

# A scalable tool for longitudinal Twitter analysis: understanding the impact of COVID-19 on public discourse

Walter Dempsey<sup>25</sup>   Alfred Hero<sup>134</sup>   Conrad Hougen<sup>3</sup>   Brandon Oselio<sup>2</sup>  
Wayne Wang<sup>4</sup>

<sup>1</sup>Department of Biomedical Engineering

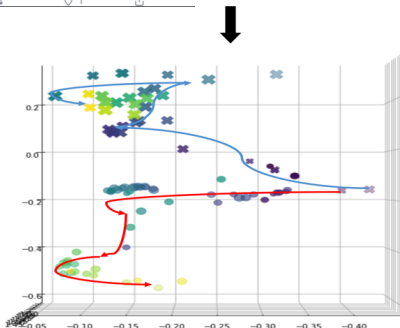
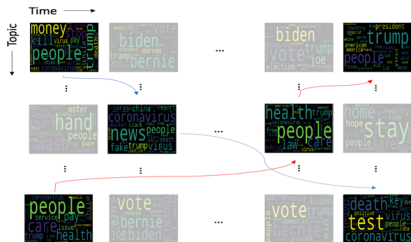
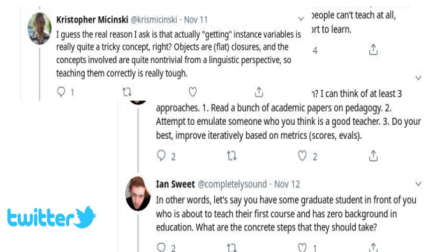
<sup>2</sup>Department of Biostatistics

<sup>3</sup>Department of EECS

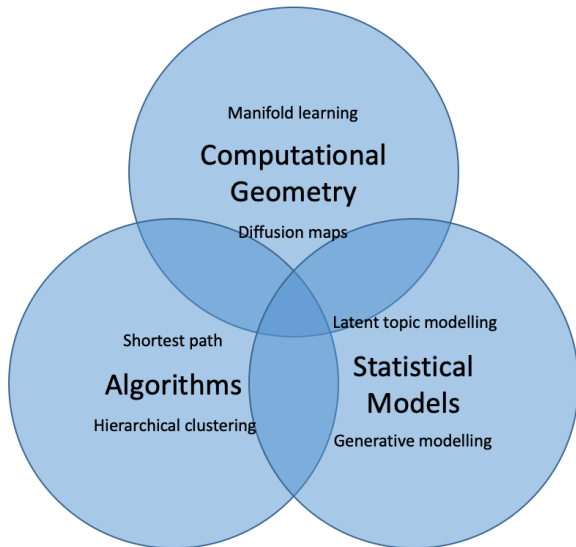
<sup>4</sup>Department of Statistics

<sup>5</sup>Institute of Social Research

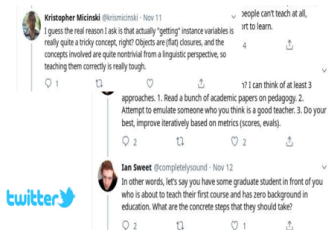
# Overview of results



# Data science tools



# Roadmap

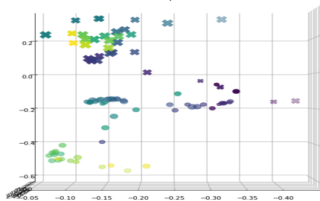


1. LDA on temporally smoothed corpus

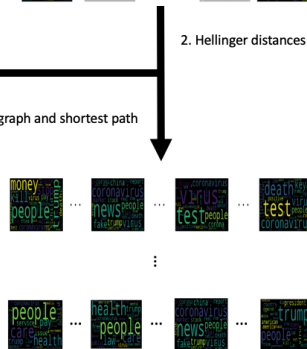


2. Hellinger distances

3a. PHATE embedding



3b. Neighborhood graph and shortest path



4. Visualize and interpret multiple paths and clusters



# Outline

- 1 Twitter Decahose Stream
- 2 Introduction to Probabilistic Topic Models
- 3 Connecting Topics Time by Time
- 4 Geometric Embedding of Structured High-dimensional Objects
- 5 Visualization and Interpretation of COVID-19 Discussions on Twitter

# Twitter decahose stream <sup>1</sup>

## Decahose stream

The Decahose delivers a  $\sim 10\%$  **random sample** of the realtime Twitter Firehose (300 to 500 million tweets per day) through a streaming connection. This is accomplished via a realtime sampling algorithm which randomly selects the data, while still allowing for the expected low-latency delivery of data as it is sent through the firehose by Twitter. One of the features available with Decahose: enhanced reliability - **geographic diversity** ( $\sim 0.1\% - 0.5\%$  of sampled tweets contain geo-location info) of backend systems<sup>a</sup>.

---

<sup>a</sup>Decahose stream. <https://developer.twitter.com/en/docs/tweets/sample-realtime/overview/decahose>. Accessed: 2020-04-20.

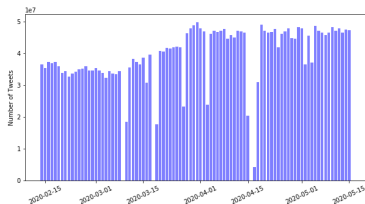
---

<sup>1</sup>Made available to us by MIDAS.

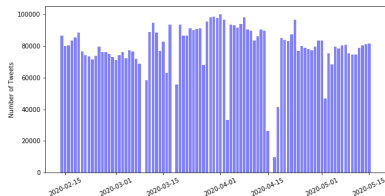
# Subsampled tweets from Feb 15 to May 15

Three levels of subsampling:

- Geotagged US tweets: to study spatial variations of tweeting behavior.
  - English tweets: to avoid confusion to topic models and to allow easy interpretation.
    - Non-retweets: to study the influence of the original tweets.



(a) Raw decahose tweets volume from Feb 15 to May 15.



(b) Geotagged US, Non-retweet, English decahose tweets volume from Feb 15 to May 15.

## Limitations of subsampled tweets

The data is noisy and subject to selection bias, and for it to be useful for general purposes one needs:

- People in the US who tweets about their opinions.
  - People in the US who tweets about their opinions + have Twitter location service on (so the data is geotagged).
- People who are truthful about their tweets.
  - People who are truthful about their tweets and use proper languages.
- ...

But we are *NOT* doing predictions, *NOR* making strong inferential decisions based on the data.

Overall, these issues do not affect our study of "public discourse" on COVID-19, and we are interested in understanding this subpopulation.



# Outline

- 1 Twitter Decahose Stream
- 2 Introduction to Probabilistic Topic Models**
- 3 Connecting Topics Time by Time
- 4 Geometric Embedding of Structured High-dimensional Objects
- 5 Visualization and Interpretation of COVID-19 Discussions on Twitter

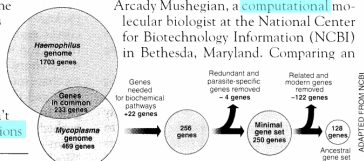
# Brief introduction to topic modelling

## Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK— How many genes does an organism need to survive? Last week at the genome meeting here,<sup>8</sup> two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



\* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

SCIENCE • VOL. 272 • 24 MAY 1996

Documents exhibit multiple topics<sup>2</sup>.

<sup>2</sup>David M Blei. "Probabilistic topic models". In: *Communications of the ACM* 55.4 (2012), pp. 77–84.

# Latent Dirichlet Allocation<sup>3</sup>

## Topics

gene 0.04  
dna 0.02  
genetic 0.01  
...

life 0.02  
evolve 0.01  
organism 0.01  
...

brain 0.04  
neuron 0.02  
nerve 0.01  
...

data 0.02  
number 0.02  
computer 0.01  
...

## Documents

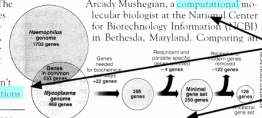
### Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many **genes** does an **organism** need to **survive**? Last week at the genome meeting here,<sup>3</sup> two genome researchers with radically different approaches presented complementary views of the basic genes needed for **life**. One research team, using **computer** analyses to compare known **genomes**, concluded that today's **organisms** can be sustained with just 250 genes, and that the earliest life forms required a mere 128 **genes**. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough. Although the numbers don't match precisely, those **predictions**

<sup>3</sup> Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

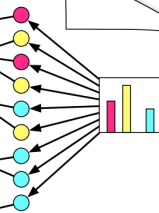
SCIENCE • VOL. 272 • 24 MAY 1996

"are not all that far apart," especially in comparison to the 75,000 **genes** in the human genome, notes Siv Andersson, a geneticist at Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a **genetic numbers** game, particularly if more and more **genomes** are completely cataloged and sequenced. "It may be a way of organizing any newly **sequenced genome**," explains Arady Muskhogian, a **computational** molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



**Stripping down.** Computer analysis yields an estimate of the minimum modern and ancient genomes.

## Topic proportions and assignments



<sup>3</sup>David M Blei, Andrew Y Ng, and Michael I Jordan. "Latent dirichlet allocation". In: *Journal of machine Learning research* 3.Jan (2003), pp. 993–1022.

# LDA and Twitter LDA<sup>4</sup> as generative processes

## LDA

For each topic  $t$

- $\phi^t \sim \text{Dir}(\beta)$

For each doc  $d$

- $\theta^d \sim \text{Dir}(\alpha)$
- For each word  $i$  in doc  $d$ 
  - $z_{d,i} \sim \text{Multi}(\theta^d)$
  - $w_{d,i} \sim \text{Multi}(\phi^{z_{d,i}})$

But, tweets (Micro-text) concentrate on *single topics*, and *aggregation/pooling* of tweets is needed.

## Twitter LDA

For each topic  $t$

- $\phi^t \sim \text{Dir}(\beta)$

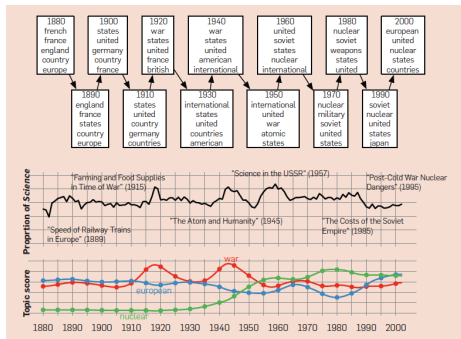
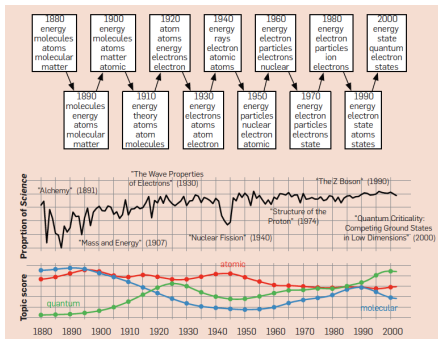
For each user  $u$

- $\theta^u \sim \text{Dir}(\alpha)$
- For each tweet  $s$ 
  - $z_{u,s} \sim \text{Multi}(\theta^u)$
  - For each word  $i$ 
    - $w_{u,s,i} \sim \text{Multi}(\phi^{z_{u,s}})$

But, we are interested in *time evolution* of the topics as well.

<sup>4</sup>Wayne Xin Zhao et al. "Comparing twitter and traditional media using topic models". In: *European conference on information retrieval*. Springer. 2011, pp. 338–349.

# Dynamic topic model



Two topics from a dynamic topic model, which was fit to the Science journal from 1880 to 2002. Top words at each decade were illustrated<sup>56</sup>.

Idea: chaining together the topic and topic proportions distributions through random walk stochastic processes and *jointly fit the model*.

<sup>5</sup>David M Blei and John D Lafferty. "Dynamic topic models". In: *Proceedings of the 23rd international conference on Machine learning*. 2006, pp. 113–120.

<sup>6</sup>David M Blei. "Probabilistic topic models". In: *Communications of the ACM* 55.4 (2012), pp. 77–84.

# Why not dynamic LDA?

## Why not dynamic LDA?

- Pros: dynamics explicitly built into the generative model.
- Cons:
  - Sensitive to model assumptions (e.g., stationarity)
  - Computationally unstable and expensive (relying on approximate inference algorithms)
  - Reliable existing software for dynamic short-document LDA

Goal: provide a simple suite of tools for general use, using robust and widely available software.

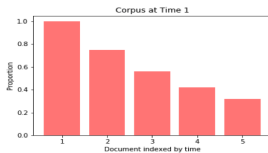
Idea: mimic the idea of dynamic topic models while taking simple, modular, and interpretable approaches.

## Emulation of dynamic model: temporally smoothed corpus

Doc 1	Doc 2	Doc 3	Doc 4	Doc 5
Tweet 1	Tweet 1	Tweet 1	Tweet 1	Tweet 1
Tweet 2	Tweet 3	Tweet 3	Tweet 3	
Tweet 3	Tweet 4	Tweet 4		
Tweet 4	Tweet 5			
Tweet 5				



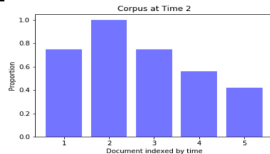
$$5 \times 1 = 5 \rightarrow 5 \times 0.75^1 \approx 4 \rightarrow \dots \rightarrow 5 \times 0.75^4 \approx 1$$



Doc 1	Doc 2	Doc 3	Doc 4	Doc 5
Tweet 1	Tweet 1	Tweet 1	Tweet 1	Tweet 1
Tweet 2	<del>Tweet 2</del>	Tweet 3	Tweet 3	<del>Tweet 2</del>
<del>Tweet 3</del>	Tweet 3	Tweet 4	<del>Tweet 4</del>	
Tweet 4	Tweet 4	<del>Tweet 5</del>		
Tweet 5	Tweet 5			



$$5 \times 0.75^1 \approx 4 \rightarrow 5 \times 1 = 5 \rightarrow \dots \rightarrow 5 \times 0.75^3 \approx 2$$



Smoothed subsampling using examples of 5 documents each containing 5 tweets. Doc 1 aggregates tweets from day 1, Doc 2 aggregates tweets from day 2, etc.

T-LDA applied *independently in parallel* to each temporally smoothed corpus from Feb 15 to May 15 ( $\approx 90$  time points) with 50 topics  $\Rightarrow$  *marginal models fitted time by time*.

# Outline

- 1 Twitter Decahose Stream
- 2 Introduction to Probabilistic Topic Models
- 3 Connecting Topics Time by Time**
- 4 Geometric Embedding of Structured High-dimensional Objects
- 5 Visualization and Interpretation of COVID-19 Discussions on Twitter

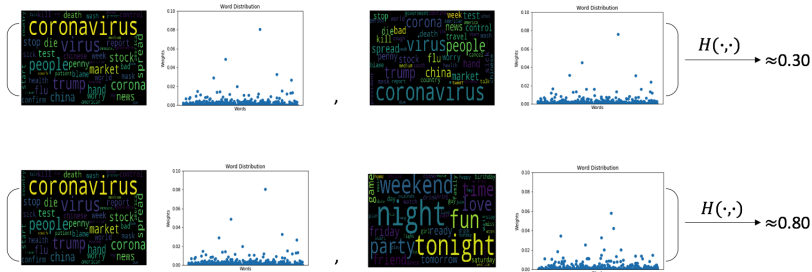




# Hellinger distances between two topics

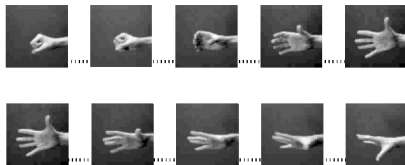
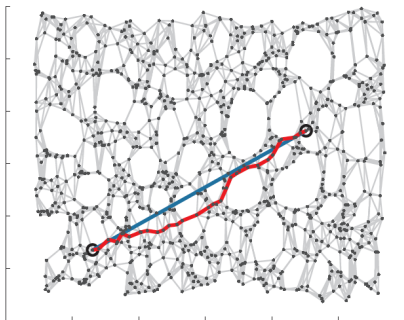
Hellinger distance <sup>7</sup> between discrete probability distributions  $P = (p_1, \dots, p_N)$  and  $Q = (q_1, \dots, q_N)$ :

$$H(P, Q) = \frac{1}{\sqrt{2}} \sqrt{\sum_{n=1}^N (\sqrt{p_n} - \sqrt{q_n})^2}, \quad 0 \leq H(\cdot, \cdot) \leq 1.$$



<sup>7</sup>A case of  $f$ -divergence that measures distance between probability distributions: <https://en.wikipedia.org/wiki/F-divergence>.

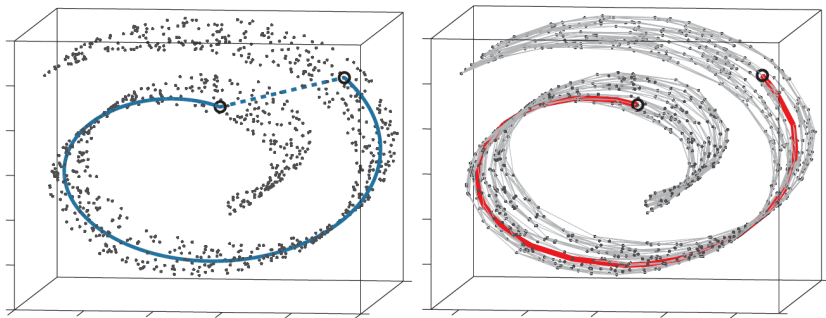
# Shortest path on neighborhood graph



Shortest distance on neighborhood graph captures perceptually natural but highly nonlinear morphs of the corresponding high-dimensional observations by transforming them approximately along geodesic paths (solid curve on the left plot)<sup>8</sup>.

<sup>8</sup>Joshua B Tenenbaum, Vin De Silva, and John C Langford. "A global geometric framework for nonlinear dimensionality reduction". In: *science* 290.5500 (2000), pp. 2319–2323.

# Full vs. neighborhood weighted graph of observations



The "Swiss roll" data set, illustrating how nearest neighborhood graph (7 nearest neighbors in this case) exploits geodesic paths for nonlinear dimensionality reduction<sup>9</sup>.

<sup>9</sup>Joshua B Tenenbaum, Vin De Silva, and John C Langford. "A global geometric framework for nonlinear dimensionality reduction". In: *science* 290.5500 (2000), pp. 2319–2323.

## Topic transformation on full vs. 10-nearest neighbor graph

10-NN graph shortest path



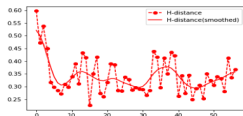
...



...



...



Full graph shortest path



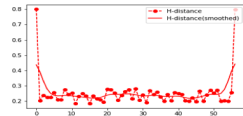
...



...



...

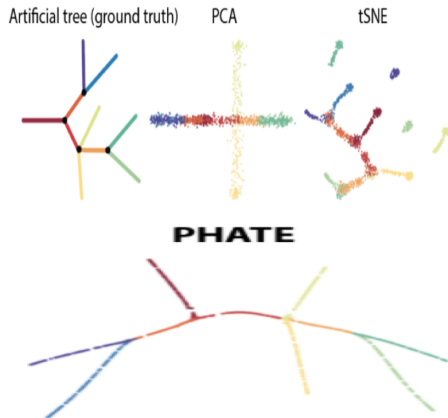


# Outline

- 1 Twitter Decahose Stream
- 2 Introduction to Probabilistic Topic Models
- 3 Connecting Topics Time by Time
- 4 Geometric Embedding of Structured High-dimensional Objects**
- 5 Visualization and Interpretation of COVID-19 Discussions on Twitter

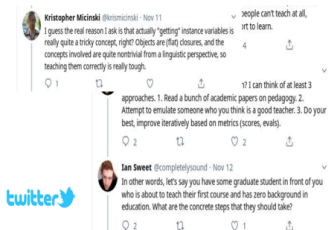
# PHATE<sup>10</sup>

- To visualize and interpret high dimensional word distributions, need *lower-dimensional embedding* that capture the intrinsic high-dimensional *trajectory structure* of the data.
- Traditional methods like PCA assumes linearity.
- Nonlinear methods like t-SNE DO NOT naturally exhibits trajectory or progression.
- PHATE (Potential of Heat-diffusion for Affinity-based Trajectory Embedding) is designed explicitly to preserve progression structure in data.

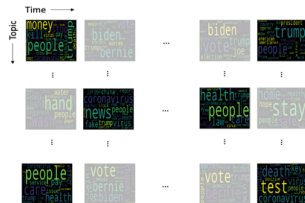


<sup>10</sup>Kevin R Moon et al. "Visualizing structure and transitions in high-dimensional biological data". In: *Nature Biotechnology* 37.12 (2019), pp. 1482–1492.

# Summary of the procedure

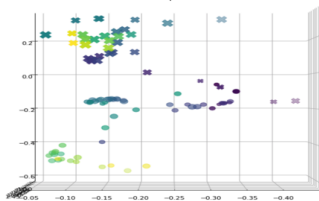


1. LDA on temporally smoothed corpus



2. Hellinger distances

3a. PHATE embedding



3b. Neighborhood graph and shortest path



4. Visualize and interpret multiple paths and clusters

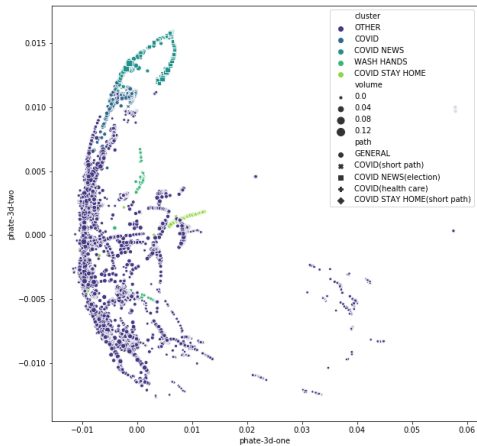




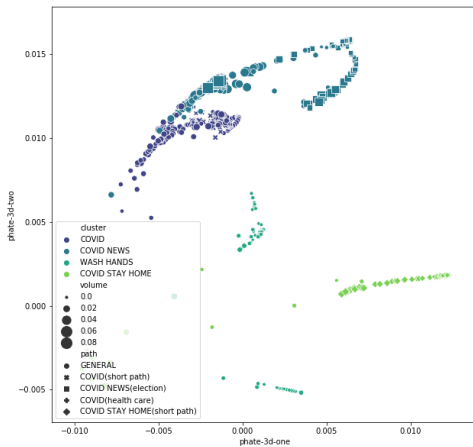
# Outline

- 1 Twitter Decahose Stream
- 2 Introduction to Probabilistic Topic Models
- 3 Connecting Topics Time by Time
- 4 Geometric Embedding of Structured High-dimensional Objects
- 5 Visualization and Interpretation of COVID-19 Discussions on Twitter**

# PHATE embedding of 4500 word distributions

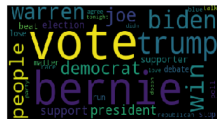
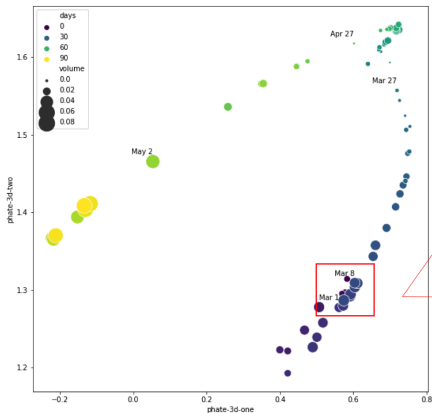


# PHATE embedding of selected COVID topics

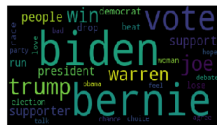


# Case study I: presidential election topic

Idea: use real events we know about to understand the trajectories.

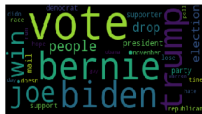
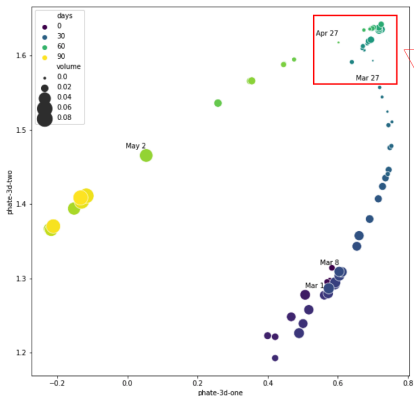


Super Tuesday, Mar 3

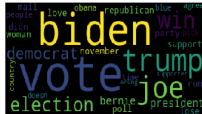


Super Tuesday week clustered with high volume.

## Case study I cont'd: April subcluster



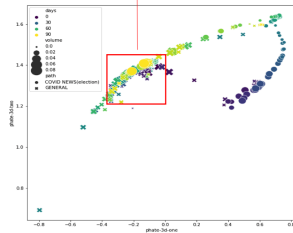
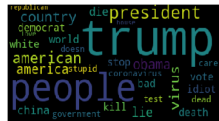
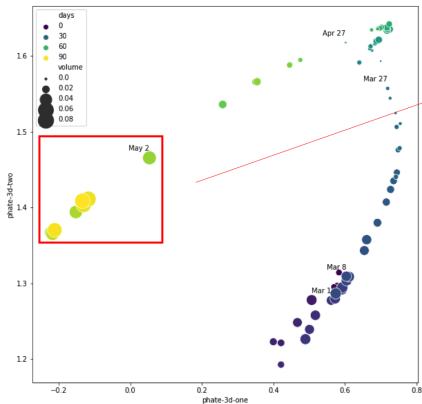
Bernie Sanders dropped out, April 8



People endorsed Biden, who wins Wisconsin &amp; Ohio, April 27

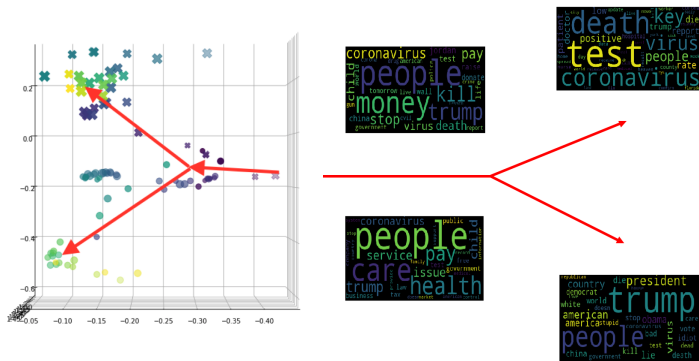
Tweets volume dies out in late March but re-surged and clustered in April because Bernie Sanders dropped out and more election-related activities.

## Case study I cont'd: May subcluster



Election tweets in May converged to general COVID-related news tweets.

# Case study II: general COVID topic



Divergence of two similar COVID-related conversations that ended up as a COVID-testing conversation and a COVID-politics conversation, respectively.

## Case study II cont'd: dynamics



# Conclusion

This talk:

- Generative statistical model & Computational geometry & Algorithms collectively provide efficient ways for understanding high-dimensional noisy Twitter data  $\Rightarrow$  visualization and interpretation of public discourse around COVID-19.
- Longitudinal analysis of Twitter data done without complicated models  $\Rightarrow$  understanding the impact of COVID-19 on various societal aspects.

Future work:

- Wrapping various tools into a single re-usable implementation for the public.
- Apply such tools globally or to tweets from other regions other than the US.
- Spatial analysis of Twitter data.