# DATA MODELING AND ANALYZING CORONA VIRUS SPREAD USING DATA SCIENCE AND PYTHON

A J-Component Report

*submitted by*

*K C S S D CHAKRADHAR 17BEC0873*

*V SAI YASHWANTH 17BEC0369*

*M SAI PRANAV 17BEC0708*

*B K HEM CHARAN 17BEC0189*

*K SURYA TEJA 17BEC0138*

*J KONG KON 17BEC0076*

*For the course*

**ECE-3999 TECHNICAL ANSWERS FOR REAL WORLD PROBLEMS**

*To*

**(Prof.) ILAVARASAN T**



**SCHOOL OF ELECTRONICS ENGINEERING**

**SLOT:** TD2

# ACKNOWLEDGEMENT

# **ABSTRACT**

In this project we tried to gather the information on how the covid-19 has spread throughout the world and how it mainly affected the countries INDIA and ITALY as ITALY is the first one to be attacked after CHINA and INDIA well is the country we live in. So we did what we can to analyze the data of the covid-19 which was available online and are able to visualize how this covid-19 has spread  and also were able to make a prediction on the new cases. we will be taking raw data from the most reliable sources and convert it into tables, graphs and to other forms of organized data for **machine learning engineers** or some other researchers which might help them a bit along their process of actively trying to fight COVID-19.We will be modelling and analyzing the data and try to predict the COVID-19 cases in the future, we will be fitting the data to a familiar model for pandemics called **SIR** (susceptible infected recovered) which is a part of compartmental model techniques. Data modeling is a set of tools and techniques used to understand and analyze how an organization should collect, update, and store data. It is a critical skill for the business analyst who is involved with discovering, analyzing, and specifying changes to how software systems create and maintain information.

# **OBJECTIVE**

- In this present situation we are not able to understand how this COVID-19 is spreading around the world. One doesn't know how many cases are rising day by day in different places and the factors that are taken by governments in different countries.
- So our main idea is to make a SIR model to divide the people who are infected and recovered and who are susceptible to COVID-19 in countries like **India** and **Italy**.
- To provide a modelled data and to help people analyze and improve the current situation of the pandemic.
- Our Project is an attempt of data modelling and analyzing Coronavirus (COVID-19) spread with the help of data science and data analytics in python.

# **INTRODUCTION**

- Presently, there are endless dashboards and measurements around the Coronavirus spread accessible everywhere on the web.
- With this so much data and master assessments, to see various countries receiving various procedures, from complete lockdown to social separating to crowd insusceptibility, Someone is left intuition concerning what the correct system is for them.
- So this is an endeavor of information displaying and dissecting Coronavirus (COVID-19) spread with the assistance of information science and information examination in python code.
- This examination will assist us with finding the premise behind basic ideas about the infection spread from absolutely a dataset point of view.

# CORONA VIRUS IN ITALY

## HOW COVID-19 STARTED IN ITALY:

The virus was first confirmed to have spread to Italy on 31 January 2020, when two Chinese tourists in Rome tested positive for the virus. One week later an Italian man repatriated back to Italy from the city of Wuhan, China, was hospitalized and confirmed as the third case in Italy. Clusters of cases were later detected in Lombardy and Veneto on 21 February, with the first deaths on 22 February. By the beginning of March, the virus had spread to all regions of Italy. On 6 March 2020, the Italian College of Anesthesia, Analgesia, Resuscitation and Intensive Care (SIAARTI) published medical ethics recommendations regarding triage protocols that needed to be employed.

## TIMELINE

January 2020: First confirmed cases

February–March 2020: Clusters in Northern Italy

March 2020: Spread to other regions

March–May 2020: Under national lockdown

May–September 2020: Reduction of cases and loosening of restrictions

September 2020–ongoing: Arrival of the second wave

## MEASURES TAKEN BY GOVERNMENT IN ITALY:

In Italy, after the shutdown of the educational system (schools and Universities will remain closed at least for one month) and the collapse of the touristic sector (90% of travels and reservations cancelled), the Government officially locked down residents of all the region of Milan (Lombardia) and other 11 provinces. To avoid the imprisonment, hundreds thousand people left those areas with any possible mean in the night of March 7th, just before the law was signed by the prime Minister Giuseppe Conte, thus turning his purpose of slowing down the epidemics exactly into the opposite. The Governors of Southern regions adopted limitations for this huge mass of potentially infected incoming people, with the risk of disseminating suspicion in the population ("hunting the greaser"). Violence exploded in several prisons, due to the cancellation of all the family visits and fear of the virus. Just 48 hours later, the Italian Government has extended these exceptional war-
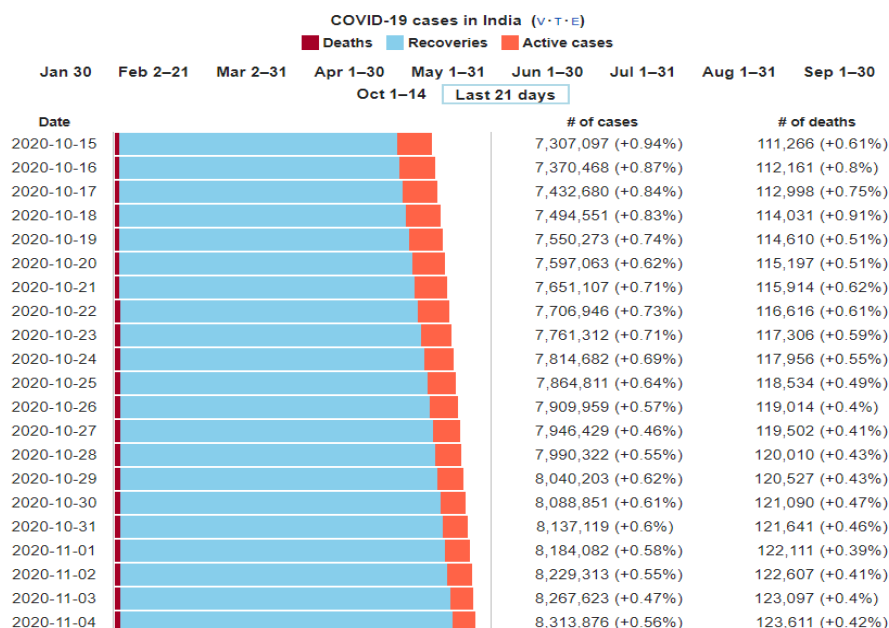
like measures to the entire nation: nobody is allowed to exit from home other than for compelling job or health reasons; museums, cinemas, theatres, sport facilities, and even churches have been closed; restaurants and bars must stop at 6 PM their activity. Further restrictions are probably going to be considered in next days. The Italian prime minister was clear in his video-message of March 5th: the emergency does not come from the lethality of the virus, but from the impossibility of the Italian healthcare system to cope with the impact of a rapid epidemic spreading of the Covid-19.

# CORONA VIRUS IN INDIA

## HOW COVID-19 STARTED IN INDIA:

The first case of COVID-19 in India, which originated from China, was reported on 30 January 2020. India currently has the largest number of confirmed cases in Asia, and has the second-highest number of confirmed cases in the world after the United States, with almost 8 million reported cases of COVID-19 infection, more than 1 lakh deaths and more than 7 million recovered. By mid of 2020, India had approached in position of conducting highest number of daily tests in the world which subsequently translated into highest number of daily new cases in world and has sustained highest number of daily cases spike since then.

## TIMELINE

COVID-19 cases in India (v·T·E)
Deaths   Recoveries   Active cases

| Jan 30 | Feb 2–21 | Mar 2–31 | Apr 1–30 | May 1–31 | Jun 1–30 | Jul 1–31 | Aug 1–31 | Sep 1–30 |

Oct 1–14   Last 21 days

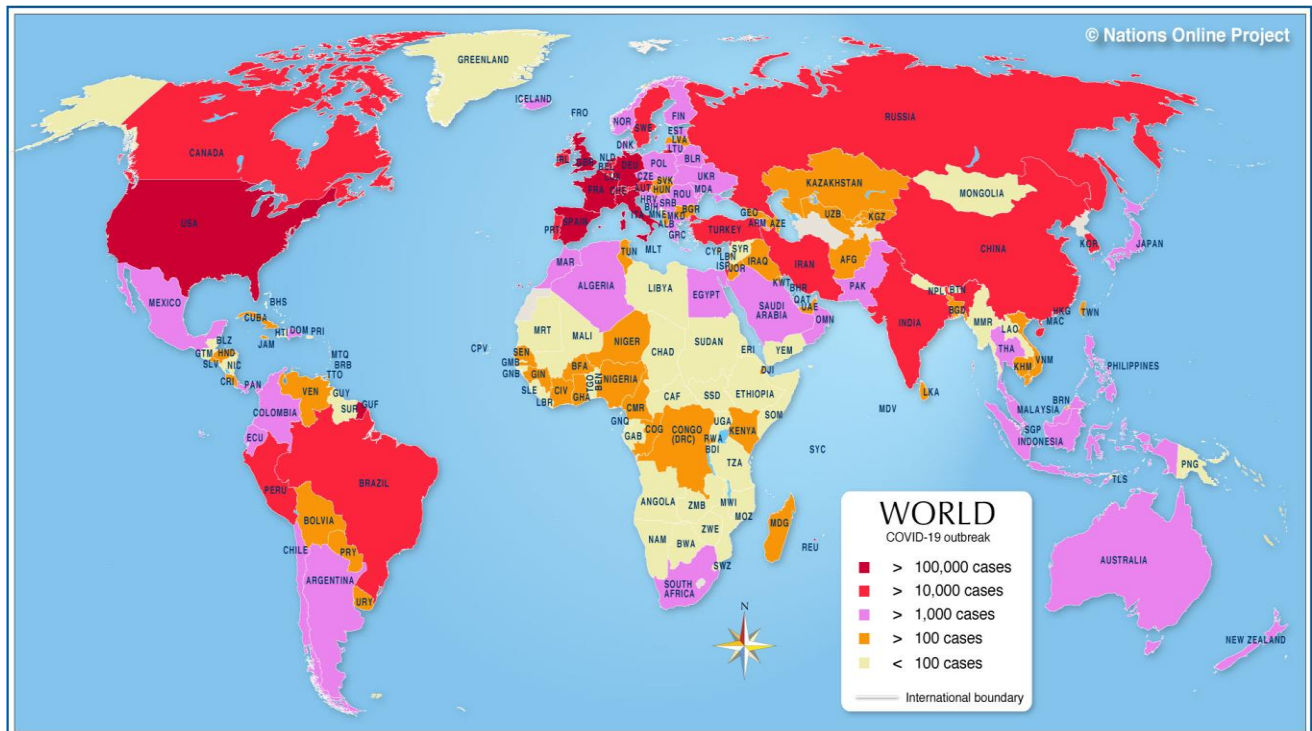| Date | # of cases | # of deaths |
| --- | --- | --- |
| 2020-10-15 | 7,307,097 (+0.94%) | 111,266 (+0.61%) |
| 2020-10-16 | 7,370,468 (+0.87%) | 112,161 (+0.8%) |
| 2020-10-17 | 7,432,680 (+0.84%) | 112,998 (+0.75%) |
| 2020-10-18 | 7,494,551 (+0.83%) | 114,031 (+0.91%) |
| 2020-10-19 | 7,550,273 (+0.74%) | 114,610 (+0.51%) |
| 2020-10-20 | 7,597,063 (+0.62%) | 115,197 (+0.51%) |
| 2020-10-21 | 7,651,107 (+0.71%) | 115,914 (+0.62%) |
| 2020-10-22 | 7,706,946 (+0.73%) | 116,616 (+0.61%) |
| 2020-10-23 | 7,761,312 (+0.71%) | 117,306 (+0.59%) |
| 2020-10-24 | 7,814,682 (+0.69%) | 117,956 (+0.55%) |
| 2020-10-25 | 7,864,811 (+0.64%) | 118,534 (+0.49%) |
| 2020-10-26 | 7,909,959 (+0.57%) | 119,014 (+0.4%) |
| 2020-10-27 | 7,946,429 (+0.46%) | 119,502 (+0.41%) |
| 2020-10-28 | 7,990,322 (+0.55%) | 120,010 (+0.43%) |
| 2020-10-29 | 8,040,203 (+0.62%) | 120,527 (+0.43%) |
| 2020-10-30 | 8,088,851 (+0.61%) | 121,090 (+0.47%) |
| 2020-10-31 | 8,137,119 (+0.6%) | 121,641 (+0.46%) |
| 2020-11-01 | 8,184,082 (+0.58%) | 122,111 (+0.39%) |
| 2020-11-02 | 8,229,313 (+0.55%) | 122,607 (+0.41%) |
| 2020-11-03 | 8,267,623 (+0.47%) | 123,097 (+0.4%) |
| 2020-11-04 | 8,313,876 (+0.56%) | 123,611 (+0.42%) |

## MEASURES TAKEN BY GOVERNMENT IN INDIA:

- The Delhi government declared coronavirus as an epidemic in the state with 6 confirmed cases. All schools, colleges and cinema halls in Delhi will remain shut till March 31 as a measure to counter the coronavirus, Chief Minister Arvind Kejriwal announced on Thursday, March 12. The Delhi government had shut the primary schools earlier this month, the secondary classes were left open in view of the exams. The Chief Minister also declared that all cinema halls will remain shut in Delhi till 31st March. Schools and colleges where exams are not being held will also remain closed. Delhi government has also made disinfecting all public places, including government, private offices and shopping malls compulsory. Furthermore, vacant flats owned by Delhi Urban Shelter Improvement Board will be used for quarantine, says CM Kejriwal.

- Over 1,500 people under observation for coming in contact with 73 positive cases of coronavirus, as revealed by the Union Health Ministry. 10.5 lakh people screened so far at 30 designated airports in India, says the ministry.

- Union Ministry of Road Transport and Highways advised states and Union Territories to take all steps for sanitisation of public transport vehicles and terminals. This is to ensure sanitation of seats, handles and bars at all bus terminals are disinfected. The ministry also suggested that the public transport should display public health messages in vehicles, bus terminals and bus stops and asked states and UTs to take expeditious action and mobilise all necessary support in this regard.

- All educational institutions, stadiums and sports clubs in Srinagar are closed from till further orders amid the coronavirus scare as a precautionary step, the city administration has said. Srinagar Mayor Junaid Azim Mattu said that its "an unavoidable decision" to allow the Srinagar Municipal Corporation (SMC) to plan, sterilize and sanitize schools and colleges. Under a special set of statutory provisions, SMC has also ordered closure of all educational institutions, public clubs, sports clubs, indoor and open stadiums, coaching centres within Srinagar city limits till further orders. The mayor said the Corporation will procure 1,000 full quarantine body kits and spraying machines and all public as well as private

hospitals will go into heightened sanitation mode. The Corporation also ordered cancellation of all sports events within its limits and phased segregation of flea markets, including weekly markets, with an immediate effect. It has also issued an advisory asking people to desist shopping, especially eatables and garments, from roadside vendors.

# COUNTRIES SUFFERING WITH COVID-19 DISEASE

# HOW TO ANALYZE THIS COVID-19 EFFECTIVELY

So everyone looks at this problem in their own way. some will think of analyzing using different charts and online tools and some will think of some mathematical way of analyzing it.

We came up with a solution of analyzing this COVID-19 by using the data modelling which is a part of data science and which helps us in viewing this problem easily.

## HOW DID WE START?

- We first gone through the tools required in order to perform the data modelling and get to know the tool as python which is used by the software **Anaconda.**
- So in order to create a model we need a IDE so we started our search for that and found the best suitable one as **Jupyter Notebook**.
- Now we require a base for analyzing that is nothing but the **data.** So we started collecting the datasets from a trusted site and started our working.
- The datasets we used are:

  https://raw.githubusercontent.com/CSSEGISandData/COVID-19/master/csse_covid_19_data/csse_covid_19_time_series/time_series_covid19_confirmed_global.csv
  https://raw.githubusercontent.com/CSSEGISandData/COVID-19/master/csse_covid_19_data/csse_covid_19_time_series/time_series_covid19_deaths_global.csv
  https://raw.githubusercontent.com/CSSEGISandData/COVID-19/master/csse_covid_19_data/csse_covid_19_time_series/time_series_covid19_recovered_global.csv
  https://raw.githubusercontent.com/CSSEGISandData/COVID-19/web-data/data/cases_country.csv

## CONFIRMED CASES:

| | Province/State | Country/Region | Lat | Long | 1/22/20 | 1/23/20 | 1/24/20 | 1/25/20 | 1/26/20 | 1/27/20 | ... | 10/25/20 | 10/26/20 | 10/27/20 | 10/28/20 | 10/29/ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | NaN | Afghanistan | 33.939110 | 67.709953 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 40833 | 40937 | 41032 | 41145 | 412 |
| 1 | NaN | Albania | 41.153300 | 20.168300 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 19157 | 19445 | 19729 | 20040 | 203 |
| 2 | NaN | Algeria | 28.033900 | 1.659600 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 56143 | 56419 | 56706 | 57026 | 573 |
| 3 | NaN | Andorra | 42.506300 | 1.521800 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 4038 | 4325 | 4410 | 4517 | 45 |
| 4 | NaN | Angola | -11.202700 | 17.873900 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 9381 | 9644 | 9871 | 10074 | 102 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 264 | NaN | Western Sahara | 24.215500 | -12.885800 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 10 | 10 | 10 | 10 | |
| 265 | NaN | Yemen | 15.552727 | 48.516388 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 2060 | 2060 | 2060 | 2061 | 20 |
| 266 | NaN | Zambia | -13.133897 | 27.849332 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 16117 | 16200 | 16243 | 16285 | 163 |
| 267 | NaN | Zimbabwe | -19.015438 | 29.154857 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 8276 | 8303 | 8315 | 8320 | 83 |
| 268 | NaN | NaN | NaN | NaN | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 | |

And it will be continued till present day

## DEATHS GLOBAL:

| | Province/State | Country/Region | Lat | Long | 1/22/20 | 1/23/20 | 1/24/20 | 1/25/20 | 1/26/20 | 1/27/20 | ... | 10/25/20 | 10/26/20 | 10/27/20 | 10/28/20 | 10/29/ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | NaN | Afghanistan | 33.939110 | 67.709953 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 1514 | 1518 | 1523 | 1529 | 15 |
| 1 | NaN | Albania | 41.153300 | 20.168300 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 477 | 480 | 487 | 493 | 4 |
| 2 | NaN | Algeria | 28.033900 | 1.659600 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 1914 | 1922 | 1931 | 1941 | 19 |
| 3 | NaN | Andorra | 42.506300 | 1.521800 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 69 | 72 | 72 | 72 | |
| 4 | NaN | Angola | -11.202700 | 17.873900 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 268 | 270 | 271 | 275 | 2 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 264 | NaN | Western Sahara | 24.215500 | -12.885800 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 1 | 1 | 1 | 1 | |
| 265 | NaN | Yemen | 15.552727 | 48.516388 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 599 | 599 | 599 | 599 | 5 |
| 266 | NaN | Zambia | -13.133897 | 27.849332 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 348 | 348 | 348 | 348 | 3 |
| 267 | NaN | Zimbabwe | -19.015438 | 29.154857 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 237 | 242 | 242 | 242 | 2 |
| 268 | NaN | NaN | NaN | NaN | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 | |

and it will be continued to the present day

## RECOVERED CASES:

| | Province/State | Country/Region | Lat | Long | 1/22/20 | 1/23/20 | 1/24/20 | 1/25/20 | 1/26/20 | 1/27/20 | ... | 10/25/20 | 10/26/20 | 10/27/20 | 10/28/20 | 10/29/ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | NaN | Afghanistan | 33.939110 | 67.709953 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 34129 | 34150 | 34217 | 34237 | 342 |
| 1 | NaN | Albania | 41.153300 | 20.168300 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 10654 | 10705 | 10808 | 10893 | 110 |
| 2 | NaN | Algeria | 28.033900 | 1.659600 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 39095 | 39273 | 39444 | 39635 | 396 |
| 3 | NaN | Andorra | 42.506300 | 1.521800 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 2729 | 2957 | 3029 | 3144 | 32 |
| 4 | NaN | Angola | -11.202700 | 17.873900 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 3508 | 3530 | 3647 | 3693 | 37 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 251 | NaN | Western Sahara | 24.215500 | -12.885800 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 8 | 8 | 8 | 8 | |
| 252 | NaN | Yemen | 15.552727 | 48.516388 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 1360 | 1360 | 1364 | 1366 | 13 |
| 253 | NaN | Zambia | -13.133897 | 27.849332 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 15179 | 15445 | 15481 | 15559 | 155 |
| 254 | NaN | Zimbabwe | -19.015438 | 29.154857 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 7797 | 7797 | 7804 | 7845 | 78 |
| 255 | NaN | NaN | NaN | NaN | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 | |

And it will be continued to the present day

## COUNTRY WIDE CASES:

| | Country_Region | Last_Update | Lat | Long_ | Confirmed | Deaths | Recovered | Active | Incident_Rate | People_Tested | People_Hospitalized | Mortality |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Afghanistan | 2020-11-04 07:24:54 | 33.939110 | 67.709953 | 41728.0 | 1544.0 | 34355.0 | 5829.0 | 107.191827 | NaN | NaN | 3.7 |
| 1 | Albania | 2020-11-04 07:24:54 | 41.153300 | 20.168300 | 21904.0 | 532.0 | 11473.0 | 9899.0 | 761.136980 | NaN | NaN | 2.4 |
| 2 | Algeria | 2020-11-04 07:24:54 | 28.033900 | 1.659600 | 58979.0 | 1980.0 | 40577.0 | 16422.0 | 134.498511 | NaN | NaN | 3.3 |
| 3 | Andorra | 2020-11-04 07:24:54 | 42.506300 | 1.521800 | 4910.0 | 75.0 | 3627.0 | 1208.0 | 6354.753122 | NaN | NaN | 1.5 |
| 4 | Angola | 2020-11-04 07:24:54 | -11.202700 | 17.873900 | 11577.0 | 291.0 | 5230.0 | 6056.0 | 35.224565 | NaN | NaN | 2.5 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 185 | West Bank and Gaza | 2020-11-04 07:24:54 | 31.952200 | 35.233200 | 55408.0 | 501.0 | 47744.0 | 7163.0 | 1086.129812 | NaN | NaN | 0.9 |
| 186 | Western Sahara | 2020-11-04 07:24:54 | 24.215500 | -12.885800 | 10.0 | 1.0 | 8.0 | 1.0 | 1.674116 | NaN | NaN | 10.0 |
| 187 | Yemen | 2020-11-04 07:24:54 | 15.552727 | 48.516388 | 2063.0 | 601.0 | 1375.0 | 87.0 | 6.916791 | NaN | NaN | 29.1 |
| 188 | Zambia | 2020-11-04 07:24:54 | -13.133897 | 27.849332 | 16661.0 | 349.0 | 15763.0 | 549.0 | 90.627937 | NaN | NaN | 2.0 |
| 189 | Zimbabwe | 2020-11-04 07:24:54 | -19.015438 | 29.154857 | 8410.0 | 246.0 | 7942.0 | 222.0 | 56.583740 | NaN | NaN | 2.9 |

- So we successfully collected the required datasets but it has to be modelled in such a way that it can be visualized.
- In order to do that we have to undergo some advanced python coding and use some of the online resources which help us reach our desired result.
- So first we tried to eliminate the empty data values which we can see as NA or NaN because they shouldn't be in a model which is going to undergo machine learning this process is called **imputing**.
- Next we tried to group all the data into one to plot the graphs required to visualize.
- Now we first tried to plot the total confirmed cases around the globe by removing the unnecessary data from the dataset we have and arrived at the graph shown below:

- Now we started working on how to create a custom plot function with the help of the online resources and also understood what is color array and how it is useful in plotting the graph.

```python
# Initializing Color Array to be used across the analysis
color_arr = px.colors.qualitative.Dark24
```

```python
def draw_plot(ts_array, ts_label, title, colors, mode_size, line_size, x_axis_title , y_axis_title, tickangle = 0, yaxis_type = '
    # initialize figure
    fig = go.Figure()
    # add all traces
    for index, ts in enumerate(ts_array):
        fig.add_trace(go.Scatter(x=ts.index,
                                 y = ts.values,
                                 name = ts_label[index],
                                 line=dict(color=colors[index], width=line_size[index]),connectgaps=True,))
    # base x_axis prop.
    x_axis_dict = dict(showline=True,
                       showgrid=True,
                       showticklabels=True,
                       linecolor='rgb(204, 204, 204)',
                       linewidth=2,
                       ticks='outside',
                       tickfont=dict(family='Arial',size=12,color='rgb(82, 82, 82)',))
    # setting x_axis params
    if x_axis_title:
        x_axis_dict['title'] = x_axis_title

    if tickangle > 0:
        x_axis_dict['tickangle'] = tickangle

    # base y_axis prop.
    y_axis_dict = dict(showline = True,
                       showgrid = True,
                       showticklabels=True,
                       linecolor='rgb(204, 204, 204)',
                       linewidth=2,)
    # setting y_axis params
    if yaxis_type != "":
        y_axis_dict['type'] = yaxis_type

    if y_axis_title:
        y_axis_dict['title'] = y_axis_title


    fig.update_layout(xaxis = x_axis_dict,
                      yaxis = y_axis_dict,
                      autosize=True,
                      margin=dict(autoexpand=False,l=100,r=20,t=110,),
                      showlegend=True,
                      )

    # base annotations for any graph
    annotations = []
    # Title
    annotations.append(dict(xref='paper', yref='paper', x=0.0, y=1.05, xanchor='left', yanchor='bottom',
                            text=title,
                            font=dict(family='Arial',size=16,color='rgb(37,37,37)'),showarrow=False))
    # adding annotations in params
    if len(additional_annotations) > 0:
        annotations.append(additional_annotations)

    #updating the layout
    fig.update_layout(annotations=annotations)

    return fig
```
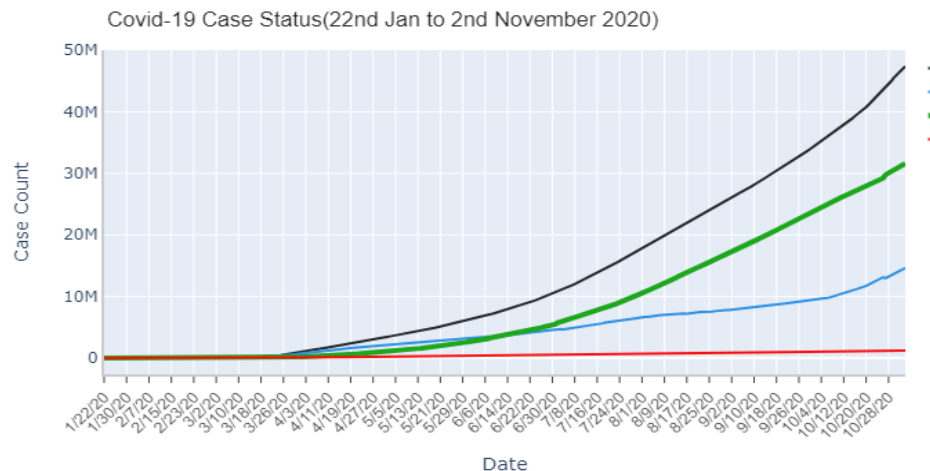
With this we arrived at creating a custom plot function.

- Now we tried to analyze the total active cases count on each day the problem that we faced here is there is no dataset of active cases in our list so we found out the active cases using the formula

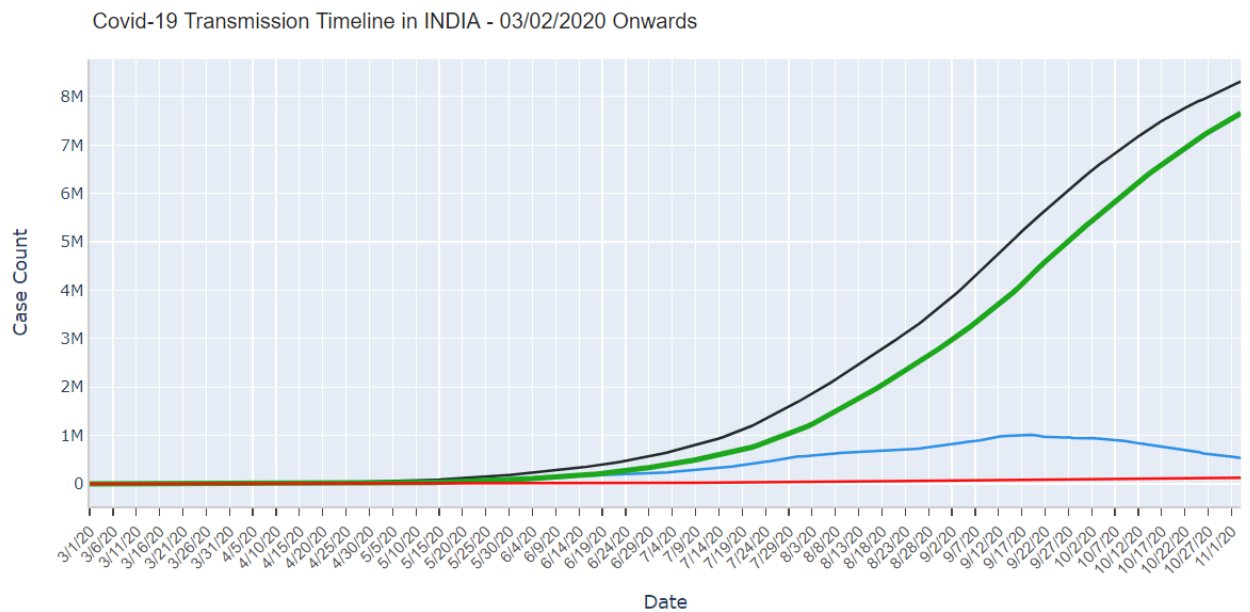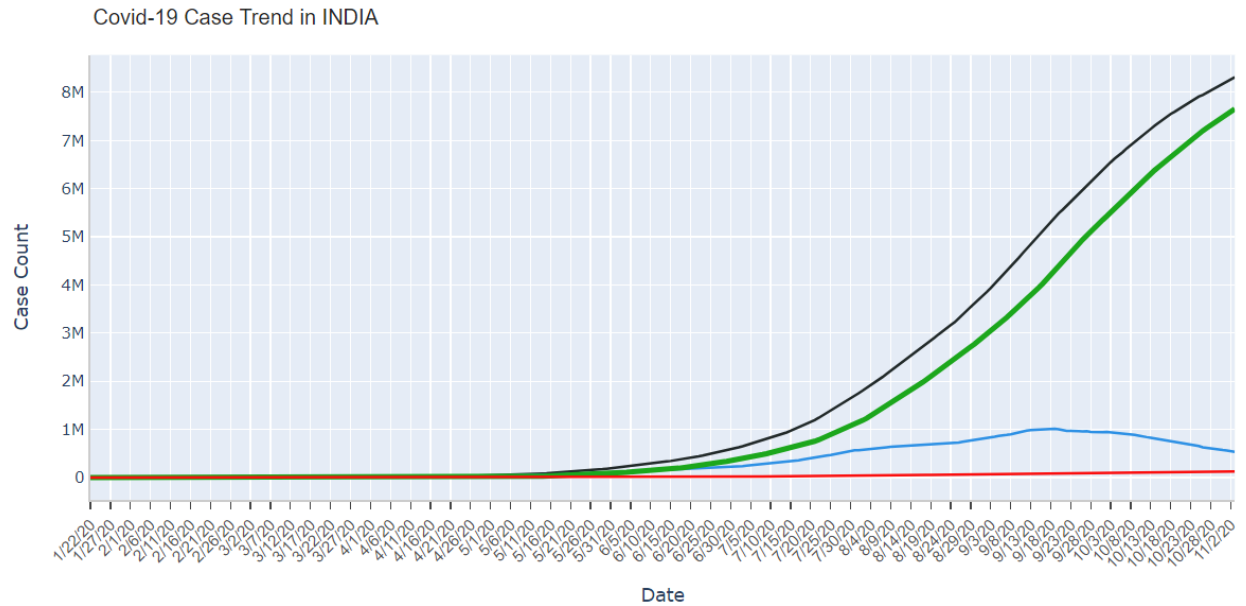  Active_cases = Confirmed_cases – deaths – recovered.
- We tried to plot these active cases using our custom plot function and you can observe the figure shown below:
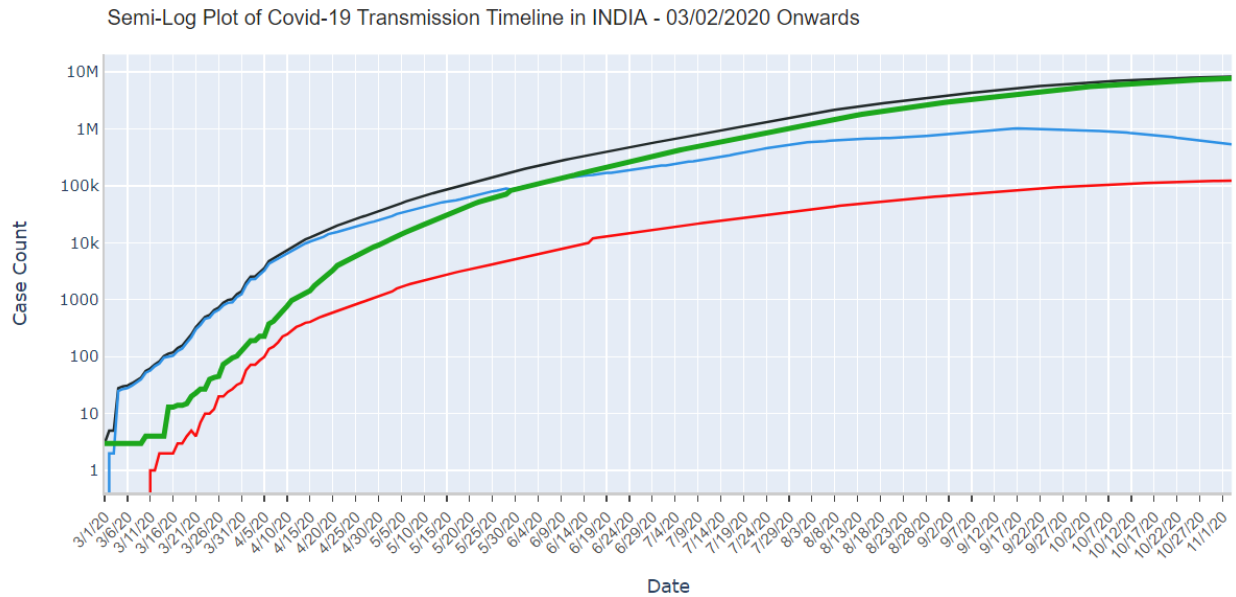


- Here we tried to give total case count of confirmed, deaths, recovered of all the countries summed up.

| | Country_Region | Confirmed | Deaths | Recovered | Active | Incident_Rate | People_Tested | People_Hospitalized | Mortality_Rate | U |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | US | 9384277.000000 | 232627.000000 | 3705130.000000 | 5446520.000000 | 2848.326971 | | | 2.478902 | 8 |
| 1 | India | 8313876.000000 | 123611.000000 | 7656478.000000 | 533787.000000 | 602.452868 | | | 1.486804 | 3 |
| 2 | Brazil | 5566049.000000 | 160496.000000 | 5060697.000000 | 344856.000000 | 2618.585094 | | | 2.883482 | |
| 3 | Russia | 1661096.000000 | 28611.000000 | 1244012.000000 | 388473.000000 | 1138.247950 | | | 1.722417 | 6 |
| 4 | France | 1461391.000000 | 37492.000000 | 123664.000000 | 1296627.000000 | 2238.872944 | | | 2.565501 | 2 |
| 5 | Spain | 1259366.000000 | 36495.000000 | 150376.000000 | 1072495.000000 | 2693.555438 | | | 2.897887 | 7 |
| 6 | Argentina | 1195276.000000 | 32052.000000 | 1009278.000000 | 153946.000000 | 2644.663018 | | | 2.681556 | |
| 7 | Colombia | 1099392.000000 | 31847.000000 | 993877.000000 | 73668.000000 | 2160.632247 | | | 2.896783 | 1 |
| 8 | United Kingdom | 1077099.000000 | 47340.000000 | 2906.000000 | 1026853.000000 | 1586.628961 | | | 4.395139 | 8 |
| 9 | Mexico | 938405.000000 | 92593.000000 | 803086.000000 | 42726.000000 | 734.320536 | | | 9.867062 | 4 |
| 10 | Peru | 902503.000000 | 34476.000000 | 832929.000000 | 35098.000000 | 2737.192816 | | | 3.820043 | 6 |

- Now we started our focus mainly on INDIA and tried to visualize the data whatever we have related to it

**Covid-19 Case Trend in INDIA**



**Covid-19 Transmission Timeline in INDIA - 03/02/2020 Onwards**

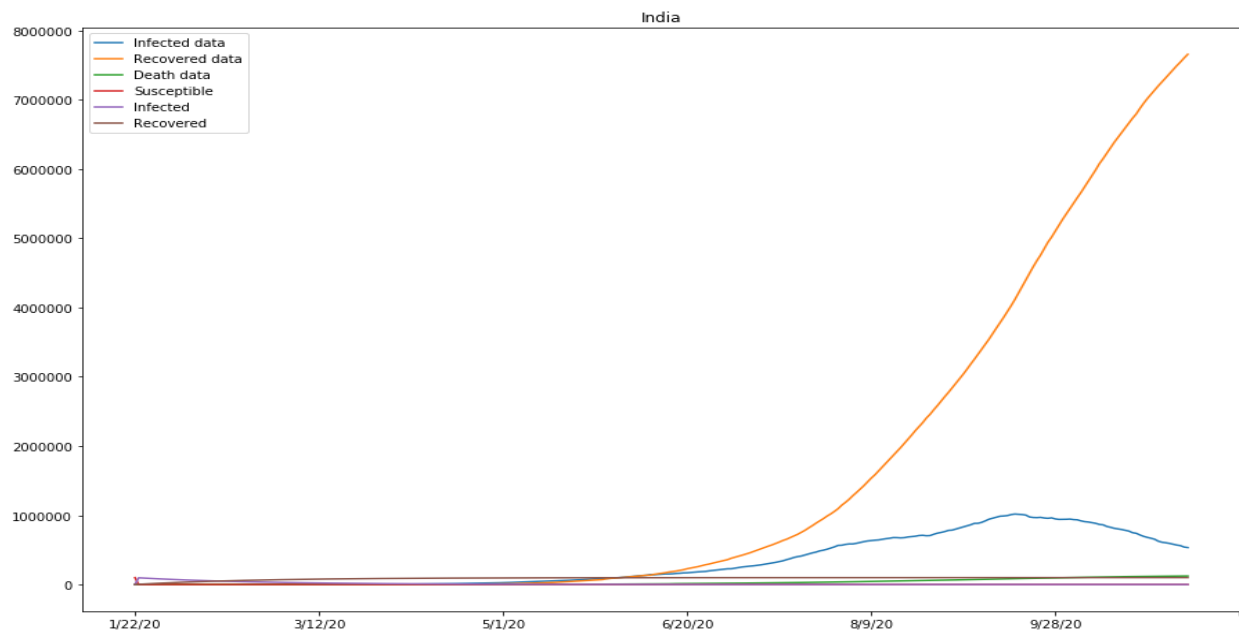Semi-Log Plot of Covid-19 Transmission Timeline in INDIA - 03/02/2020 Onwards

- Now we started making a predictive SIR model with using L-BFGS-B(large scale bound constrained optimization) for comparing different variables and the SIR model that is used is differential equation model.
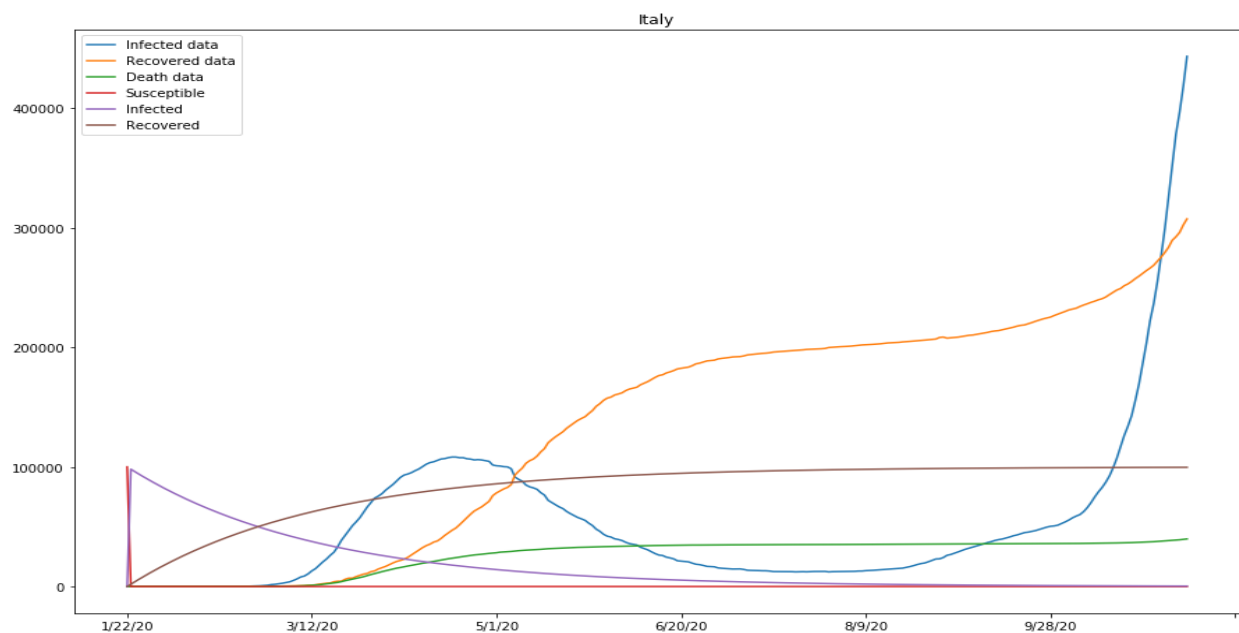
   ### What is a SIR model?
   The **SIR** model is a simple mathematical model of epidemics. An epidemic is when the number of people infected with a disease is increasing in a population. S, I, and R stand for: S – susceptible, I – infected, R – recovered.

## SIR MODEL PREDICTION FOR INDIA

```
        fun: 2454550.0801889934
   hess_inv: <2x2 LbfgsInvHessProduct with dtype=float64>
        jac: array([0.3259629 , 1.95577741])
    message: b'CONVERGENCE: REL_REDUCTION_OF_F_<=_FACTR*EPSMCH'
       nfev: 81
        nit: 8
     status: 0
    success: True
          x: array([0.39999942, 0.03051931])
country=India, beta=0.39999942, gamma=0.03051931, r_0:13.10643882
```

## SIR MODEL PREDICTION FOR ITALY



Italy

```
        fun: 90702.35410836339
  hess_inv: <2x2 LbfgsInvHessProduct with dtype=float64>
       jac: array([ 0.27066562, -0.56024874])
   message: b'CONVERGENCE: REL_REDUCTION_OF_F_<=_FACTR*EPSMCH'
      nfev: 54
       nit: 6
    status: 0
   success: True
         x: array([0.39844962, 0.01956381])
country=Italy, beta=0.39844962, gamma=0.01956381, r_0:20.36666971
```

# CONCLUSION

As our model has shown an increase in the rise of cases in INDIA the country has to take utmost measures in order to stop this pandemic COVID-19 or else there will be a huge risk ahead. However as our plots shown the number of recovered people is rising gradually so we can understand that the people are starting to develop immunity in their bodies which helps in fighting the COVID-19 as we can also see that the number of deaths are decreasing so we can predict that the world is getting back to normal state but this has to be continued till the arrival of the corona virus vaccine.

# REFERENCES

https://www.thelancet.com/journals/eclinm/article/PIIS2589-5370(20)30064-X/fulltext

https://swachhindia.ndtv.com/coronavirus-here-are-the-steps-taken-by-india-to-control-the-spread-of-covid-19-42304/

https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html

https://mathworld.wolfram.com/SIRModel.html

https://en.wikipedia.org/wiki/Compartmental_models_in_epidemiology

https://in.springboard.com/blog/data-modelling-covid/##

**THANK YOU.**