



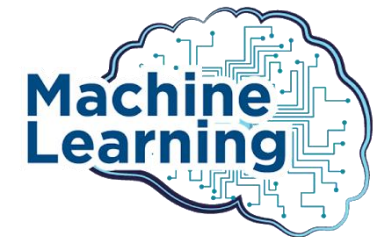
Discussing the Covid19 Pandemic

Capstone – Final Presentation

By:
Shreyas Chitransh

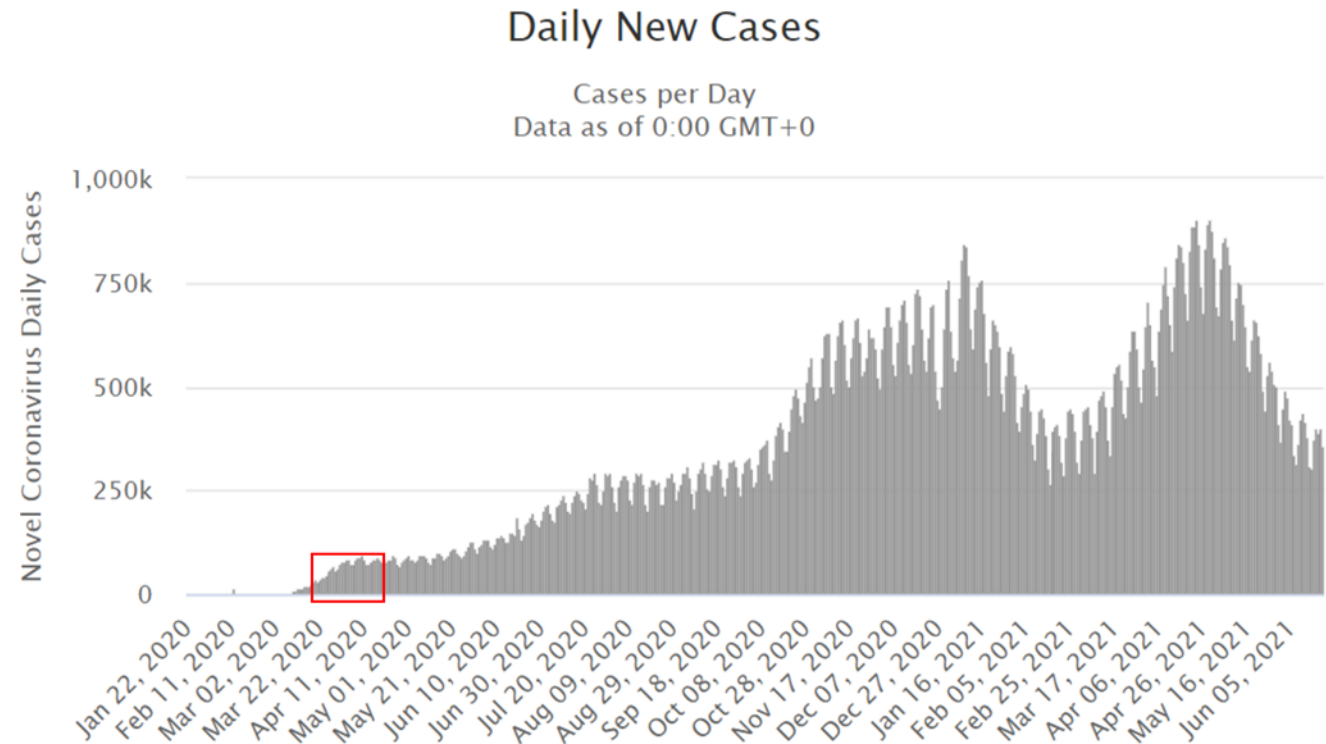
Agenda

- Background
 - Objectives
 - Methodology
- Process:
 - Data gathering and Cleaning
 - Model Optimization
- Results
- Conclusions & Further Work



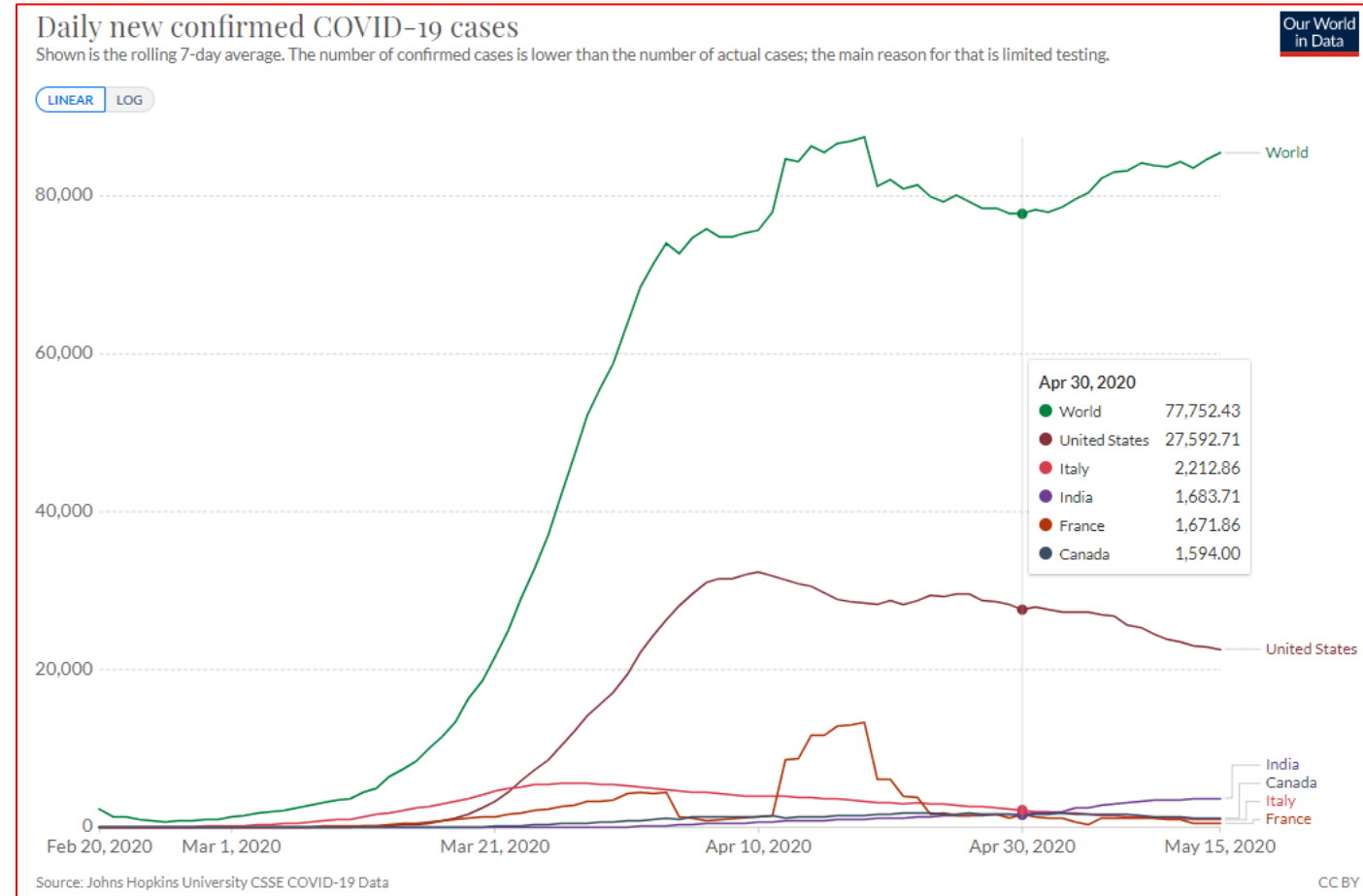
Background

- Covid-19 culminated in fear and grief but also some of the best solidarity amongst the global citizens;
- Goal:
 - To automatically identify the underlying topics being discussed with context to Covid19;
- Methodology:
 - Analyze Twitter data from the first global surge of Covid19 cases using Unsupervised ML;
 - Topic Modeling - Latent Dirichlet Allocation (LDA)
 - 28th March – 30th April 2020



Background

- Covid-19 culminated in fear and grief but also some of the best solidarity amongst the global citizens;
- Goal:
 - To automatically identify the underlying topics being discussed with context to Covid19;
- Methodology:
 - Analyze Twitter data from the first global surge of Covid19 cases using Unsupervised ML;
 - Topic Modeling - Latent Dirichlet Allocation (LDA)
 - 28th March – 30th April 2020



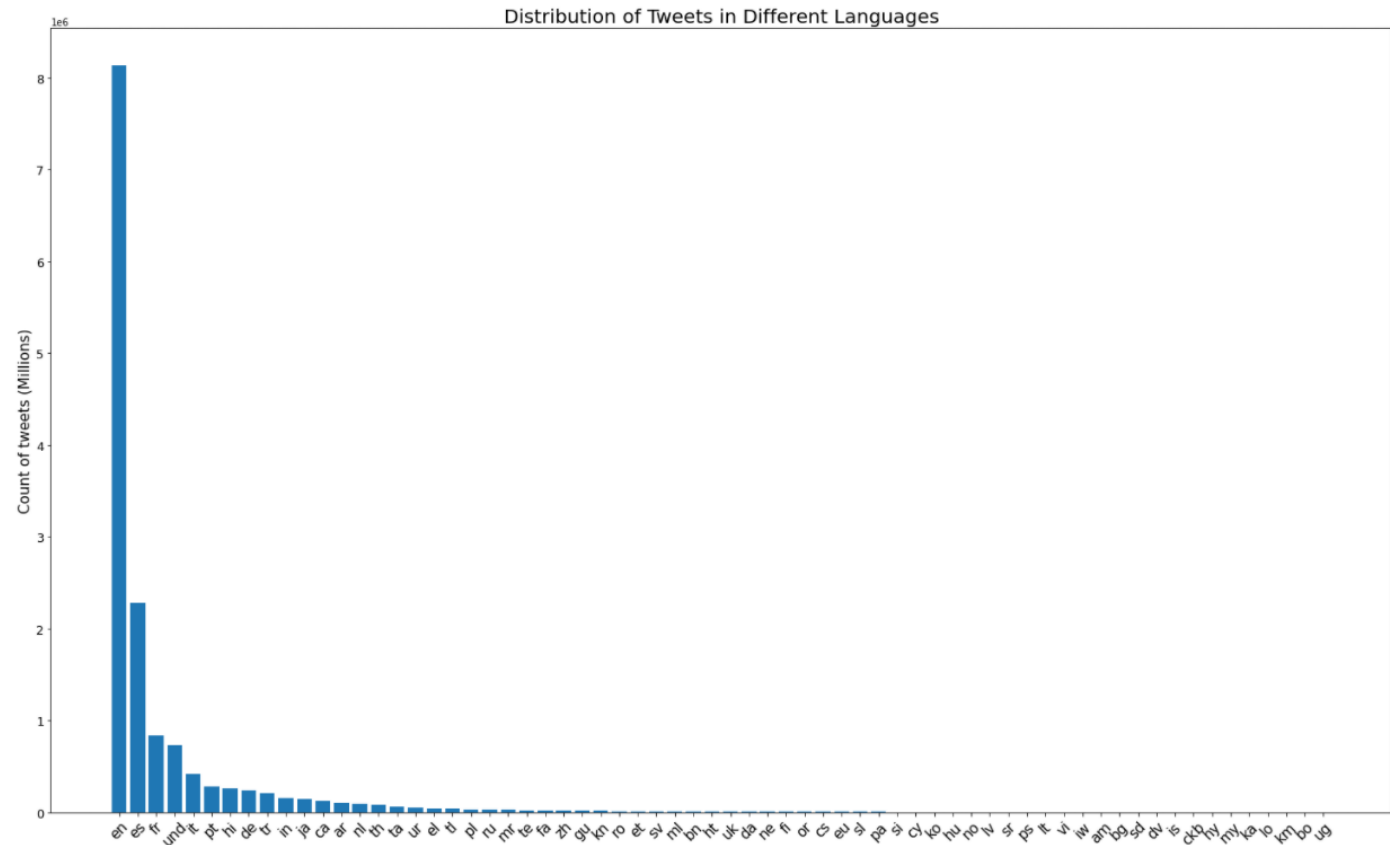
Process - Data Gathering

- Kaggle dataset >4GB of Tweets spanning 28th March – 30th April 2020;
 - CSV Files
- >14M Rows and 22 Columns
 - Rows : Separate Tweets
 - Columns : Tweet and User Metadata
- Metadata:
 - Tweet:
 - Status ID, Date, Time, Favorites, Retweets , Autodetected language etc.
 - User:
 - User ID, Screen Name, Account Creation Date etc.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14607013 entries, 0 to 14607012
Data columns (total 22 columns):
#   Column                Dtype
---  -
0   status_id             int64
1   user_id               int64
2   created_at            object
3   screen_name           object
4   text                  object
5   source                object
6   reply_to_status_id    float64
7   reply_to_user_id      float64
8   reply_to_screen_name  object
9   is_quote              bool
10  is_retweet             bool
11  favourites_count       int64
12  retweet_count          int64
13  country_code           object
14  place_full_name        object
15  place_type             object
16  followers_count        int64
17  friends_count          int64
18  account_lang           float64
19  account_created_at     object
20  verified               bool
21  lang                   object
dtypes: bool(3), float64(3), int64(6), object(10)
memory usage: 2.1+ GB
```

Data Clean-Up

- Clean up and Preprocessing most important steps for Unsupervised NLP
- Cleaning:
 - Removed Duplicates
 - Removed columns containing null values (over 80% null)
 - For remaining data, removed rows containing null values (884 rows, <0.05%)
 - Removed rows containing non-English Tweets
 - Final English Tweets --> 8.1M



Analysis Dataset – Improve Manageability

Original Dataset

14.6M Rows with 22 Columns

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14607013 entries, 0 to 14607012
Data columns (total 22 columns):
#   Column                Dtype
---  -
0   status_id             int64
1   user_id               int64
2   created_at            object
3   screen_name           object
4   text                  object
5   source                object
6   reply_to_status_id    float64
7   reply_to_user_id      float64
8   reply_to_screen_name  object
9   is_quote              bool
10  is_retweet             bool
11  favourites_count       int64
12  retweet_count          int64
13  country_code           object
14  place_full_name        object
15  place_type             object
16  followers_count        int64
17  friends_count          int64
18  account_lang           float64
19  account_created_at     object
20  verified               bool
21  lang                   object
dtypes: bool(3), float64(3), int64(6), object(10)
memory usage: 2.1+ GB
```

---> 8.1M
English --->
Tweets with
9 Columns

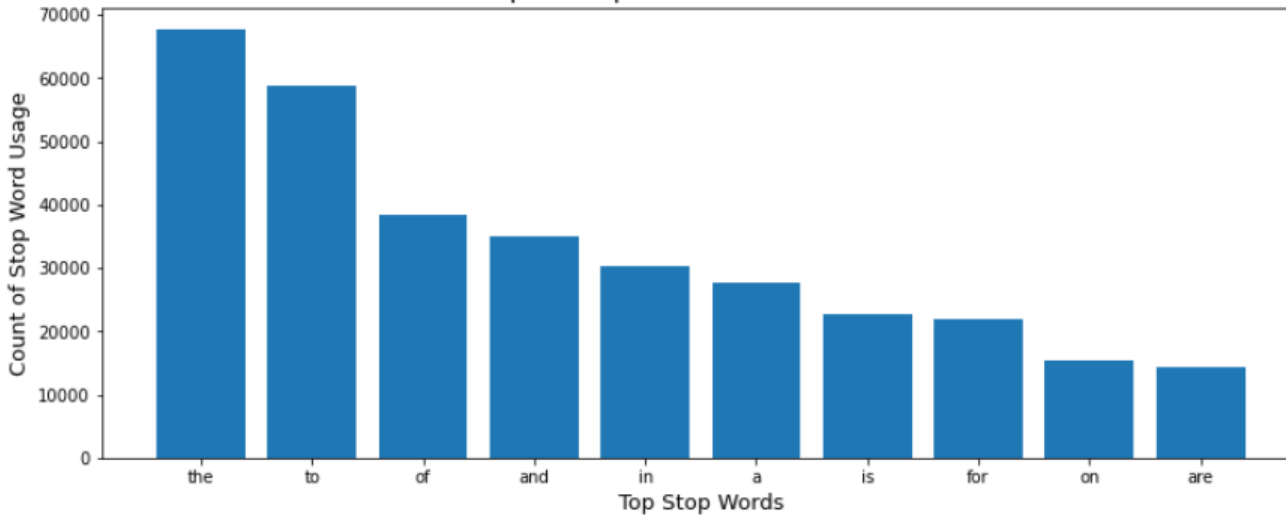
Analysis Dataset

81.3K Rows with 9 Columns

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 81333 entries, 0 to 81332
Data columns (total 9 columns):
#   Column                Non-Null Count  Dtype
---  -
0   created_at            81333 non-null  datetime64[ns, UTC]
1   screen_name           81333 non-null  object
2   text                  81333 non-null  object
3   is_retweet            81333 non-null  bool
4   favourites_count       81333 non-null  int64
5   retweet_count         81333 non-null  int64
6   followers_count       81333 non-null  int64
7   friends_count         81333 non-null  int64
8   verified              81333 non-null  bool
dtypes: bool(2), datetime64[ns, UTC](1), int64(4), object(2)
memory usage: 4.5+ MB
```

Preprocessing (Crucial Step) & EDA

Top 10 Stop Words and their Count



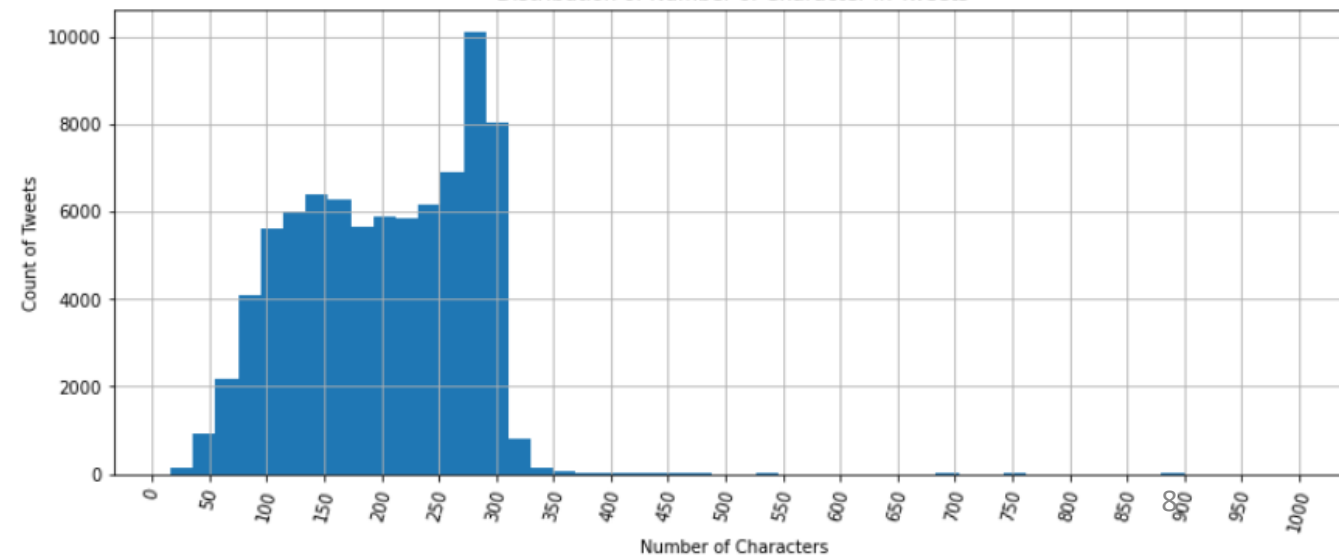
Preprocessing:

- Removing URLs, mentions, symbols and stopwords;
- Tokenization and Lemmatization;
- Changed with evolving models.

EDA:

- Identifying data trends to improve analysis methodology;
- Changed with evolving models.

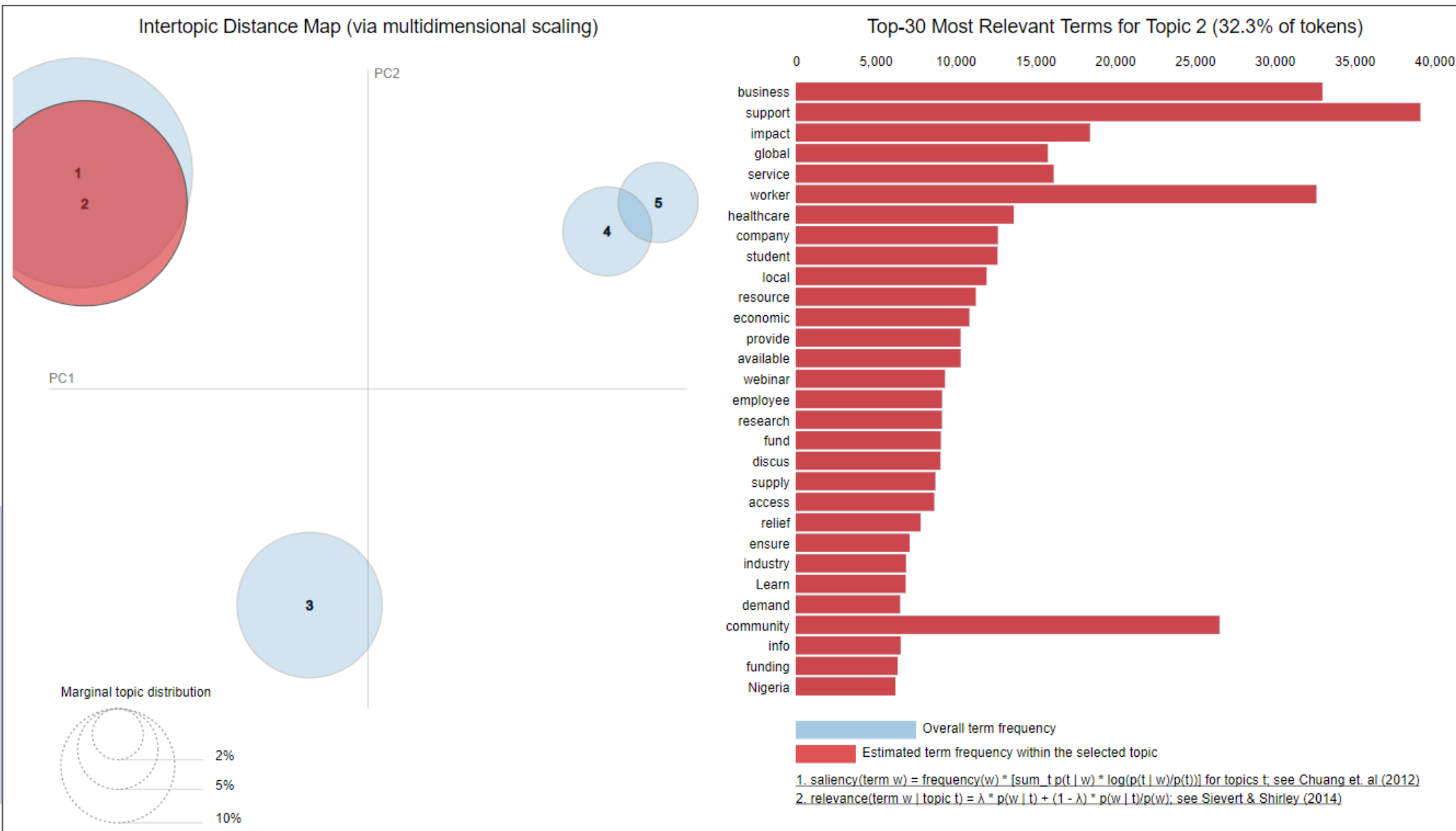
Distribution of Number of Character in Tweets



Preprocessed Corpus: 1.5M



LDA Model 1 - 5 Topics & 35% Coherency



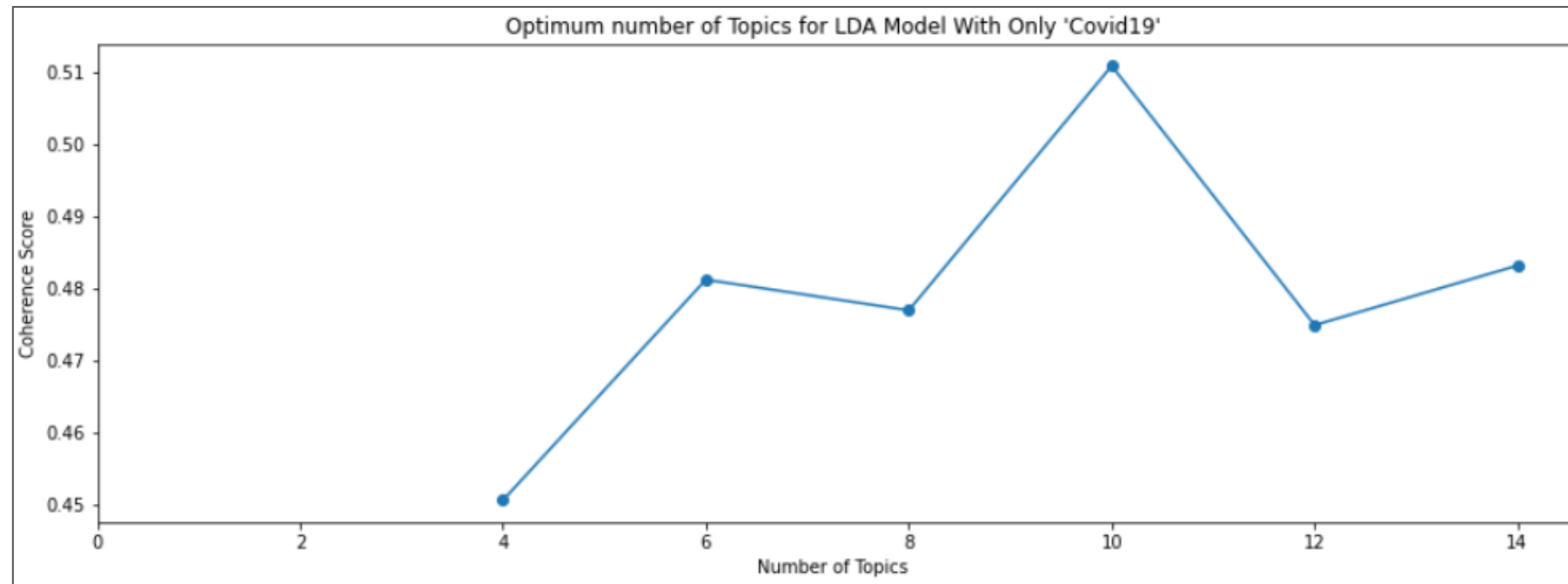
Topics:

- American Politics:**
 - 'Trump' and 'realDonaldTrump'
- Business/Economic Impact:**
 - 'business', 'support', 'impact', 'global' and 'economic'
- International Outlook:**
 - 'case', 'India', 'total', 'confirmed', 'Spain' and 'France'
- Social Media with American Politics:**
 - 'Youtube', 'MAGA', 'IngrahamAngle', 'Joe', 'PressBriefing'
- Random:**
 - 'support', 'last', 'you'

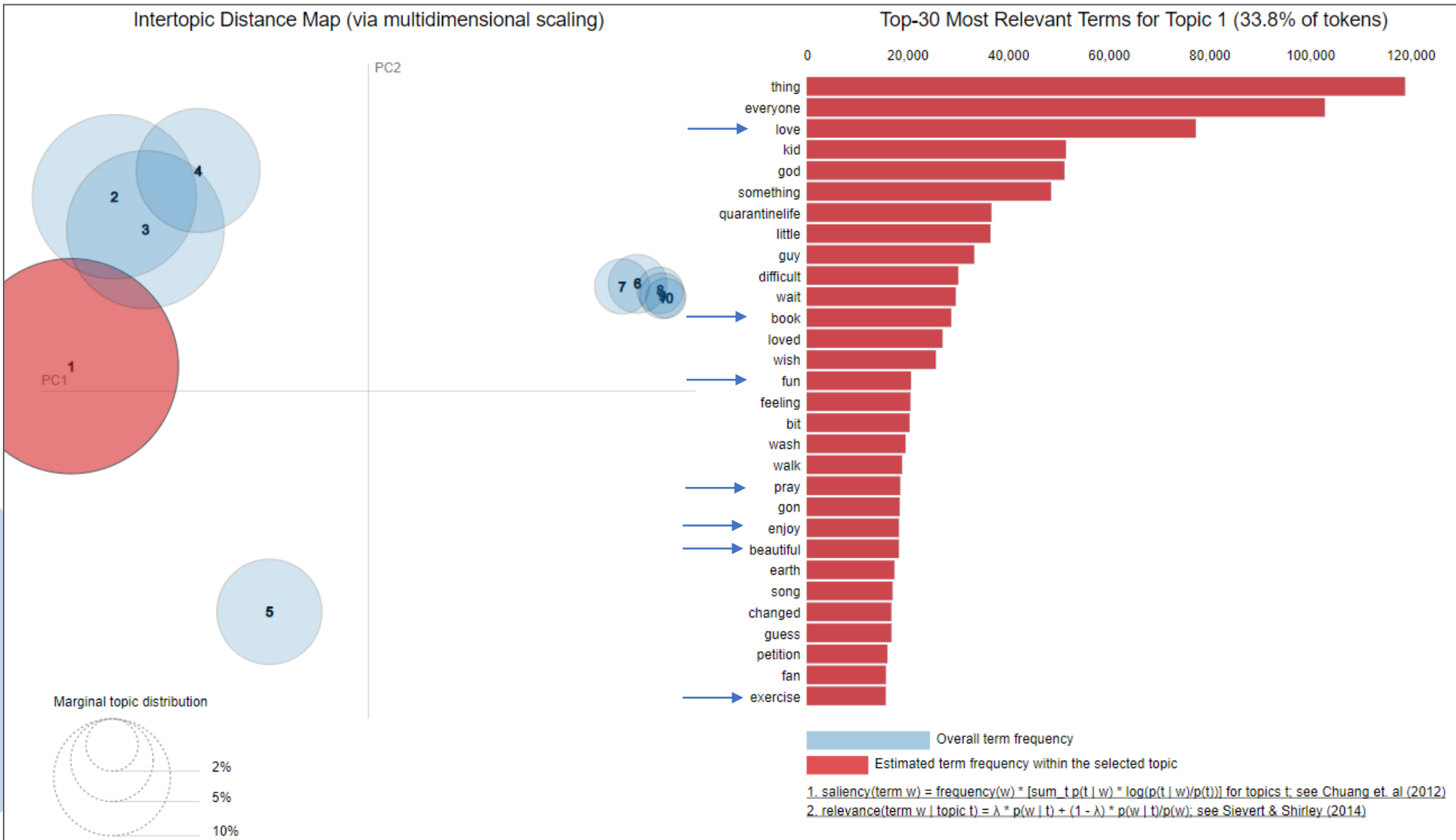
LDA Model Optimization

2 Intermediate Models (2 & 3) Created with Successive Improvements

- Key Preprocessing Changes (Based on Learning):
 - Covid19 Synonym Unification;
 - All tokens turned lowercase;
 - Balance between number of Topics and Coherency Score;



Results: LDA Model 4 - 10 Topics & 51% Coherency



Topics:

- Quarantine Silver Lining:
 - 'loved', 'beautiful', 'feeling', 'enjoy' and 'fun';
- Business and Financial Impacts:
 - 'business', 'impact', 'industry' and 'financial';
- Covid19 Tracking;
- Covid19 Medical Advances;
- American Republican Politics;
- Indian Politics w/ Randomness;
- Democratic and Canadian Politics;
- Random:
 - 'taiwan', 'housing', 'tired', 'togetherathome', 'imf' and 'novadairy';
- Negative Lockdown Aspect;
- American Conspiracy Theories:
 - 'agenda', 'qanon', 'immigration', 'wwg1wga' and 'russian';

Conclusions & Potential Further Work

Conclusions

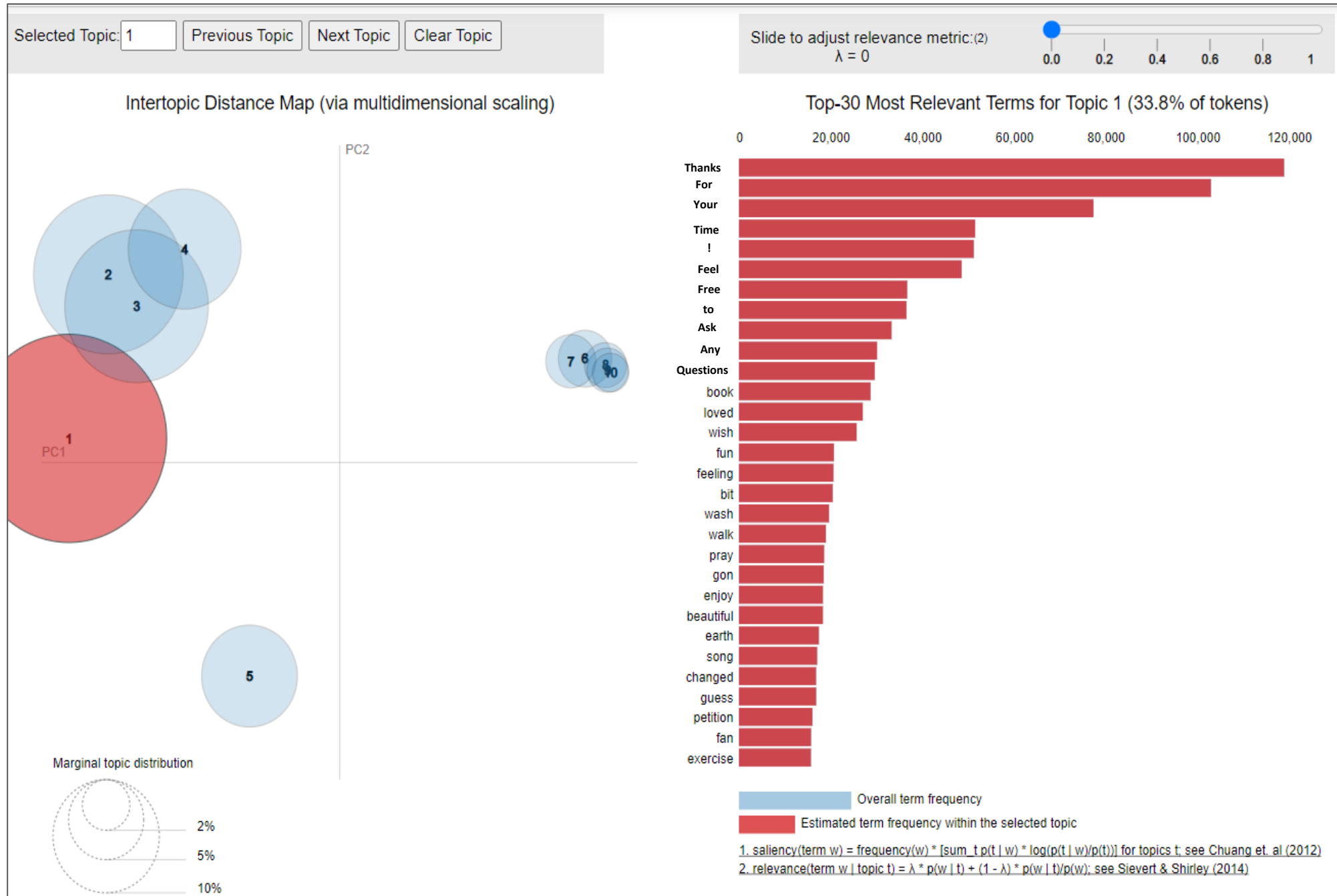
- Goal:
 - To automatically identify the underlying topics being discussed with context to Covid19;
- Generated 4 Latent Dirichlet Allocation Models
 - Preprocessing Changes
 - Data Comprehension
 - Consistent improvements
- Outcome:
 - 9 Coherent Topics and 1 Random Topic

Model Type (name)	Corpus Characteristics	Number of Topics	Coherence Score (%)
LDA Model 1	corpus_LDA_with_Covid19_synonyms	5	~35.0
LDA Model 2	corpus_LDA_with_Covid19_synonyms	20	~46.0
LDA Model 3	corpus_LDA_with_Only_Covid19	5	~47.5
LDA Model 4	corpus_LDA_with_Only_Covid19	10	~51.5

Further Work

- Use different vectorizers:
 - Tf-Idf;
- Bigger proportion of dataset leveraging cloud computing;
- Tweets from other months of Covid19;
- Entity Recognition;

Questions



Backups

Issues with Stemming

Stemming and Lemmatization

```
words = ["connects", "connected", "strange", "is", "am"]
```

```
stemmed = ["connect", "connect", "strang", "is", "am"]
```

```
lemmatized = ["connect", "connect", "strange", "be", "be"]
```

Saliency and Relevance

Saliency is not only a measure of how frequently the term is used but also how important it is in the corpus for detecting a topic.

1. saliency(term w) = frequency(w) * [sum t $p(t | w) * \log(p(t | w)/p(t))$] for topics t ; see Chuang et. al (2012)
2. relevance(term w | topic t) = $\lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$; see Sievert & Shirley (2014)

Increased Characters

```
Out[10]: 179      889
          1475     872
          2448     482
          9465     650
          11917    464
          ...
          79575    749
          79750    921
          80767    485
          80851    858
          81178    917
          Name: text, Length: 76, dtype: int64
```

We will check one of the Tweets below. We have selected index 79750 as it shows a character count of 921.

```
In [11]: analysis_df.iloc[79750]
```

```
Out[11]: Unnamed: 0              79750
          created_at      2020-04-30 07:15:17+00:00
          screen_name      PhillyComptonMW
          text      @laurasessions10 @SusanBordo @Sammysgranny @su...
          is_retweet      False
          favourites_count      31740
          retweet_count      0
          followers_count      1966
          friends_count      1034
          verified      False
          Name: 79750, dtype: object
```

It's apparent that the Tweet contains a high number of mentions, however we are unable to see the full text of the Tweet. Let's isolate that and check the text below.

```
In [12]: analysis_df['text'].iloc[79750]
```

```
Out[12]: '@laurasessions10 @SusanBordo @Sammysgranny @suewashko @bannerite @grammy4lphhl @kjoerwin @SophieInCT @WPalmerCurl @messengerjs @P
uestoLoco @MooPersists @Retinalia @doctordill3 @letat_lechat @NWMouzer @KayTweetTweet @last_person_on_ @morgfair @veedubyoo @m3
3gs @workingtrucker @goprabebuster @PodcastObsessed @BrindaStar @ginadem @Bellarealness @AHamiltonSpirit @co_rapunzel4 @NaphiSo
c @EricWolfson @TheBaxterBean @Pandeism @Bvweir @TheWomensWatch @RonSupportsYou @admiralmpj @linksteroh @tedlieu @SenWarren @ju
dapeters @ColtSTaylor @kalpenn @JohnWesleyShipp @jonfavs @ProudResister @JuliaLikesFrogs @TheDailyShow @crzyfkinworld @Starz_Wa
yne Trump & Republicans: What virus? https://t.co/W1QmC9Lp5F\n#Trump #coronavirus #Georgia \n#TrumpPandemic #GOPGenocide #G
OP #TrumpIsTheWORSTPresidentEVER #TrumpGenocide #Republicans #TrumpliedPeopleDied #TrumpliesAmericansDie #RepublicansAreKilling
Us #MoscowMitch #COVID19 #Wisconsin'
```

Final Topic WordClouds Topics 1 & 2

('Quarantine Silver Lining',)

[illegible]

('Business and Financial Impacts',)

A word cloud of terms related to the COVID-19 pandemic. The words are arranged in a circular pattern, with some words appearing larger than others. The words include: public, community, free, people, global, student, work, fund, care, support, business, today, webinar, effort, time, impact, service, read, food, great, resource, latest, thank, need, health, pandemic, fight, team, healthcare, covid19, worker, response, protect, online, working, learn, please, join, company, take, new, help, government, crisis, share, mask, information, call.

Final Topic WordClouds Topics 3 & 4

('Covid19 Tracking',)

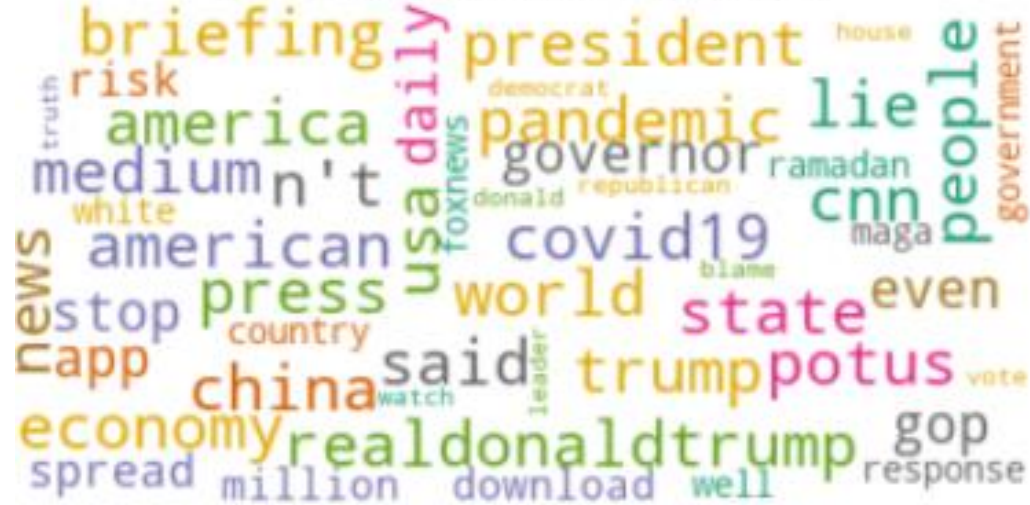


('Covid19 Medical Advances',)



Final Topic WordClouds Topics 5 & 6

('USA Republican Politics',)



('Indian Politics w/ Randomness',)



Final Topic WordClouds Topics 7 & 8

('USA Democratic and Canadian Politics',)



('Random',)



Final Topic WordClouds Topics 9 & 10

('Negative Lockdown Aspects',)



('American Conspiracy Theories',)

