

# **CoVis: COVID-19 Twitter Sentiment Analysis & Extraction Visualizer**

## **1. Introduction**

### **1.1 Overview**

In this era of flourishing technology, Social Media has become a powerful platform for the public to voice their concerns and beliefs. Among them one such platform is Twitter. Twitter has been a popular platform for microblogging in the past few years. In this context, Sentiment Analysis is extremely useful in social media monitoring as it allows us to gain an overview of the wider public opinion behind certain topics. Across the past few years, as the organizations and governments across the world start to adopt the ability to extract insights from social data, the applications of sentiment analysis are broad and powerful. There has been a clear implication that shifts in sentiment on social media correlate with shifts in the economics of a country and also the common notion among the public.

Due to the recent COVID-19 pandemic, there has been a wide change in sentiments of various sectors of the Indian public towards the government policies/actions. Studying the sentiment of the people on the epidemic and government decisions is very important as it acts as a sanity check for the effectiveness of the adopted government policies. This study also provides insight into the business models required to be adopted to suit this new age of post-COVID-19 where people's sentiments have widely changed. In this context 'Sentiment Analysis of COVID-19 Tweets' is a very important problem statement.

### **1.2 Purpose**

The outbreak of COVID-19 caused heavy disruption to the everyday lives of people across the globe. In a country like with a large, diverse population like India, there are bound to be instances of mass hysteria and panic which are further fuelled by unreliable and sometimes inaccurate data. Gauging the feelings/emotions of the citizens would provide insights into the public mindset and would pave the way for the government and many organizations to address these situations by providing them with the right data and information, eradicating fake news, thereby helping in suppressing unnecessary panic among the people. Social media acts as the bridge between the people, the government, and such organizations. The scope of this project lies in the application of sentiment analysis to the views expressed by people on social media, twitter, in this case, to analyze the trends in the dynamic mood of the population. Usually, the terms "fight" and "positive" are used in a negative and positive context respectively, but we observe a role reversal in this situation. The identification of such terms and their usage according to the context would be an essential part of the project. Also, the scope of the project can be found in stopping the spread of fake news related to the pandemic, *creating an interactive dashboard that delivers information about*

*the current situation, real-time sentiment analysis of tweets, trend analysis of various COVID-19 related hashtags, engagement on Twitter, overall sector-wise polarity score of the tweets and the public emotion charts.*

## **2. LITERATURE SURVEY**

### **2.1 Existing Problem**

The whole world is dealing with the coronavirus pandemic right now. But another crisis has also manifested itself along with the virus, which is as detrimental as the virus itself. That is the large flow of information regarding the virus in the form of tweets, blogs, news, and too little analysis of this tremendous amount of information flowing in and out of the system every moment. This information could sometimes be fake, mixed news or just some opinion about the pandemic. The spreading of such incomplete, inaccurate news would be the cause of mass hysteria and fear among the public and hence the need of the hour is to address and better understand this information crisis regarding the pandemic. The government, companies, and many organizations form their decisions to cater to the needs of the people, and hence understanding the right sentiment and opinions of the people towards the pandemic plays a crucial role in policymaking and creating appropriate business models that are of necessity to the people.

### **2.2 Proposed Solution**

The exponential rise in the social media usage by the people in the last decade to express their opinions, perspectives and also as a source of news rather than the traditional news has laid emphasis on the usage of Deep Learning and Artificial Intelligence methods in gauging information from these sites to analyze and extract sentiments that act as valuable sources of insights for corporate companies and the government alike in making policies and appropriate business models to match the trend of the public. Hence, the main purpose of this project lies in the creation of a Public Sentiment Analysis Dashboard that depicts the sentiment trend among the Indian public towards the pandemic. Twitter is a microblogging site that has gained popularity in recent years, for the governments, organizations, and people alike, as a platform to make official announcements, expressing moods, opinions towards current, ongoing issues. Hence the tweets, that are the short posts that are made on this platform, act as our data set for creating this dashboard. The application of Deep Learning methods to analyze these tweets, gauge sentiments from them, and then representing these results on a live interactive dashboard is the main aim of this project.

A sentiment extraction model that analyses this data on a large scale and provides valuable insights into the sentiment trends of the public. Recent language models such as BERT, XLNET, T5, Roberta, Electra have shown a great contextual understanding of language. So a transfer learning approach based on these models was used. A caveat that appeared in the context of COVID-19 is that few words portray a negative sentiment whereas originally in a normal context they convey a positive sentiment, for instance, the word 'positive'. Instances like these were looked upon and it was found out after a certain amount of

analysis that these words were not frequently used out of context as they were initially assumed to be. Also, another interesting challenge that was taken up was the process of extraction of common notions of public regarding COVID-19 actions taken by the Indian government and business corporations. To facilitate features to overcome these challenges during sentiment analysis, language models capable of simultaneous sentiment classification and extraction were developed. Also, the problem statement requires a good graphical user interface which provides keen insights into the public opinion on government and business decisions during COVID. For this, a user-friendly interactive dashboard that delivers information about the current situation, real-time sentiment analysis of tweets, trend analysis of various COVID-19 related hashtags, engagement on Twitter, overall sector-wise polarity score of the tweets and the public emotion charts was created. Thus, in this context, the sentiment extraction approach and the post prediction analysis have facilitated the easy extraction of the climate of opinion surrounding the COVID-19 pandemic.

### 3. Theoretical Analysis

#### 3.1 Block Diagram



## 3.2 Hardware/ Software Design

Initially, the IEEE Coronavirus (COVID-19) Tweets Data set was downloaded from their website. Upon inspection, it was found that many tweets did not have geo-location tags, and also many were in different languages apart from English. Due to this challenge in obtaining proper data, a new data set named Geo-Tagged Coronavirus (COVID-19) Tweets Data set was obtained from the same website. These tweets were then hydrated using the “Hydrator” software and also a few python commands. Then, tweets in “English” and tweets from “India” were randomly chosen and a new dataset was created. Further, other data set containing COVID-19 related tweets from India were obtained from Kaggle. This data set was then cleaned and normalized to make it useful for further analysis.

This cleaned data were then subjected to further analysis by extracting bigrams, trigrams, and plotting frequency bar graphs, Word Clouds, Relationship Nexus, Boxplots, etc. This was done using both Python and R. Some Interactive Plots were also plotted. This complete process is termed as Exploratory Data Analysis.

After Exploratory Data Analysis is completed, the tweets are then Tokenized and are made in a format suitable for the Language Model. In this Step, two models were used:

- 1) Roberta Model: Transfer learning methods were implemented to carry out sentiment analysis. Sentiment Analysis of Tweets was carried out by integrating and using both the Huggingface Transformer Library and FastAI. Further **Slanted Triangular Learning Rates**, **Discriminate Learning Rate** and even **Gradual Unfreezing** were used, as a result of which, state-of-the-art results were obtained rapidly without even tuning the parameters. The tokenized data was then passed through the RoBERTa model to perform Sentiment Analysis. This yielded a model with an accuracy of 97% over the data set. The Tweepy API was used to scrape tweets in real-time which were then passed through the model to obtain the sentiments.
- 2) RoBERTa-CNN Sentiment Extractor: After the completion of the sentiment analysis the data was further explored for the sentiment triggers in the tweets. HuggingFace transformers don't have a TFRobertaForQuestionAnswering, for this purpose, a TFRobertaModel was created to convert trained data into arrays that the Roberta model can interpret. While training the Sentiment Extractor model, 5 stratified KFoldes were used in such a way that, in each fold, the best model weights were saved and these weights were reloaded before carrying out testing and predictions. Roberta with CNN head was used for Twitter Sentiment Extraction. Thus after passing the data through this model we obtained a new column of the extracted text for the sentiments which was also used to plot certain graphs.

Now the entire process pertaining to the data and Model building is completed. Now, the Flask APP is built for the purpose of Deployment. First, the application is deployed on the localhost and debugged and then we move on to deploying on the WebServers.

A flask app was used for setting up website routing. It is used to integrate the back end machine learning models with the dashboard. Then Socketio (web sockets) were used for dynamic implementations on the website, namely the Real-Time Plot Generators and Twitter live feed. The basic functionality of the Flask Socketio lies in running background threads when the client is not connected to the website thereby enabling dynamic plotting. The

above built Dashboard was deployed on the Local Machine and debugged for any possible errors. The scraping rate and other parameters were monitored and corrected accordingly.

**Project Flask/Socketio APP with Sentiment Analyzer & Extractor Code:**

07/15/20 09:22:37 C:\Users\Nikhil Keetha\Downloads\app\main.py

```

1 import numpy as np # linear algebra
2 import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
3 from pathlib import Path
4 from sklearn.metrics import confusion_matrix, classification_report, accuracy_score
5 from sklearn.model_selection import StratifiedKFold
6 import os
7
8 # pytorch
9 import torch
10 import torch.optim as optim
11
12 import random
13
14 # fastai
15 from fastai import *
16 from fastai.text import *
17 from fastai.callbacks import *
18
19 # transformers
20 from transformers import *
21 from transformers import PreTrainedModel, PreTrainedTokenizer, PretrainedConfig
22 from transformers import RobertaForSequenceClassification, RobertaTokenizer, RobertaConfig
23
24 # tensorflow
25 import tensorflow as tf
26 import tensorflow.keras.backend as K
27
28 import tokenizers
29 import pickle
30 import math
31 import re
32 import string
33 import seaborn as sns
34 color = sns.color_palette()
35 import matplotlib.pyplot as plt
36 from nltk import bigrams
37 import nltk
38 nltk.download('stopwords')
39 from nltk.corpus import stopwords
40 eng_stopwords = stopwords.words('english')
41 import collections
42 from wordcloud import WordCloud
43 from textwrap import wrap
44 import plotly.graph_objects as go
45 import plotly.express as px
46 import networkx as nx
47
48 from tweepy import OAuthHandler
49 #from tweepy.streaming import StreamListener
50 import tweepy
51 import csv
52 import time
53 import datetime
54
55 from flask import Flask, request, jsonify, make_response
56 from flask import render_template, url_for, flash, redirect, copy_current_request_context
57 from flask_socketio import SocketIO, emit
58 from threading import Thread, Event
59
60 from flask_restful import reqparse, abort, Api, Resource
61 import json
62 from flask import jsonify
63
64 app = Flask(__name__)
65 app.config['SECRET_KEY'] = '5791628bb0b13ce0c676dfde280ba012'
66 app.config['DEBUG'] = True
67 app.config['SEND_FILE_MAX_AGE_DEFAULT'] = 0
68
69 socketio = SocketIO(app, async_mode=None, logger=True, engineio_logger=True)
70
71 # Twitter credentials
72 consumer_key = 'rm2bLDjA2Bzlj0A0GomL5o6W7'
73 consumer_secret = 'xiFBG4VKWpuQts1v3uqAesllpDp36y44YkFnzBtezSbSYW9dBV'
74 access_key = '935519854064418816-s0BxmFMAyGAX3FQXRBJH0drpZ20XpB'
75 access_secret = 'Gb0Tefzapdet9vpmR3H90BRuJN3Ns1cIAdh5HrkIYPJz'
76
77 # Pass your twitter credentials to tweepy via its OAuthHandler
78 auth = OAuthHandler(consumer_key, consumer_secret)
79 auth.set_access_token(access_key, access_secret)
80 t_api = tweepy.API(auth, wait_on_rate_limit=True, wait_on_rate_limit_notify=True)
81
82 class TransformersBaseTokenizer(BaseTokenizer):
83     """Wrapper around PreTrainedTokenizer to be compatible with fast.ai"""
84     def __init__(self, pretrained_tokenizer: PreTrainedTokenizer, model_type = 'bert', **kwargs):
85         self.pretrained_tokenizer = pretrained_tokenizer

```

```

86         self.max_seq_len = pretrained_tokenizer.max_len
87         self.model_type = model_type
88
89     def __call__(self, *args, **kwargs):
90         return self
91
92     def tokenizer(self, t:str) -> List[str]:
93         """Limits the maximum sequence length and add the special tokens"""
94         CLS = self._pretrained_tokenizer.cls_token
95         SEP = self._pretrained_tokenizer.sep_token
96         if self.model_type in ['roberta']:
97             tokens = self._pretrained_tokenizer.tokenize(t, add_prefix_space=True)[:self.max_seq_len - 2]
98             tokens = [CLS] + tokens + [SEP]
99         else:
100             tokens = self._pretrained_tokenizer.tokenize(t)[:self.max_seq_len - 2]
101             if self.model_type in ['xlnet']:
102                 tokens = tokens + [SEP] + [CLS]
103             else:
104                 tokens = [CLS] + tokens + [SEP]
105         return tokens
106
107     class TransformersVocab(Vocab):
108         def __init__(self, tokenizer: PreTrainedTokenizer):
109             super(TransformersVocab, self).__init__(itos = [])
110             self.tokenizer = tokenizer
111
112         def numericalize(self, t:Collection[str]) -> List[int]:
113             "Convert a list of tokens 't' to their ids."
114             return self.tokenizer.convert_tokens_to_ids(t)
115             #return self.tokenizer.encode(t)
116
117         def textify(self, nums:Collection[int], sep=' ') -> List[str]:
118             "Convert a list of 'nums' to their tokens."
119             nums = np.array(nums).tolist()
120             return sep.join(self.tokenizer.convert_ids_to_tokens(nums)) if sep is not None else
121             self.tokenizer.convert_ids_to_tokens(nums)
122
123     def __getstate__(self):
124         return {'itos':self.itos, 'tokenizer':self.tokenizer}
125
126     def __setstate__(self, state:dict):
127         self.itos = state['itos']
128         self.tokenizer = state['tokenizer']
129         self.stoi = collections.defaultdict(int,{v:k for k,v in enumerate(self.itos)})
130
131     # defining our model architecture
132     class CustomTransformerModel(nn.Module):
133         def __init__(self, transformer_model: PreTrainedModel):
134             super(CustomTransformerModel,self).__init__()
135             self.transformer = transformer_model
136
137         def forward(self, input_ids, attention_mask=None):
138             # attention_mask
139             # Mask to avoid performing attention on padding token indices.
140             # Mask values selected in ``[0, 1]``:
141             # ``1`` for tokens that are NOT MASKED, ``0`` for MASKED tokens.
142             attention_mask = (input_ids!=pad_idx).type(input_ids.type())
143
144             logits = self.transformer(input_ids,
145                                     attention_mask = attention_mask)[0]
146             return logits
147
148     def predict_sentiment(learner, text):
149         sentiment = learner.predict(text)[1].item()
150         return sentiment
151
152     def sentiment_label (Sentiment):
153         if Sentiment == 2:
154             return "positive"
155         elif Sentiment == 0:
156             return "negative"
157         else :
158             return "neutral"
159
160     def replace_url(string): # cleaning of URL
161         text = re.sub(r'http\S+', 'LINK', string)
162         return text
163
164     def replace_email(text):#Cleaning of Email related text
165         line = re.sub(r'[\w\.-]+@[\w\.-]+\S+', 'MAIL',str(text))
166         return "".join(line)
167
168     def rep(text):#cleaning of non standard words
169         grp = text.group(0)
170         if len(grp) > 3:
171             return grp[0:2]

```



[illegible]



```

258     lst=[]
259     for h in hashtags:
260         lst.append(h['text'])
261     try:
262         text = tweetretweeted_status.full_text
263     except AttributeError: # Not a Retweet
264         text = tweet.full_text
265
266     itweet =
[username,acctdesc,location,following,followers,totaltweets,usercreatedts,tweetcreatedts,retweetcount,text,lst]
267     db_tweets.loc[len(db_tweets)] = itweet
268
269     noTweets += 1
270
271     #filename = "tweets.csv"
272     #with open(filename, "a", newline='') as fp:
273     #    wr = csv.writer(fp, dialect='excel')
274     #    wr.writerow(itweet)
275
276     if i+1 != numRuns:
277         time.sleep(920)
278
279     filename = "static/analysis.csv"
280     db_tweets['text'] = db_tweets['text'].apply(change)
281     db_tweets = db_tweets[['retweetcount', 'text']]
282     # Store dataframe in csv with creation date timestamp
283     db_tweets.drop_duplicates(subset="text", keep = 'first', inplace = True)
284     db_tweets.to_csv(filename, index = False) #
285     print('Scrapping Done')
286
287 # Functions for Sentiment Extractor
288 def save_weights(model, dst_fn):
289     weights = model.get_weights()
290     with open(dst_fn, 'wb') as f:
291         pickle.dump(weights, f)
292
293
294 def load_weights(model, weight_fn):
295     with open(weight_fn, 'rb') as f:
296         weights = pickle.load(f)
297         model.set_weights(weights)
298     return model
299
300 def loss_fn(y_true, y_pred):
301     # adjust the targets for sequence bucketing
302     ll = tf.shape(y_pred)[1]
303     y_true = y_true[:, :ll]
304     loss = tf.keras.losses.categorical_crossentropy(y_true, y_pred,
305         from_logits=False, label_smoothing=LABEL_SMOOTHING)
306     loss = tf.reduce_mean(loss)
307     return loss
308
309 #Global Constants for TF sentiment extractor
310 MAX_LEN = 310
311 PAD_ID = 1
312 num_splits = 1
313 SEED = 88888
314
315 PATH = 'static/Tf-Roberta/'
316 tokenizer = tokenizers.ByteLevelBPETokenizer(
317     vocab_file=PATH+'vocab-roberta-base.json',
318     merges_file=PATH+'merges-roberta-base.txt',
319     lowercase=True,
320     add_prefix_space=True
321 )
322
323 def build_model():
324     ids = tf.keras.layers.Input((MAX_LEN,), dtype=tf.int32)
325     att = tf.keras.layers.Input((MAX_LEN,), dtype=tf.int32)
326     tok = tf.keras.layers.Input((MAX_LEN,), dtype=tf.int32)
327     padding = tf.cast(tf.equal(ids, PAD_ID), tf.int32)
328
329     lens = MAX_LEN - tf.reduce_sum(padding, -1)
330     max_len = tf.reduce_max(lens)
331     ids_ = ids[:, :max_len]
332     att_ = att[:, :max_len]
333     tok_ = tok[:, :max_len]
334
335     config = RobertaConfig.from_pretrained(PATH+'config-roberta-base.json')
336     bert_model = TFRobertaModel.from_pretrained(PATH+'pretrained-roberta-base.h5', config=config)
337     x = bert_model(ids_, attention_mask=att_, token_type_ids=tok_)
338
339     x1 = tf.keras.layers.Dropout(0.1)(x[0])
340     x1 = tf.keras.layers.Conv1D(768, 2, padding='same')(x1)
341     x1 = tf.keras.layers.LeakyReLU()(x1)
342     x1 = tf.keras.layers.Dense(1)(x1)
343     x1 = tf.keras.layers.Flatten()(x1)
344     x1 = tf.keras.layers.Activation('softmax')(x1)

```

```

345 x2 = tf.keras.layers.Dropout(0.1)(x[0])
346 x2 = tf.keras.layers.Conv1D(768, 2, padding='same')(x2)
347 x2 = tf.keras.layers.LeakyReLU()(x2)
348 x2 = tf.keras.layers.Dense(1)(x2)
349 x2 = tf.keras.layers.Flatten()(x2)
350 x2 = tf.keras.layers.Activation('softmax')(x2)
351
352 model = tf.keras.models.Model(inputs=[ids, att, tok], outputs=[x1,x2])
353 optimizer = tf.keras.optimizers.Adam(learning_rate=3e-5)
354 model.compile(loss=loss_fn, optimizer=optimizer)
355
356 # this is required as `model.predict` needs a fixed size!
357 x1_padded = tf.pad(x1, [[0, 0], [0, MAX_LEN - max_len]], constant_values=0.)
358 x2_padded = tf.pad(x2, [[0, 0], [0, MAX_LEN - max_len]], constant_values=0.)
359
360 padded_model = tf.keras.models.Model(inputs=[ids, att, tok], outputs=[x1_padded,x2_padded])
361 return model, padded_model
362
363 def generate_wordcloud(data,title):
364     wc = WordCloud(width=400, height=330, max_words=150, colormap="Dark2", background_color='white',
365     collocations=False).generate_from_frequencies(data)
366     plt.figure(figsize=(10,8))
367     plt.imshow(wc, interpolation='bilinear')
368     plt.tight_layout(pad=0)
369
370 @app.route("/", methods=['GET', 'POST'])
371 @app.route("/dashboard", methods=['GET', 'POST'])
372 def home():
373     return render_template('dashboard.html')
374
375 @app.route("/live_case_count", methods=['GET', 'POST'])
376 def live_count():
377     return render_template('live_case_count.html')
378
379 @app.route("/about_us", methods=['GET', 'POST'])
380 def about():
381     return render_template('about_us.html')
382
383 HOUR = 3600;
384
385 thread = Thread()
386 thread_stop_event = Event()
387
388 model_type = 'roberta'
389 pretrained_model_name = 'roberta-base'
390
391 model_class, tokenizer_class, config_class = RobertaForSequenceClassification, RobertaTokenizer, RobertaConfig
392
393 transformer_tokenizer = tokenizer_class.from_pretrained(pretrained_model_name)
394 transformer_base_tokenizer = TransformersBaseTokenizer(pretrained_tokenizer = transformer_tokenizer, model_type =
395 model_type)
396 fastai_tokenizer = Tokenizer(tok_func = transformer_base_tokenizer, pre_rules=[], post_rules=[])
397
398 pad_idx = transformer_tokenizer.pad_token_id
399
400 p = 'static/Roberta_Model'
401 learner = load_learner(p, 'transformer.pkl')
402
403 def plotGenerator():
404     """
405     Generates real time plots every 1 day.
406     """
407     while not thread_stop_event.isSet():
408         # Initialise these variables:
409         print('Code is Running!!!!')
410
411         search_words = "(#India AND #COVID-19) OR #COVID19India"
412         yesterday = datetime.datetime.now() - datetime.timedelta(days = 1)
413         date_since = yesterday.strftime("%Y-%m-%d")
414         date_until = datetime.datetime.today().strftime('%Y-%m-%d')
415         numTweets = 2000
416         numRuns = 1
417         # Call the function scraptweets
418         program_start = time.time()
419         scraptweets(search_words, date_since, date_until, numTweets, numRuns)
420         program_end = time.time()
421
422         path = 'static/analysis.csv'
423         predictions = pd.read_csv(path)
424
425         print('Start Prediction')
426
427         predictions['Prediction'] = predictions['text'].apply(lambda x: predict_sentiment(learner, x))
428         predictions['Prediction'] = predictions['Prediction'].apply(sentiment_label)
429         class_names = ['negative','positive','neutral']
430
431         predictions.rename(columns={'Prediction':'sentiment'}, inplace=True)

```

```

431     print('Predictions Done')
432
433     path = 'static/analysis.csv'
434     predictions.to_csv(path,index=False)
435     test = pd.read_csv(path)
436
437     MAX_LEN = 310
438     PAD_ID = 1
439     num_splits = 1
440     SEED = 88888
441
442     PATH = 'static/Tf-Roberta/'
443     tokenizer = tokenizers.ByteLevelBPETokenizer(
444         vocab_file=PATH+'vocab-roberta-base.json',
445         merges_file=PATH+'merges-roberta-base.txt',
446         lowercase=True,
447         add_prefix_space=True
448     )
449
450     test['len'] = test['text'].str.len()
451     test = test[test['len']<=310]
452     test.drop("len",axis=1,inplace=True)
453     test.reset_index(drop=True, inplace=True)
454
455     ct = test.shape[0]
456     input_ids_t = np.ones((ct,MAX_LEN),dtype='int32')
457     attention_mask_t = np.zeros((ct,MAX_LEN),dtype='int32')
458     token_type_ids_t = np.zeros((ct,MAX_LEN),dtype='int32')
459     sentiment_id = {'positive': 1313, 'negative': 2430, 'neutral': 7974}
460
461     for k in range(test.shape[0]):
462
463         # INPUT_IDS
464         text1 = " "+" ".join(test.loc[k,'text'].split())
465         enc = tokenizer.encode(text1)
466         s_tok = sentiment_id[test.loc[k,'sentiment']]
467         input_ids_t[k,:len(enc.ids)+3] = [0, s_tok] + enc.ids + [2]
468         attention_mask_t[k,:len(enc.ids)+3] = 1
469
470     DISPLAY=1 # USE display=1 FOR INTERACTIVE
471     preds_start = np.zeros((input_ids_t.shape[0],MAX_LEN))
472     preds_end = np.zeros((input_ids_t.shape[0],MAX_LEN))
473
474     # for fold in range(0,5):
475     K.clear_session()
476     model, padded_model = build_model()
477     path = 'static/R_CNN_weights/'
478     weight_fn = path+'v0-roberta-0.h5'
479
480     print('Loading model...')
481     # model.load_weights('%s-roberta-%i.h5'%(VER,fold))
482     load_weights(model, weight_fn)
483
484     print('Predicting Test...')
485
486     preds = padded_model.predict([input_ids_t,attention_mask_t,token_type_ids_t],verbose=DISPLAY)
487     preds_start += preds[0]/num_splits
488     preds_end += preds[1]/num_splits
489
490     all = []
491     for k in range(input_ids_t.shape[0]):
492         a = np.argmax(preds_start[k,:])
493         b = np.argmax(preds_end[k,:])
494         if a>b:
495             st = test.loc[k,'text']
496         else:
497             text1 = " "+" ".join(test.loc[k,'text'].split())
498             enc = tokenizer.encode(text1)
499             st = tokenizer.decode(enc.ids[a-2:b-1])
500         all.append(st)
501
502     test['selected_text'] = all
503     test.to_csv('static/analysis.csv',index=False)
504
505     print('Extraction Done')
506
507     #Plots start
508     data=pd.read_csv("static/analysis.csv")
509     data['selected_text'] = data['selected_text'].astype(str)
510     df = data.sentiment.value_counts()
511     size = list(df.values)
512     names = list(df.index)
513     fig = plt.figure(figsize=(10,10))
514     plt.xlabel("Sentiment",FontSize = 16)
515     plt.ylabel("Frequency",FontSize = 16)
516     sns.barplot(names,size,alpha = 0.8)
517     fig.savefig("static/images/realtime/bar.png")
518

```

```

519 df_new = pd.DataFrame(dict(
520     r=list(df.values),
521     theta=list(df.index)))
522 plt.figure(figsize=(10,10))
523 fig = px.line_polar(df_new, r='r', theta='theta', line_close=True)
524 fig.update_traces(fill='toself')
525 fig.write_image("static/images/realtime/radar_plot.png")
526
527 # Pie chart
528 labels = list(df.index)
529 sizes = list(df.values)
530 # only "explode" the 2nd slice
531 explode = (0.1, 0.1, 0.1)
532 #add colors
533 colors = ['#ff9999', '#66b3ff', '#99ff99']
534 fig1, ax1 = plt.subplots()
535 ax1.pie(sizes, explode=explode, labels=labels, colors=colors, autopct='%1.1f%%',
536        shadow=True, startangle=90)
537 # Equal aspect ratio ensures that pie is drawn as a circle
538 ax1.axis('equal')
539 plt.tight_layout()
540 plt.savefig('static/images/realtime/pie_chart.png')
541
542 for i in range(3):
543     Data = data[data["sentiment"]==df.index[i]]
544     Word_frequency = pd.Series(' '.join(Data.selected_text).split()).value_counts()[:20]#Calculating the words
frequency
545     plt.figure(figsize=(25,10))
546     plt.ylabel("Frequency",fontsize=16)
547     plt.title("Sentiment Triggers")
548     sns.barplot(Word_frequency.index,Word_frequency.values,alpha=0.8)
549     plt.savefig("static/images/realtime/wordfrequency_"+df.index[i]+".png")
550
551     for i in range(3):
552         Analysis_Data = data
553         data["selected_text"]=data["selected_text"].apply(lambda x: ' '.join([word for word in x.split() if word
not in (eng_stopwords)]))
554         Sentiment = Analysis_Data[Analysis_Data['sentiment'] == df.index[i]]#Creating the dataframe of having
same sentiment
555         Word_frequency = pd.Series(' '.join(Sentiment.selected_text).split()).value_counts()[:50]#Calculating
the words frequency
556         generate_wordcloud(Word_frequency.sort_values(ascending=False),data.index[i])
557         plt.savefig("static/images/realtime/Wordcloud_" +df.index[i]+".png")
558
559         data["text"]=data["text"].apply(lambda x: ' '.join([word for word in x.split() if word not in
(eng_stopwords)]))
560         bigrams = [b for l in data.text for b in zip(l.split(" ")[:-1], l.split(" ")[1:])]
561         bigram_counts = collections.Counter(bigrams)
562         bigram_df = pd.DataFrame(bigram_counts.most_common(10),
563                                columns=['bigram', 'frequency'])
564         x = bigram_df.bigram
565         y = bigram_df.frequency
566
567         fig, ax = plt.subplots(1, 1, figsize = (20, 15), dpi=300)
568         sns.barplot(x,y,alpha=0.8)
569         plt.ylabel("Frequency",fontsize=16)
570         ax.set_xlabel('')
571         plt.savefig('static/images/realtime/bigram_freq.png')
572
573         ext_data_negative = data[data["sentiment"]=="negative"]
574         ext_data_positive = data[data["sentiment"]=="positive"]
575         bigrams = [b for l in ext_data_positive.selected_text for b in zip(l.split(" ")[:-1], l.split(" ")[1:])]
576         bigram_counts = collections.Counter(bigrams)
577         bigram_df_positive = pd.DataFrame(bigram_counts.most_common(60),
578                                         columns=['bigram', 'frequency'])
579         bigrams = [b for l in ext_data_negative.selected_text for b in zip(l.split(" ")[:-1], l.split(" ")[1:])]
580         bigram_counts = collections.Counter(bigrams)
581         bigram_df_negative = pd.DataFrame(bigram_counts.most_common(80),
582                                         columns=['bigram', 'frequency'])
583
584         # Create network plot
585         G=nx.grid_2d_graph(2,2)
586
587         pos = nx.fruchterman_reingold_layout(G,k=10,iterations=100)
588         fig,ax = plt.subplots(figsize=(50,30))
589         d = bigram_df_negative.set_index('bigram').T.to_dict('records')
590         for k, v in d[0].items():
591             G.add_edge(k[0], k[1], weight=(v * 10))
592         pos = nx.fruchterman_reingold_layout(G,k=10,iterations=100)
593
594         nx.draw_networkx(G, pos,
595                        font_size=16,
596                        width=4,
597                        edge_color='#e25a4b',
598                        node_size=500,
599                        title = "Negative Sentiment",
600                        with_labels = False,
601                        ax=ax)

```



```

602 x_values, y_values = zip(*pos.values())
603 x_max = max(x_values)
604 x_min = min(x_values)
605 x_margin = (x_max - x_min) * 0.25
606 plt.xlim(x_min - x_margin, x_max + x_margin)
607
608 for key, value in pos.items():
609     x, y = value[0]+.135, value[1]+.045
610     ax.text(x, y,
611             s=key, bbox=dict(facecolor='ffcd94', alpha=0.4),
612                             horizontalalignment='center', fontsize=35)
613 plt.savefig("static/images/realtime/ext_negative.png")
614 fig, ax = plt.subplots(figsize=(50,30))
615 d = bigram_df_positive.set_index('bigram').T.to_dict('records')
616 for k, v in d[0].items():
617     G.add_edge(k[0], k[1], weight=(v * 10))
618 pos = nx.fruchterman_reingold_layout(G, k=10, iterations=100)
619
620 nx.draw_networkx(G, pos,
621                  font_size=16,
622                  width=4,
623                  edge_color='#999894',
624                  node_size=500,
625                  with_labels = False,
626                  title = "Positive Sentiment",
627                  ax=ax)
628 x_values, y_values = zip(*pos.values())
629 x_max = max(x_values)
630 x_min = min(x_values)
631 x_margin = (x_max - x_min) * 0.25
632 plt.xlim(x_min - x_margin, x_max + x_margin)
633
634 # Create offset labels
635 for key, value in pos.items():
636     x, y = value[0]+.135, value[1]+.045
637     ax.text(x, y,
638             s=key, bbox=dict(facecolor='#7c99d0', alpha=0.4),
639                             horizontalalignment='center', fontsize=35)
640 plt.savefig("static/images/realtime/ext_positive.png")
641
642 data["text"] = data["text"].apply(lambda x: ' '.join([word for word in x.split() if word not in
643 (eng_stopwords)]))
644 bigrams = [b for l in data.text for b in zip(l.split(" ")[:-1], l.split(" ")[1:])]
645 bigram_counts = collections.Counter(bigrams)
646 bigram_df = pd.DataFrame(bigram_counts.most_common(60),
647                           columns=['bigram', 'frequency'])
648
649 d = bigram_df.set_index('bigram').T.to_dict('records')
650 # Create network plot
651 G = nx.Graph()
652 for k, v in d[0].items():
653     G.add_edge(k[0], k[1], weight=(v * 10))
654
655 fig, ax = plt.subplots(figsize=(20,20))
656 pos = nx.spring_layout(G, dim=2, k=5)
657
658 # Plot networks
659 nx.draw_networkx(G, pos,
660                  font_size=12,
661                  width=4,
662                  edge_color='grey',
663                  node_color='#4a4140',
664                  node_size=500,
665                  with_labels = False,
666                  ax=ax)
667 x_values, y_values = zip(*pos.values())
668 x_max = max(x_values)
669 x_min = min(x_values)
670 x_margin = (x_max - x_min) * 0.25
671 plt.xlim(x_min - x_margin, x_max + x_margin)
672
673 # Create offset labels
674 for key, value in pos.items():
675     x, y = value[0]+.135, value[1]+.045
676     ax.text(x, y,
677             s=key,
678             bbox=dict(facecolor='ffcd94', alpha=0.4),
679                     horizontalalignment='center', fontsize=25)
680
681 fig.savefig('static/images/realtime/network.png')
682
683 fig = px.box(data, y="retweetcount", points="all")
684 fig.update_layout(
685     yaxis_title="Retweet Count",
686     font=dict(
687         family="Courier New, monospace",
688         size=18,

```

```

689         color="#7f7f7f"
690     )
691 )
692 fig.write_image('static/images/realtime/retweet_count_boxplot.png')
693
694 fig = px.box(data, y="sentiment", points="all")
695 fig.update_layout(
696     yaxis_title="Sentiment",
697     font=dict(
698         family="Courier New, monospace",
699         size=18,
700         color="#7f7f7f"
701     )
702 )
703 fig.write_image('static/images/realtime/sentiment_boxplot.png')
704
705 print('Plotting Done.')
706
707 socketio.sleep(24*HOUR)
708
709 @app.route('/real_time_analysis', methods=['GET', 'POST'])
710 def realtime():
711     #only by sending this page first will the client be connected to the socketio instance
712     return render_template('real_time_analysis.html')
713
714 @socketio.on('connect', namespace='/test')
715 def test_connect():
716     # need visibility of the global thread object
717     global thread
718     print('Client connected')
719     #Start the plot generator thread only if the thread has not been started before.
720     if not thread.isAlive():
721         print("Starting Thread")
722         thread = socketio.start_background_task(plotGenerator)
723
724 @socketio.on('disconnect', namespace='/test')
725 def test_disconnect():
726     print('Client disconnected')
727
728 twitter_thread = Thread()
729 twitter_thread_stop_event = Event()
730
731 def gettweets():
732
733     while not twitter_thread_stop_event.isSet():
734         print('Tweet feed starts')
735         for tweet in tweepy.Cursor(t_api.search,q="#" + "COVID19India" + " -filter:retweets",rpp=5,lang="en",
736 tweet_mode='extended').items(20):
737             temp = {}
738             text = change(tweet.full_text)
739             temp["text"] = text
740             temp["username"] = tweet.user.screen_name
741             text = change(tweet.full_text)
742             prediction = predict_sentiment(learner,text)
743             prediction = sentiment_label(prediction)
744             temp["label"] = prediction
745             socketio.emit('tweets', {'Text': temp['text'], 'Username': temp['username'], 'Sentiment':
746 temp['label']}, namespace='/twitter')
747             print('Twitter feed end')
748             socketio.sleep(300)
749
750 @app.route("/twitter_live_feed", methods=['GET', 'POST'])
751 def twitter():
752     return render_template('twitter_live_feed.html')
753
754 @socketio.on('connect', namespace='/twitter')
755 def twitter_connect():
756     # need visibility of the global thread object
757     global twitter_thread
758     print(' Twitter Client connected')
759
760     #Start the tweet feed thread only if the thread has not been started before.
761     if not twitter_thread.isAlive():
762         print("Starting Twitter Thread")
763         thread = socketio.start_background_task(gettweets)
764
765 @socketio.on('disconnect', namespace='/twitter')
766 def twitter_disconnect():
767     print('Twitter Client disconnected')
768
769 #Run APP
770
771 if __name__ == '__main__':
772     socketio.run(app)

```

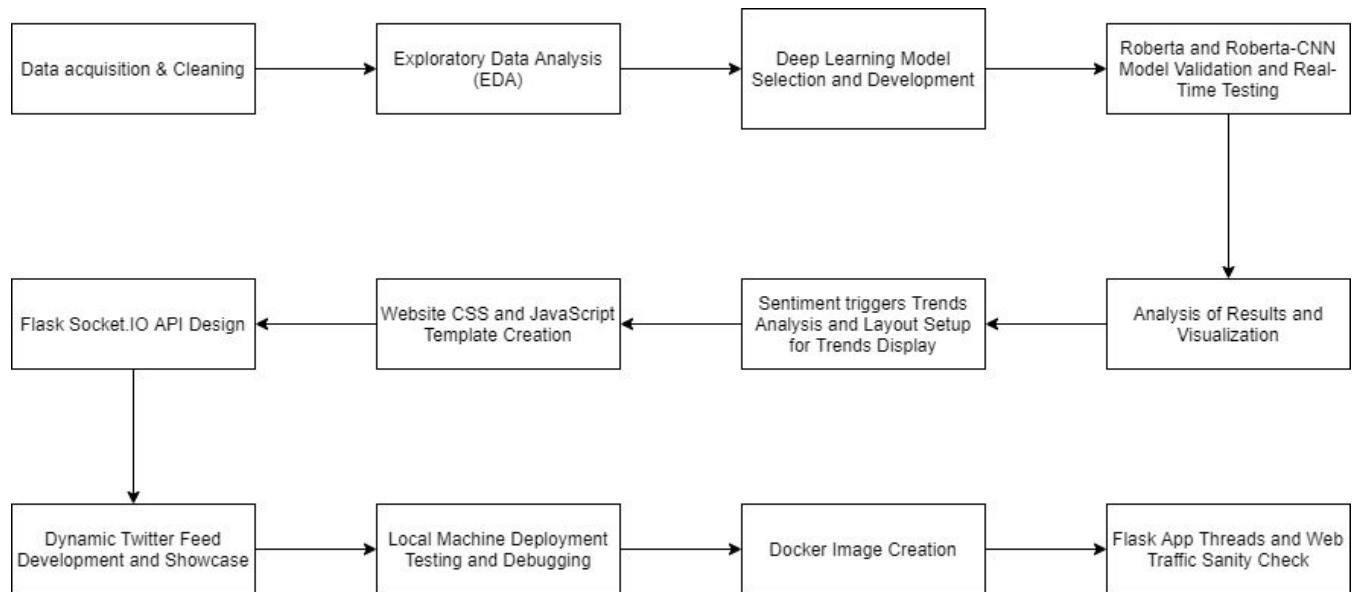
## 4. Experimental Investigations

During Exploratory Data Analysis (EDA), various types of graphs based on the sentiments and sentiment triggers were plotted to gain valuable insights from the data. The frequency distribution graphs give us a good idea about the data and also gives us an insight into predicting the model's generalization capability. From the frequency graphs and trends, it was clearly evident that the most frequently used terms by Twitterati in India and the rest of the world were "Corona" and "COVID-19" when compared to the other words. "rt", the most frequently used term in the world dataset, means "retweet". The high usage of the word in this pandemic could be attributed to the fact that people showed great enthusiasm in retweeting the tweets related to coronavirus. Also, the trigger "link" had many appearances in the graphs, this is because the trigger "link" was the code word used by the model in place of the URLs present in the tweets, this "link" has many appearances in the data set. From this, it can be understood that the model should not generalize and classify tweets based on the presence of the word "Corona", "COVID-19" and "rt". The trends from the graphs plotted further strengthen the belief that the Indian dataset is similar to the world dataset in terms of sentiments expressed.

The initial assumption about terms like positive and negative conveying the opposite meaning as opposed to the observed trends after the analysis and hence the assumption was dropped. Words like "fat" and "sick" have frequent occurrences in tweets, showing that people started paying attention to their health more than they previously did. "Love" was the most frequently used positive term and "nasty" was the most used negative term. By the nexus relationship graphs, we received better insights into the relationships between the sentiment triggers.

## 5. Flow Chart





## 6. Result

**CoVis** - The COVID-19 Twitter Sentiment Analysis Dashboard is the final product of the project. The dashboard contains many interactive tabs depicting various features like “Public Sentiment Dashboard”, “Real-Time Analysis Dashboard”, “Twitter Live Feed & Analysis” and “Live Case Counts” with live, robust, interactive graphs.

- *Public Sentiment Dashboard* - This tab contains an analysis of the tweets made by Twitterati in India in the past 90 days on dates when announcements like “Nation-wide Lockdown”, “Unlock 1.0”, etc. were made by the Government of India. These tweets were collected on the basis of the hashtags used in the tweets. These tweets were then passed to the Roberta model created to extract sentiments. Finally, these sentiments were represented in the form of various graphs on the dashboard.
- *Real-Time Analysis Dashboard* - Tweets are scraped from Twitter based on hashtags given as input by using the Tweepy API. A total of 2,500 tweets are scraped (which is the maximum scraping limit per day for the Tweepy API), their sentiments are extracted and live, interactive graphs are plotted which are updated every 24 hours.
- *Twitter Live Feed & Analysis* - A few tweets are scraped using the Tweepy API every 5 minutes and their Sentiment Analysis is done. This page has a live scrolling of the tweets extracted, an interactive 3-D plot depicting the sentiment frequency of those tweets, and a block that displays the twitter user ID, the tweet body, and their extracted sentiment. This page gets updated on its own, every 5 minutes.
- *Live Case Counts* - This dashboard contains an overview of the live COVID-19 statistics of India in a state-wise manner. It has categories of “Active cases”, “Recovered”, “Deaths” and “Number of tests”, which gets updated periodically.

## **7.ADVANTAGES & DISADVANTAGES**

### Advantages

The effect of COVID-19 pandemic is visible all over the world. National healthcare systems are facing the contagion with incredible strength, but concern regarding psychosocial and economic effects is critically growing. In a fast-moving crisis, as information swarms in from every direction, citizens look to their governments for information, guidance, and leadership. Sentimental Analysis is only the option in this current situation to understand the psychological condition/mental condition of the public. By Sentimental Analysis, the public opinion on COVID-19, regime policies, and actions can be understood. After Analysis, amendments can be made to the decisions taken by the regime policies, and the public can be fortified in such a way so as to enhance the sentiment towards a positive outlook. Not only this but also sentiment analysis will help NGOs and various organizations to come forward to help the people. Businesses can adapt their products and services to match the requirements of the people based on the real-time trending mood of the public, which will not only help businesses to grow but will also help the public meet their need of the hour. Also, this will enable the government to make business and people-friendly rules and laws to help in the betterment of the economy and the market in these untested times.

### Disadvantages

Natural Language Processing models, in general, face a problem in recognizing human aspects of a language like irony, sarcasm, negotiations, exaggerations, and jokes - the sorts of things humans wouldn't face many problems in understanding. Machines sometimes fail in recognizing these aspects, which leads to skewed and incorrect results.

## **8. APPLICATIONS**

The application of Sentiment Analysis lies in taking the sentiments of the people into consideration and creating appropriate business strategies and government policies so as to meet their needs. For example, many food delivery companies like Swiggy and online retail stores like Amazon are now opting for "Contactless Delivery" as a preventive step towards controlling the spread of the virus. This step taken by such companies would not only increase their profits but also would provide the essential services to the people while taking their sentiments into consideration in this pandemic. Not only this but also based on the sentiments analyzed from the tweets, rehabilitation or positive suggestions can be made to users who have consistently expressed negative sentiments. This way, public health can also be monitored using sentiment analysis.

## **9. CONCLUSION**

Sentiment analysis or opinion mining is a hot topic in deep learning. There is still a long way to go before sentiments can be accurately detected from texts, because of the complexity involved in the English language, and even more when other languages like Hindi are considered. Though the Roberta model developed as a part of this project has predicted and classified the sentiments of the test data set into positive, negative and neutral categories with an accuracy of 97%, by making necessary modifications and additions to the model, sentiment analysis can be done with greater accuracy by taking the language complexities into consideration.

## 10. FUTURE SCOPE

Sentiment Analysis still has many aspects that can still be worked upon. This project can be further enhanced by training the model about the human aspects of a language and making it more accurate in cases where sarcasm, irony, and other aspects are used. Taking the actions of all the ministries into consideration while gauging the sentiments of the public can also make the analysis more detailed and sector-specific, which would help in the analysis of the area of development required in those sectors. The classifier can be further improved by trying to extract more features from the tweets, trying different kinds of features, tuning the Hyperparameters, and also by making it work on various Indian languages.

## 11. BIBLIOGRAPHY

1. <https://towardsdatascience.com/twitter-data-collection-tutorial-using-python-3267d7cfa93e> (Data collection)
2. <https://ieee-dataport.org/open-access/coronavirus-covid-19-tweets-dataset> (Datasets)
3. <https://www.kaggle.com/ragnisah/text-data-cleaning-tweets-analysis> (Data cleaning)
4. <https://www.youtube.com/watch?v=ZhyLqPnOeh0> (Normalization)
5. <https://cognitiveclass.ai/courses/data-visualization-with-python> (Data Visualization)
6. <https://github.com/DocNow/hydrator/blob/master/README.md> (Hydrator)
7. <https://dphi.tech/data-science-bootcamp-day-15-exploratory-data-analysis/> (Exploratory Data Analysis)
8. <https://medium.com/swlh/understand-tweets-better-with-bert-sentiment-analysis-3b054c4b802a> (BERT)
9. <https://towardsdatascience.com/fasttext-sentiment-analysis-for-tweets-a-straightforward-guide-9a8c070449a2> (FastText)
10. <https://medium.com/swlh/using-xlNet-for-sentiment-classification-cfa948e65e85> (XLNet(This given paper is only for Binary Classification.))
11. <https://www.creative-tim.com/product/black-dashboard> (Dashboards)
12. <https://link.medium.com/ZZvca1Mf27> (Twitter live feed Sentimental Analysis)
13. <https://flask-socketio.readthedocs.io/en/latest/> (Flask Socket.IO)

14. <https://github.com/covid19india/api> (Live case count)

## APPENDIX

### A. Source code

<https://github.com/SmartPracticeschool/SBSPS-Challenge-2700-Twitter-Sentiment-Analysis-Extraction-for-COVID-19>