# Tweet Topics and Sentiments Analysis

**Analyze tweet topics and sentiment using Unsupervised machine learning**

Cindy Su 11.12.2021

# Motivation

Twitter with 192 million daily active users worldwide.

It contains rich source of data to analyze and understand social behavior.

# Goal

- How do people think about the the overall economic environment during the pandemic?

- Whether the sentiment in social media relates to the stock market?

NLP PIPELINE

1. TEXT INFORMATION
Twitter

2. SEGMENTATION AND TOKENIZATION
Spacy

3. TEXT CLEANING
RegExr

7. INTERPRETATION OF THE RESULT
pyLDAvis, WordCloud, textacy

6. MACHINE LEARNING ALGORITHMS
NMF, TruncatedSVD **LatentDirichletAllocation**

5. TEXT LEMMATIZATION AND STEAMING
Spacy Lemm SnowballStemmer

4. VECTORIZATION AND FEATURE ENGINEERING
TfidfVectorizer CountVectorizer

4

# Exploratory Data Analysis



- 2020.01-2021.10 : 69,060 data
- LikeCount > 2
- Keyword/hashtag: economics, business, financ

# Latent Dirichlet Allocation – Topic 11



- Politics/ coronavirus/ country'

# Latent Dirichlet Allocation – Topic 7



○ Education /research /school /study

# Latent Dirichlet Allocation – Topic 5



○ People /Talk /mind
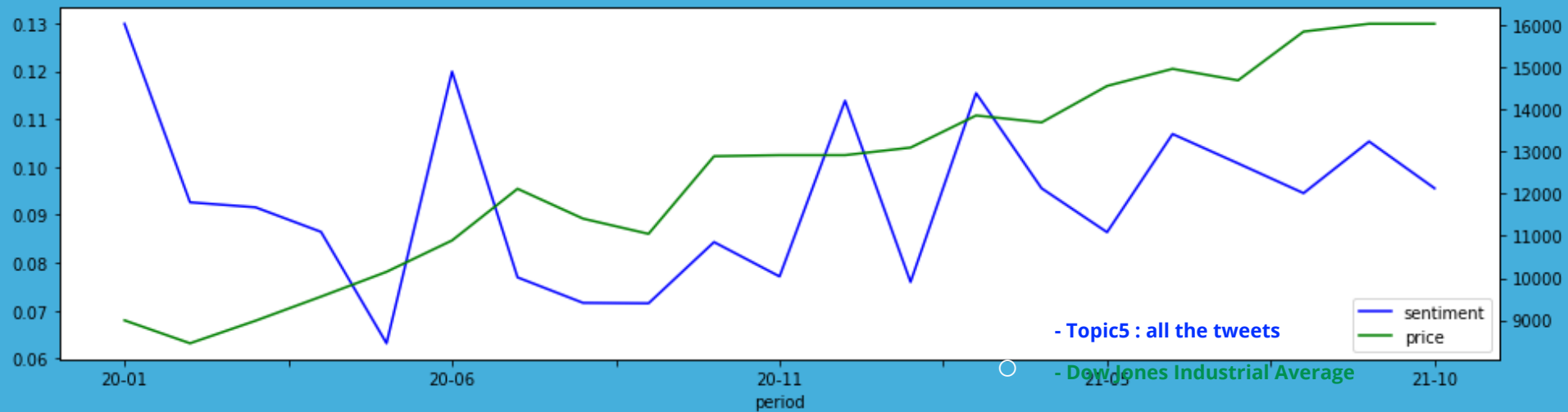
# Latent Dirichlet Allocation

- **Categories:**
- Education/ research/ school/ study: 4, 7, 10
- Wealth/ work: 2, 9
- People/talk/mind: 1, 5
- Politics/coronavirus/country : 11
- Market, tax, government: 3

# Sentiment Analysis



- Topic5 : People/ Talk/ mind

- Topic11: Politics/ coronavirus/ country'

# Sentiment Analysis



- Topic5 : all the tweets
- Dow Jones Industrial Average

11

# Future work

- Explore more on top tweets in each topics
- May need more data/months to observe better result (300-500/months)

# THANKS!

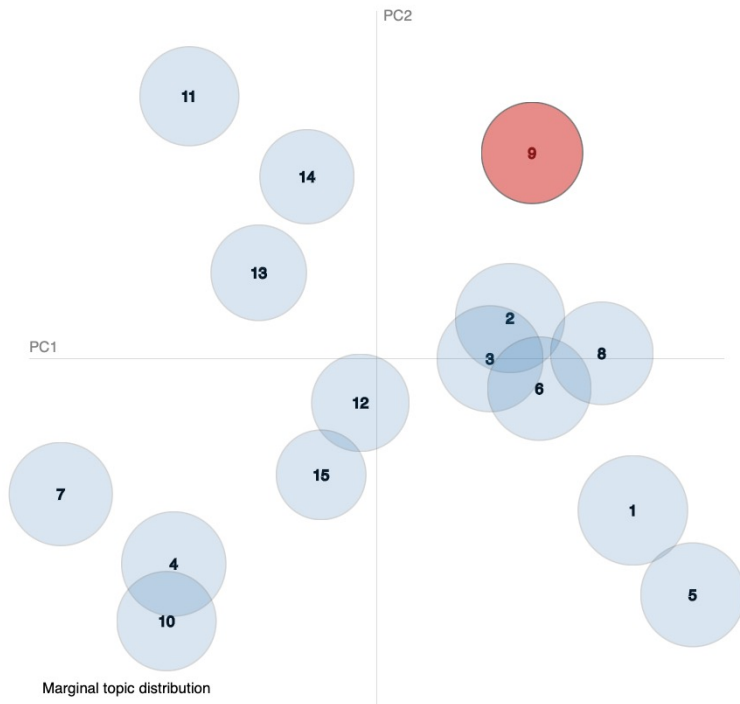Any questions?

# Latent Dirichlet Allocation – Topic 9



○ Wealth/ Tax/ Trickle-down economics