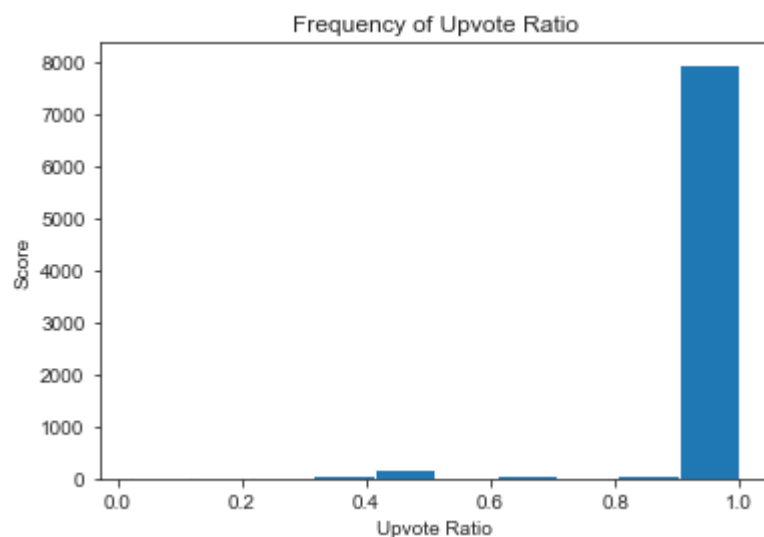


What we saw from the above graph showed us where there were some highly downvoted posts. However, the biggest peaks were on January 22 and February 6. This visualization did not give us

as much information as we were hoping for, as there were no notable price changes in Silver stock in the 24 hours around these peaks.

We then chose to graph the distribution of posts with different scores. Our next visualization puts into perspective that there were lots of posts which had close to a 1.0 ratio of upvoted posts.

```
In [39]: 1 #upvote ratio distribuibtion
2
3 subs_df.upvote_ratio.plot.hist()
4 plt.xlabel("Upvote Ratio")
5 plt.ylabel("Score")
6 plt.title("Frequency of Upvote Ratio")
7 plt.show()
```



Additional Data Cleaning and Sentiment Analysis

Add column showing if the selftext has been removed

```
In [40]: 1 is_self_text_removed = []
2
3 for post in final_df['selftext']:
4     if post == "[removed]":
5         is_self_text_removed.append('yes')
6     else:
7         is_self_text_removed.append('no')
8
9 final_df['is_self_text_removed'] = is_self_text_removed
10 final_df
```

Use textblob to add a column that calculates the polarity of the post title

```
In [41]: 1 from textblob import TextBlob
2
3 #blob = TextBlob('testing this is great')
4
5
6 blob_sent_polarity = []
7 for title in final_df['title']:
8     blob = TextBlob(title)
9     blob_sent_polarity.append(blob.polarity)
10
11 final_df['blob_sent_polarity'] = blob_sent_polarity
12
13 final_df.head()
```

Add a column tokenizing the title by word

```
In [42]: 1 import nltk
2 nltk.download('punkt')
3 from nltk.tokenize import sent_tokenize, word_tokenize#, #WordPunctTokenizer
4
5
6 final_df['word_tokenize_title'] = final_df['title'].apply(word_tokenize)
7 final_df
```

Import the empath library to run sentiment analysis on tokenized text.

```
In [43]: 1 #!pip install empath
2 from empath import Empath
3
4 lexicon = Empath()
5
6 empath_list = []
7 empath_list_title = []
8 for title in final_df['word_tokenize_title']:
9     categ = lexicon.analyze(title, normalize = True)
10     for key, value in categ.items():
11         if value != 0:
12             empath_list_title.append(key)
13     empath_list.append(empath_list_title)
14     empath_list_title = []
15
16
17 final_df['empath_categories'] = empath_list
18
19 final_df
```

Remove stopwords

```
In [44]: 1 # import nltk
2 # #nltk.download('stopwords')
3 # #!python -m nltk.downloader stopwords
4 from nltk.corpus import stopwords
5
6
7 from nltk.tokenize import sent_tokenize, word_tokenize, WordPunctTokenizer,
8 from nltk.tokenize.treebank import TreebankWordDetokenizer
9
10 stopeng = set(stopwords.words('english'))
11
12 stop_list = []
13 for word in final_df['title']:
14     tokens = word_tokenize( word.lower() )
15     tokens_nostop = [w for w in tokens if w not in stopeng]
16     tokens_nostop_updated = TreebankWordDetokenizer().detokenize(tokens_nostop)
17     stop_list.append(tokens_nostop_updated)
18
19 stop_list
20
21 final_df['title_tokenized'] = stop_list
```

Get the 20 most used words. Buy, squeeze and short seem to be some of the most popular words.

```
In [46]: 1 final_df = pd.read_csv('updated_final_df.csv')
2
3 from collections import Counter
4 from collections import defaultdict
5
6 def top20(thislist):
7     # First make a string out of the entire list
8     BIGstr = " ".join(thislist)
9     wordlist = BIGstr.split(" ")
10    wordcount = Counter(wordlist)
11    return(wordcount.most_common(20))
12
13 print(top20(final_df['title_tokenized']))
```

```
[('silver', 7093), ('.', 2284), ('buy', 1169), (''', 971), ('gme', 719), ('squeeze', 697), ('short', 557), ('buying', 447), ('-', 425), ('physical', 406), ('silver?', 380), ('silver,', 335), ('slv', 318), ('gold', 304), ('hold', 298), ('like', 296), ('amc', 296), ('news', 276), ('reddit', 275), ('silver!', 262)]
```

Look at the amount of posts with self text removed vs not removed

Visualizations and Analysis

Post Deletion

T test

Null Hypothesis: deleted post would have same upvote ratio as a undeleted post.

Alternative Hypothesis: Deleted post and Undeleted post would have different upvote ratio.

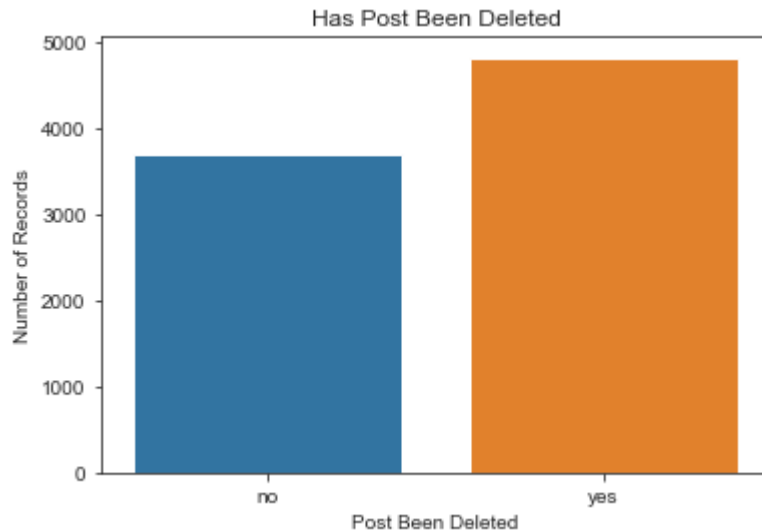
Based on the test, there was a significant difference between the samples.

```
In [47]: 1 #Import dataset
2 df = pd.read_csv('updated_final_df.csv')
3
4 #Array of delete post
5 options = ['[removed]', '[deleted]']
6
7 #create deleted only dataset
8 deleted_df = df.loc[df['selftext'].isin(options)]
9
10 #create non-deleted only dataset
11 notdeleted_df = df.loc[~df['selftext'].isin(options)]
12
13 #Sample 1000 rows from each dataset
14 deleted_df_sample = deleted_df.sample(n = 1000)
15 notdeleted_df_sample = notdeleted_df.sample(n = 1000)
16
17 #T test
18 from scipy import stats as st
19 #to each upvote_ratio to array
20 a = deleted_df_sample[['upvote_ratio']].to_numpy()
21 b = notdeleted_df_sample[['upvote_ratio']].to_numpy()
22 st.ttest_ind(a=a, b=b, equal_var=True)
```

```
Out[47]: Ttest_indResult(statistic=array([2.46741533]), pvalue=array([0.01369252]))
```

Of our dataset, over 60% of posts ended up being deleted. This shows that r/WallStreetBets moderators didn't want anyone to see most of the posts about silver, implying these posts contained misinformation.

```
In [48]: 1 #import seaborn as sns
2 final_df = pd.read_csv('updated_final_df.csv')
3
4 sns.countplot(x='is_self_text_removed', data = final_df)
5 plt.xlabel("Post Been Deleted")
6 plt.ylabel("Number of Records")
7 plt.title("Has Post Been Deleted")
8 plt.show()
```



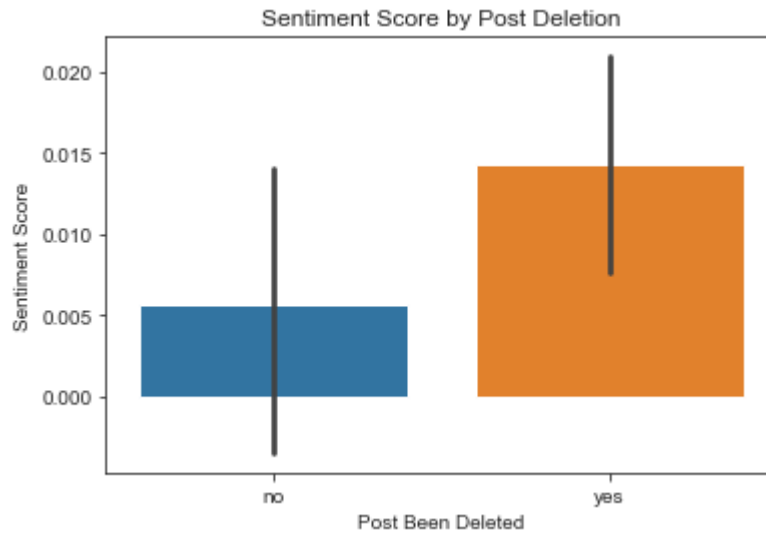
What the next visualization does is shows the comparison between the number of comments and the karma that the comments get on a given post.

```
In [49]: 1
2
3 final_df = pd.read_csv('updated_final_df.csv')
4 temporary_df = final_df.copy()
5 temporary_df.rename(columns= {'is_self_text_removed': 'Is Post Deleted'}, in
6
7 fig, ax = plt.subplots()
8 # the size of A4 paper
9
10 fig.set_size_inches(8.7, 8.27)
11
12 #sns.boxplot(x= "Is Post Deleted", y= "comment_karma_list", data= temporary_
13 sns.scatterplot(x='num_comments', y='comment_karma_list', data= temporary_df
14 plt.xlabel("Number of Comments")
15 plt.ylabel("Comment Karma")
16 plt.title("Number of Comments vs Karma")
17 plt.show()
```

From the results, we can see that there were several deleted posts with a lot of comment karma.

Look at the sentiment between removed posts and not removed posts. Posts which ended up being deleted had higher sentiment than posts which did not. This may be because people advocating for silver used more positive words to describe silver.

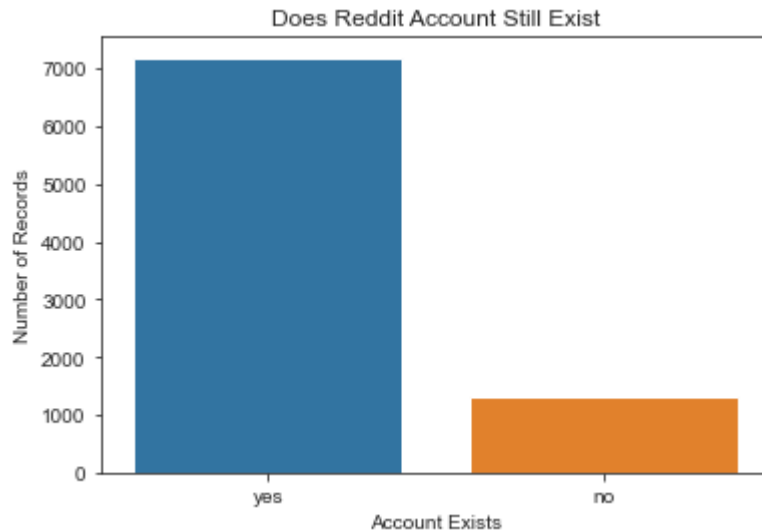
```
In [50]: 1 sns.barplot(x="is_self_text_removed", y="blob_sent_polarity", data=final_df)
2 plt.xlabel("Post Been Deleted")
3 plt.ylabel("Sentiment Score")
4 plt.title("Sentiment Score by Post Deletion")
5 plt.show()
6
```



Account Deletion

Look to see how many reddit accounts have been removed since the original post. Over 12% of accounts have been deleted since their original post relating to silver. Although we don't currently have a baseline to compare this to, this number seems very high.

```
In [51]: 1 import seaborn as sns
2
3 sns.countplot(x='account_exists', data = final_df)
4
5 plt.xlabel("Account Exists")
6 plt.ylabel("Number of Records")
7 plt.title("Does Reddit Account Still Exist")
8 plt.show()
```



Posts which have been deleted have slightly lower average comment karma than posts which have not been deleted.

Look at the 20 most used words for empath categories for account deleted and not deleted. They had very similar results. The account deleted had slightly more occurrences of economic related posts than the accounts which have not been deleted.


```

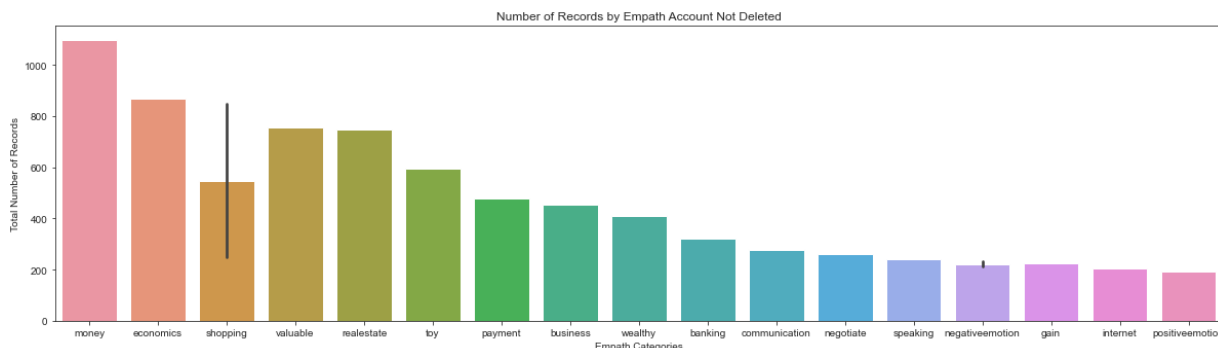
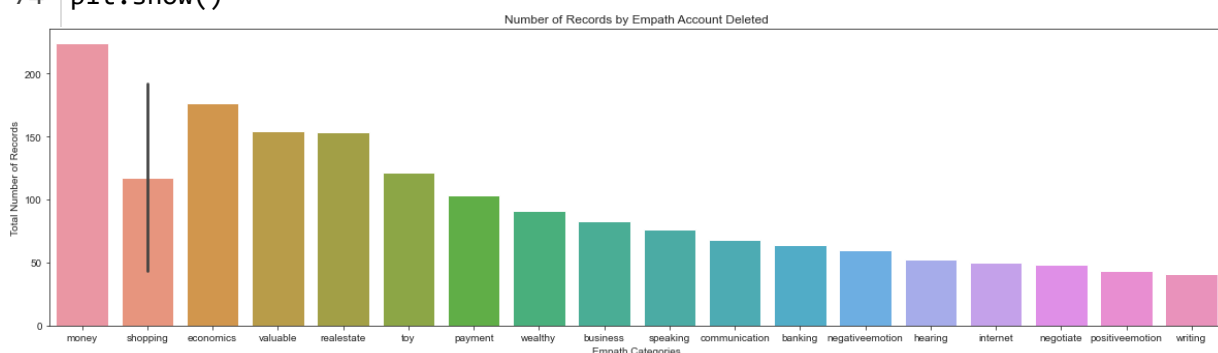
In [53]: 1 account_no_exist = final_df[final_df.account_exists == 'no']
2 account_yes_exist = final_df[final_df.account_exists == 'yes']
3
4 #split account exists vs account doesn't exist
5
6
7 word_cloud_list_key = []
8 word_cloud_list_value = []
9 word_cloud_thing = top20(account_no_exist['empath_categories'])
10 for thing in word_cloud_thing:
11     the_key = thing[0]
12     import re
13     the_key = re.sub(r'^A-Za-z', '', the_key)
14     word_cloud_list_key.append(the_key)
15     the_value = thing[1]
16     #the_value = re.sub(r'^A-Za-z', '', the_value)
17     word_cloud_list_value.append(the_value)
18
19 word_cloud_list_key[0] = "blank"
20 word_cloud_list_key
21
22 from pandas import DataFrame
23
24 word_cloud_df = DataFrame (word_cloud_list_key,columns=['Key'])
25 word_cloud_df['value'] = word_cloud_list_value
26
27 word_cloud_df = word_cloud_df.drop(labels= 0, axis=0)
28
29 word_cloud_df
30
31 sns.set_style('ticks')
32 fig, ax = plt.subplots()
33 # the size of A4 paper
34 fig.set_size_inches(20.4, 5.27)
35 sns.barplot(x="Key", y = "value", data = word_cloud_df)
36 plt.xlabel("Empath Categories")
37 plt.ylabel("Total Number of Records")
38 plt.title("Number of Records by Empath Account Deleted")
39 plt.show()
40
41
42 word_cloud_list_key = []
43 word_cloud_list_value = []
44 word_cloud_thing = top20(account_yes_exist['empath_categories'])
45 for thing in word_cloud_thing:
46     the_key = thing[0]
47     import re
48     the_key = re.sub(r'^A-Za-z', '', the_key)
49     word_cloud_list_key.append(the_key)
50     the_value = thing[1]
51     #the_value = re.sub(r'^A-Za-z', '', the_value)
52     word_cloud_list_value.append(the_value)
53
54 word_cloud_list_key[0] = "blank"
55 word_cloud_list_key
56

```

```

57 from pandas import DataFrame
58
59 word_cloud_df = DataFrame (word_cloud_list_key,columns=['Key'])
60 word_cloud_df['value'] = word_cloud_list_value
61
62 word_cloud_df = word_cloud_df.drop(labels= 0, axis=0)
63
64 word_cloud_df
65
66 sns.set_style('ticks')
67 fig, ax = plt.subplots()
68 # the size of A4 paper
69 fig.set_size_inches(20.4, 5.27)
70 sns.barplot(x="Key", y = "value", data = word_cloud_df)
71 plt.xlabel("Empath Categories")
72 plt.ylabel("Total Number of Records")
73 plt.title("Number of Records by Empath Account Not Deleted")
74 plt.show()

```



Account Creation

Look at account creation over time. The amount of created accounts showed a similar pattern to the total number of comments peaking around the end of January. Over 1,000 accounts were created and one of their first posts had to do with silver.

```

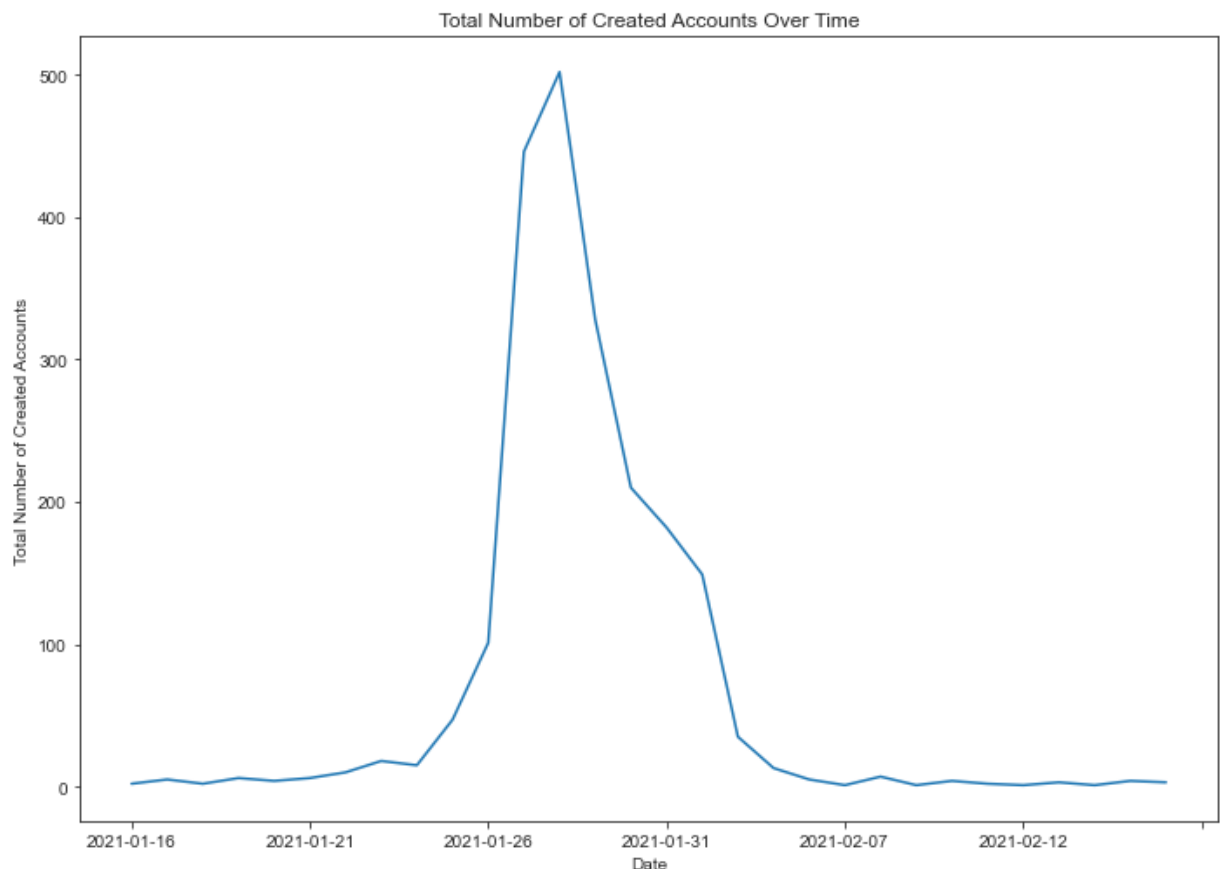
In [54]: 1 import datetime
2 import time
3
4 created_df = final_df.copy()
5 #created_df = final_df[['created_utc_list']]
6 created_df = created_df.dropna()
7 created_df['created_pst'] = created_df['created_utc_list'].apply(lambda t: t
8 filter_created_df = created_df[created_df.created_pst > "2021-01-15"]
9
10 import numpy as np
11 import pandas as pd
12 import matplotlib.pyplot as plt
13 %matplotlib inline
14 s = filter_created_df['created_pst'].value_counts().sort_index()
15 #print(s)
16
17 #plt.tick_params(axis='x', which='major', labelsize=3)
18 fig, ax = plt.subplots()
19 # the size of A4 paper
20 fig.set_size_inches(11.7, 8.27)
21 plt.xlabel("Date")
22 plt.ylabel("Total Number of Created Accounts")
23 plt.title("Total Number of Created Accounts Over Time")
24 s.plot()

```

```

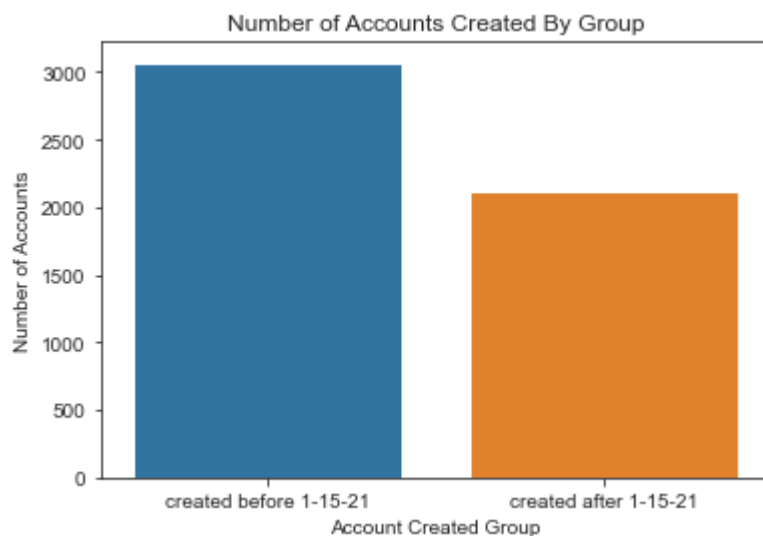
Out[54]: <AxesSubplot:title={'center':'Total Number of Created Accounts Over Time'}, xla
bel='Date', ylabel='Total Number of Created Accounts'>

```



Next, we split our data into accounts that were created before and after January 15th. We saw that nearly 40% of accounts which posted about silver were posted after January 15, 2021.

```
In [55]: 1 #downvoted posts
2
3 #subs_df.dtypes
4
5 #as_int(sub_df['upvote_ratio'])
6
7 created_df_list = []
8
9 for x in created_df['account_exists']:
10     created_df_list.append('created before 1-15-21')
11
12 created_df['date_filter'] = created_df_list
13
14
15 created_after_1_15_21 = (created_df.created_pst > "2021-01-15")
16
17 created_df.loc[created_after_1_15_21, 'date_filter'] = 'created after 1-15-21'
18
19 sns.countplot(x='date_filter', data= created_df)
20 plt.xlabel("Account Created Group")
21 plt.ylabel("Number of Accounts")
22 plt.title("Number of Accounts Created By Group")
23 plt.show()
```



Summary of Key Findings

1. Many posts have been deleted by moderators (60%)

2. A high percentage of reddit accounts have been deleted
3. There are some notable differences in sentiment of deleted posts and non deleted posts
4. There was a spike in activity on January 27-28, concurrently with the relative drop in price in silver stock
 - This activity included mass downvoting of posts, surges in creating and deleting accounts, and comments

Overall Summary

Overall when we were looking at Reddit data on r/WallStreetBets we started with the raw data and transformed that into data which we would be able to analyze. We layed out our definitions for bots this allowed us to look at specific characteristics in which we could determine if there was significant bot activity which could influence the prices of silver by influencing users to pursue it the same way that they bought up Gamestop stocks.

What our visualizations have showed us is that there has been a significant amount of accounts which were created just before the price of silver went up and also deleted days later, to us this proves that there were bots influencing users of this reddit and causing the price of silver to increase more than it had in many years.