

Covid19 Twitter Project

Summary

=====

This Twitter dataset is crawled from ids provided in the TweetsCOV19 dataset <<https://data.gesis.org/tweetscov19/>>.

The data are contained in the files `tweets.txt`, `users.txt`, `places.txt` and `media.txt`. More details about the contents and use of all these files follows.

For questions, please contact Mohamed Reda Bouadjenek <rbouadjenek@gmail.com>.

Content and Use of Files

=====

Formatting and Encoding

The dataset files are formatted using the JSON file format (<https://en.wikipedia.org/wiki/JSON>) where each row is one instance. These files are encoded as UTF-8.

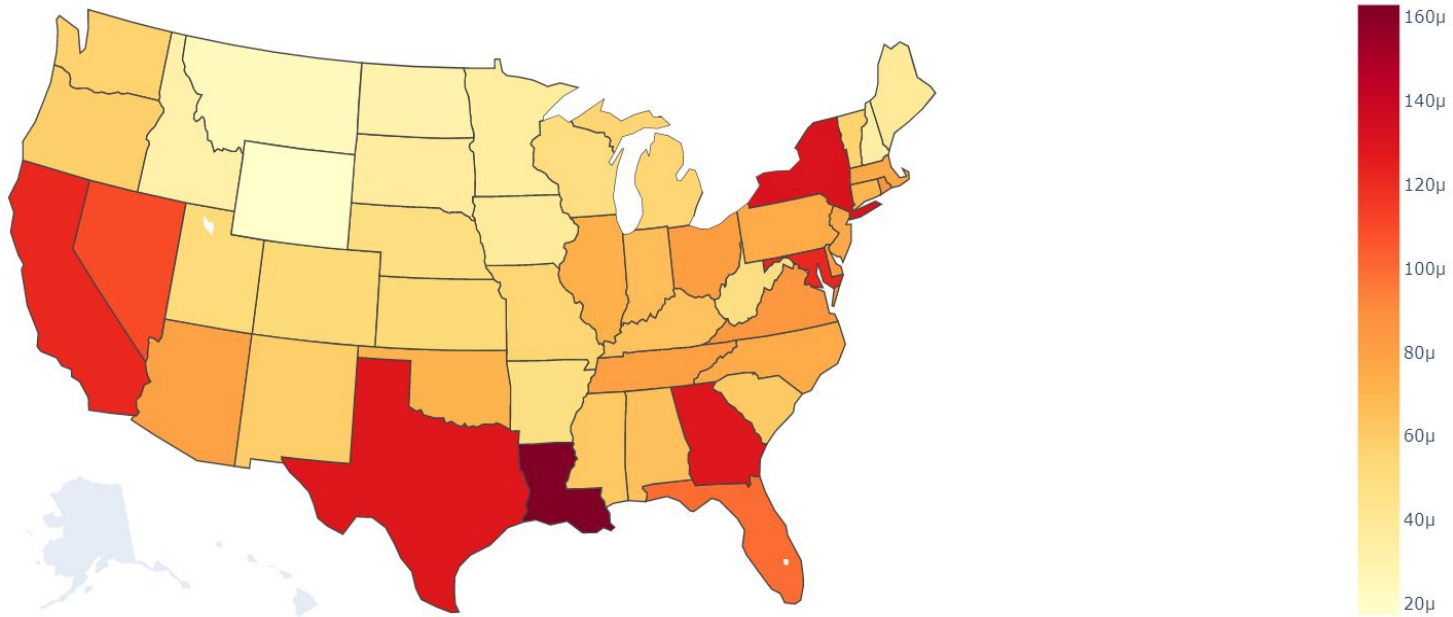
Tweets File Structure (tweets.txt)

All tweets are contained in the file `tweets.txt`. Each line of this file represents the information crawled for one tweet id, and has ****mainly**** the following JSON keys:

- id: the tweets id
- author_id: the id of the user who posted the tweet
- created_at: creation date of the tweet
- text: the content of the tweet
- public_metrics: public metrics like retweet count, reply count, like count, and quote count
- entities: entities mentionned in the text including mentions and hashtags
- geo: which gives the id of the place that can be retrieved from 'places.txt'.

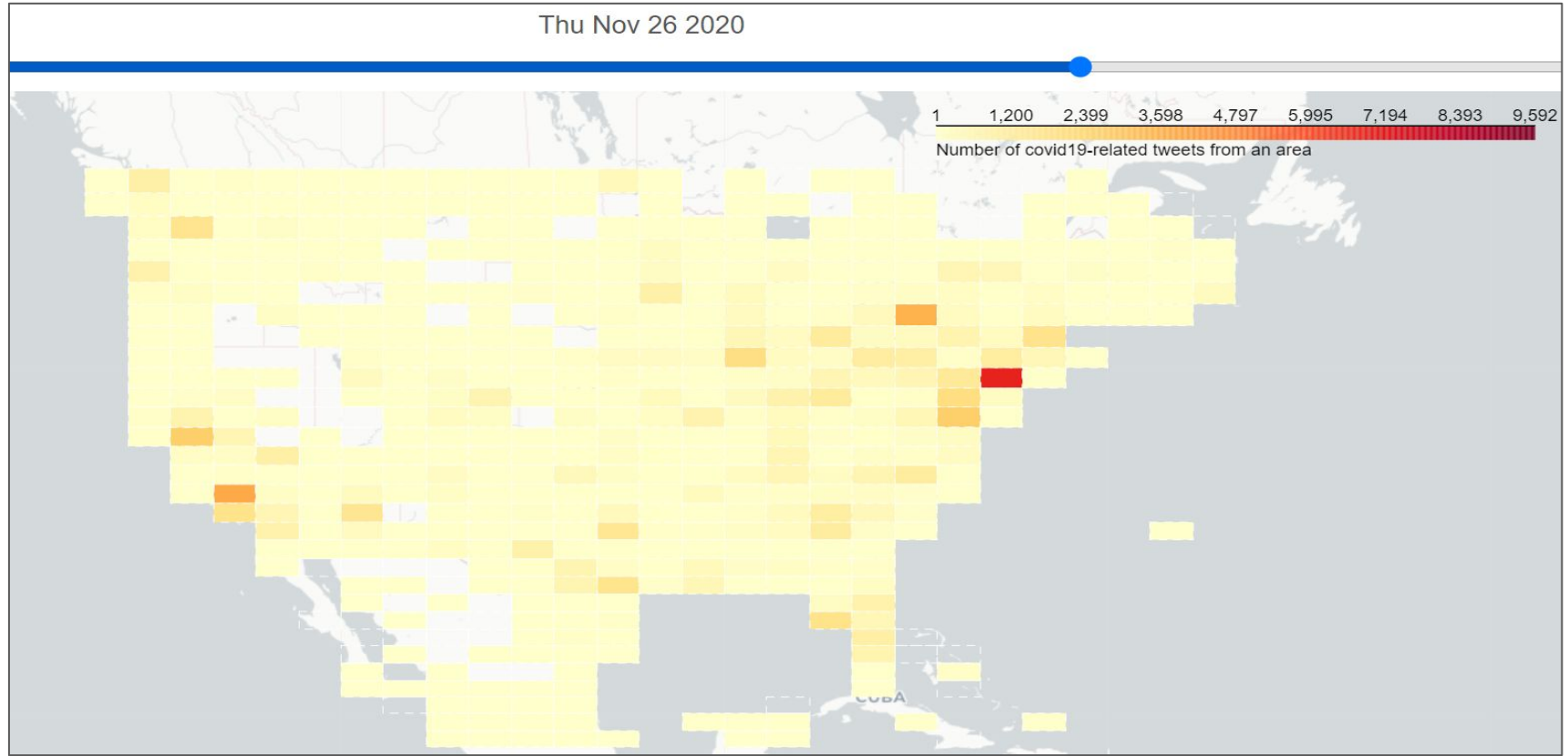
**Hshmat Sahak
Litao (John) Zhou**

Time Series Visualizations



Weekly Twitter Count by State (Normalized by Population)

Grid Map with Time Slider



Same data as previous slide, but different visualization approach. This is grid-wise; not state-wise

Sentiment Analysis

Current Strategy:

Remove election-related tweets:

- Trump|Biden|Election|democratic|republican|party|President|campaign|elector|candidate

Then divide into:

Family-related Tweets

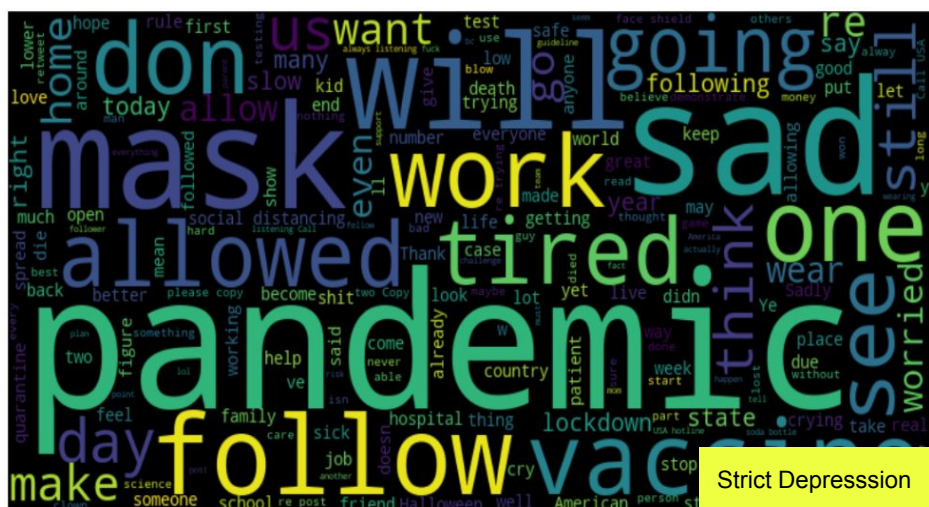
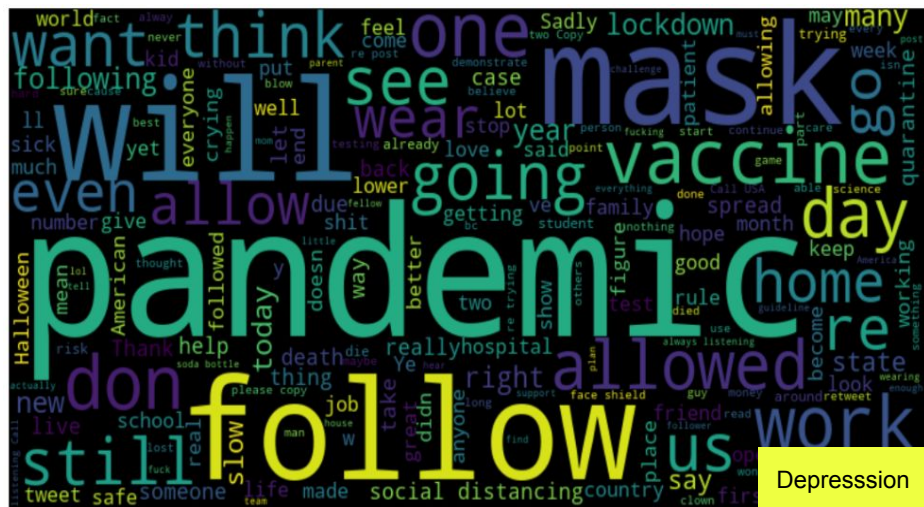
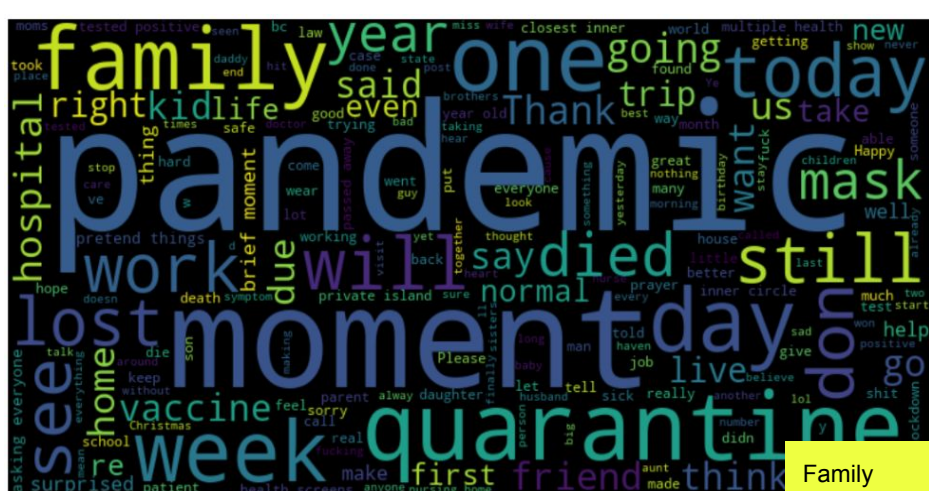
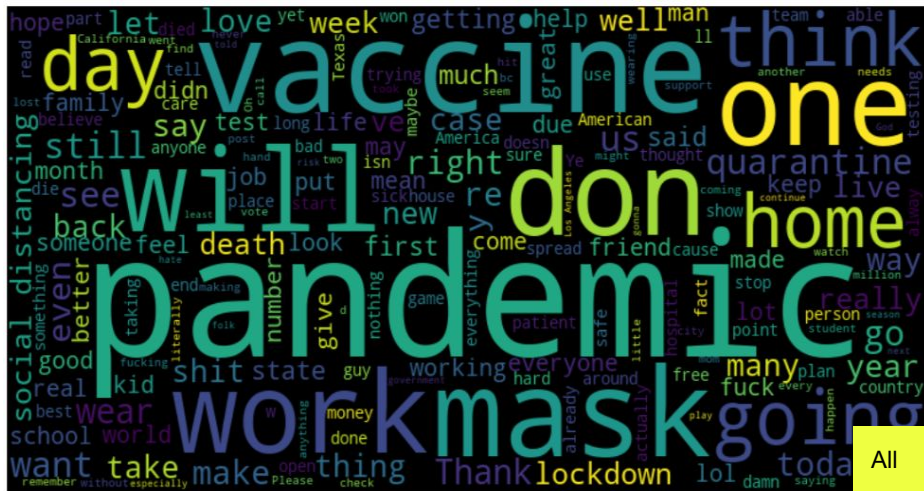
- sister|brother|mother|father|grandpa|grandma|grandparents|grandmother|grandfather|cousin|dad|mom

Depression-related Tweets

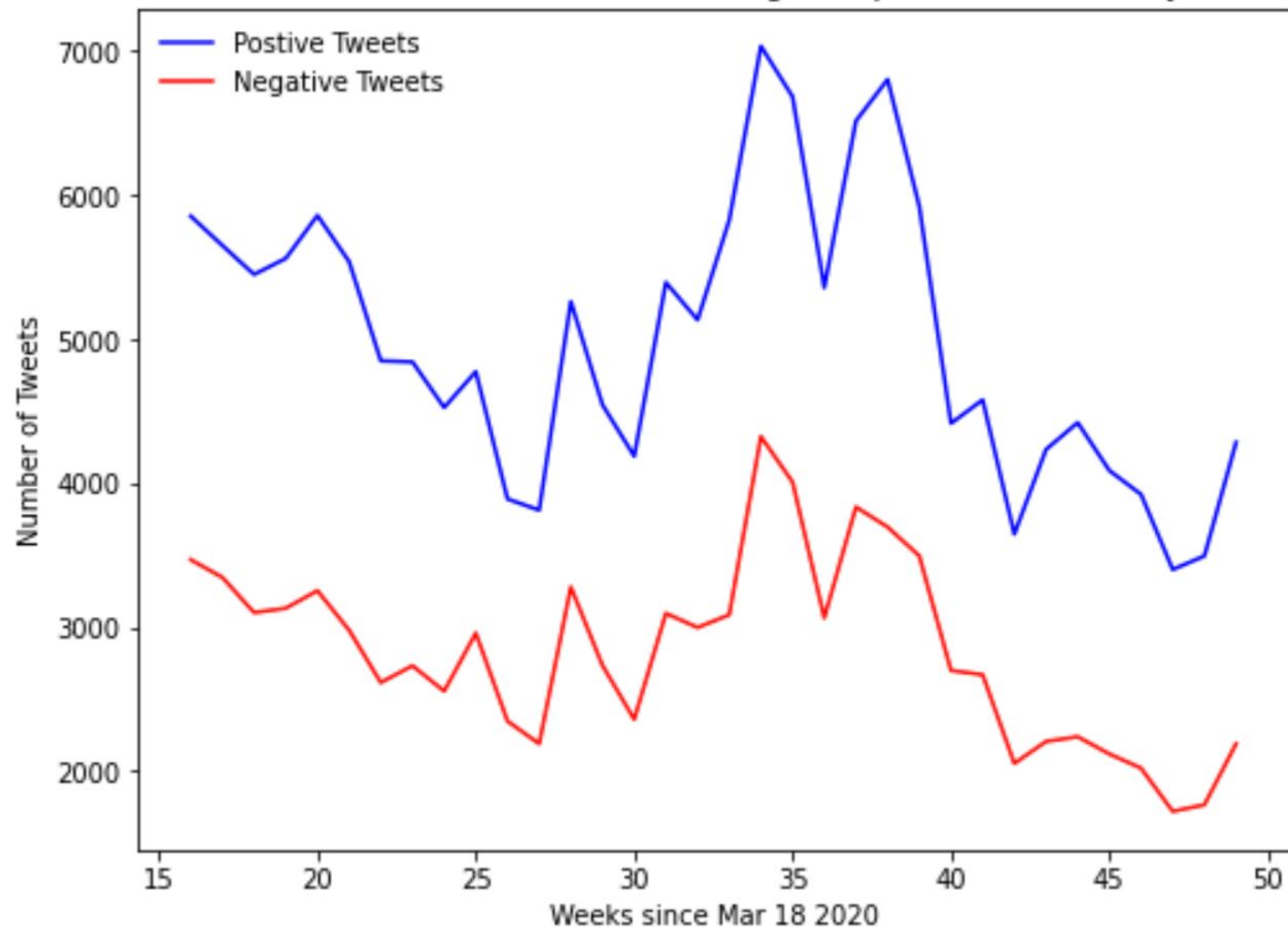
- overwhelmed|exhausted|distressed|anxiety|anxious|tired|low|depression|discouraged|desperate|demotivated|insomnia|cry|nervous|worried|lonely|sad|empty|suicide|antidepressant|hopeless

Strict Depression-related Tweets

- overwhelmed|exhausted|distressed|anxiety|anxious|depression|discouraged|demotivated|insomnia|lonely|empty|suicide|antidepressant|hopeless



US Number of Tweets of Postive and Negative polarities on weekly basis



Filters for topics of family,covid and us election

```
[3]: #use regex to clean and filter out the useful info from tweets related to family
      #also filter out the word related to us election
      regex_election = re.compile(r'(?i)Trump|Biden|Election|democratic|republican|party|President|campaign|elector|candidate')
      regex_family = re.compile(r'(?i)sister|brother|mother|father|grandpa|grandma|grandparents|grandmother|grandfather|cousin|child|dad|mom|uncle|aunt|nephew|niece')
      regex_covid = re.compile(r'(?i)corona|covid|virus|pandemic|epidemic|quarantine|lockdown|social distancing|isolation|mask|infectious|formite|vanccine|vaccine')

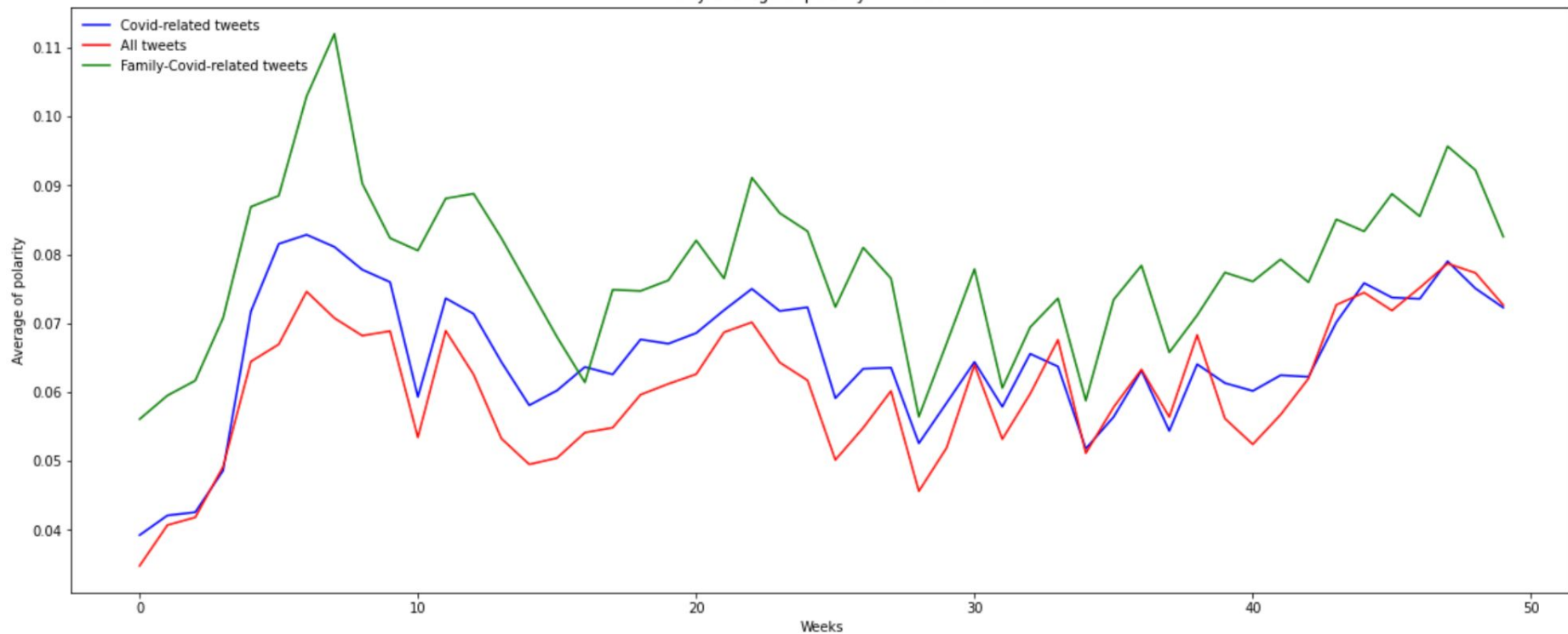
      def covid_related_text(text):
          if regex_covid.search(text):
              return True
          return False

      def family_related_text(text):
          if regex_family.search(text):
              return True
          return False

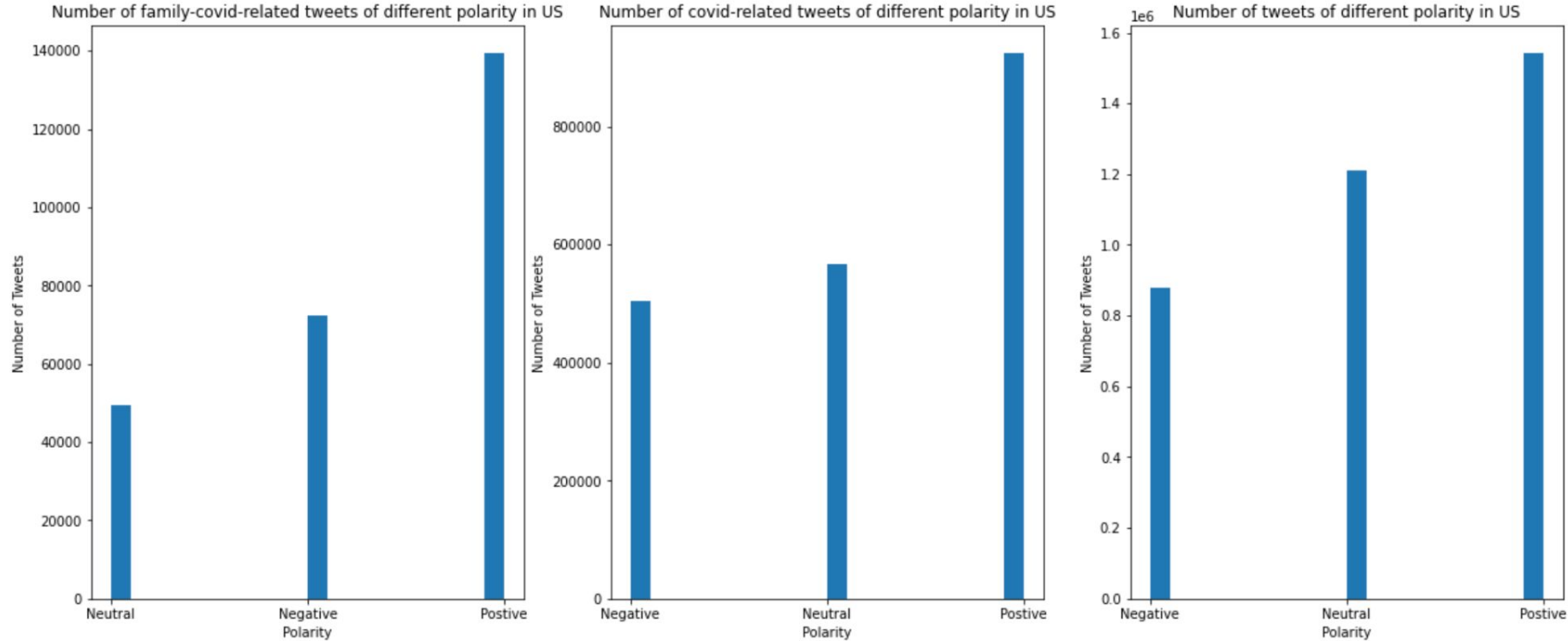
      def election_related_text(text):
          if regex_election.search(text):
              return True
          return False

      total_df['election_related'] = total_df['processed_text'].apply(election_related_text)
      total_df['covid_related'] = total_df['processed_text'].apply(covid_related_text)
      total_df['family_related'] = total_df['processed_text'].apply(family_related_text)
```

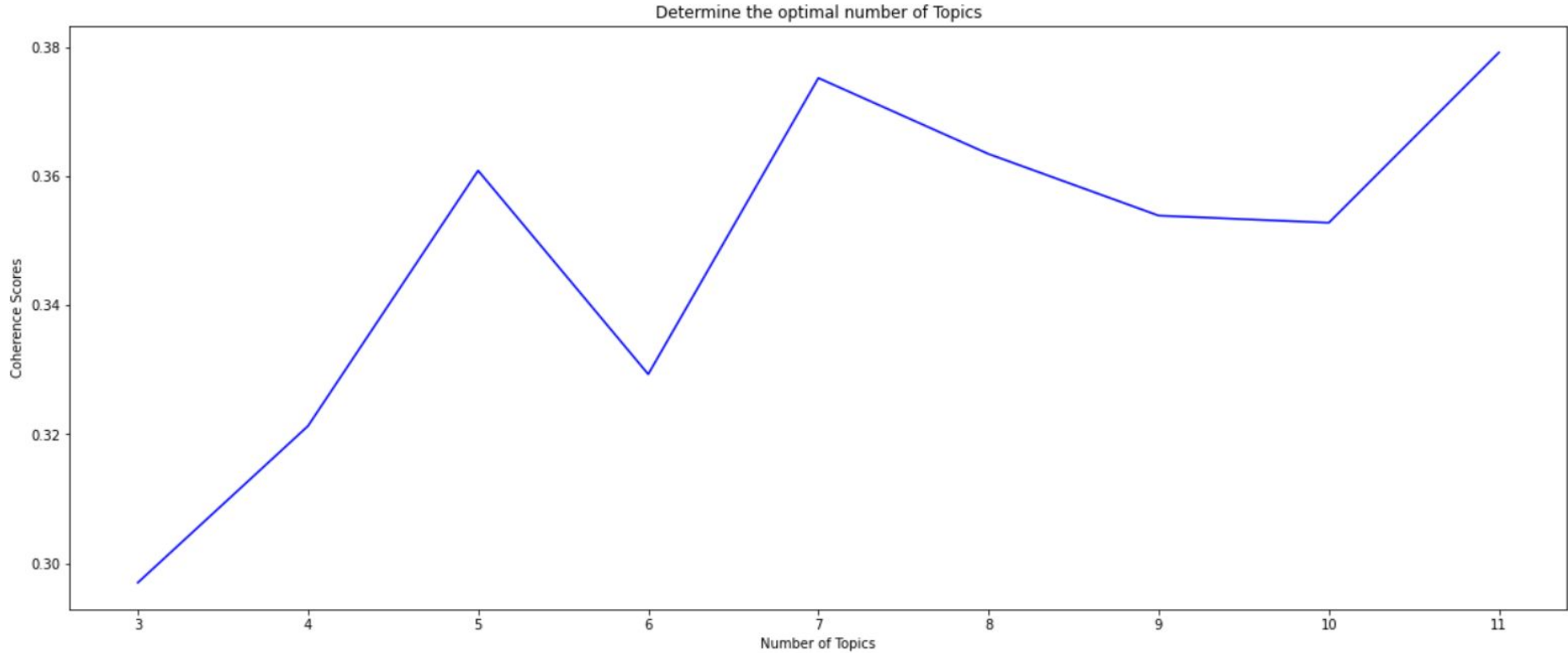
weekly average of polarity of tweets in US



Textblob toolkit for tweets polarity analysis



Optimal Number of Topics



Pick 7 as optimal number of topics

Topic Modelling - Family covid-related

Selected Topic:

Slide to adjust relevance metric:(2)

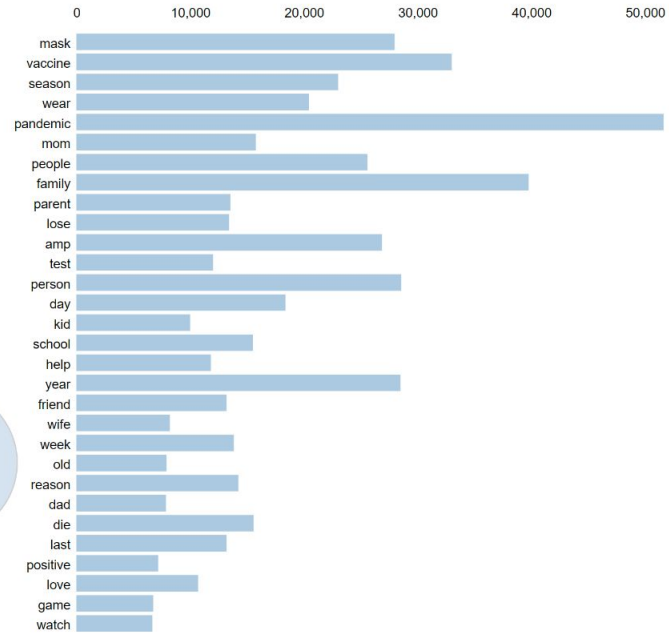
$\lambda = 1$

0.0 0.2 0.4 0.6 0.8 1.0

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Salient Terms¹



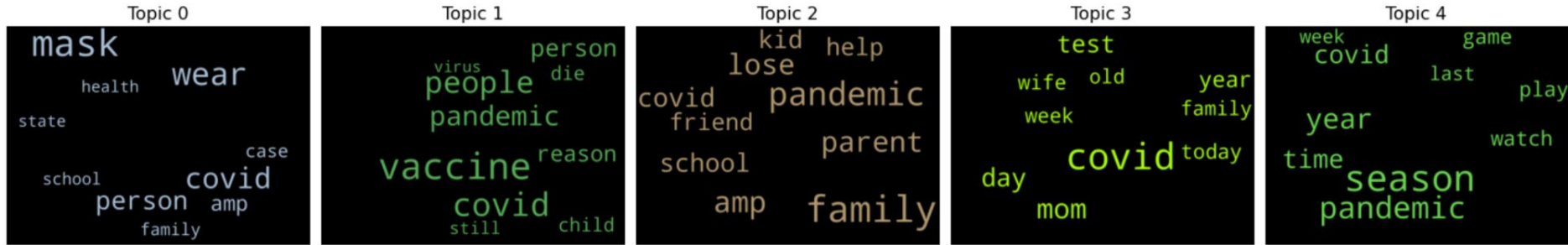
Overall term frequency

Estimated term frequency within the selected topic

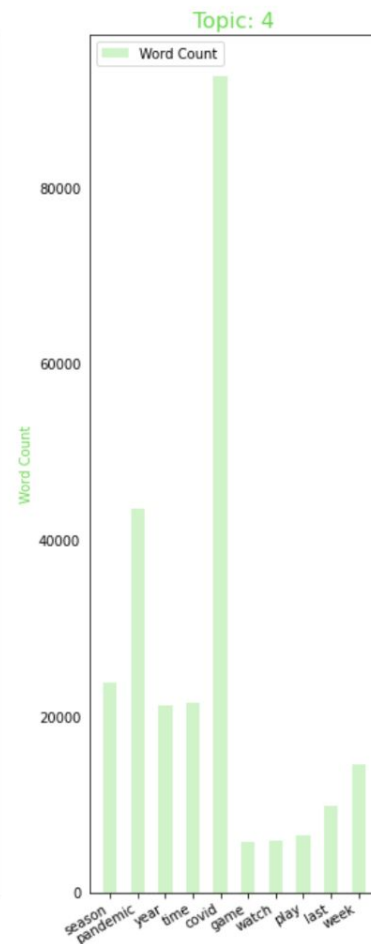
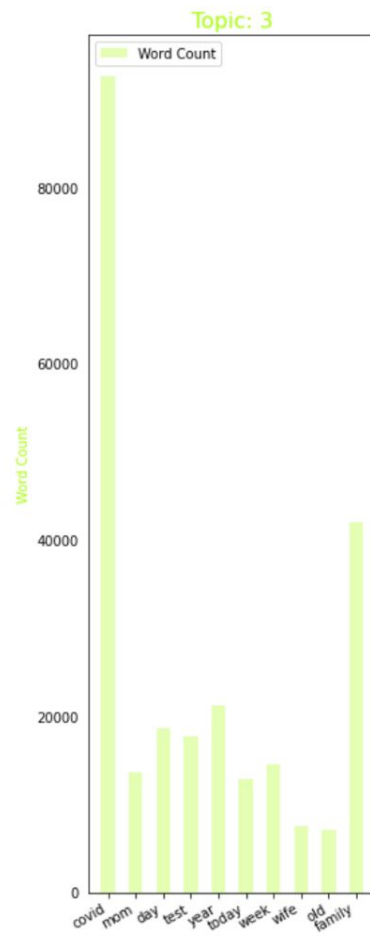
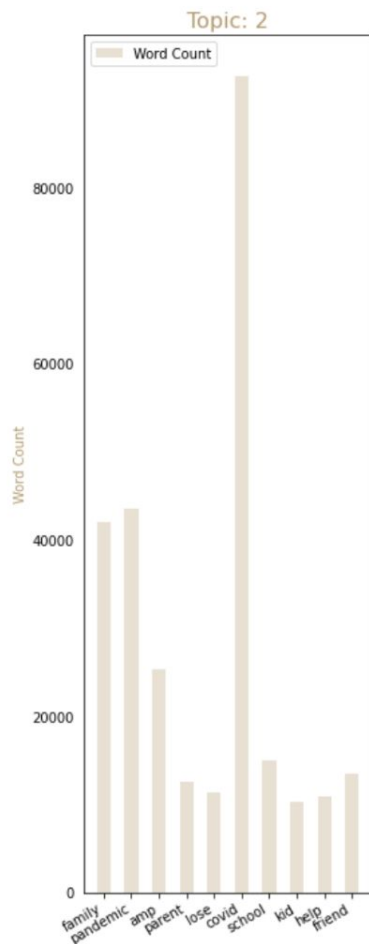
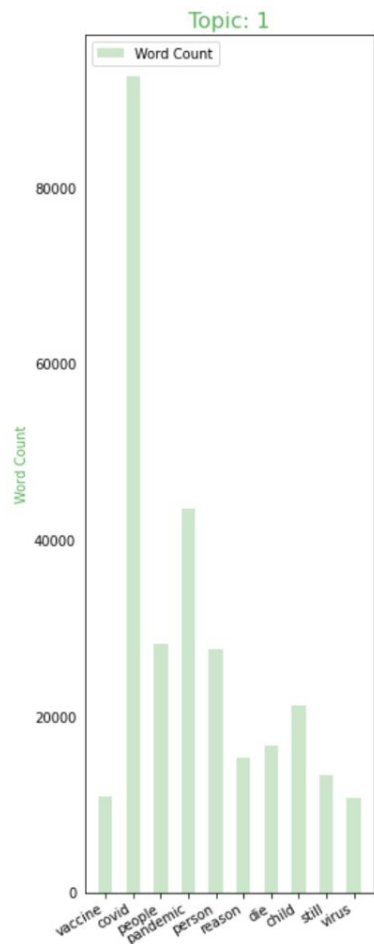
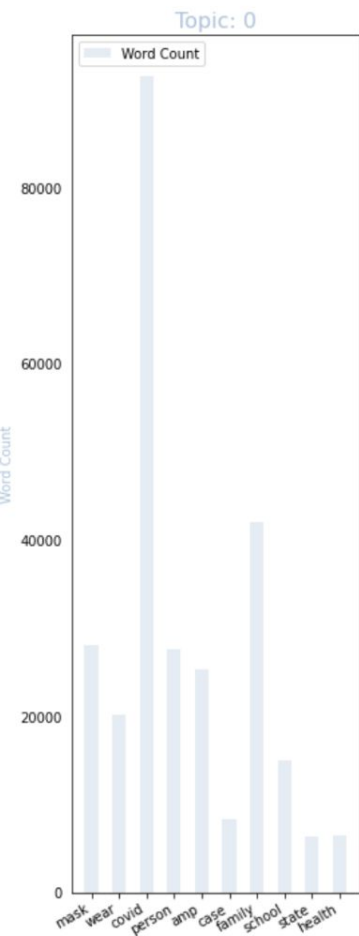
1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))]; see Chuang et. al (2012)

2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (2014)

Word Cloud for Topic Modelling - Family-covid-related

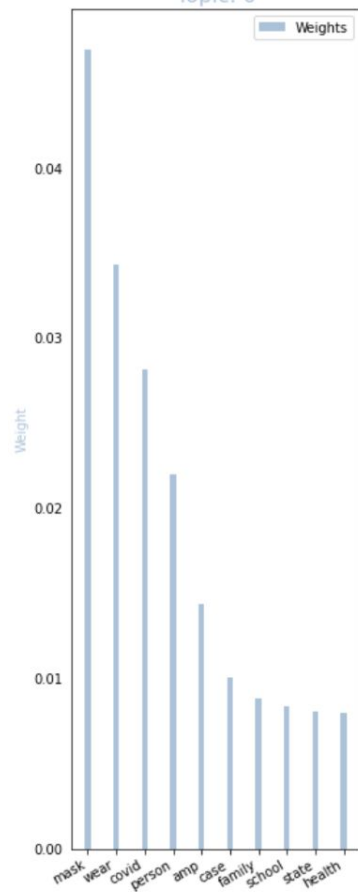


Word Count of Topic Keywords for the family-covid-related tweets

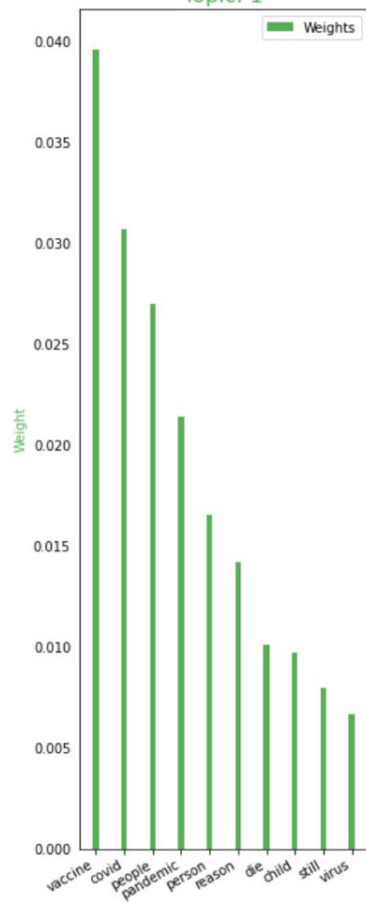


Word weight of Topic Keywords for the family-covid-related tweets

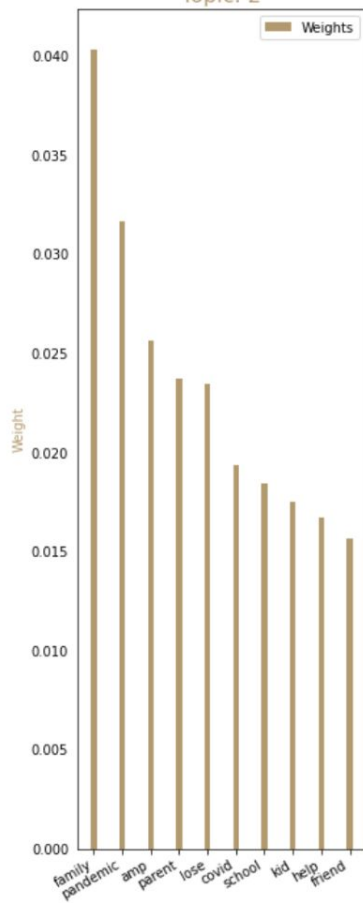
Topic: 0



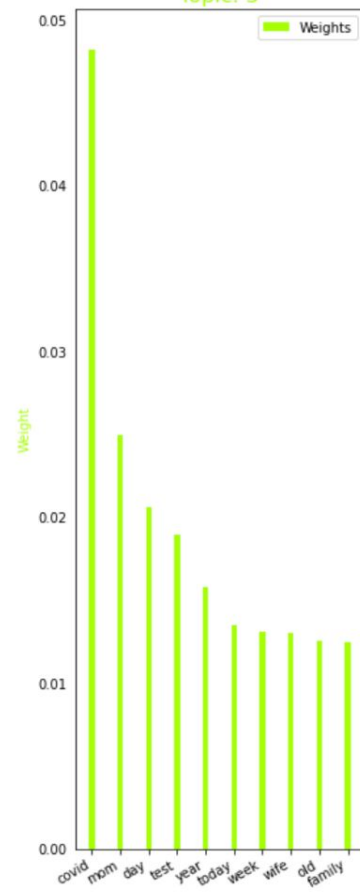
Topic: 1



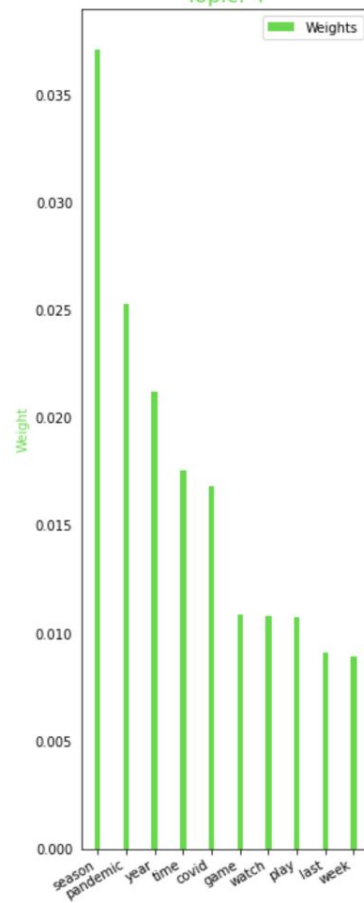
Topic: 2



Topic: 3

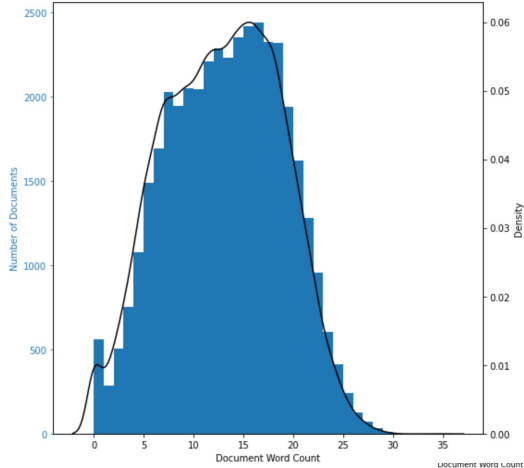


Topic: 4

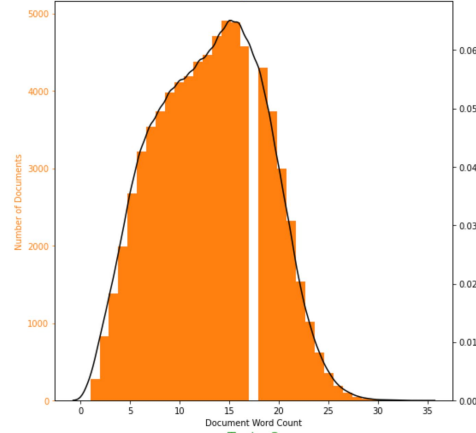


Frequency of word count in each documents

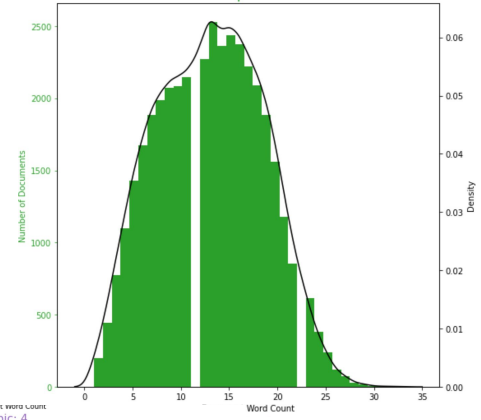
Topic: 0



Topic: 1

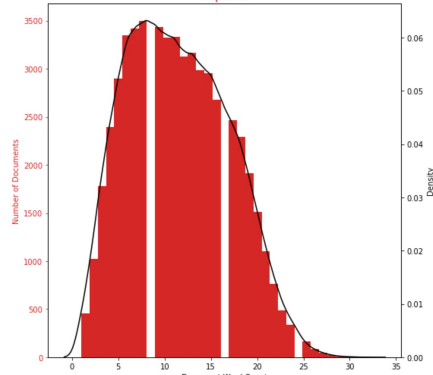


Document Word Count
Topic: 2



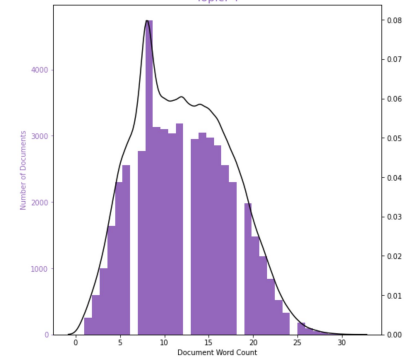
Document Word Count

Topic: 3



Document Word Count

Topic: 4



Topic Modelling- Strict Depression

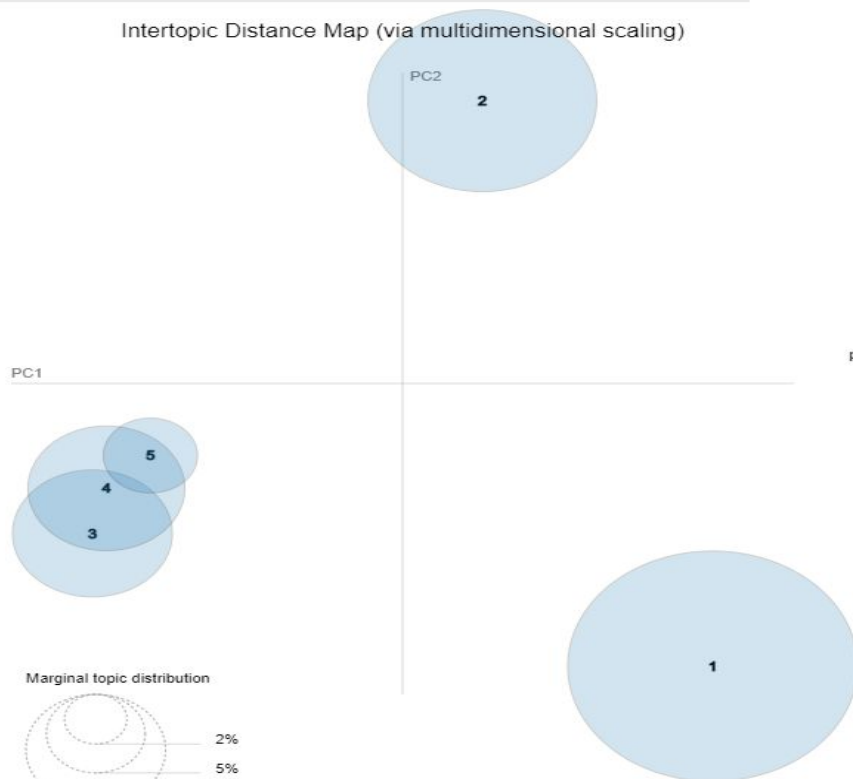
Selected Topic: Previous Topic Next Topic Clear Topic

Slide to adjust relevance metric:(2)

$\lambda = 1$

0.0 0.2 0.4 0.6 0.8 1

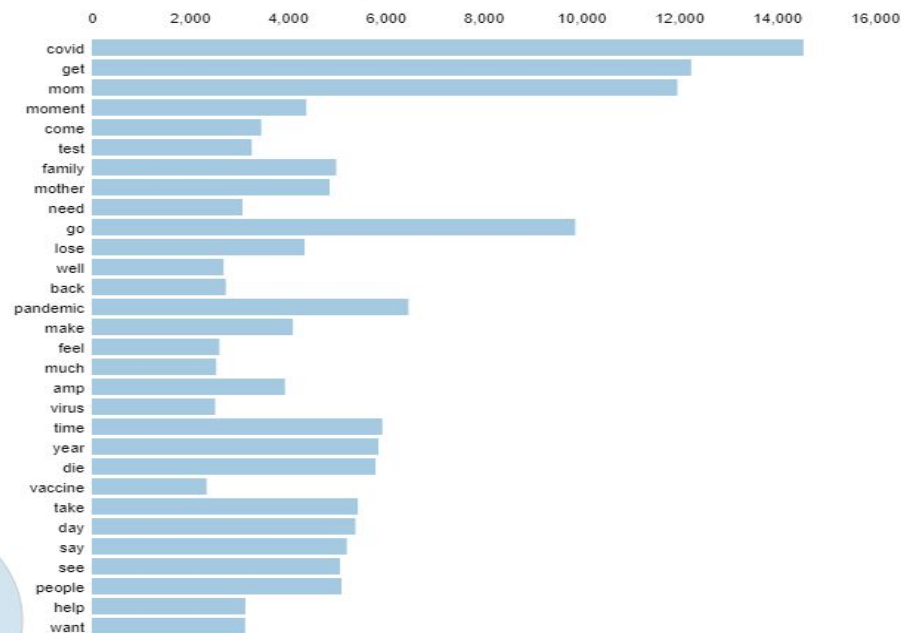
Intertopic Distance Map (via multidimensional scaling)



Marginal topic distribution



Top-30 Most Salient Terms¹



Overall term frequency

Estimated term frequency within the selected topic

1. $\text{saliency}(\text{term}, w) = \text{frequency}(w) * [\sum_t p(t|w) * \log(p(t|w)/p(t))]$ for topics t . see Chuang et. al (2012)

2. $\text{relevance}(\text{term}, w | \text{topic}, t) = \lambda * p(w|t) + (1 - \lambda) * p(w|t)/p(w)$; see Sievert & Shirley (2014)

What we need help with

- What particular problem should we try to approach that hasn't been covered by others? (or, that we can do better)
 - Publication is a goal here...
- What better ways can we detect depression
 - Something better than just keywords...
 - Reference papers use keywords or train ML algorithms. We could get advice from a mental health expert...