

# **Covid19 Misinformation Trends**

Student Name: Roshan Jangid, Rishabh Gupta

Roll Numbers: MT21031, MT21070

Capstone report submitted in partial fulfillment of the requirements  
for the Degree of M.Tech. in Computer Science & Engineering  
on May 16, 2022

**Capstone Advisor**

Dr. Tavpritesh Sethi

Indraprastha Institute of Information Technology  
New Delhi

## **Abstract**

The COVID19 pandemic, a health care catastrophe that began in December 2019, has placed us all in an unusual scenario where physical and social interaction is strictly limited, and everyone is confined to their houses. The home confined condition has raised the relevance of social media & online support in our daily lives. As social media is a crucial platform for expressing emotions and feelings, it is critical to understand the sort of emotions associated with that tweet and the change in individuals' behavior due to the significant events that occurred during this epidemic.

**Keywords** - *COVID19, Emotions, Twitter, Social media, Pandemic Events*

## **Acknowledgments**

First and foremost, we sincerely thank our Capstone project advisor, Dr. Tavpritesh Sethi, for providing valuable suggestions, guidance, and support throughout the project. We extend our gratitude to Ridam Pal, a Ph.D. student at IIITD, and Sargun Nagpal, a Research Assistant at TavLab, for constantly guiding us throughout the project, helping us, and keeping track of our progress. We would also like to acknowledge IIITD for giving us this opportunity to work on a meaningful project.

## **Work Distribution**

All members have equally contributed to this project.

## **Winter 2022**

During Jan-mar, we analyze requirements, find tools and data sources, and collect all the required data mentioned in chapters 1, 2, & 3. In April-May, we study behavioral indicators and events as shown in chapters 4, 5, and 6.

# Contents

## **1. Introduction**

- 1.1. Related Work

## **2. Data**

- 2.1. Tools for Data Extraction
- 2.2. Source for Dataset Extraction
- 2.3. Dataset Descriptions

## **3. Experiments- Understanding the data**

- 3.1. Tweet Analysis and Exploration
- 3.2. Geographic Trends

## **4. Indicators and Events**

- 4.1. Emotion
  - 4.1.1. Models
  - 4.1.2. Analysis
- 4.2. Events
  - 4.2.1. Identification and Exploration

## **5. Measuring Emotions**

- 5.1. Cross-Correlation between Emotions and behavioral indicators
- 5.2. Change Point analysis

## **6. Summary**

- 6.1. Discussion
- 6.2. Limitations
- 6.3. Future work

## **References**

# Chapter 1

## Introduction

COVID19 rapidly spread worldwide, starting from Wuhan, China, from December 2019 onwards. WHO announced COVID19 as a pandemic and forced multiple countries to strict lockdown, where people were constrained to remain in their homes. This lockdown resulted in more usage of social media platforms where people started expressing their concerns and COVID stories and also shared medical advice with those who tested positive. Social media platforms played a significant role during this pandemic, and at the same time, it has also affected people psychologically.

We focused on Twitter as our social media platform, a network of people where they can express their views, ideas, and opinions. In particular, we focused on COVID19-related misinformation, where false belief is shared among singles, groups, or in the community related to COVID19.

Multiple observations were carried out like which tweets get retweeted the most, how legitimate the information is, and what are the emotions associated with those tweets. This gives rise to the motivation of research questions like, What type of COVID topics are people talking about, What emotions are associated with the tweets related to COVID-19 misinformation, Changes caused by major pandemic events, and how are the emotions associated with these changes. We build our dataset based upon panacea lab's data tweets id and fetching complete data using these ids by Twitter APIs; we further collect covid19 events and behavioral indicators data. We comprehensively described each piece of data and documented it summarized in this report for future references.

Further, We Studied Geographic and popularity trends of tweets on Twitter to get an overall data analysis. To do depth analysis of user behavior, we did emotions analysis using empath and plutchik and compared both the models; we map covid19 events with plutchik emotions strength monthly and on 15 days moving average daily basis, to support events emotions analysis, we do change point analysis which gives us a clear picture of peaks in time series emotions strength data. Next, we do a cross-correlation analysis of emotional strength with behavioral indicator data.

Our analysis and study will better understand users' emotions and covid19 events' effects on users' emotions.

## 1.1 Related Work

Our contributions are based on already existing work to identify misinformation spread on social media(especially in tweets), where the author discusses the demographics of persons who actively participate in the distribution of COVID misinformation and claims to disprove those disseminating misinformation. The user's tweets were evaluated, and for each tweet that fell into the informed group, the user was awarded a +1 valence. In comparison, tweets that fell into the uninformed category were given a -1 valence. The total valence of a user would be the sum of the valences of each tweet. If the final valence is more than zero, the user is placed in the informed category. If the final valence is less than zero, the user is classified as uninformed. Network analysis is carried out, and it is observed that misinformed sub-communities are denser than informed communities.

Vaccine Beliefs on Twitter with Lexical Embeddings talk about how the changes in an emotional category are reflected when Vaccine Rollout, Misinformation, Health Effects, and Inequities are performed. The experiment is carried out by extracting the corpus from Twitter posts related to COVID-19 vaccination for countries like India, the United States of America, Australia, the UK, and Brazil. A community detection algorithm is used to find positive correlation networks. According to the data, tweets expressing doubt regarding vaccinations contain the most references to health-related impacts. Results mentioned by the paper include the change in linear trends of hesitation and contentment category before. After vaccination, the importance of negative emotions in the alluvial diagram and the link between Emotions and Contributing Factors were discovered to differ among countries.

## Chapter 2

# Data

### 2.1 Tools for Data Extraction

#### 1. Twint

It is a python based free tool for extracting data from Twitter without using Twitter API.

**Attributes returned by Twint** - These are the information we get while fetching tweets using Twint.

1. Id - Tweet unique id for differentiating each tweet
2. conversation\_id - It is the Tweet id of the original tweet, it will show reply threads on the tweets i.e all reply threads in the tweet have the same conversation id
3. created\_at,- it represents data and time when the tweet is first created in UTC format
4. date - represent the date of the tweet created
5. timezone, - Timezone of the place from where the tweet is posted
6. user\_id, - unique id for every user assigned,
7. username, - username of the user used in Twitter account
8. name, - Personal name of the user.
9. tweets, - Original tweet text posted by the user
10. link, - Tweet link
11. retweet - boolean attributes which specified whether the tweet is retweeted or not
12. nreplies - specify the number of replies counts on the tweet
13. Nretweet - specify the number of retweet counts on the tweet
14. likes\_count - specify the number of likes counts on tweets
15. hashtags - give hashtags mentioned in the tweet
16. quote url - Url link which is quoted in the tweet.
17. near - it specifies the radius within the location to limit the collection data within that radius
18. geo(lat, long, radius), - it specifies geolocation features of the tweet like latitude, longitude, and radius within the tweet is fetched
19. reply to - original tweet in which the represented tweet is replied
20. retweet\_date - the date when the tweet is retweeted if the retweet is true.

**Some important points for this tool -**

1. We can obtain location-based tweets(hardcoded) with timelines, but ISP ban after a specific limit, so it is tough to collect extensive data using it.
2. Language is not filtering correctly (sometimes it works and sometimes not).

**Different Fetching techniques for tweets using Twint-**

1. Using hashtags/ keywords (can be used along with username)
2. Using followers & followings
3. Using links

## **2. Snsrape**

It is a python scraping tool to scrape data from various social media sites like Facebook, Instagram, Twitter, and Reddit.

**Attributes returned by Snsraspe -**

1. url- Twitter Tweets URL
2. date- date and time of the tweets when it is first created
3. content - the full text of the tweet
4. id - Unique id of the tweet
5. username - username of the user-posted the tweet.

**Limitations of This tool -**

1. Does not get the location with the tweets in the data
2. we have to define a location to fetch particular tweets under that location - can define a city with a radius predefined and coordinates
3. We could fetch 500 tweets in 2 minutes location defined - “Delhi” as a city and radius of 1000km - got all tweets with all languages (mainly English and Hindi) - we need to filter language.
4. We were able to fetch 50000 tweets in 34 minutes.

## **3. Twitter API**

Twitter API is the most efficient and official twitter tool to extract data from Twitter, but it comes with some restrictions. Twitter divided its data extraction account into four-level, namely, Essential, Elevated, Elevated+, and Managed. There are two versions of Twitter API v1 and v2. Using the Essential account, we can fetch 500K tweets per month, one application environment per account, and support for v2 is only available with this account. In Elevated access, we can fetch 2M tweets per month, and three app

environments and support for v1 and v2, Essential and Elevated accounts are freely available. In the Elevated+ access, we can fetch 10M tweets per month, three app environments, and support for v1 and v2 APIs, but this account is not entirely free and made available to researchers across different countries on special request by filling out a google form, Managed is paid version of Twitter API, and it comes with dedicated account management and support by Twitter it is enterprise Level API access account.

When we request data from Twitter, it returns a Twitter object, User Object, which consists of different information about tweets posted and information about the user who interacted with a tweet.

### **Attributes returned by Tweet Object -**

1. `created_at` - UTC when this Tweet was created
2. `id/id_str` - unique identifier for the tweet, `id_str` represent id greater than 53 bits
3. `text` - Tweet text posted by the user
4. `source` - device and platform source by which tweet is posted, e.g., Twitter, Web client
5. `truncated` - true if the tweet text length is more significant than 140 characters
6. `in_reply_to_status_id/ in_reply_to_status_id_str` - if the representative tweet is replying than this will contain the id of the original tweet
7. `in_reply_to_user_id/in_reply_to_user_id_str` - if the representative tweet is a reply, then it will contain the user id of the original Tweet author profile
8. `In_reply_to_screen_name` - if the representative tweet is a reply, then it will contain the screen name of the original Tweet author profile.
9. `coordinates` - represents the point type of geographic location of the tweet reported by a user or client application, in the format of longitude and latitude, respectively
10. `Place` - represent polygon type of geographic location with country code, country name, country full name, country code, and Place type like city/country.
11. `quoted_status_id/ quoted_status_id_str` - this represents the quote id of the tweet is the representative tweet is quoted
12. `is_quote_status` - True if the tweet is quoted; else, False
13. `quoted_status` - contains Tweet Object of the original tweet if `is_quote_status` is true.
14. `retweeted_status` - contains a representation of the *original* Tweet that was retweeted.
15. `quote_count` - Indicates approximately how many times this Tweet has been quoted by Twitter users.
16. `reply_count` - represents the number of times this Tweet has been replied to.



17. `retweet_count` - represents the number of times this Tweet has been retweeted.
18. `favorite_count` - it indicates how many times the user likes the tweet.
19. `entities` - contains URLs, hashtags, user mentions, media, symbols, and polls of the tweet
20. `extended_entities` - contains media URL when a tweet contains more than one media
21. `lang` - indicate the language of the tweet
22. `User` - The user object contains all the information about the user like username, user id, follower count, friends count, favorite count, statuses count(total tweets), is the user verified or not, location of the user, Bio of the user, date when the account is created, profile, and background profile banner.

## **4. Tweepy**

It is a python library to access tweets using Twitter API. It is open-source and easy to use to access data from Twitter. It has built-in classes and methods/functions for the Twitter Tweet Model and handles data encoding and decoding.

## **2.2 Sources for Data Extraction -**

### **1. Twitter Dataset - Panacea Lab dataset:-**

This dataset consists of the tweet id of the tweets related to covid-19. It uses seven keywords(COVID19, CoronavirusPandemic, COVID-19, 2019nCoV, CoronaOutbreak, Coronavirus, WuhanVirus) to fetch covid-19 tweets from Twitter with tweet id; in addition to tweet id, it also contains language of the tweets. The timeline of the dataset available is Jan 2020 to April 2022. Dataset count from Jan 2020 to 11th March 2020 is low but has high counts of tweet after 11th march 2020. This dataset only consists of tweet ids and language to fetch entire attributes, text, and location and needs to hydrate tweets using Twitter API. The dataset is updated regularly.

### **2. Our World in data:-**

This dataset consists of Behavioural indicators of covid-19 like No. of COVID-19 Tests, Vaccines, Stringency index, No. of ICU patients, No. of hospitalized patients, Cases, and Deaths. The dataset was updated daily until March 2022, and now weekly updates are available.

### 3. Covid-19 popular events dataset:-

Think Global Health dataset for Covid-19 events and major headlines happening worldwide, this data is updated frequently, and the data is available from January 2020, the beginning of the covid-19 pandemic.

Business Standard Covid-19 timeline events since the lockdown proposed in India. It covers major Covid-19 events of India only and their timeline.

## 2.3 DataSet Description:-

### 1. Twitter Dataset:-



We hydrated/fetched ~80k tweets per month from Jan 2020 to Dec 2021. Total data is 2.2 million with location and English language. In the below tables, we fetch some statistics about data like how many tweets contain media, as media in tweets is vastly used to spread misinformation. We fetched around 10k tweets per day from Jan 2020 to Dec 2021, totaling around 7.2 million tweets, of which 2.2 million contain locations.

Duration	Jan - Dec, 2021	Jan - Dec 2020
Total Tweets collected	~ 1 million	~ 1.2 million
Total Tweets with media content	~ 0.17 million (17%)	~0.17 million(13%)
Tweets per month	70 - 80k	1 Lakh
Number of countries	46	46

## Metadata about the User:-

Hashtags, URLs, reply\_count, user\_followers, user\_friends\_count, user\_statuses\_count

1. In tweet Object, there is a column quoted status which contains Tweet object of the quoted tweets we created quoted\_status dataset with language filter English and location might not be in the data this will be fetched around 1 million extra tweets. Below are some of the quote tweets mentioned -

Inauthentic christian Talibangelical money lender in the Temple says what???? 🤔🤔🤔🤔 <a href="https://t.co/b2ynIFg8qz">https://t.co/b2ynIFg8qz</a>
FINALLY the experts have spoken... @GotabayaR over to you Sir ! <a href="https://t.co/vUHgK43LiM">https://t.co/vUHgK43LiM</a>
Is the bar for an Order of Ontario that low that they can't find good candidates?! <a href="https://t.co/VXYLu0KxFz">https://t.co/VXYLu0KxFz</a>
Nope just a cold 🤧 <a href="https://t.co/f9GdWRc9Xv">https://t.co/f9GdWRc9Xv</a>

## 2. Replies / Comments:

- 12% of the data.
- A comment that contains a COVID-19 keyword. The original post may or may not be on COVID.
- The ID of the original post in API response.

Samples for replied tweets:-

@MiddleMolly @gvravel @andyllassner @gtconway3d You already have sheriffs and mayors defying #COVID19 rules even against GOP guys like DeWine and Scott never mind Dems 2/2 <https://t.co/lz1v0v1pty>

So many I know gone due to COVID-19 today another family friend, and the oldest yet 92 at home with freakin COVID-19 gee whiz!! <https://t.co/niNk5oR528>

#BlackLivesMatter is the largest Social Justice protest not just here in the USA but it United the World.  
#COVID19 as deathly has amplifies the social economic inequality not just here in USA but worldwide.  
#MedicareForAll  
#SocialJustice  
#IncomeInequality

1. Tweets sharing an external article
  - Original tweet may not have enough context.
  - ~9% tweets: No keyword in a tweet, but only in the shared URL.

Sample:-

Oh look. The UK government ignoring science and gambling with lives.  
Who'd have thought it <https://t.co/3DnvmivDJe>  
<https://www.nytimes.com/2021/01/01/health/coronavirus-vaccines-britain.html>

2. Retweets contribute 0.2% of the data. Truncated Tweets < 140 chars were truncated. Full text retrieved for truncated tweets

## 2. Our World In data:-

Only countries with specific conditions are taken into account:

1. populations greater than 5 million
2. countries in which at least 21 days had passed since the 100th confirmed case

Important columns -These are the attributes we are relating with misinformation available in datasets:-

Attributes	Available for - No. of countries
Total_cases, new cases	228
total_deaths,new_deaths	219
Icu_patients	36
Hosp_patients	38
stringency_index	182
reproduction_rate	188
Total_tests	144
new_tests	137
positive_rate - rolling 7-days averaged, tests_per_case	142
total_vaccinations	213
People_vaccinated, People_fully_vaccinated	231
total_boosters	151
new_vaccinations	185
Hospital_beds_per_thousand	137

In rare cases where our source for confirmed cases & deaths reports a negative daily change due to a data correction, we set the corresponding metric (new\_cases or new\_deaths) to NA. This also means that rolling metrics (7-day rolling average, weekly rolling sum, biweekly rolling sum) are set to NA until this missing value leaves the rolling window.

Stringency\_index indicates how strict the government is based on nine indicators on a scale of 0 to 100. (school closures, workplace closures, cancellation of public events, restrictions on public gatherings, closures of public transport, stay-at-home requirements, public information campaigns, restrictions on internal movements, and international travel controls), and this data is available for 182 countries.

The location also consists of the word "world," which means the date is for all over the world, continents, and the Islands.

"excess\_mortality" indicates the % difference between projection deaths and actual deaths.

As datasets are country-wise, we have to analyze data for each country. For each column, data is available for some countries and not for others (which is none in this analysis datasets).

We find the start\_date(starting date from which data is available), end\_date(up to when data is available), and no\_of\_days(for how many days between start and end date data is available ) for each row values which are not null mean (the row and corresponding column which indicate None/null value means that column value for that particular country is not available in dataset while for other countries it is available ), datasets has some negative values also(to decide either drop or replace techniques should use).

### **Sample for Country -**

**India** - total\_cases are available from 30-01-2020 upto 28-02-2022 for 761 no. of days

attribute	date
total_cases_start_date	2020-01-30
total_cases_End_date	2022-02-28
total_cases_no_of_days	761

**Afghanistan** - new vaccinations are available from 27-05-2021 to 27-01-2022 for only 3 days.

attribute	date
new_vaccinations_start_date	2021-05-27
new_vaccinations_end_date	2022-01-27
new_vaccinations_no_of_days	3

### 3. Events dataset:-

Web Scraping a list of ~1.6k COVID-19 related events across the globe with their dates and used NER to extract the associated location names.

#### **Approach:**

Tried to identify the location name using some predefined locations list available to us.

#### **Challenges -**

- City names were missing
- UN/ United Nations, United Nations are not captured properly
- Some locations might be missing due to non-availability

Using Spacy library to overcome the above problems, we find city/country/state associated with events using named entities.

## Chapter 3

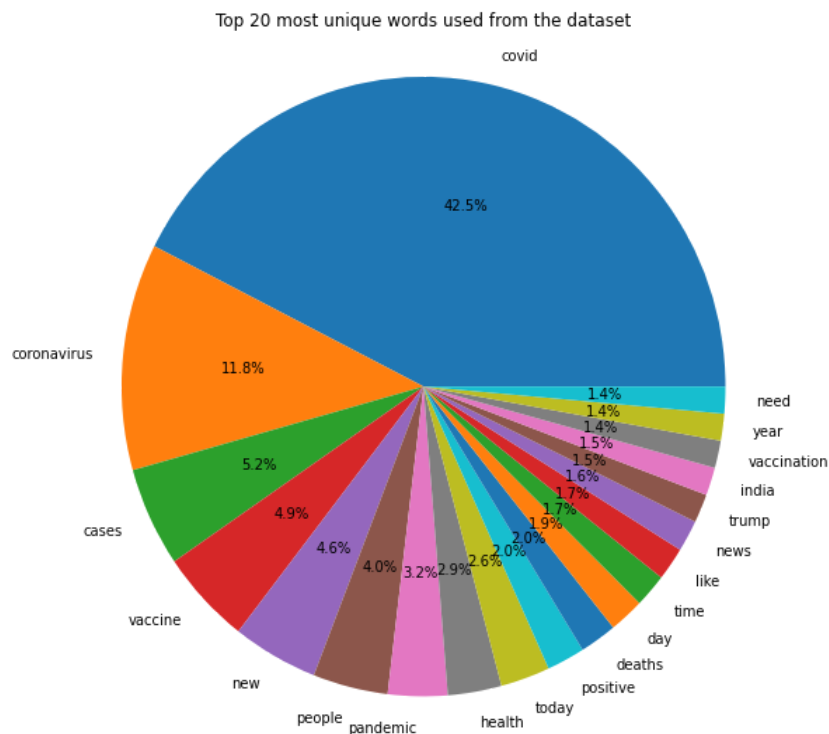
# Experiments- Understanding the data

### 3.1 Tweet Analysis and Exploration-

In the context of observing the content of the tweets, we are analyzing;

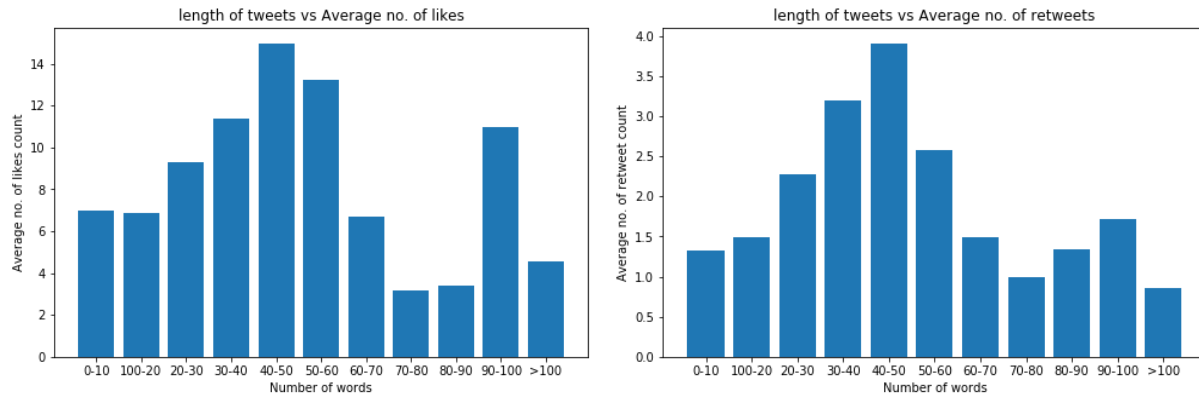
1. How the distribution of the unique words are spread across all the text,
2. What is the distribution of popularity measurement based on the length of words?

Tweets are a massive variety of words in which words of all categories are included directly or indirectly, and identifying the behavior of those tweets is not handy all the time. We observe only some of the most popular tweet words to find which type of words are more frequently used than others during this covid-19 pandemic and come up with the top 20 most frequently used unique words from our dataset. The pie chart below depicts how the words like covid, coronavirus, cases, vaccine, and pandemic cover more than 75% of the tweets, which shows the relevance of our dataset. The words like people, deaths, covid times, trump, and India are some of the hot words used most frequently for discussion related to covid.





Below is the popularity measurement like average likes count vs. the number of words in a tweet and the average number of retweet count vs. the number of words in a tweet. The histograms below show that the average length of tweets between 40-50 gets the highest number of average counts for both likes and retweets. And least for the range of words between 70-90, or more than 100.

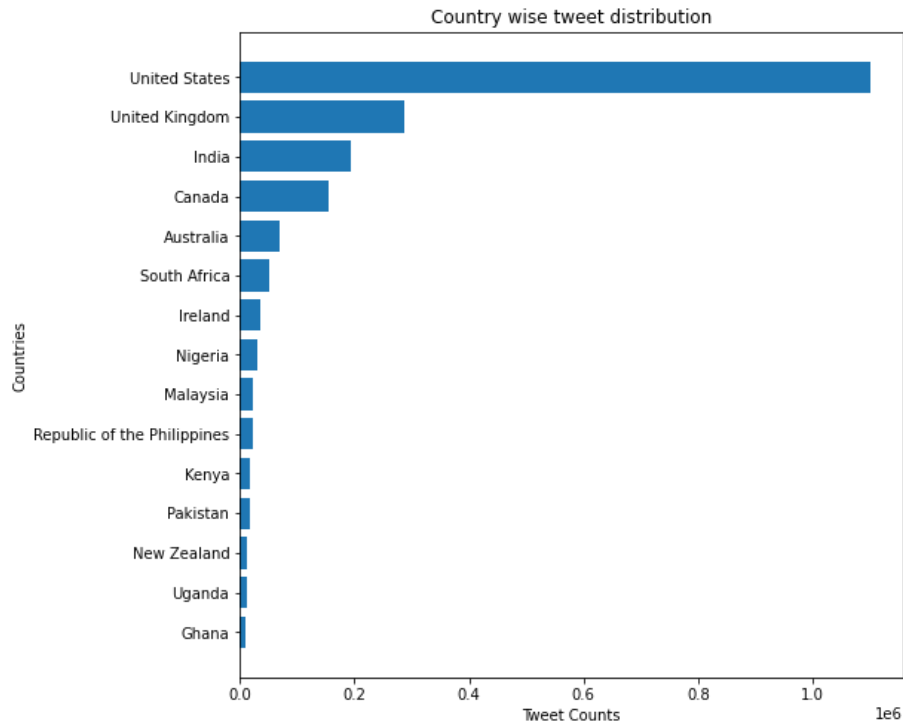


## 3.2 Geographic Trends

We create basic plots to understand how the tweet data is represented. The following questions are answered in this section:

1. What is the distribution of the tweets among different geographic regions?
2. Which hashtags are used the most during the COVID timeline for Geographic trends of the United States, India, and the United Kingdom?

The below bar plot shows how the geographic distribution of posts on Twitter takes place for different countries, where the United States has the most frequent tweets and contributes more than 50% of our dataset(which is around 11 lakhs tweets), followed by the United Kingdom, India, Canada, and Australia(naming a few topmost countries).



x-axis represents the Tweet Counts with values ranging from 0 to 12 lakhs.  
y-axis represents the Country wise distribution of posts on Twitter.

We observe the data according to geographic location and analyze which hashtags are most frequently used in countries like the United States, India, and United Kingdom. The bar plot for these countries is given in fig 3.2.1, fig 3.2.2, & fig 3.2.3. In all the three graphs given below, the y-axis represents the most frequent hashtags used in the tweets of that country, and the x represents the frequency of the hashtag used.

The first figure represents the hashtag frequency for the United States. The words like pandemic, vaccine, trump, trumpvirus, breaking, bidenharris2020 are some of the hot topics for discussion related to misinformation. The second and third figure represents the hashtag frequency for India and UK, where words like "coronavirusupdate," "nifty," "unitingpeoplewiththepossibilities," "news," and "community," etc. are all relatable words that can be used for misinformation.

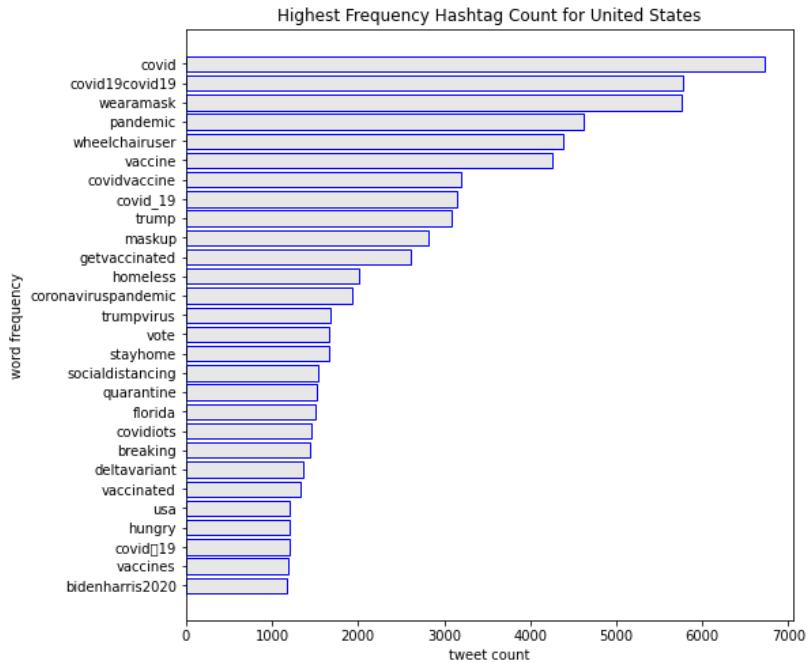


Fig 3.2.1 :- hashtags counts for United states

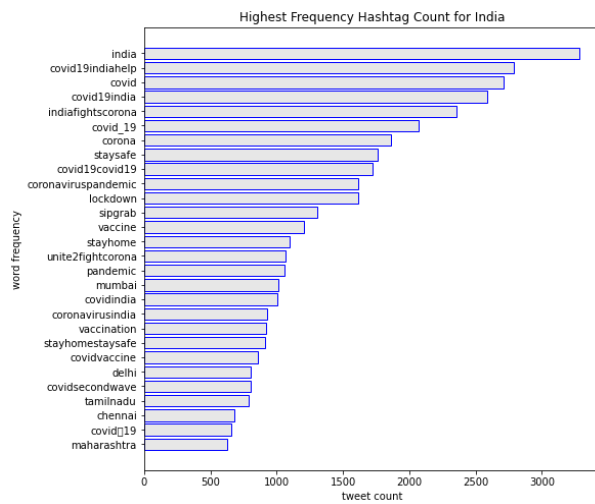


Fig 3.2.2 :- hashtags counts for India

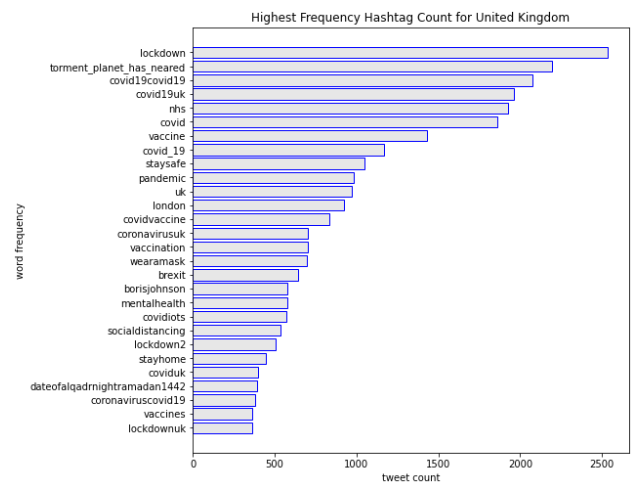


Fig 3.2.3 :- hashtags counts for United Kingdom

## Chapter 4

# Indicators and Events

In this chapter we analyze emotions strength with covid19 events, we applied to methodology to find emotions in the tweets first is empath based and another is plutchik based. Later in this chapter we analyzed covid19 events with Emotions strength.

1. How do misinformation impact people's emotions and behavior? Useful information for policy-makers.
2. Knowing what kind of misinformation follows a pandemic event can help design strategies to counter it.  
For example, What happens on social media when there is news of a new SARS-CoV-2 variant?

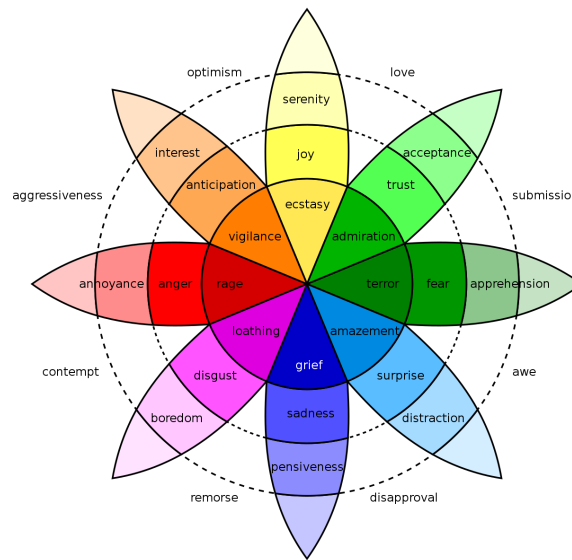
### 4.1 Emotion

Emotions are mental states produced by neurophysiological changes related to varied ideas, sensations, ability to respond, and a degree of pleasure or dissatisfaction.



#### 4.1.1 Models

**Plutchik's Wheel of Emotions:** In one of the basic emotion models, it is observed that some emotions are more essential than others. Plutchik's wheel of emotions serves as the foundation for our multidimensional emotion classification. This taxonomy, which has been in use, seeks to define human emotions as a mixture of four categories. They are represented in duality as **Joy - Sadness, Anger - Fear, Trust - Disgust, and Surprise - Anticipation**. All classification scores for the Fear, Disgust, Joy, Anger, and Sadness emotions are present in Plutchik's wheel given in the figure below.



Plutchik's wheel of emotions (Plutchik 1979).

**Empath:** It is a tool that can develop and evaluate new lexical categories on request based on a limited collection of seed words (such as "bleed" and "punch" to produce class violence). There are 197 words in empath for emotions which can be used to identify the emotions associated with the tweets. However, this model is not found to be appropriate for our dataset. The brief for the same is provided in the below section.

#### 4.1.2 Analysis

1. **Empath-based:** Using the Empath Reddit model on our extracted tweet data, the number of tweets classified was only 18%, as it contains only some of the root emotions.

Trump says up to 100,000 Americans may die from coronavirus	Medical_emergency, Kill, death
---	--------------------------------

2. **Plutchik Transformer:** Using this model, multi-label classification was possible and helpful to classify more tweets, along with the availability of the probability for each emotion in the parameter of strength. This model can classify 99.35% of the tweets on our dataset, but some tweets are problematic for humans to verify for that particular emotion.

Anyone else feel like half the people talking about the coronavirus have no idea what a virus is?	Anger, fear, disgust
---	----------------------

By comparing Plutchik and Empath models, it was found that there is only a 0.6% overlap between both the results (where the result of both the emotion model matches the same tweet). Also, among the 0.6% similarity, there is less percentage of tweets showing similarity for disgust emotion (almost 64% coverage), and most for the surprise emotion (94% coverage) going from left to right.

disgust	trust	anger	sadness	anticipation	fear	joy	surprise
64	53	62	64	64	71	71	94

As a concluding factor, we are using the Plutchik wheel (because of high coverage for unknown tweets) for our dataset to identify the emotions held in each of the tweets.

## 4.2 Events

Events have had a profound impact on how the Coronavirus pandemic has been viewed and understood worldwide. As possible vectors for the virus's spread have had a tremendous impact on the events industry in India and worldwide. Web scraping is carried out to extract significant events for our objective, and the website linked <https://www.thinkglobalhealth.org/article/updated-timeline-coronavirus> is referred to, where all the major events all around the world have been taken into account. This event's timeline is from January 2020 to December 2021 (the Time frame for two years of a pandemic). Below are a few of the samples depicting the events in India that have taken place.

30-01-2020	India confirms first first case.
12-03-2020	India records its first coronavirus death.
19-03-2020	The Indian governmentÂ bans exportsÂ of masks, ventilators, as well as certain medications and supplements.
08-06-2020	India lifts lockdown restrictions, despite fears of a surge.
20-10-2020	India reports fewer than 50,000 new daily COVID-19 cases for the first time in three months.
04-03-2021	India's Covaxin COVID-19 vaccine is shown to be 81 percent effective in new study.
19-05-2021	India sets global record for daily COVID-19 deaths, surpassing the U.S. record.

#### 4.3.1. Identification and Exploration

In this section, we will try to answer the question like;

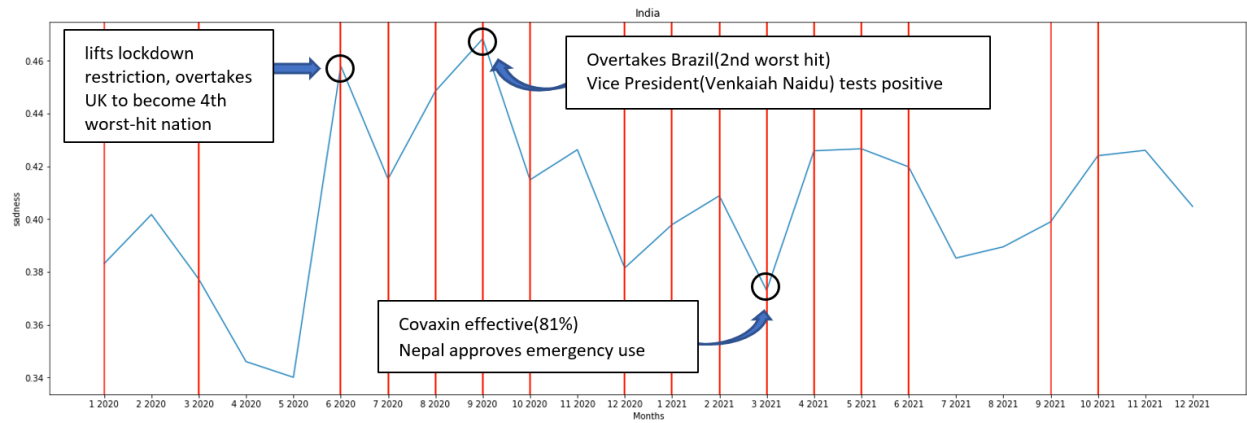
1. How emotion changes due to significant pandemic events
  - a. How does the news of a new variant affect misinformation on different topics?
2. Monthly and daily time series of Emotion strengths?

The study of different emotions associated with tweets is observed against the time series to know the change in emotional strength when an event takes place. The below graph shows how the strength of emotions like fear, sadness, trust, and anger is considered over two years in India. Few significant events are annotated in the graph representing the event in that month. The x-axis represents the time series ranging from January 2020 to December 2021, and the y-axis represents the strength of emotion.

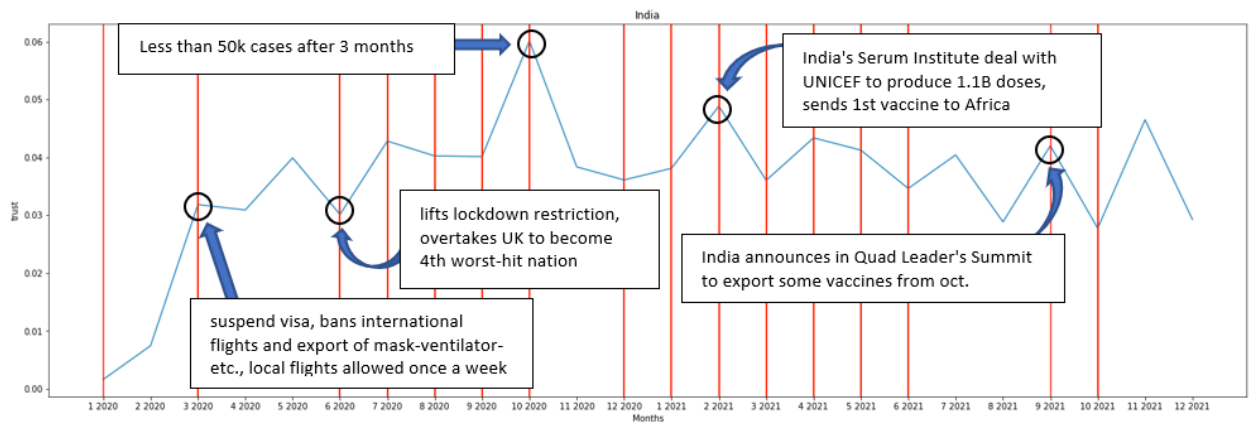
By looking at the figure(4.2.1), the pandemic event by observing the “1st covid case in India” shows how the strength of fear is at its peak in January 2020. Little spikes in the emotion of fear were observed when the lockdown restriction was lifted in June 2020 or when a new covid strain patient was discovered in India in December 2020. Similarly, the change in emotional strength is observed for other figures.



figure(4.2.1) fear vs time series

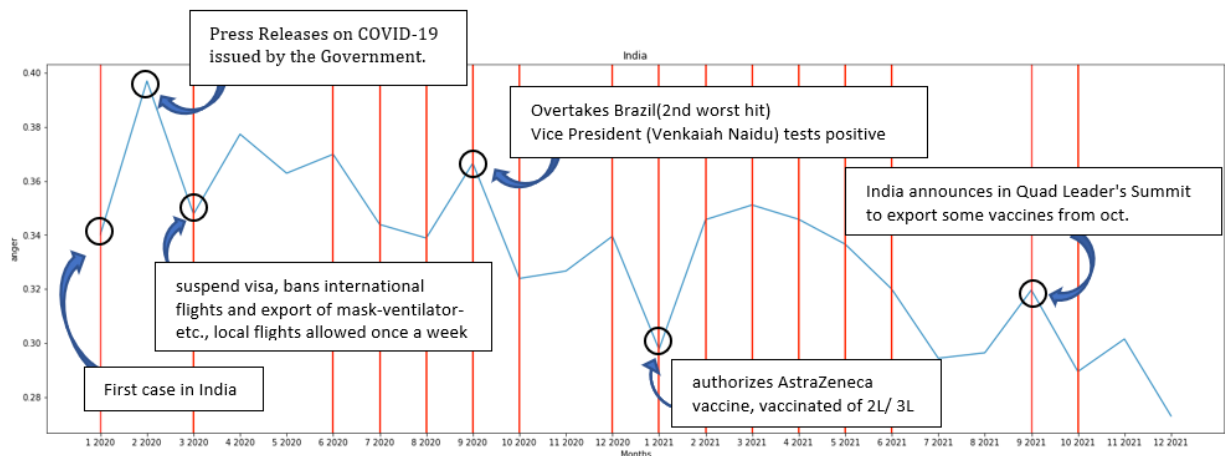


figure(4.2.2) sadness vs time series



figure(4.2.3) trust vs time series

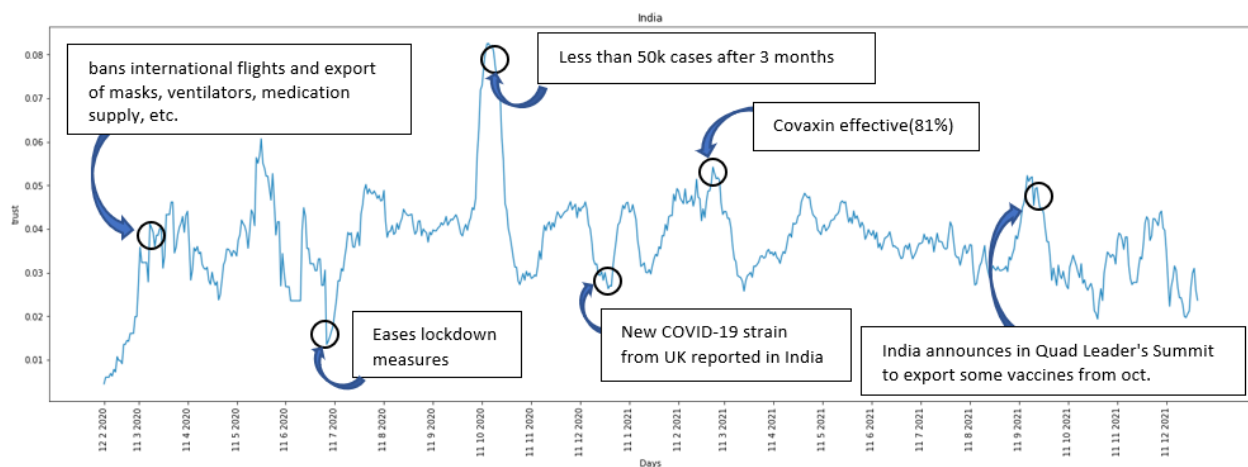




figure(4.2.4) anger vs time series

Along with the month-wise observations of emotion vs. time series, we have carried out how the graph looks for the event if we consider the time series day-wise. 15-days smoothing is applied to the emotional strength to eliminate the noise.

The figure below depicts how the emotional strength “trust” changes along with the events when the time series is plotted day-wise.



trust vs time series(15 days smoothing)

## Chapter 5

# Measuring Emotions

To measure Emotions and prove our above results, we have done two analyses first is a cross-correlation between emotions, namely, fear, anger, sadness, trust, anticipation, joy, disgust, and surprise, with behavioral indicators of new cases, new tests, new vaccination, stringency index, hospitalized and ICU patients, and favorable rate. The second one is change point analysis on the above event-emotions graph.

### 5.1 Cross Correlation between Emotions and behavioral indicators

1. Cross-correlation is the method to find relativeness between two-time series data. It measures the correlation between a time series and lagged version of another time series. More importantly, it tells us whether one-time series data can predict another time series. Here we measure the daily emotional strength of India with behavioral indicators like new cases, deaths, testing, vaccination, positive rate, and stringency index, so we can detect whether fear in people can be a leading signal for these behavioral indicators or not?, we use 'R' language for this analysis; specifically we use CCF(cross-correlation function ) of R, Below are the graph for same and inferences from the graph:-
2. **Stringency index vs. fear** - Section 2.3 talks about the Stringency index. It represents how strict the government was based on nine factors: school closure, international travel restrictions, and other lockdown measures. Here we plot a graph (fig. 5.1.1) between stringency index vs. fear; it shows the leading negative correlation of stringency index with fear, which indicates that as fear increases, the government will take more measures like lockdowns and other restrictions lead to a feeling of safety amongst the population.

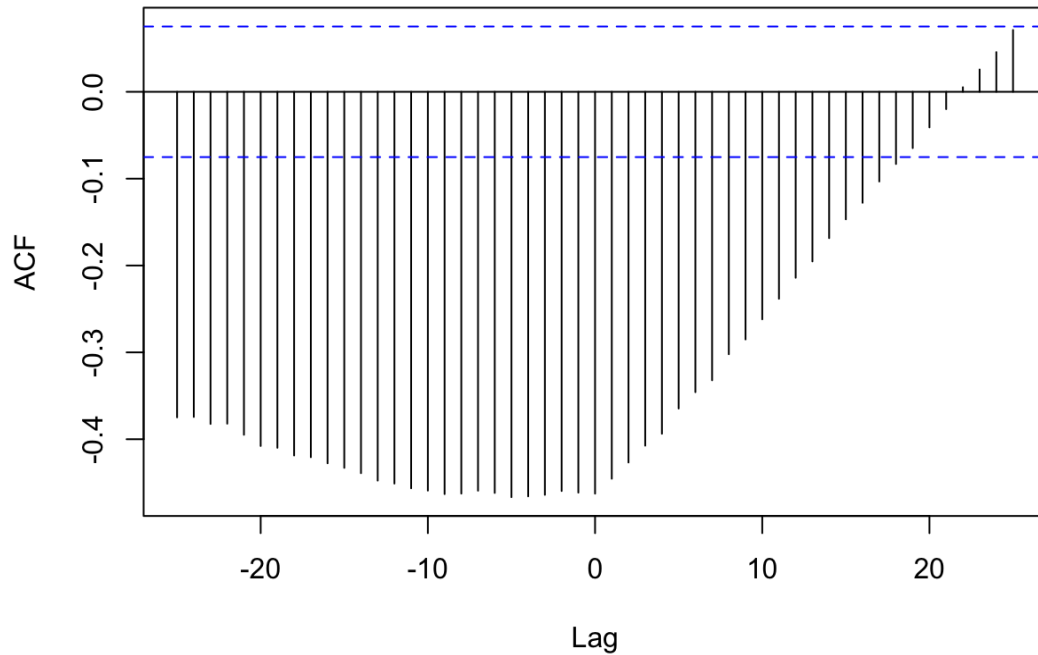


Fig 5.1.1:- Stringency index vs fear

3. **Daily new vaccination vs fear** - In this graph we plotted cross correlation graph between fear and daily new testing, fear has negative correlation with stringency index that is more fear less people going for testing. Which is also expected behaviour amongst population.

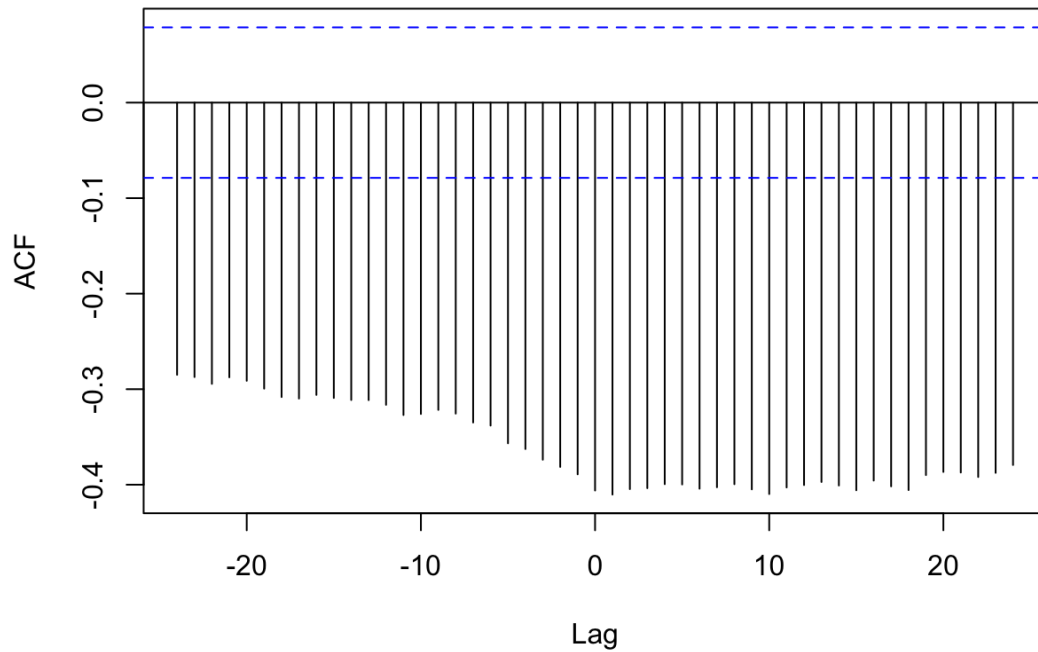


Fig 5.1.2 : - Daily testing vs fear

4. **Daily new vaccination vs. fear** - We plotted a cross-correlation graph between fear and daily new testing. Fear has a negative correlation with the stringency index, which is more fear-less people going for testing, which is also expected behavior.

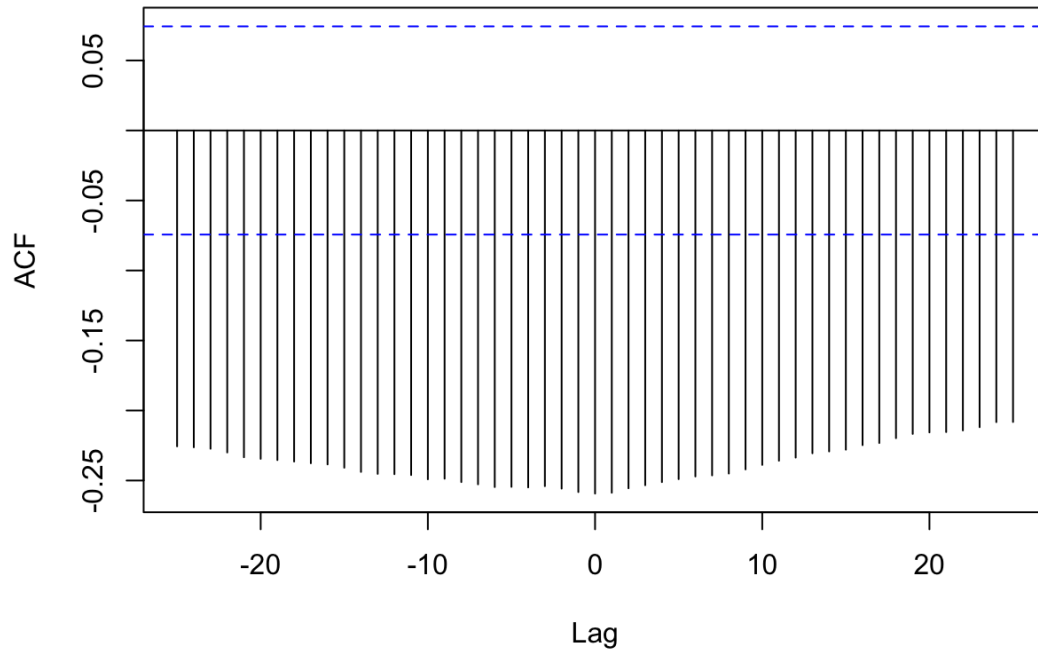


Fig 5.2.3 :- Daily cases vs fear

## 5.2 Change Point analysis

In chapter 4 we performed events emotions analysis, There are various graphs where we are indicating major events whenever there is high or low peak in emotions strength, to back this analysis we had done change point analysis with same data , this analysis determine, did the changes really occurred on the peaks, and where in the graph point the change occur and whether there is one or more changes how confidently we can say these are the real changes, for measuring these functionalities we measure mean, variance and mean variance change points on the graphs. We perform this analysis on daily 15 day moving average events- emotions graph.

1. Fig 5.2.1 shows 15 day moving average change point analysis of Emotions “trust”
2. Fig 5.2.2 shows 15 day moving average change point analysis of Emotions “Anger”.

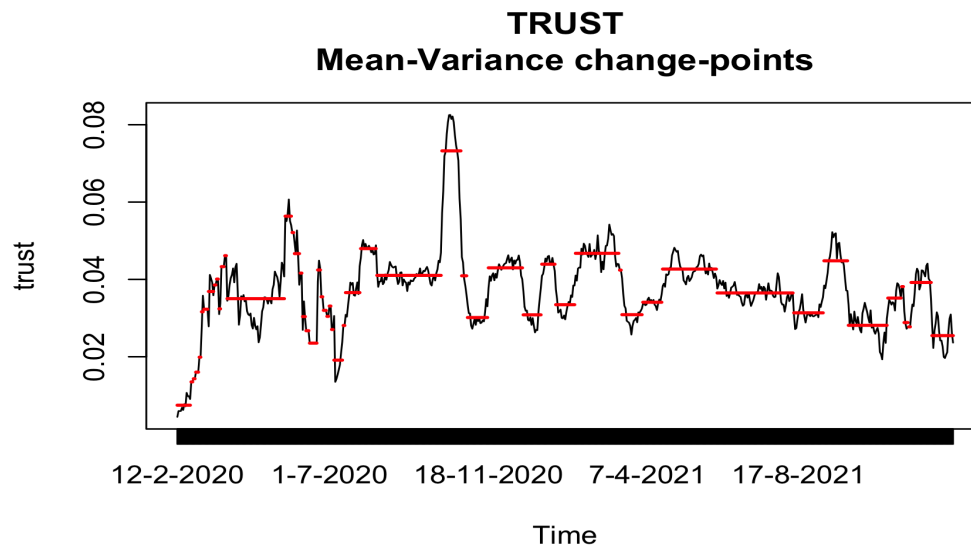


Fig 5.2.1 :- Trust change point analysis

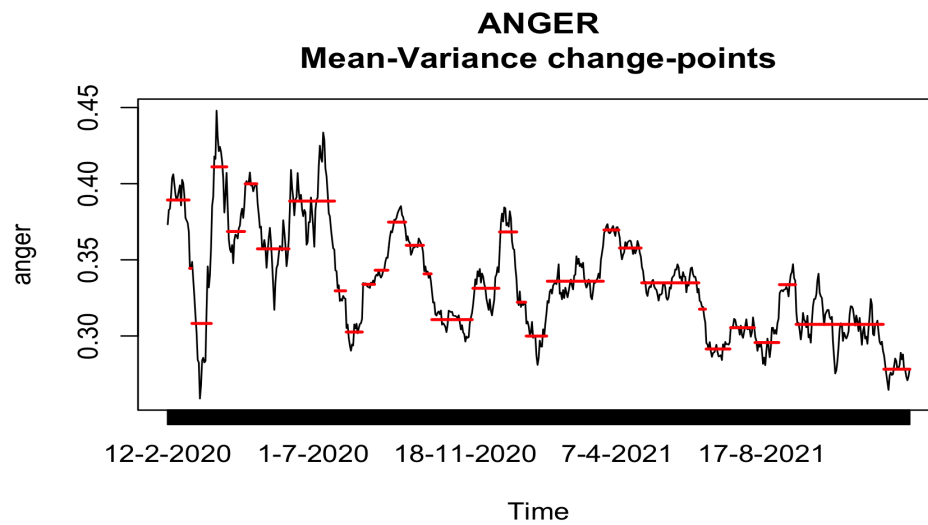


Fig 5.2.2 :- Anger change point Analysis

## Chapter 6

# Summary

### 6.1 Discussion

In this study, we are analyzing covid 19 related tweets from the perspective of misinformation. We collected a large set of datasets consisting of the location of each tweet. Besides this, we also have data collection of tweets related to covid-19 without location, behavioral indicator data, and significant pandemic events. We analyzed various tools for collecting data and sources of dataset collection and provided an easy-to-use overview of Tweets and User objects returned by Twitter API while fetching data.

In going forward in the study, we analyzed the dataset comprehensively, like an overview of data and platform, measuring the popularity of tweets and measuring words used while tweeting covid-19 tweets. We provide counts of hashtags used geographically. To come up with some inside properties in the dataset, we did Emotions detection using Empath and Plutchik and compared the output of both. We relate emotions strength-wise and daily with events during the covid-19 timeline from Jan 2020 to Dec 2021. We have further done the most significant task of cross-correlation between emotion strength and behavioral indicators. We also did a change point analysis of the events- emotions graph to support our findings.

This type of study is beneficial in analyzing the emotional strength of netizens with covid-19 chitchat happening on social media as people are talking more on social media, and sometimes wrong information may lead to considerable loss. Therefore we thought this type would support online communities.

### 6.2 Limitations

1. Panacea data is not an exhaustive data source.
2. Limited keywords in the panacea data set.
3. It does not contain locations anymore in their recent data upload of 2022 in Panacea lab data.
4. Panacea data only provides a unique id for tweets, but the content needs to be fetched from Twitter API.

5. Events data is limited and sometimes does not provide city attributes leading to manual validation of the location of events.
6. We only consider eight plutchik emotions, and some tweets might not be classified into these eight emotions.
7. In Cross correlation of emotions strength we have seen some counter intuitive result which is opposite to real word scenario like fear has negative correlation with daily new cases.

## **6.3 Future work**

1. Further extending the timeline from Jan 2022 and updating the study and analysis.
2. Extending this study to find significant changes in behavior after covid-19 events like covid-19 wave surge or after lockdown.
3. Analyzing this study for other parts of the world or countries where covid-19 has continuously spread.
4. We can further consider more emotions from the plutchik wheel and classify tweets.

# References

1. Panacea Lab data - [https://github.com/thepanacealab/covid19\\_twitter](https://github.com/thepanacealab/covid19_twitter)
2. Our world in data - <https://github.com/owid/covid-19-data/tree/master/public/data>
3. World events data - <https://www.thinkglobalhealth.org/article/updated-timeline-coronavirus>
4. India Events data - [https://www.business-standard.com/article/current-affairs/here-s-a-timeline-of-events-since-lockdown-was-imposed-in-india-120070201413\\_1.html](https://www.business-standard.com/article/current-affairs/here-s-a-timeline-of-events-since-lockdown-was-imposed-in-india-120070201413_1.html)
5. Twitter API docs - <https://developer.twitter.com/en/docs/twitter-api>
6. Tweepy tools - <https://docs.tweepy.org/en/stable/>
7. Snsrape tools - <https://github.com/JustAnotherArchivist/snsrape>
8. Twint tools - <https://github.com/twintproject/twint>
9. Plutchik emotion wheel - <https://www.6seconds.org/2022/03/13/plutchik-wheel-emotions/>
10. Change point analysis - <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5464762/>
11. Cross correlation - <https://www.sciencedirect.com/topics/earth-and-planetary-sciences/cross-correlation>
12. Plutchick model - <https://arxiv.org/abs/1812.01207>