

Ryan Downing, Stefan Obradovic, Sahil Goel, Candace Austin, Rohan Samy, Kafana Ouattara

Dataset Description

We collected tweets from individuals on Twitter regarding three separate industries in the economy by grabbing tweets with explicit ticker mentions

- Twitter posts made between March 01, 2020 and March 12, 2022
- At most 10 tweets per analyst

Technology:

 AAPL, MSFT, GOOG, AMZN, TSLA, FB, NVDA, AVGO, CSCO, ADBE, CRM, ORCL, ITC, AMD, NFLX, TXN, INTU, PYPL, NOW, IBM, QQQ, TQQQ, ARKK

Financials:

V, MA, AXP, PYPL, JPM, BAC, WFC, C, COF, PNC, USB, BRK-B, BRK-A, SCHW, MS, GS, MET, PRU, BK, TROW

Energy:

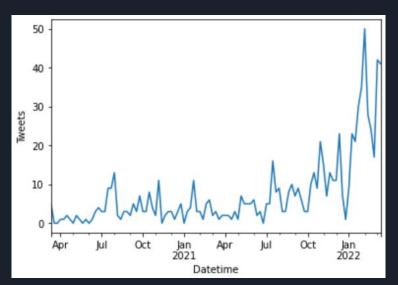
XOM, CVX, SHEL, PTR, TTE, COP, EQNR, BP, PBR, ENB, CNQ, SNP, EOG, SLB, PXD, EPD, OXY, TRP, E, MPC, SU, KMI, VLO, WMB, ICLN, SPWR, ENPH, CWEN, SHFT, AMRC, XLE

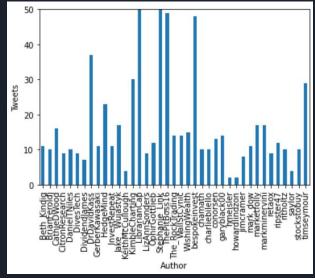


Collecting the Dataset

Used Python's tweepy Twitter API

- Collected list of twitter usernames from active users in the financial analyst community.
- Obtained all tweets which contained relevant cashtags (i.e. \$V, \$MA).
- Subsampled twitter users who had many tweets for later analysis to remove bias.





Punctuation, Numbers, and Syntax

Cleaning text before applying analysis

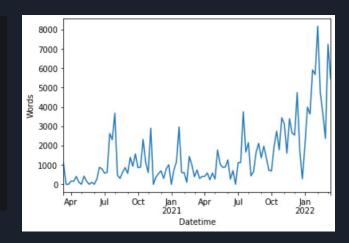
- Removed all characters which are not alphabetical or whitespace (including emojis and cashtags).
- Use NLTK's porter stemmer to get the roots of all words.
- Split on whitespace to get tokens, remove stopwords, and recreate sentence.

```
stemmer = PorterStemmer()
sw = set(stopwords.words("english"))

def stem_and_remove_sw(tokens):
    return [stemmer.stem(token) for token in tokens if stemmer.stem(token) not in sw]

data["cleaned_text"] = data["Text"].str.replace("[^A-Za-z\s]", "", regex=True).str.lower()

data['tokens'] = data["cleaned_text"].str.split(',')
data['tokens'] = data['tokens'].apply(stem_and_remove_sw)
data["tokenized_sentence"] = data["tokens"].apply(lambda tokens: ' '.join(tokens))
```

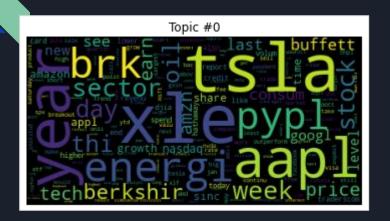


Gensim

```
import pandas as pd
import re
from nltk.stem.porter import *
from nltk.corpus import stopwords
from pandas.tseries import offsets
from wordcloud import WordCloud
from gensim.corpora.dictionary import Dictionary
from gensim.models.ldamodel import LdaModel
from pathlib import Path
import matplotlib.pyplot as plt
%matplotlib inline
```

	Author	Text	Date	Stock(s) Mentioned	Tweet Link	Misc.	Topic
0	Stephanie_Link	Didn't do the wage increase situation justice	2021-05-24 19:34:08+00:00	FB	https://twitter.com/twitter/statuses/139691241	NaN	Technology
1	bespokeinvest	Apple \$AAPL now up 14% over the last month.	2021-07-14 14:35:44+00:00	AAPL	https://twitter.com/twitter/statuses/141531910	NaN	Technology
2	HedgeMind	\$AMD Unstoppable momentum, up 11%+ after repor	2021-11-08 19:49:26+00:00	AMD	https://twitter.com/twitter/statuses/145779742	NaN	Technology
3	jimcramer	critical reversal day and then \$AMD announces	2022-02-24 21:22:06+00:00	AMD	https://twitter.com/twitter/statuses/149695863	NaN	Technology
4	Beth_Kindig	Today we recap this week's biggest tech earnin	2021-05-21 14:25:35+00:00	CSCO	https://twitter.com/twitter/statuses/139574760	NaN	Technology

Word Cloud



Topic 0: Energy Stocks

Topic 1: Tech Stocks

Topic 2: Financial Stocks





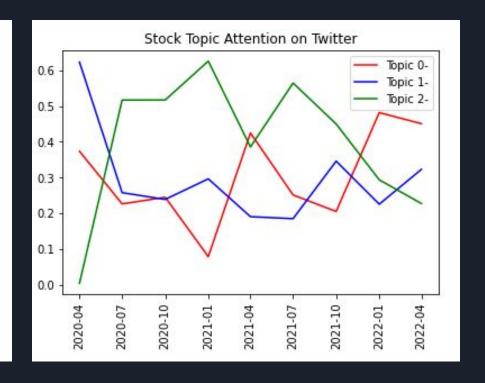
Document Topic Matrix

2	1	0	
(2, 0.0039743935)	(1, 0.62279147)	(0, 0.37323412)	0
(2, 0.5167407)	(1, 0.25737366)	(0, 0.22588567)	1
(2, 0.5168945)	(1, 0.23853293)	(0, 0.24457252)	2
(2, 0.6256254)	(1, 0.2960779)	(0, 0.07829677)	3
(2, 0.38526794)	(1, 0.19016759)	(0, 0.42456445)	4
(2, 0.5644213)	(1, 0.1846899)	(0, 0.25088882)	5
(2, 0.44960126)	(1, 0.3456513)	(0, 0.20474742)	6
(2, 0.29333073)	(1, 0.22509378)	(0, 0.48157552)	7

(2, 0.22691186)

(0, 0.45047432) (1, 0.32261384)

Topics Over Time



Topic Model with K>3

Topic Model was re-ran with K=15



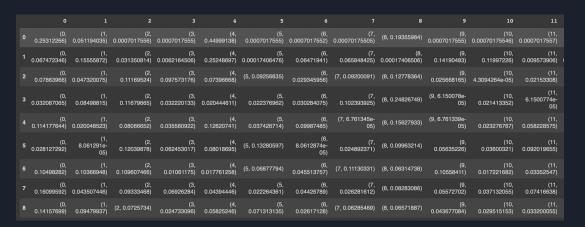


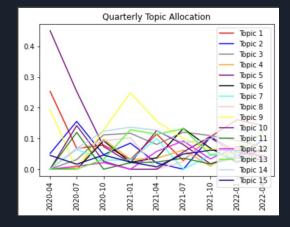
Noticeable Changes:

- 1. Topics are more defined and easier to interpret
- 2. Topics include subcategories of the original three topics









Topic Model Prediction



<u>Predictions of Which Topic Should have the Highest Weight</u>

Tweet 1: Technology
Tweet 2: Finance
Tweet 3: Finance
Tweet 4: Energy

Topic Model Weight Predictions

Tweet 5: Finance

	Unnamed:	User	Tweet	clean_tokens	0		2
0		ARKInvest	ARK's updated open-source Tesla model yields a	[ark, updat, open, sourc, tesla, model, yield,	0.025014	0.612669	0.362318
		xxxx_chen	\$PYPL bears getting louder \n\nhttps://t.co/0u	[pypl, bear, get, louder, uhxdvcdau, thi, doe,	0.016240	0.071141	0.912618
2		bishnuvardhan	\$PYPL waking up and outperforming \$QQQ. Poten	[pypl, wake, outperform, qqq, potenti, someth,	0.413838	0.266361	0.319800
3		clpirtle25	\$XOM NEW ARTICLE : Exxon sees carbon capture m	[xom, new, articl, exxon, see, carbon, captur,	0.090699	0.635759	0.273543
4		MYFORTPIERCE	I had someone tell me they are convinced that	[someon, tell, convinc, run, bank, point, info	0.281458	0.318613	0.399929

Topic Model's Highest Weighted Topic

Tweet 1: Technology
Tweet 2: Finance

Tweet 3: Energy

Tweet 4: Energy
Tweet 5: Finance

Paper Summary - Background

 Climate change presents risks to investors that cannot be easily hedged with futures or insurance contracts due to the widespread implications of climate disasters

 The authors propose an approach to hedge the risk of climate change in the long term by utilizing text mining techniques on climate change news



Paper Summary - High Level Strategy Overview

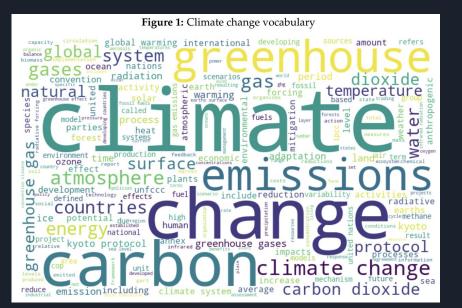
• Main idea is to construct a portfolio of publicly traded U.S. equities that hedges short-term climate change news in order to be protected against the dynamic nature of climate change in the long run

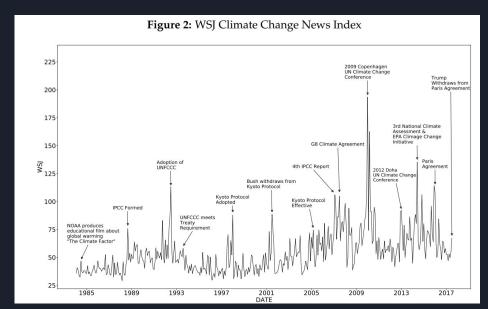
- Strategy follows a dynamic hedging approach created by Black, Scholes and Merton
- Two indices are created which measure the level and sentiment of climate related news based on the Wall Street Journal and data analytics vendor Crimson Hexagon (compilation of >1000 news outlets)
- From these indices, the portfolio is constructed by detecting stocks with returns that are negatively correlated with the level of negative news pertaining to climate

• This portfolio is frequently updated to account for recent news, thus protecting the investor overtime

Paper Summary - Climate Change Risk Indexes

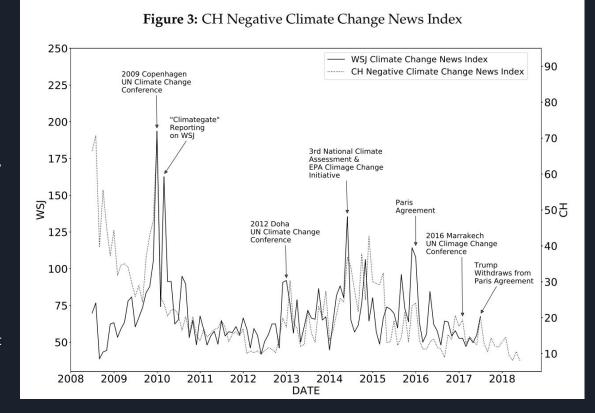
- Climate Change Vocabulary (CCV) Corpus
 - 19 Climate change White-papers from government agencies and conferences
 - 55 climate change word glossaries
 - Constructed by counting the frequency of stemmed words in the corpus (tf-idf matrix calculated using corpus as a single document)
- WSJ Index
 - Constructed by computing cosine similarity between tf-idf of CCV corpus and Wall Street Journal Articles
 - No separation between positive and negative sentiment regarding climate change news





Paper Summary - Climate Change Risk Indexes

- Crimson Hexagon (CH) Negative
 Climate Change News Index
 - Constructed using proprietary tool that calculates sentiment regarding a particular topic across many news outlets
 - Computed magnitude of negative sentiment regarding "climate change"
- Hedging Portfolios
 - Long companies that perform well when climate change news is up
 - Short companies that perform poorly when climate change news is up
 - Constructed using correlation between company ESG metrics
 - Reduce exposure to other risks to allow for climate change hedging that does not affect other market risks.



Paper Summary - Conclusions

- While the authors argue that their hedging approach is successful, they recognize that their strategy isn't definitively the "best" and that their findings should be used as a starting point for further research
- More research should be done to account for the fact that there are different types of climate risk. For example, this methodology does not distinguish between news pertaining to physical climate risk and regulatory climate risk
- Additionally, more types of assets could be considered in the hedging portfolio such as ETFs and international equities

Thank You!!