

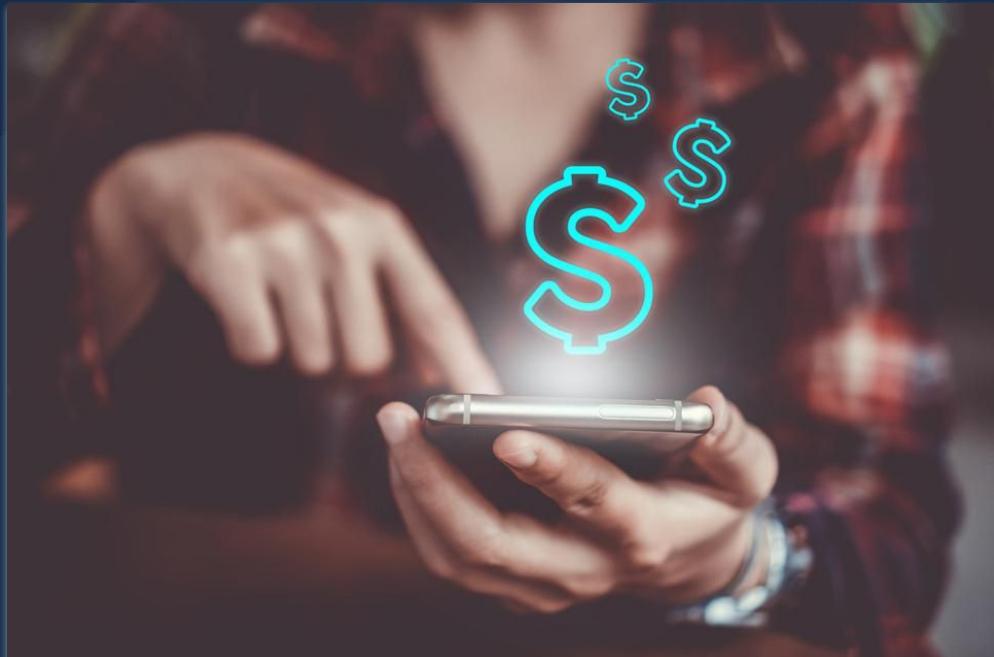
UC San Diego

**Social Media Talks Stocks
2021-12-10**

Amit Sachan, Mustafa Burny, Sam Courtney

Overview

- Every day, millions of people discuss stock prices and financial markets on social media.
- Recently, social media has started to significantly impact stock prices.
- We wanted to investigate the relationship between financial markets and social media activity.
- We wanted to investigate the engagement period.



Problem

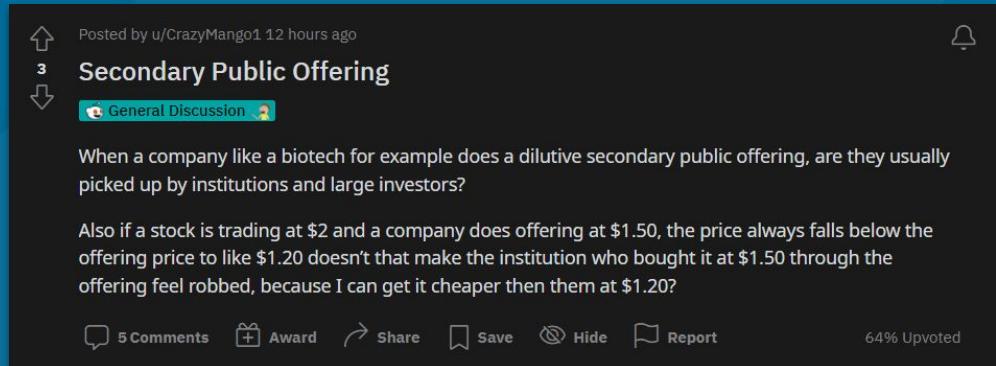


1. Can we extract what companies and stocks are mentioned in Reddit discussions?
2. Can we correlate the Reddit discussions to stock price movements?
3. Can we understand the topics being mentioned in Reddit discussions?
4. Can we understand relation between posts and comments lifetime?

DATASET OVERVIEW

Reddit

- Anonymous users can create Reddit submissions in sub-reddits associated with specific topics.
- Each Reddit submission has an author, a title, and (optional) body.
- Each Reddit submission can have a collection of comments and sub-comments.
- Reddit submissions and Reddit comments have a score determined by how many Reddit users “liked” or “disliked” the content. It’s called “karma”



Posted by u/CrazyMango1 12 hours ago

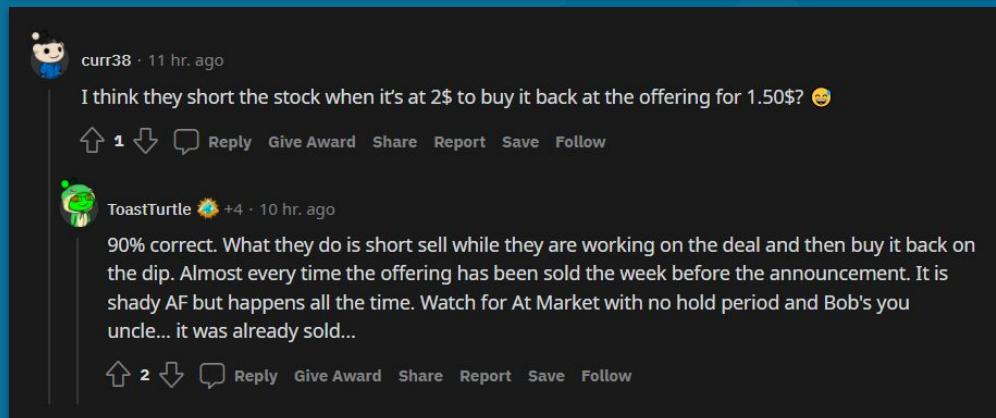
Secondary Public Offering

General Discussion

When a company like a biotech for example does a dilutive secondary public offering, are they usually picked up by institutions and large investors?

Also if a stock is trading at \$2 and a company does offering at \$1.50, the price always falls below the offering price to like \$1.20 doesn't that make the institution who bought it at \$1.50 through the offering feel robbed, because I can get it cheaper then them at \$1.20?

5 Comments Award Share Save Hide Report 64% Upvoted



curr38 · 11 hr. ago

I think they short the stock when it's at 2\$ to buy it back at the offering for 1.50\$? 😊

1 Reply Give Award Share Report Save Follow

ToastTurtle +4 · 10 hr. ago

90% correct. What they do is short sell while they are working on the deal and then buy it back on the dip. Almost every time the offering has been sold the week before the announcement. It is shady AF but happens all the time. Watch for At Market with no hold period and Bob's your uncle... it was already sold...

2 Reply Give Award Share Report Save Follow

Extracting Data from Reddit - Reddit API & PRAW

- Reddit provides a REST API to retrieve submissions, comments, and user information
- Requires creating application/bot credentials via Reddit web interface
- praw (Python Reddit API Wrapper) is a Python package that facilitates Reddit API access
- API is rate-limited to about 1 request per second

<https://www.reddit.com/dev/api/>
<https://praw.readthedocs.io/>

```
# Initialize Reddit instance
reddit = praw.Reddit(
    client_id = REDDIT_CLIENT_ID,
    client_secret = REDDIT_CLIENT_SECRET,
    user_agent = REDDIT_USER_AGENT
)
```

```
# Get submission details
submission = reddit.submission(id='<SUBMISSION ID>')
print(f' Title = {submission.title}')
print(f' Author = {submission.author.name}')
print(f' Created = {submission.created}')
print(f' Score = {submission.score}')

# Get comment details
comment = reddit.comment(id='<COMMENT ID>')
print(f' Created = {comment.created}')
print(f' Author = {comment.author.name}')
print(f' Body = {comment.body}')
print(f' Score = {comment.score}'')
```

Extract Data from Reddit - PushShift API & PSAW

- PushShift is a third-party archive of Reddit data
- Provides useful features that are not available via official Reddit API:
 - Find top-scoring submissions or comments in a specific date range
 - Search multiple subreddits in a single API call
 - Filter submissions and comments by a simple wild-card search (e.g. `$AAPL|AAPL|"apple inc"`)

<https://pushshift.io/api-parameters/>
<https://psaw.readthedocs.io>

```
# Initialize Pushshift API instance
api = PushshiftAPI()

# Search submissions that include the search query `q`
submissions = api.search_submissions(
    after = start_epoch,
    before = end_epoch,
    q = '$AAPL|AAPL|"apple inc"',
    subreddit = 'wallstreetbets,wallstreetbetsnew,Investing',
    sort_type = 'num_comments',
    sort = 'desc',
    filter = [
        'id', 'url', 'subreddit', 'created', 'author',
        'num_comments', 'num_crossposts', 'title', 'selftext'
    ],
    limit = 10
)

for submission in submissions:
    pass # Do something with each submission
```

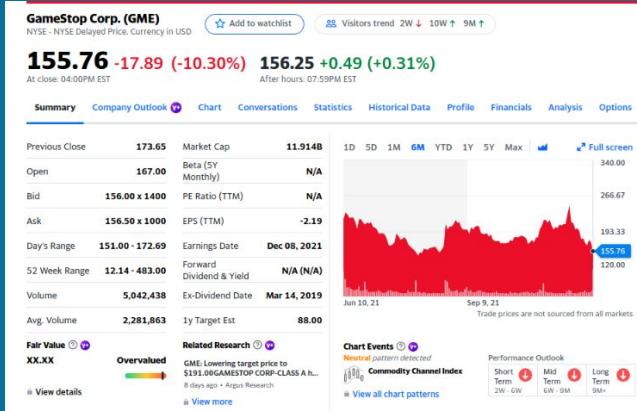
Extracting Data from Reddit - Kaggle

- Ref: <https://www.kaggle.com/pavellexyr/wallstreetbets-posts-and-comments-for-august-2021>
- Contains 2 csv files:
 - **wsb-aug-2021-posts** ['type', 'id', 'subreddit.id', 'subreddit.name', 'subreddit.nsfw', 'created_utc', 'permalink', 'domain', 'url', 'selftext', 'title', 'score']
shape - (25,751 rows, 12 columns)
 - **wsb-aug-2021-comments** ['type', 'id', 'subreddit.id', 'subreddit.name', 'subreddit.nsfw', 'created_utc', 'permalink', 'body', 'sentiment', 'score']
shape - (1,001,160 rows, 10 columns)

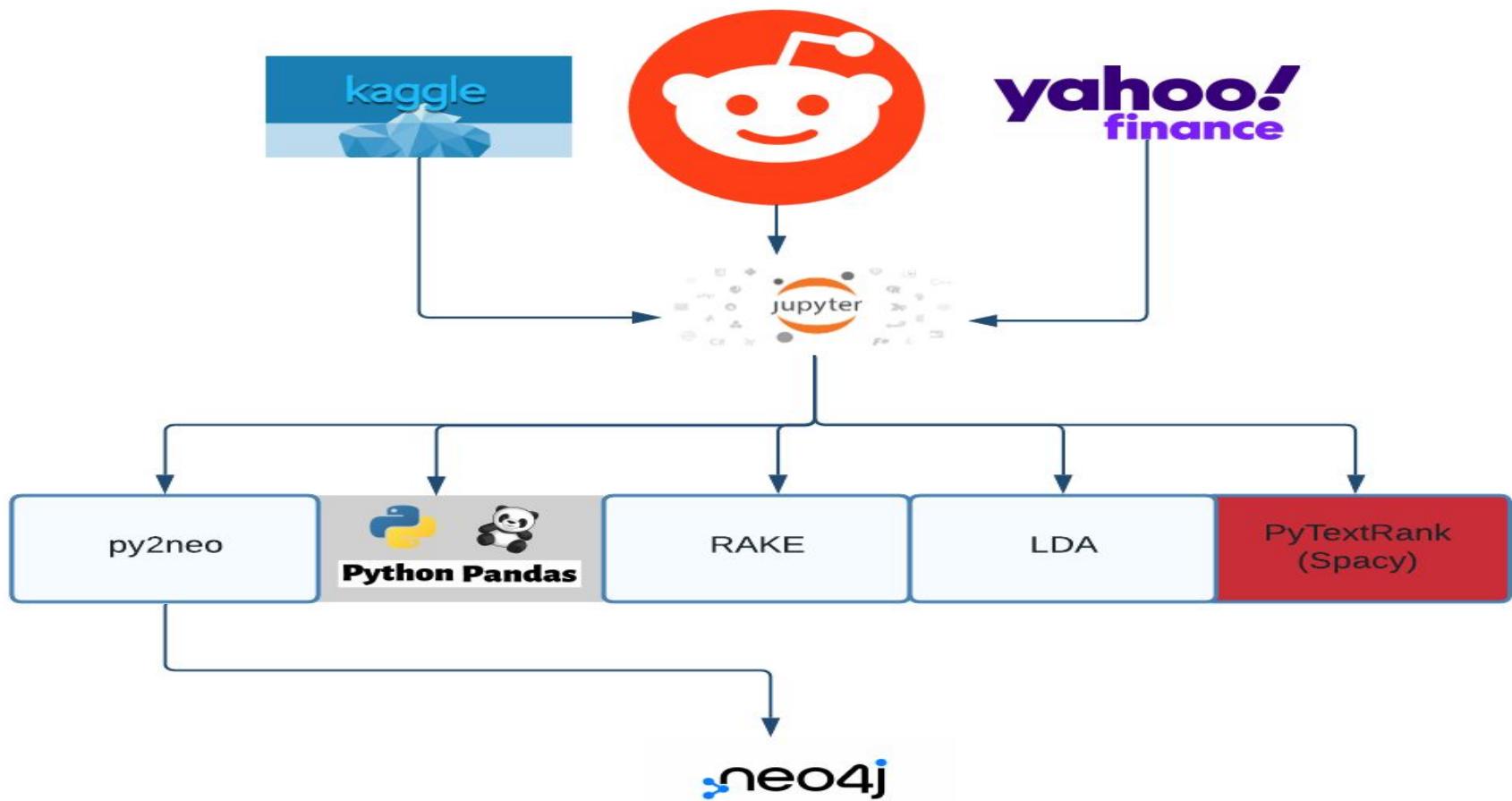
Yahoo! Finance

- yfinance is a python package that provides information about publicly traded companies
- This is used to replace an API that was decommissioned by Yahoo!
- Provides up to 151 fields of information - including opening, closing, and daily delta
- Used to match and validate legitimate tickers extracted via RAKE

<https://pypi.org/project/yfinance/>



THE ANALYSIS



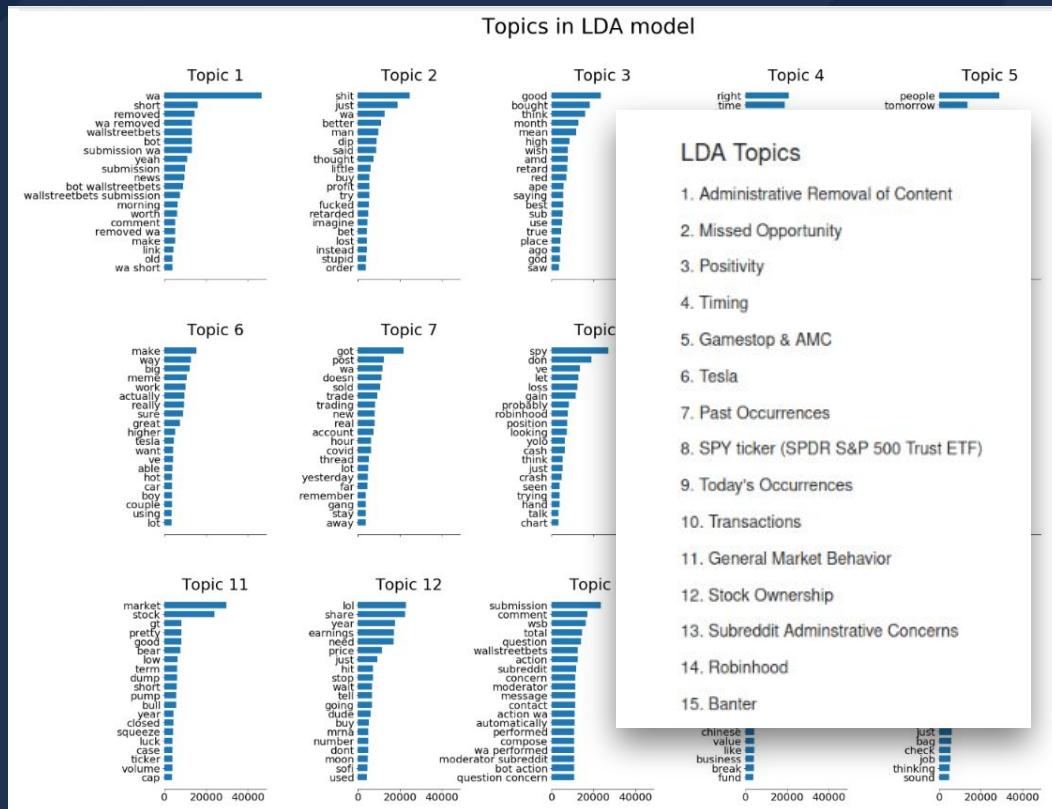
The #1 Database for Connected Data

Keywords extraction

- To keep the extraction output manageable, we extracted all keywords from posts and comments and then ran a regex to only look for keywords with \$ prefix.
- 1456 company keywords extracted through RAKE.
- 294 stock symbols remained after validating against yahoo finance.

```
['sdc', 'leg', 'adm', 'crm', 'grwg', 'd', 'tt', 'coty', 'cpop',  
'open', 'gdrx', 'abnb', 'amd', 'cvna', 'rgbp', 'pcg', 'fubo',  
'jblu', 'dq', 'nio', 'mj', 'u', 'rrgb', 'alb', 'wmt', 'crsr',  
'low', 'zim', 'psfe', 'dfen', 'roku', 'tup', 'panw', 'gme', 'dnut',  
'clov', 'cvs', 'paya', 'wix', 'dash', 'zm', 'spce', 'snndl', 'bark',  
'aso', 'unh', 'bud', 'ms', 'coin', 'base', 'mrna', 'lmt', 'ecvt',  
'spot', 'nkla', 'tal', 'lc', 'nvda', 'vxx', 'tr', 'corn', 'sony',  
'cscs', 's', 'astr', 'dkng', 'uvxy', 'save', 'mnst', 'x', 'eat',  
'dis', 'pypl', 'snps', 'sq', 'igt', 'crnt', 'carr', 'coke', 'rblx',  
'me', 'clf', 'hood', 'azn', 'grom', 'play', 'luv', 'bili', 'nflx',  
'fvrr', 'maxr', 'ejh', 'geni', 'mdlz', 'fsm', 'amc', 'khc', 'ater',  
'b', 'baba', 'cat', 'ge', 'jmia', 'jack', 'su', 'vzio', 'amzn',  
'xl', 'ed', 'azz', 'uwmc', 'webr', 'spy', 'tsla', 'qs', 'zbra',  
'afrm', 'pbya', 'pton', 'stla', 'manu', 'clne', 'root', 'gaymf',  
'gld', 'mmnd', 'asml', 'mu', 'et', 'atvi', 'gm', 'hyln', 'ftch',  
'pltr', 'viac', 'jd', 'car', 'tdoc', 'alto', 'you', 'hog', 'cvsi',  
'arw', 'm', 'ibkr', 'axsm', 'y', 'vym', 'hut', 'fas', 'spxu',  
'jepi', 'simp', 'gt', 'cmp', 'pins', 'wish', 'ha', 'crsp', 'goog',  
'ship', 'sklz', 'mo', 'xlf', 'bnrg', 'mrvl', 'rgr', 'aal', 'twlo',  
'lcid', 'cron', 'ba', 'mrk', 'dow', 'dude', 'job', 'rtx', 'fxlv',  
'cook', 'flws', 'posh', 'nvax', 'ftnt', 'snap', 'se', 'msft',  
'cyxt', 'tch', 'tgt', 'ndaq', 'rklb', 'hexo', 'sgmd',  
'edit', 'crwd', 'sens', 'wfc', 'muds', 'bntx', 'sbux', 'shop',  
'etsy', 'sli', 'ally', 'chwy', 'land', 'hd', 'ebet', 'yinn',  
'dooo', 'bb', 'abt', 'arkk', 'asts', 'dal', 'dole', 'sesn', 'dats',  
'goev', 'four', 'wkhs', 'gis', 'ride', 'fnma', 'body', 'dnmr',  
'bbig', 'riot', 'net', 'dia', 'mara', 'jnj', 'v', 'rkt', 'didi',  
'upst', 'sava', 'woof', 'avpt', 'sofi', 'wen', 'ptr', 'tm',  
'tlry', 'tell', 'wing', 'aapl', 'ida', 'zy', 'slqt', 'tcehy',  
'rdnt', 'f', 'ntnx', 'xpo', 'celh', 'wwe', 't', 'apps', 'bac',  
'nakd', 'ayx', 'nok', 'ely', 'indi', 'gold', 'pdd', 'amat', 'view',  
'antm', 'zg', 'doyu', 'ap', 'googl', 'el', 'ccl', 'aci', 'bidu',  
'ebay', 'pfe', 'j', 'negg', 'evgo', 'fb', 'vrns', 'tmf', 'mvst',  
'gotu', 'pg', 'path', 'penn', 'znga', 'mcfe'], dtype=object)
```

Analysis Details - LDA - Latent Dirichlet Allocation



- LDA is a dimensionality reduction algorithm
- It is used to generate synthetic topics from a corpus of documents and a dictionary of terms
- The end result is a set of topics that can be used to classify documents based on how well they matches the topics in the trained set

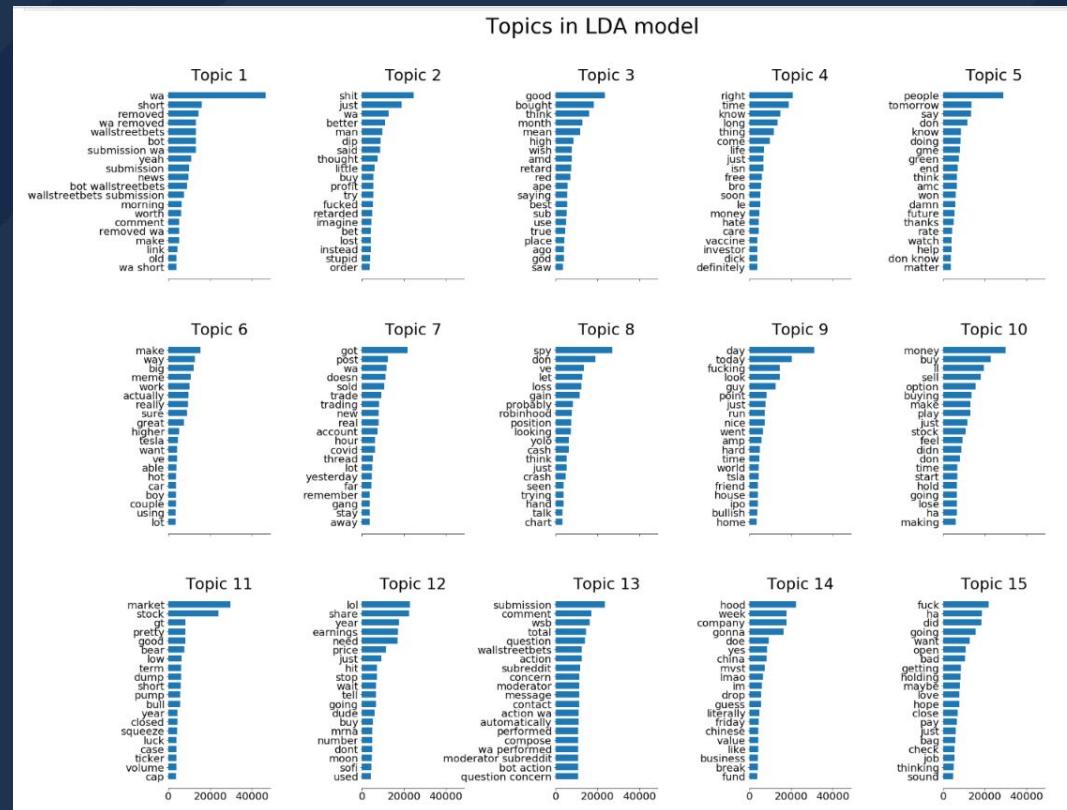
Analysis Details - LDA (cont.)

Text Pre-Processing

- Removed missing and deleted posts
 - Remove numbers
 - Tokenize
 - Lemmatization
 - Remove words occurring in >95% of documents
 - Remove stop words
 - Vectorize
 - Removed html code
 - Lower case
 - Limit parts-of-speech
 - Remove single letter words
 - Remove words that only occur in one document
 - Create bi-grams

Analysis Details - LDA - Latent Dirichlet Allocation

- Processed vector fed into Latent Dirichlet Allocation model
- 15 component LDA model
- 1000 features used
- Trained on 814k posts
- Trained model used to classify posts
- Packages used - natural language toolkit (`nltk`), sci kit learn (`sklearn`)



Analysis Details - LDA - Latent Dirichlet Allocation

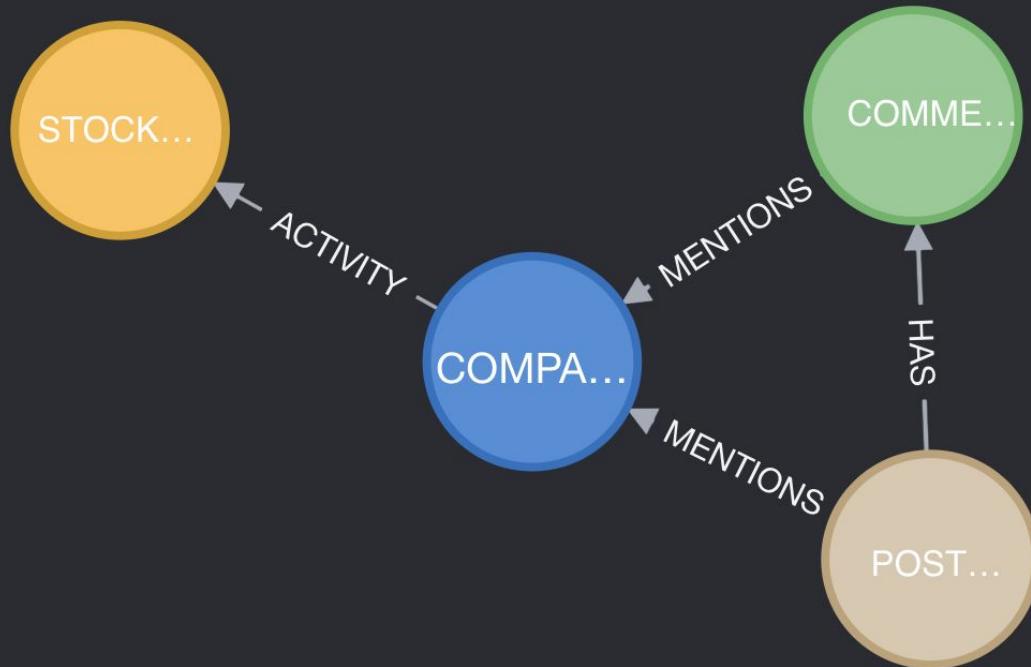
LDA Topics

1. Administrative Removal of Content
2. Missed Opportunity
3. Positivity
4. Timing
5. Gamestop & AMC
6. Tesla
7. Past Occurrences
8. SPY ticker (SPDR S&P 500 Trust ETF)
9. Today's Occurrences
10. Transactions
11. General Market Behavior
12. Stock Ownership
13. Subreddit Administrative Concerns
14. Robinhood
15. Banter

1	14.960851
2	6.792915
3	5.931684
4	5.483008
5	6.613543
6	3.787083
7	5.312113
8	5.854775
9	6.882847
10	8.630618
11	5.394674
12	7.046002
13	2.887152
14	7.880081
15	6.542654

"c.topic"	"num"
"Administrative Removal of Content"	19579
"Robinhood"	10621
"Transactions"	9568
"Missed Opportunity"	9202
"Todays Occurrences"	8749
"Gamestop & AMC"	8495
"Mixed Banter"	8089
"Stock Ownership"	7933
"Positivity"	7220
"Timing"	7153
"SPY ticker"	7070

KNOWLEDGE GRAPH



Overview



Node labels

* (4) POST (1)

STOCKPRICE (1) COMMENT (1)

COMPANY (1)

Relationship Types

* (4) HAS (1) MENTIONS (2)

ACTIVITY (1)

Displaying 4 nodes, 4 relationships.



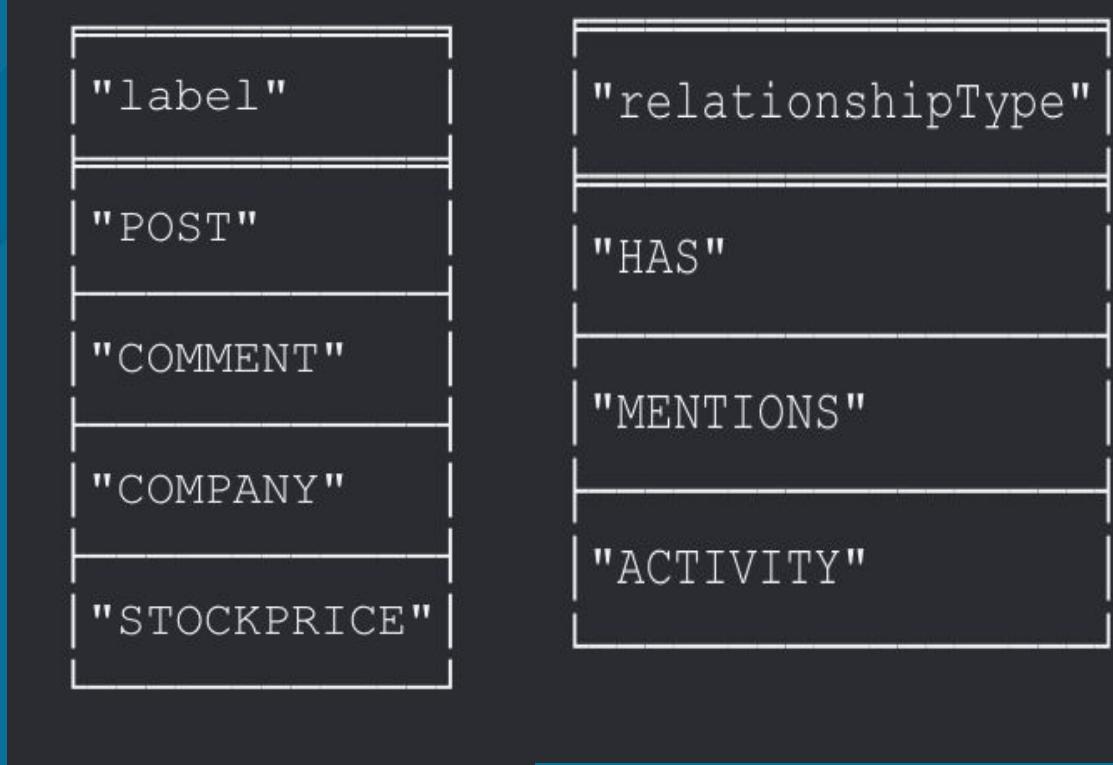
Knowledge Graph

"labels(n)"	"SampleSize"	"Avg_PropertyCount"	"Min_PropertyCount"	"Max_PropertyCount"	"Avg_RelationshipCount"	"Min_RelationshipCount"	"Max_RelationshipCount"
["COMMENT"]	12347	6.0	6	6	1.0091520207337765	1	4
["COMPANY"]	33	1.0	1	1	23.57575757575757	22	53
["STOCKPRICE"]	596	5.0	5	5	1.0	1	1
["POST"]	1	4.0	4	4	12606.0	12606	12606

- On average, each post has ~12k comments
- On average, each comment mentions 1 company.
- On average, each company has 23 relationships owing to number of open market days.
- No null property values across nodes.

Graph Stats

- Total Number of nodes - 128936
- Total Number of relationships - 129910



THE RESULTS

Most popular post

p.created	p.title	p.score
1 "24"	"Daily Discussion Thread for August 24, 2021"	"502"

Most popular comment

	p.c_id	p.body	p.score
1	"halme9b"	"Hit half mil net worth today, thank you JPOW and WSB."	"98"

Most commented post

	p.p_id	p.title	num
1	"p3sv76"	"Weekend Discussion Thread for the Weekend of August 13, 2021"	19259

Top 10 mentioned companies

	Name	num
1	"hood"	199
2	"spy"	78
3	"wish"	69
4	"gme"	56
5	"pfe"	41
6	"nvda"	34
7	"root"	32
8	"baba"	31
9	"clov"	23

Lifetime of posts

	p.title	daysdiff	num
1	"Daily Discussion Thread for August 19, 2021"	0	13272
2	"Daily Discussion Thread for August 24, 2021"	0	12550
3	"Daily Discussion Thread for August 04, 2021"	0	12277
4	"What Are Your Moves Tomorrow, August 19, 2021"	0	10341
5	"Weekend Discussion Thread for the Weekend of August 13, 2021"	1	9127
6	"Weekend Discussion Thread for the Weekend of August 06, 2021"	1	8506
7	"Weekend Discussion Thread for the Weekend of August 20, 2021"	1	8000
8	"Weekend Discussion Thread for the Weekend of August 27, 2021"	1	7199
9	"Weekend Discussion Thread for the Weekend of August 06, 2021"	0	6802

Lifetime of posts

	po.title	numComments	Company	CreatedPost	StockDate	val
1	"Weekend Discussion Thread for the Weekend of August 27, 2021"	5	"bbig"	"27"	"27"	50.14164967962419
2	"Weekend Discussion Thread for the Weekend of August 27, 2021"	1	"ater"	"27"	"27"	47.14423931778589
3	"Daily Discussion Thread for August 04, 2021"	164	"hood"	"04"	"04"	29.27456088771682
4	"Daily Discussion Thread for August 24, 2021"	44	"gme"	"24"	"24"	26.459800214369984
5	"Daily Discussion Thread for August 24, 2021"	11	"amc"	"24"	"24"	19.01048657137274
6	"Daily Discussion Thread for August 24, 2021"	1	"bbig"	"24"	"24"	17.60000228881836
7	"Daily Discussion Thread for August 24, 2021"	3	"gotu"	"24"	"24"	16.24999379118308
8	"Daily Discussion Thread for August 24, 2021"	2	"tal"	"24"	"24"	14.98973380344713
9	"Daily Discussion Thread for August 19, 2021"	3	"dnut"	"19"	"19"	14.692306518554691

DEMO

Takeaways

- Social activity goes up when there is a big change in stock price.
- Whether the volume of social activity influences stock price is inconclusive.
- Engagement with posts is only active for $t + 1$ days. Moderators need to frequently create new posts to get traffic.
- Blue chip companies are frequently mentioned alongside meme stocks.

Lessons Learned & Future Work

LESSONS LEARNED

- Power of keywords extraction.
- Power of nltk for text preprocessing.
- py2neo over neo4j driver.
- Graph tuning.
- Node and relationship creation in batch.
- Remove bot activity.
- There are different “degrees” to unstructured data!
- Everything about GRAPH!!

POTENTIAL FUTURE WORK

- Switching to a dataset (e.g. from /r/Investing_Advice) that contains more grammatically well-formed unstructured data may allow us to refine our topic modeling and keyword extraction efforts.
- Focusing down on “Daily Discussion Threads” may help to better understand time-specific trends in topics and sentiment.

QUESTIONS?

Amit Sachan
Mustafa Burny
Sam Courtney