



Case Study 1

By: Marcela Vasconcellos, Stokley Voltmer, Sirshendu Ganguly,



Search Twitter



For you

Trending

COVID-19

News

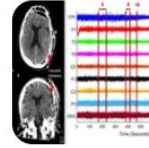
Sports

Entertainment



Daily Mail Online · Last night

Brain activity of dying man suggests our lives really do flash before our eyes as we die



Bloomberg Wealth · 3 hours ago

Million-dollar home listings dry up for wealthy suburban buyers



World news · LIVE

Ukraine: Zelenskyy calls up reservists and urges citizens to leave Russia 'immediately'



HuffPost Life · February 10, 2022

Signs Of A Toxic Job You Can Spot During Your Interview



News on twitter:

- Fast-moving information
- Demand-oriented content
- Marketplace economy:
high competition

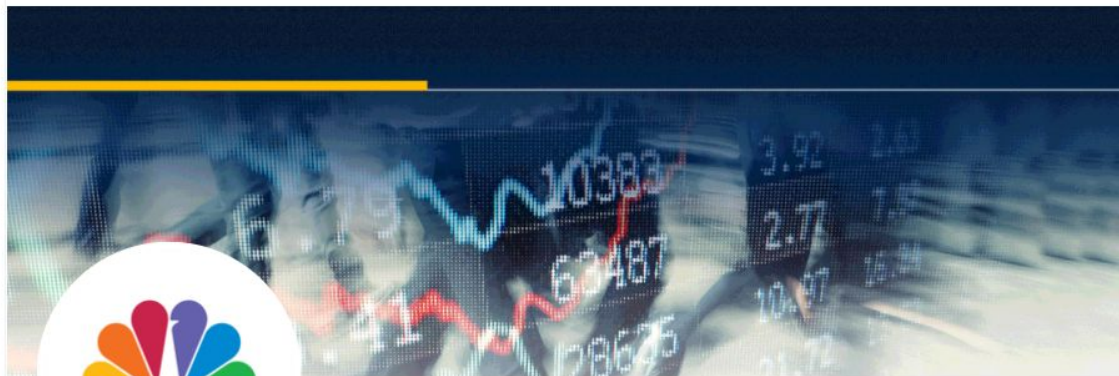
↑ Users

↑ Subscribers and Advertisers

↑ Profits



CNBC ✓
[481.2K Tweets](#)



Follow

CNBC ✓
@CNBC

First in business worldwide.

📍 Englewood Cliffs, NJ 🌐 [cnbc.com](#) 📅 Joined February 2009

851 Following [4.5M Followers](#)

Not followed by anyone you're following



Search Twitter



You might like



Bloomberg ✓
@business



Financial Times ✓
@FT



Reuters Business ✓
@ReutersBiz

Business question

- What types of news articles should we produce in order to engage users and increase profits?

Using Twitter

- Highly engaged customers: we know which articles were most popular
- Historic data: we can use Machine Learning to find patterns
- Objective: use unsupervised learning to group news posts by similarity and investigate which groups are most popular

Data Source

Attributes	Characteristics
text	The actual text from the tweets.
retweet_count	The number of times the tweets were retweeted

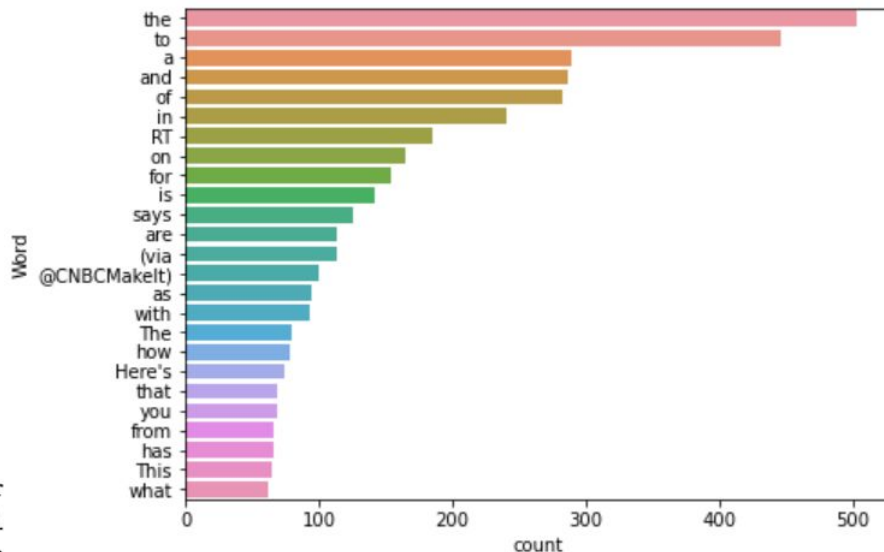
Methodology

- Exploratory Data Analysis
- Preprocessing
- TF-IDF indexing
- Clustering
- Analyzing Popularity

EDA

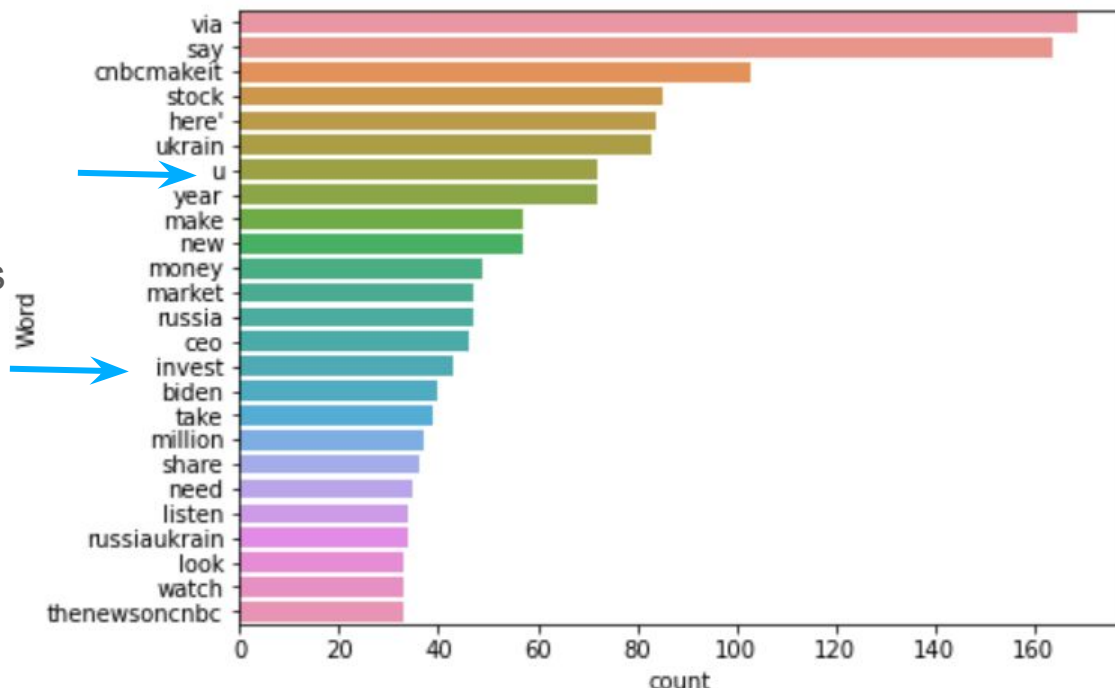
- Visualizing tweets
- 1000 tweets by @CNBC
- Word frequencies

```
0 phuck elonmusk this is a nice distraction for
1 rt marketrebels elon musk in email to cnbc bi
2 trump says he and capitol rioters wanted same
3 rt schwartzbcnbc exclusive elonmusk speaks to ...
4 rt danpriceseattle big mac prices are up in th...
5 rt jasuja how much money popular startups are ...
6 rt repmaloney i m requesting that the gsa cons...
7 rt fisa end ccp from cnbc international tv htt...
8 cnbc frankcnbc cnbc should start thinking abou...
9 aeharshada cnbc awaaz can nifty surpass
10 rt bruceduck jrubinblogger sanctions on russia...
11 did i miss something bond yieldspread cnbc htt...
12 rt sawyermerritt news musk says he would do th...
13 rt sawyermerritt news musk says he would do th...
14 rollingstone press traitortrump unforgivable a...
```



Preprocessing

- All lowercase
- Remove symbols and stopwords
- Lemmatizing
- Stemming
- New word count
- Example:



['RT', '@CNBCTechCheck:', 'Despite', 'posting', 'Q4', 'results', 'that', 'mostly', 'beat', 'the', 'Street,', 'shares', 'of', 'cloud-based', 'platform', 'm', '@mondaydotcom', 'still', 'plunged', 'to', 'ne...'],



['cnbctechcheck despit post q4 result mostli beat street share cloudbas platform mondaydotcom still plung ne ',

TF-IDF

- Tweets x Vocabulary matrix
- Term Frequency
- Inverse Document Frequency

$$tf(t, d) \cdot idf(t, D)$$

	Word1	Word2	...
Tweet 1		(Tfidf value)	
Tweet 2			

Number of features (words): 212

Sample indexed tweet:

(0, 168) 0.6355892105048007

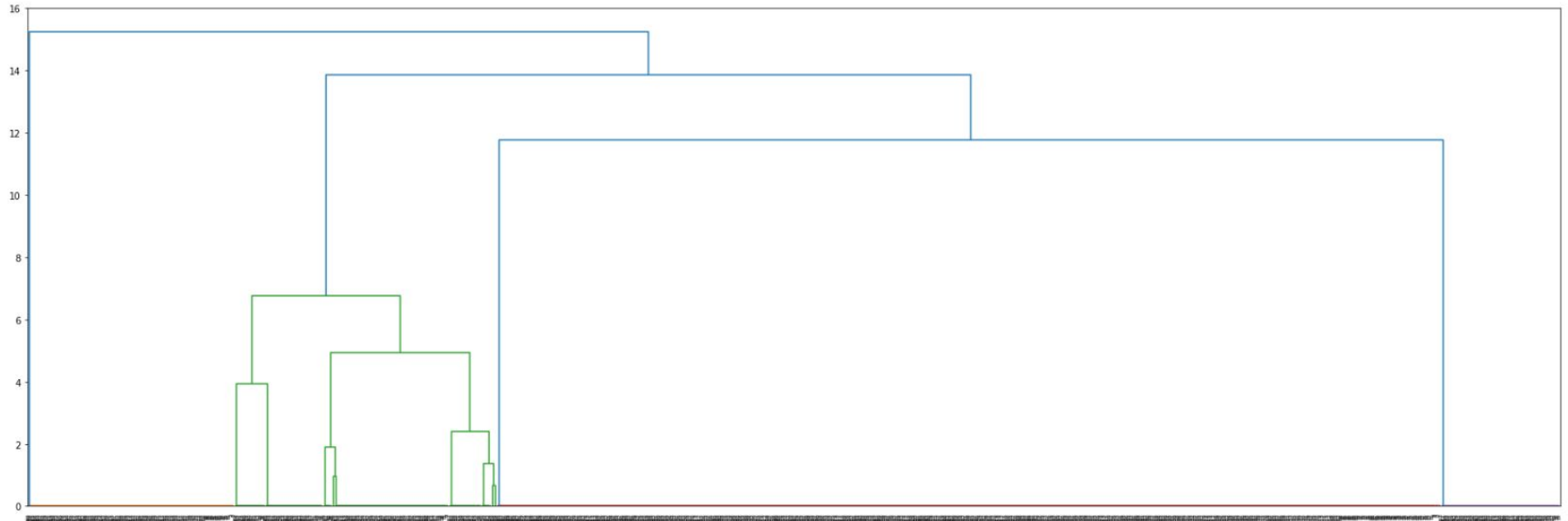
(0, 162) 0.5111822990026225

(0, 31) 0.57854905814138

Features in sample tweet: [array(['share', 'street', 'cnbctechcheck'], dtype='<U14')]

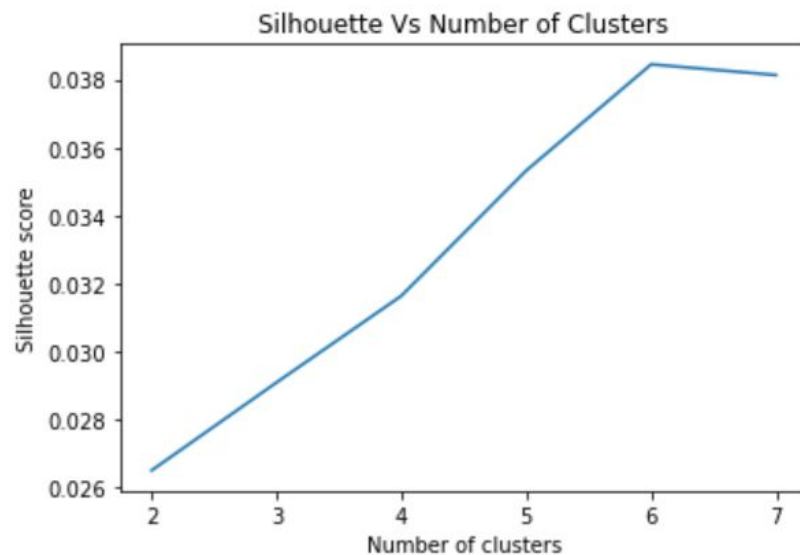
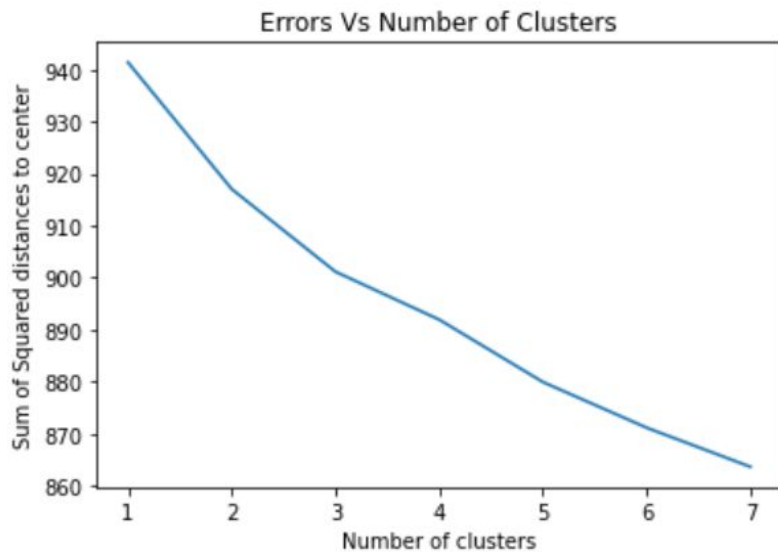
Clustering

- Dendrogram



Clustering

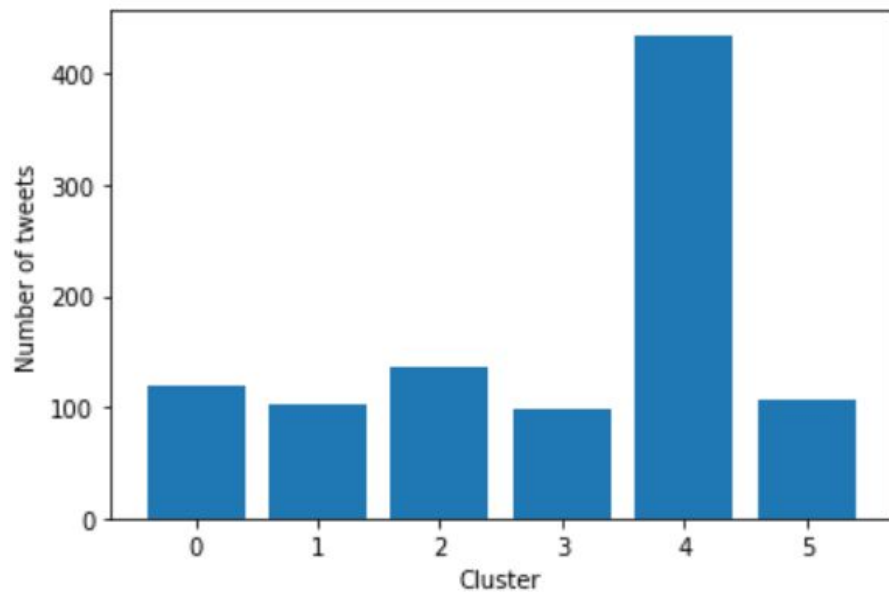
- K-means
 - Tested 1 to 7 clusters



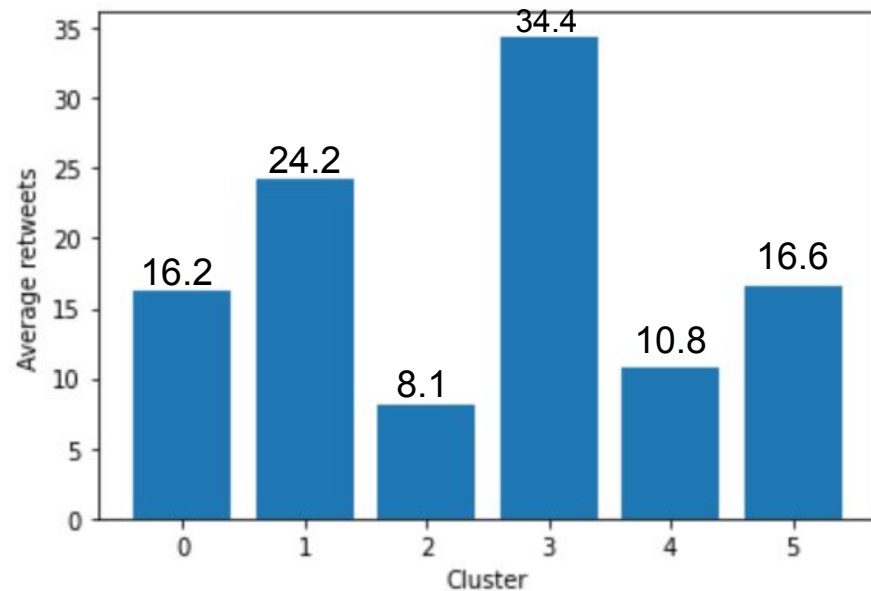
Top Words Cluster 0	Count	Top Words Cluster 1	Count	Top Words Cluster 2	Count
via	119	stock	82	here	58
cnbcmakeit	103	market	26	invest	36
here	22	say	23	million	28
year	21	russiaukrain	20	year	25
money	19	tension	19	much	21
make	13	week	9	money	19
home	12	fed	9	pay	17
top	12	listen	8	spend	17
21	9	ukrain	8	save	16
biggest	9	biggest	8	need	16
Top Words Cluster 3	Count	Top Words Cluster 4	Count	Top Words Cluster 5	Count
ukrain	73	new	33	say	103
russia	51	ceo	26	via	22
putin	34	tesla	26	expert	15
sanction	26	watch	24	harvard	13
biden	26	busi	21	make	13
say	23	via	21	avoid	12
russian	20	here	20	ceo	8
presid	20	report	20	covid	7
troop	15	take	18	food	7
invas	14	year	18	worst	7

Clustering

How many tweets are in each of the clusters



How many average retweets per cluster



Conclusion

- Our analysis will be helpful to the organization in categorizing the tweets before uploading so that they can tag their tweets with suitable keywords i.e. politics, stock market, etc.
- Thus allowing the readers of the tweet to easily scroll through the topics of their interest.
- This in return will prevent CNBC from losing their Twitter followers.