

India's COVID-19 Exploratory analysis



About

This document contains basic exploratory data analysis of COVID-19 Disease in India. This notebook serves to analyze and visualize the progress of the pandemic from various perspectives.

Introduction

The first signs of **COVID-19** in India was reported in some towns of Kerala, among three Indian medical students who had returned from Wuhan. After that, the Government of India had announced lockdown on **25 March 2020**. India faced its **first wave** from May 2020 to January 2020 with an Amplitude of around **90,000** new infections a day. As of now India is going under second wave which has proved to be more deadlier than previous one.

1. Cases, Deaths and Recovery

```
In [1]: import pandas as pd
from matplotlib import pyplot as plt
from matplotlib import dates as mldates
```

```
In [2]: ind_covid_df = pd.read_csv('https://api.covid19india.org/csv/latest/case_time_series.csv')
```

```
In [3]: ind_covid_df
```

| | Date | Date_YMD | Daily Confirmed | Total Confirmed | Daily Recovered | Total Recovered | Daily Deceased | Total Deceased |
|-----|-----------------|------------|-----------------|-----------------|-----------------|-----------------|----------------|----------------|
| 0 | 20 January 2020 | 2020-01-30 | 1 | 1 | 0 | 0 | 0 | 0 |
| 1 | 31 January 2020 | 2020-01-31 | 0 | 1 | 0 | 0 | 0 | 0 |
| 2 | 1 February 2020 | 2020-02-01 | 0 | 1 | 0 | 0 | 0 | 0 |
| 3 | 2 February 2020 | 2020-02-02 | 1 | 2 | 0 | 0 | 0 | 0 |
| 4 | 3 February 2020 | 2020-02-03 | 1 | 3 | 0 | 0 | 0 | 0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 499 | 12 June 2021 | 2021-06-12 | 80525 | 29438859 | 132964 | 28035743 | 3300 | 369816 |
| 500 | 13 June 2021 | 2021-06-13 | 71001 | 29509860 | 119574 | 28155317 | 3922 | 373738 |
| 501 | 14 June 2021 | 2021-06-14 | 60008 | 29569868 | 117376 | 28272693 | 2733 | 376471 |
| 502 | 15 June 2021 | 2021-06-15 | 62214 | 29632082 | 107767 | 28380460 | 2540 | 379011 |
| 503 | 16 June 2021 | 2021-06-16 | 67289 | 29699371 | 103900 | 28484360 | 2329 | 381340 |

504 rows × 8 columns

```
In [4]: ind_covid_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 504 entries, 0 to 503
Data columns (total 8 columns):
 #   Column      Non-Null Count  Dtype
---  ---
 0   Date        504 non-null    object
 1   Date_YMD    504 non-null    object
 2   Daily Confirmed  504 non-null    int64
 3   Total Confirmed  504 non-null    int64
 4   Daily Recovered  504 non-null    int64
 5   Total Recovered  504 non-null    int64
 6   Daily Deceased  504 non-null    int64
 7   Total Deceased  504 non-null    int64
dtypes: int64(6), object(2)
memory usage: 31.6+ KB
```

```
In [5]: ind_covid_df.isnull().sum()
```

```
Out[5]: Date          0
Date_YMD         0
Daily Confirmed   0
Total Confirmed   0
Daily Recovered   0
Total Recovered   0
Daily Deceased    0
Total Deceased    0
dtype: int64
```

```
In [6]: ind_covid_df['Date_YMD'] = pd.to_datetime(ind_covid_df['Date_YMD'])
```

```
In [7]: ind_covid_df.tail(1)
```

| | Date | Date_YMD | Daily Confirmed | Total Confirmed | Daily Recovered | Total Recovered | Daily Deceased | Total Deceased |
|-----|--------------|------------|-----------------|-----------------|-----------------|-----------------|----------------|----------------|
| 503 | 16 June 2021 | 2021-06-16 | 67289 | 29699371 | 103900 | 28484360 | 2329 | 381340 |

```
In [8]: total_cases = ind_covid_df['Total Confirmed']
dates = ind_covid_df['Date_YMD']
```

```
In [9]: curr_date = dates.max()
curr_total_cases = int(total_cases.tail(1))
```

```
In [10]: dates.max()
```

```
Out[10]: Timestamp('2021-06-16 00:00:00')
```

```
In [11]: filt = ind_covid_df.Date_YMD==dates.max()
today_cases = int(ind_covid_df.loc[filt, 'Daily Confirmed'])
today_recovered = int(ind_covid_df.loc[filt, 'Daily Recovered'])
today_deaths = int(ind_covid_df.loc[filt, 'Daily Deceased'])
curr_total_deaths = int(ind_covid_df.loc[filt, 'Total Deceased'])
curr_total_recovered = int(ind_covid_df.loc[filt, 'Total Recovered'])
```

```
In [12]: plt.style.use('fivethirtyeight')
# total_cases.plot(figsize=(10,6))
plt.figure(figsize=(10,6))
plt.plot(dates, total_cases.values/10**6, color='#000000')
# plt.plot(total_deaths.index, total_deaths.values, linewidth=1)

plt.gca().autofmt_xdate()
date_format = mldates.DateFormatter('%d %b, %y')
plt.gca().axis.set_major_formatter(date_format)
```

```
plt.xlabel('')
plt.ylabel('Count (in Millions)', fontsize=16)
plt.xticks(fontsize=14)
plt.yticks(fontsize=13)
plt.suptitle('Total Cases by Time', fontsize=20)

plt.annotate(text=str(curr_total_cases), xy=(curr_date, curr_total_cases/10**6),
             xycoords='data', xytext=(-80,1), textcoords='offset points', fontsize=14)
```

```
Out[12]: Text(-80, 1, '29699371')
```



The logarithmic rise of total cases was observed from July end 2020 till December 2020 which seemed to saturate in January 2021. But April 2020 onwards, cases started to increase at much higher rate than before

```
In [13]: daily_cases = ind_covid_df['Daily Confirmed']
```

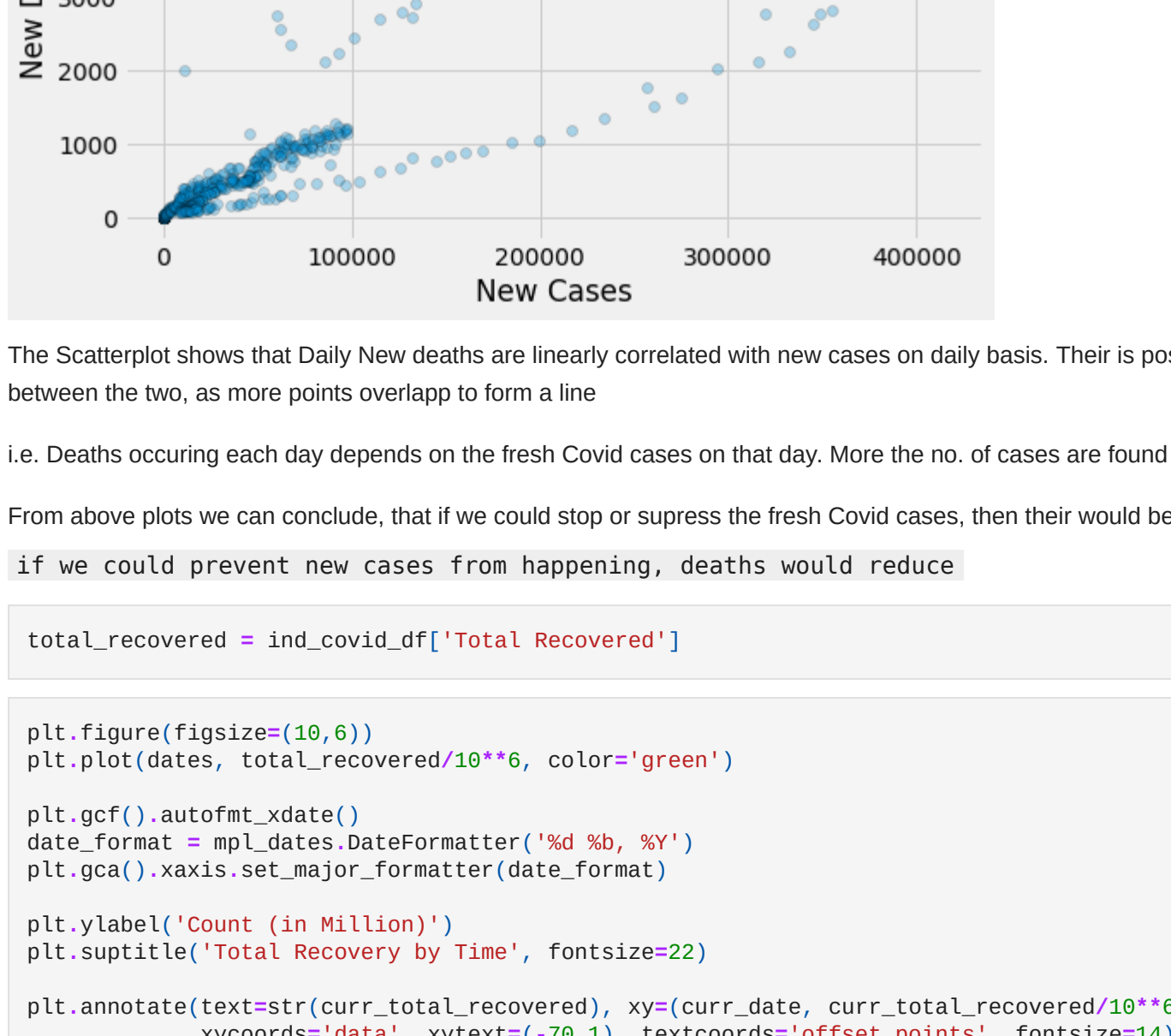
```
In [14]: plt.figure(figsize=(10,6))
plt.plot(dates, daily_cases, '-', linewidth=2)

plt.gca().autofmt_xdate()
date_format = mldates.DateFormatter('%d %b, %y')
plt.gca().axis.set_major_formatter(date_format)

plt.ylabel('Count')
plt.suptitle('Daily New Cases by Time', fontsize=22)

plt.annotate(text=str(today_cases), xy=(curr_date, today_cases),
             xycoords='data', xytext=(-56,1), textcoords='offset points', fontsize=14)
```

```
plt.savefig('Images/daily_cases.png')
```



From July 2020 Onwards infection rate started to increase and reached its first peak at September 2020 with over 90,000 cases reported per-day. Cases began to decline from October 2020 and were reported below 15,000 in January 2021 which was a good sign.

A second wave beginning in March 2021 was much larger than first, with shortages of vaccines, hospital beds, oxygen cylinders and other medicines such as remdesivir in parts of the country. By April end daily infection count reached over 400,000 which was new record

```
In [15]: total_deaths = ind_covid_df['Total Deceased']
```

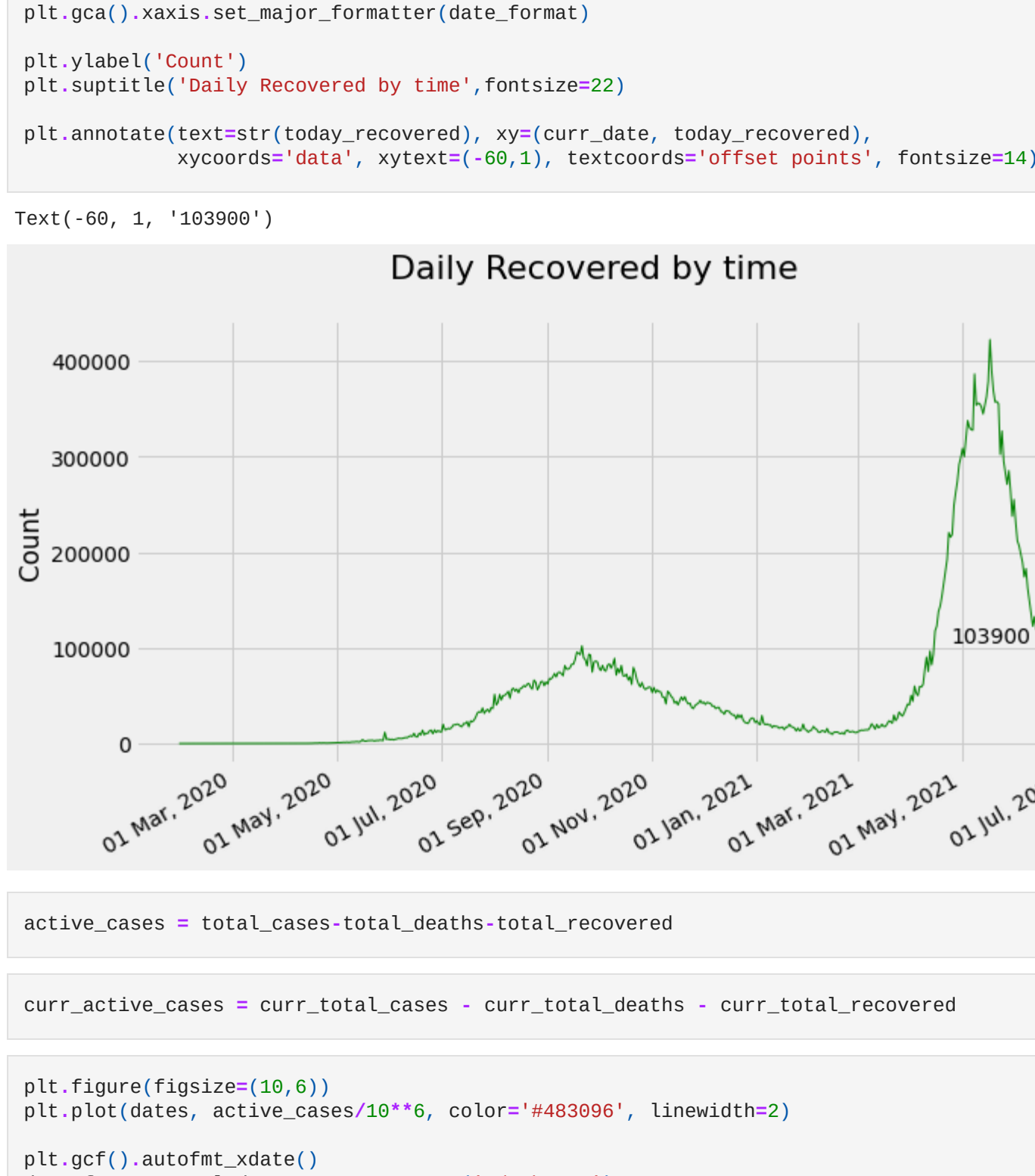
```
In [16]: plt.figure(figsize=(10,6))
plt.plot(dates, total_deaths, color='red')

plt.gca().autofmt_xdate()
date_format = mldates.DateFormatter('%d %b, %y')
plt.gca().axis.set_major_formatter(date_format)

plt.ylabel('Count')
plt.suptitle('Total Deaths by Time', fontsize=22)

plt.annotate(text=str(curr_total_deaths), xy=(curr_date, curr_total_deaths),
             xycoords='data', xytext=(-56,1), textcoords='offset points', fontsize=14)
```

```
Out[16]: Text(-56, 1, '381340')
```



```
In [17]: daily_deaths = ind_covid_df['Daily Deceased']
```

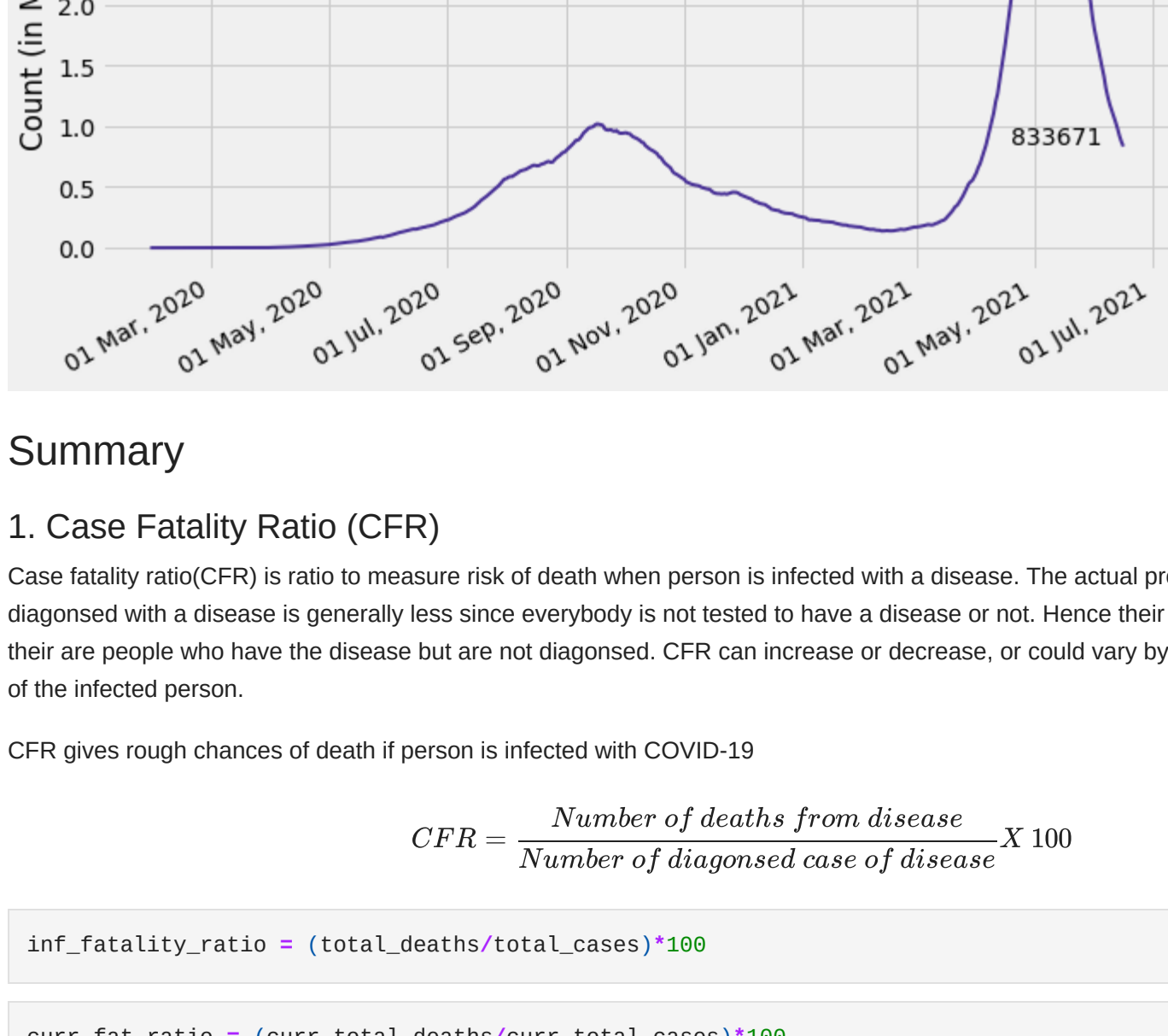
```
In [18]: plt.figure(figsize=(10,6))
plt.plot(dates, daily_deaths, 'r', linewidth=1)

plt.gca().autofmt_xdate()
date_format = mldates.DateFormatter('%d %b, %y')
plt.gca().axis.set_major_formatter(date_format)

plt.ylabel('Count')
plt.suptitle('Daily New Deaths by Time', fontsize=22)

plt.annotate(text=str(today_deaths), xy=(curr_date, today_deaths),
             xycoords='data', xytext=(-40,1), textcoords='offset points', fontsize=14)
```

```
Out[18]: Text(-40, 1, '2329')
```

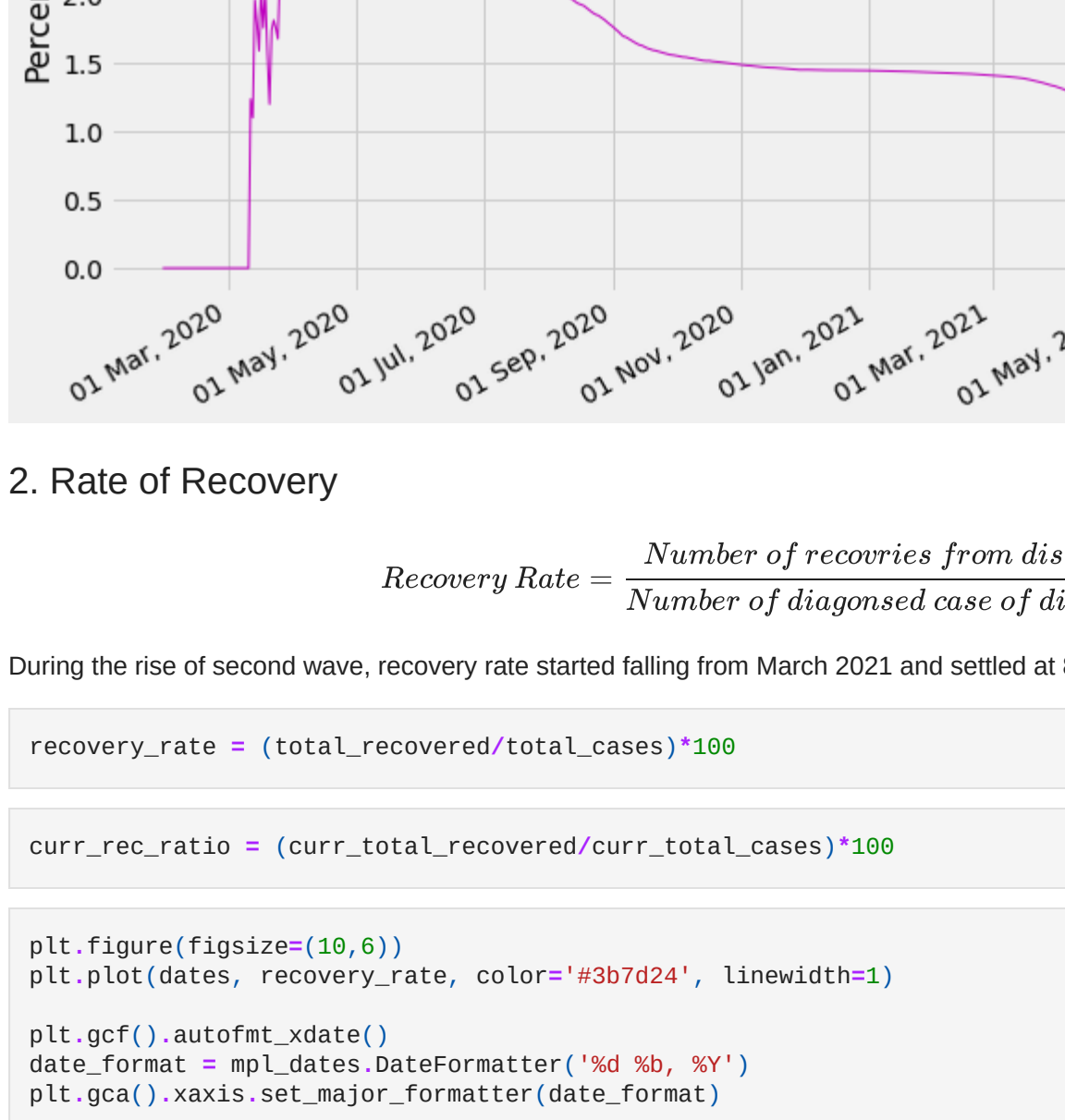


Above plot depicts that there were large no. of deaths in August, September and October months of year 2020. Sudden spike of deaths was seen in mid-June month. In Second wave the deaths are 4 to 5 times more than the previous wave

Let us see if there is any correlation between new cases and new deaths on daily basis

```
In [19]: plt.figure(figsize=(8,5))
plt.scatter(daily_cases, daily_deaths, edgecolor='black', alpha=.3)
plt.xlabel('New Cases')
plt.ylabel('New Deaths')
```

```
Out[19]: Text(0, 0.5, 'New Deaths')
```



The Scatterplot shows that Daily New deaths are linearly correlated with new cases on daily basis. Their is positive, strong relation between the two, as more points overlap to form a line

i.e. Deaths occurring each day depends on the fresh Covid cases on that day. More the no. of cases are found more deaths will occur.

From above plots we can conclude, that if we could control or suppress the fresh Covid cases, then their would be less deaths. If we could prevent new cases from happening, deaths would reduce.

```
In [20]: total_recovered = ind_covid_df['Total Recovered']
```

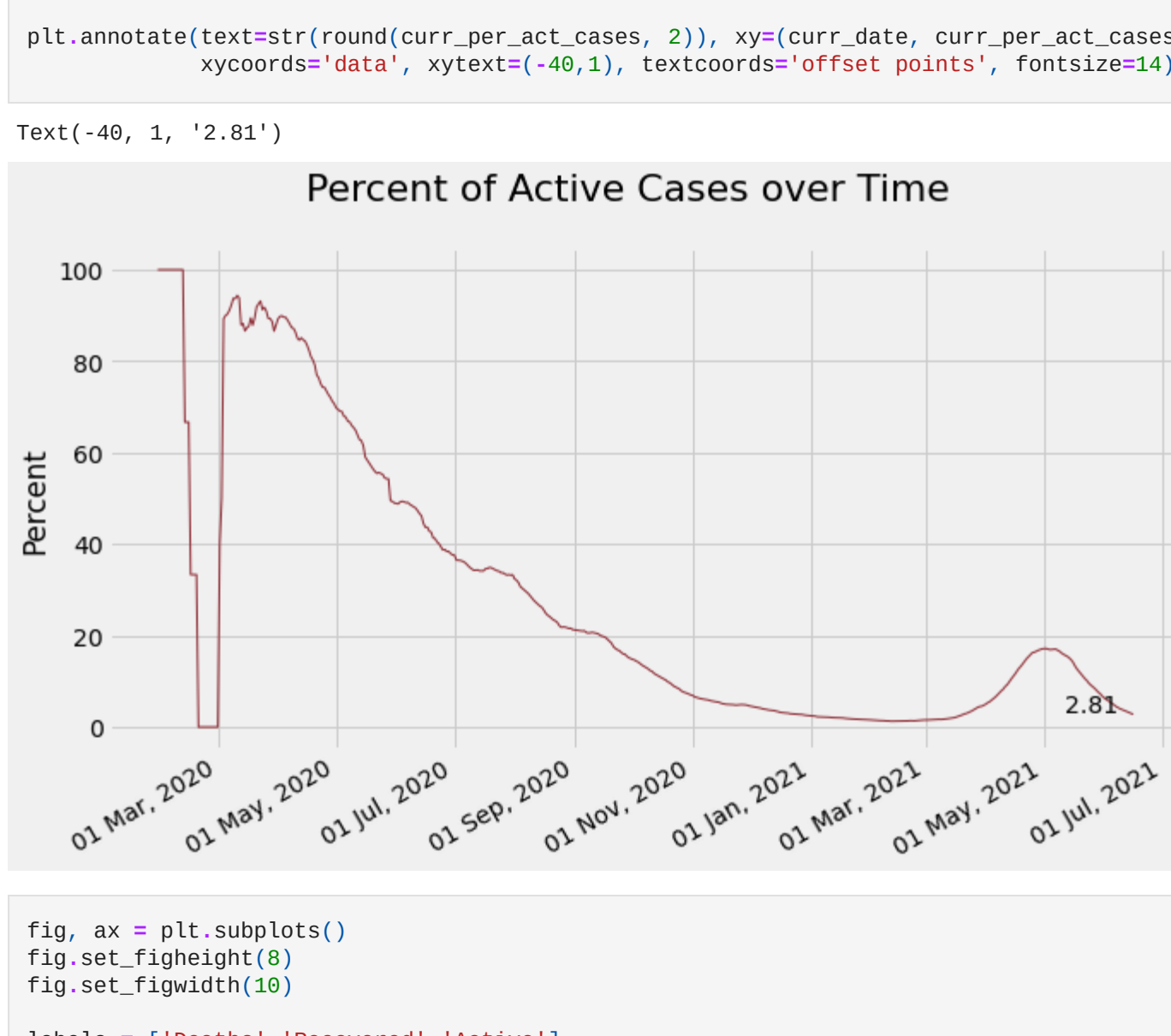
```
In [21]: plt.figure(figsize=(10,6))
plt.plot(dates, total_recovered/10**6, color='green')

plt.gca().autofmt_xdate()
date_format = mldates.DateFormatter('%d %b, %y')
plt.gca().axis.set_major_formatter(date_format)

plt.ylabel('Count (in Million)')
plt.suptitle('Total Recovery by Time', fontsize=22)

plt.annotate(text=str(curr_total_recovered), xy=(curr_date, curr_total_recovered/10**6),
             xycoords='data', xytext=(-56,1), textcoords='offset points', fontsize=14)
```

```
plt.savefig('Images/total_recov.png')
```



```
In [22]: daily_recovered = ind_covid_df['Daily Recovered']
```

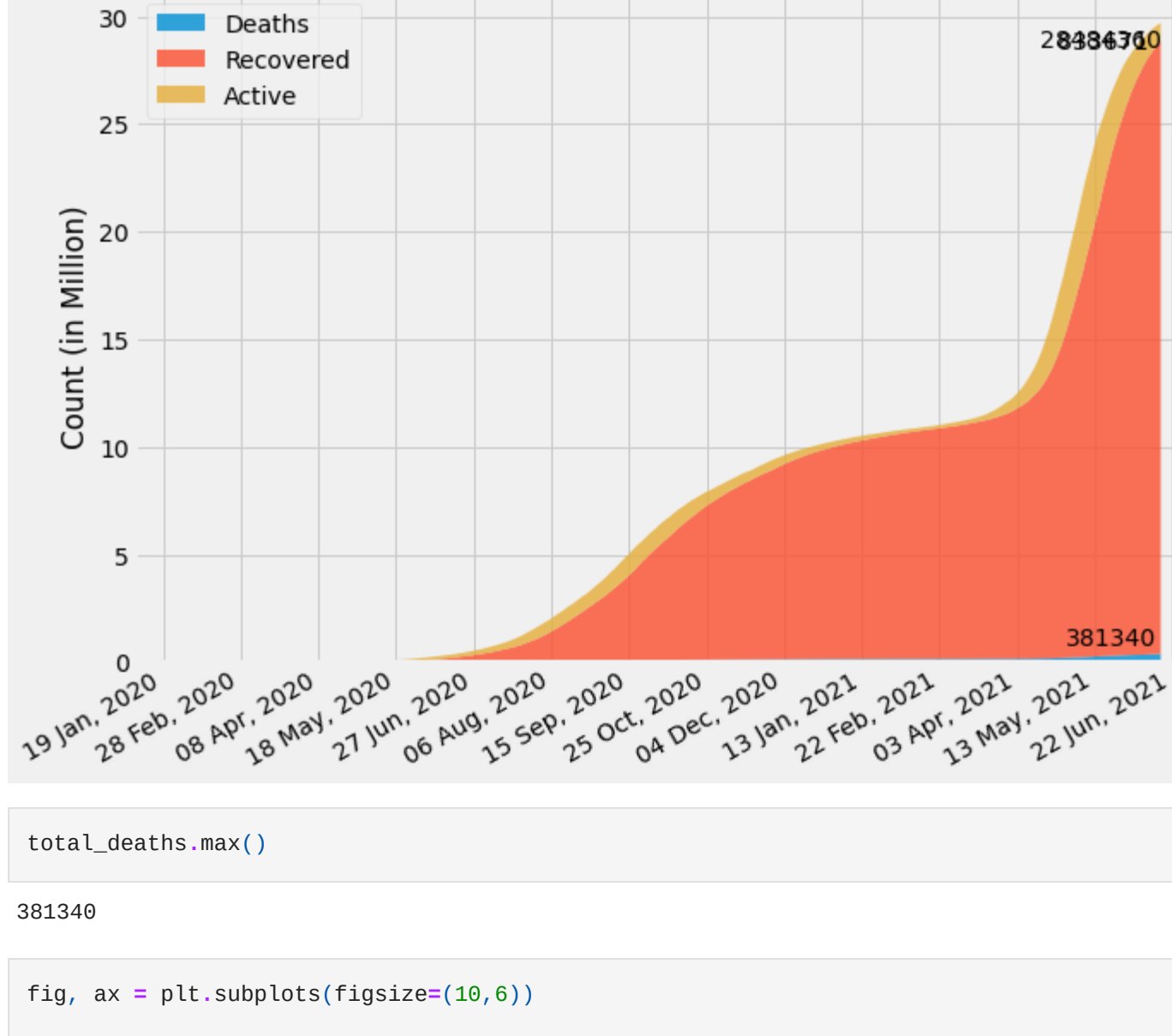
```
In [23]: plt.figure(figsize=(10,6))
plt.plot(dates, daily_recovered, 'g', linewidth=1)

plt.gca().autofmt_xdate()
date_format = mldates.DateFormatter('%d %b, %y')
plt.gca().axis.set_major_formatter(date_format)

plt.ylabel('Count')
plt.suptitle('Daily Recovered by Time', fontsize=22)

plt.annotate(text=str(today_recovered), xy=(curr_date, today_recovered),
             xycoords='data', xytext=(-60,1), textcoords='offset points', fontsize=14)
```

```
Out[23]: Text(-60, 1, '103900')
```



```
In [24]: active_cases = total_cases-total_deaths-total_recovered
```

```
In [25]: curr_active_cases = curr_total_cases - curr_total_deaths - curr_total_recovered
```

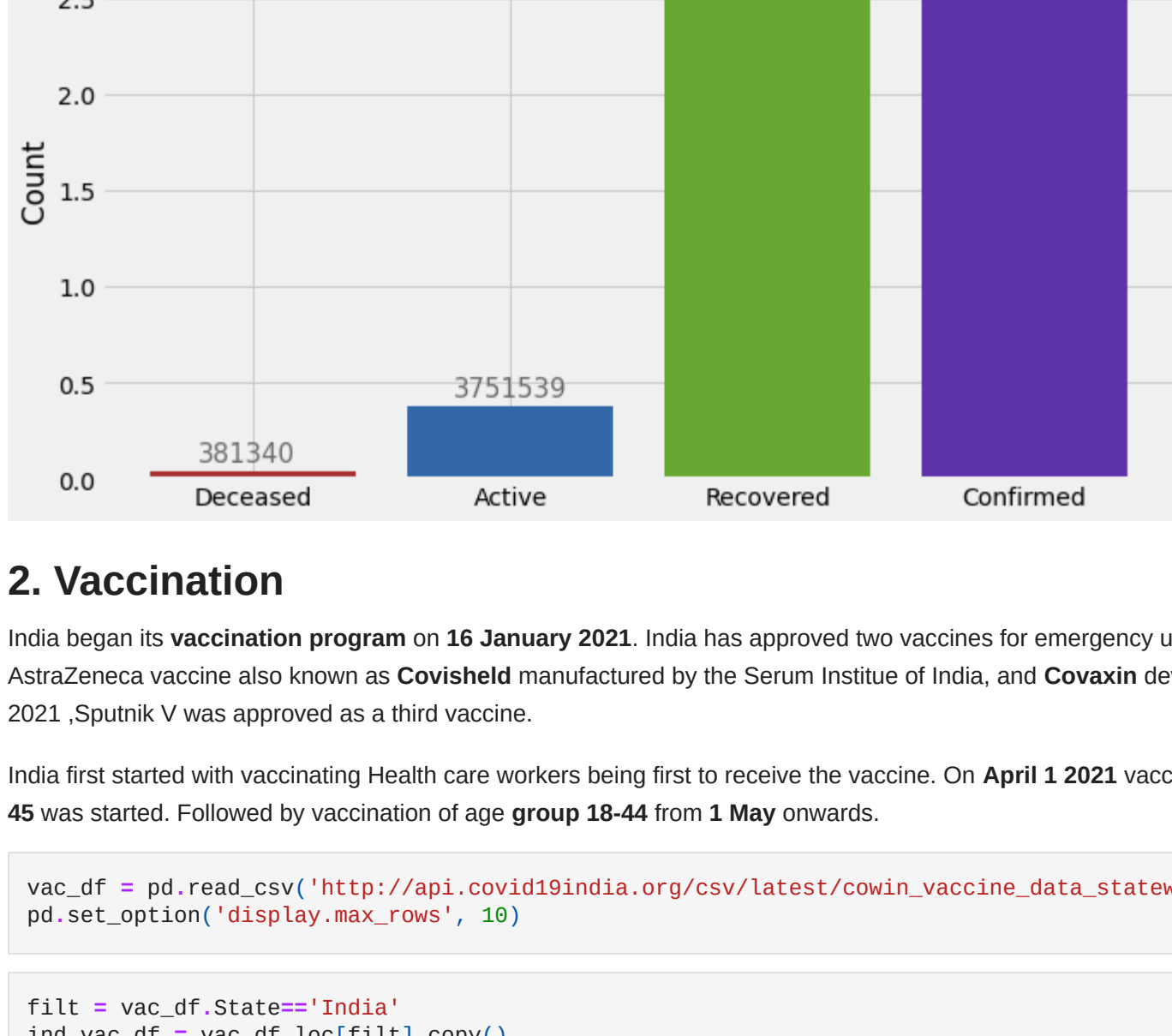
```
In [26]: plt.figure(figsize=(10,6))
plt.plot(dates, active_cases/10**6, color='#483986', linewidth=2)

plt.gca().autofmt_xdate()
date_format = mldates.DateFormatter('%d %b, %y')
plt.gca().axis.set_major_formatter(date_format)

plt.ylabel('Count (in Millions)')
plt.suptitle('Active Cases over Time', fontsize=22)

plt.annotate(text=str(curr_active_cases), xy=(curr_date, curr_active_cases/10**6),
             xycoords='data', xytext=(-66,1), textcoords='offset points', fontsize=14)
```

```
Out[26]: Text(-66, 1, '833671')
```



Summary

1. Case Fatality Ratio (CFR)

Case fatality ratio(CFR) is ratio to measure risk of death when person is infected with a disease. The actual probability of death of person diagnosed with a disease is generally less since everybody is not tested to have a disease or not. Hence their would be a scenario where there are people who have the disease but are not diagnosed. CFR can increase or decrease, or could vary by location and characteristics of the infected person.

CFR gives rough chances of death if person is infected with COVID-19

$$CFR = \frac{\text{Number of deaths from disease}}{\text{Number of diagnosed case of disease}} \times 100$$

```
In [27]: inf_fatality_ratio = (total_deaths/total_cases)*100
```

```
In [28]: curr_fat_ratio = (curr_total_deaths/curr_total_cases)*100
```

```
In [29]: plt.figure(figsize=(10,6))
plt.plot(dates, inf_fatality_ratio, 'm', linewidth=1)

plt.gca().autofmt_xdate()
date_format = mldates.DateFormatter('%d %b, %y')
plt.gca().axis.set_major_formatter(date_format)

plt.ylabel('Percent')
plt.suptitle('Infection Fatality Ratio over Time', fontsize=22)
# plt.title('Chances of Death')

plt.annotate(text=str(round(curr_fat_ratio, 3)), xy=(curr_date, curr_fat_ratio),
             xycoords='data', xytext=(-48,-10), textcoords='offset points', fontsize=14)
```

```
Out[29]: Text(-48, -10, '1.284')
```


2. Rate of Recovery

$$\text{Recovery Rate} = \frac{\text{Number of recoveries from disease}}{\text{Number of diagnosed case of disease}} \times 100$$

During the rise of second wave, recovery rate started falling from March 2021 and settled at 80% after which has started to grow again

```
In [30]: recovery_rate = (total_recovered/total_cases)*100
```

```
In [31]: curr_rec_ratio = (curr_total_recovered/curr_total_cases)*100
```

```
In [32]: plt.figure(figsize=(10,6))
plt.plot(dates, recovery_rate, color='#3070D4', linewidth=1)

plt.gca().autofmt_xdate()
date_format = mldates.DateFormatter('%d %b, %y')
plt.gca().axis.set_major_formatter(date_format)

plt.ylabel('Percent')
plt.suptitle('Recovery rate over Time', fontsize=22)

plt.annotate(text=str(round(curr_rec_ratio, 2)), xy=(curr_date, curr_rec_ratio),
             xycoords='data', xytext=(-50,-5), textcoords='offset points', fontsize=14)
```

```
Out[32]: Text(-50, -5, '95.91')
```



```
In [33]: per_act_cases = (active_cases/total_cases)*100
```

```
In [34]: curr_per_act_cases = (curr_active_cases/curr_total_cases)*100
```

```
In [35]: plt.figure(figsize=(10,6))
plt.plot(dates, per_act_cases, color='#8C2730', linewidth=1)

plt.gca().autofmt_xdate()
date_format = mldates.DateFormatter('%d %b, %y')
plt.gca().axis.set_major_formatter(date_format)

plt.ylabel('Percent')
plt.suptitle('Percent of Active Cases over Time', fontsize=22)

plt.annotate(text=str(round(curr_per_act_cases, 2)), xy=(curr_date, curr_per_act_cases),
             xycoords='data', xytext=(-40,1), textcoords='offset points', fontsize=14)
```

```
Out[35]: Text(-40, 1, '2.81')
```



```
In [36]: fig, ax = plt.subplots()
fig.set_figheight(8)
fig.set_figwidth(10)
```

```
labels = ['Deaths', 'Recovered', 'Active']
ax.stackplot(dates, total_deaths/10**6, total_recovered/10**6, active_cases/10**6, alpha=.8, labels=labels)
ax.set_ylabel('Count (in Millions)')
ax.legend(loc='upper left')
```

```
fig.autofmt_xdate()
date_format = mldates.DateFormatter('%d %b, %y')
```

```
ax.xaxis.set_major_locator(plt.MaxNLocator(10))
ax.xaxis.set_major_formatter(date_format)
```

```
plt.annotate(text=str(curr_total_deaths), xy=(curr_date, curr_total_deaths/10**6),
             xycoords='data', xytext=(-55,5), textcoords='offset points', fontsize=14)
```

```
plt.annotate(text=str(curr_total_recovered), xy=(curr_date, curr_total_recovered/10**6),
             xycoords='data', xytext=(-70,1), textcoords='offset points', fontsize=14)
```

```
plt.annotate(text=str(curr_active_cases), xy=(curr_date, curr_active_cases/10**6 + curr_total_recovered/10**6),
             xycoords='data', xytext=(-60,-10), textcoords='offset points', fontsize=14)
```

```
Out[36]: Text(-60, -10, '833671')
```



```
In [37]: total_deaths.max()
```

```
Out[37]: 381340
```

```
In [38]: fig, ax = plt.subplots(figsize=(10,6))
```

```
labels = ['Deceased', 'Active', 'Recovered', 'Confirmed']
values = [total_deaths.max(), active_cases.max(), total_recovered.max(), total_cases.max()]
ax.bar(labels, values, color=['#A83232', '#3267A8', '#67A832', '#5D32A8'])
ax.set_ylabel('Count')
```

```
# create a list to collect the plt.patches data
totals = []
```

```
# find the values and append to list
for i in ax.patches:
    totals.append(i.get_height())
```

```
# get individual bar labels using above list
for i in ax.patches:
    # set_x pulls left or right; get_height pushes up or down
    ax.text(i.get_x()+.10, i.get_height()+900000, \
            round(i.get_height(),1), fontsize=15,
            color='dimgrey')
```

```
plt.savefig('Images/cases_summary.png')
```


2. Vaccination

India began its vaccination program on **16 January 2021**. India has approved two vaccines for emergency use, including Oxford-AstraZeneca vaccine also known as **Covishield** manufactured by the Serum Institute of India, and **Covaxin** developed by Biotech. In vaccine 2021. Sputnik V was approved as a third vaccine.

India first started with vaccinating health care workers being first to receive the vaccine. On **April 1 2021** vaccination of people above **age 45** was started. Followed by vaccination of age group **18-44** from **1 May** onwards.

```
In [39]: vac_df = pd.read_csv('https://api.covid19india.org/csv/latest/covid_vaccine_data_statewise.csv')
pd.set_option('display.max_rows', 10)
```

```
In [40]: filt = vac_df.State=='India'
ind_vac_df = vac_df.loc[filt].copy()
```

```
In [41]: ind_vac_df
```

| | Updated On | State | Total Individuals Administered | Total Sessions Conducted | Total Sites | First Dose Administered | Second Dose Administered | Male(Individuals Vaccinated) | Female(Individuals Vaccinated) | Transgender(Indiv Vaccinated) |
|-----|------------|-------|--------------------------------|--------------------------|-------------|-------------------------|--------------------------|------------------------------|--------------------------------|-------------------------------|
| 0 | 16/01/2021 | India | 48276.0 | 3455.0 | 2957.0 | 48276.0 | 0.0 | 23757.0 | 24517.0 | 0.0 |
| 1 | 17/01/2021 | India | 58604.0 | 8532.0 | 4954.0 | 58604.0 | 0.0 | 27348.0 | 31252.0 | 0.0 |
| 2 | 18/01/2021 | India | 99449.0 | 13611.0 | 6583.0 | 99449.0 | 0.0 | 41361.0 | 58083.0 | 0.0 |
| 3 | 19/01/2021 | India | 195525.0 | 17855.0 | 7951.0 | 195525.0 | 0.0 | 81901.0 | 113613.0 | 0.0 |
| 4 | 20/01/2021 | India | 251280.0 | 25472.0 | 10504.0 | 251280.0 | 0.0 | 98111.0 | 153145.0 | 0.0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 149 | 14/06/2021 | India | 20727441.0 | 14048206.0 | 46824.0 | 20727441.0 | 47378599.0 | 111868784.0 | 95369557.0 | 0.0 |
| 150 | 15/06/2021 | India | 209816439.0 | 11834657.0 | 41484.0 | 209816439.0 | 47789099.0 | 113266182.0 | 96513504.0 | 0.0 |
| 151 | 16/06/2021 | India | 212945352.0 | 14252163.0 | 44447.0 | 212945352.0 | 48173116.0 | 114978703.0 | 97929256.0 | 0.0 |
| 152 | 17/06/2021 | India | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| 153 | 18/06/2021 | India | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN |

154 rows × 11 columns

```
In [42]: total_population = 13380894385
```

```
In [43]: ind_vac_df.columns
```

```
Index(['Updated On', 'State', 'Total Individuals Vaccinated', 'Total Sessions Conducted', 'Total Sites', 'First Dose Administered', 'Second Dose Administered', 'Male(Individuals Vaccinated)', 'Female(Individuals Vaccinated)', 'Transgender(Individuals Vaccinated)', 'Total Covaxin Administered', 'Total Covishield Administered', 'Total Sputnik V Administered', 'AGE1', '18-45 years (Age)', '45-60 years (Age)', '60+ years (Age)', 'Total Doses Administered'], dtype='object')
```

```
In [44]: ind_vac_df.drop(['State', 'Total Sessions Conducted', 'Total Sites', 'Male(Individuals Vaccinated)', 'Female(Individuals Vaccinated)'], axis=1, inplace=True)
```

```
In [45]: ind_vac_df
```