# CMSC-6950 PROJECT WORK

Muhammad Ifthkar, Prasun Thapa, Adnan Shafiq                    06/22/2020

Github url: https://github.com/miaftab/finalproject

## Introduction

The project is about the covid-19 Canada's data including all the provinces. Our main agenda of the project was to analyze, implement, append and plot the data in the visualized form. We are taking all the data associated with covid-19 from Canada's Health website:
https://health-infobase.canada.ca/src/data/covidLive/covid19.csv
The project was the perfect combo, to learn how to manage the source code using Github and python commands. We also learned to append, edit the csv files and create a plot diagrams based on the given data. Furthermore, it gave us a sense of testing the code, debugging and an environment to work on a team basis.

## Source Code

```
1  import matplotlib.pyplot as plt
2  from matplotlib.dates import date2num
3  import matplotlib.patches as mpatches
4  import datetime
5  import requests
6  import csv
7  import os
```

Importing Libraries

## Function responsible for plotting the data over days

```
1  def plotChartFunction(x_dates,setData,title):
2
3      # plot chart on SetData
4      x_dates = date2num(x_dates)
5      ax = plt.subplot(111)
6      legends = []
7
8      # Colors Array
```

```
 9    colors = ['orange',
10              'green',
11              'red',
12              'blue',
13              'purple',
14              'yellow',
15              'grey',
16              'magenta',
17              'black',
18              'cyan',
19              'brown',
20              'indigo',
21              'olive',
22              'navy',
23              'orchid'
24              ]
25
26    count = 0
```

## DataSet have key as each provence name and it's data list as value

```
 1  for key,value in setData.items():
 2      ax.plot(x_dates, value, color=colors[count])
 3      a=mpatches.Patch(color=colors[count],linestyle='--',label=key)
 4      legends.append(a)
 5      count = count + 1
 6  ax.xaxis_date()
 7  plt.xticks(rotation=70)
 8  plt.legend(handles=legends)
 9  plt.title(title,size=10,color='Green')
10  plt.show()
```

Global plotting function created, which is responsible for plotting data over days. It contain three arguments xdates, setData, title. Xdates will be on the X-axis while title will be the title of graph plot. The most important parameter is setData which contains key value pairs. Keys are the name of provinces and value will be the corresponding data against particular key. An array of different colors used to differentiate between all provinces in sub plotting. At the end with the help of iteration loop all graphs plotted one by one.

## Automating the download of latest Covid data

```python
1  try:
2      os.remove("covid19.csv")
3  except:
4      print("File not exist")
5
6
7  # Download new covid19 csv from HIC website
8  req = requests.get("https://health-infobase.canada.ca/src/data/covidLive/←
       covid19.csv")
9  url_content = req.content
10 csv_file = open('covid19.csv', 'wb')
11 csv_file.write(url_content)
12 csv_file.close()
13
14 # open the file and read
15 with open('covid19.csv', 'r') as infile:
16     # read the file as a dictionary for each row ({header : value})
17     reader = csv.DictReader(infile)
18     data = {}
19     for row in reader:
20         for header, value in row.items():
21             try:
22                 data[header].append(value)
23             except KeyError:
24                 data[header] = [value]
25
26 # extract the variables you want unique
27 provences = list(set(data['prname']))
28 totalCases = data['numtotal']
29
30 # Select only unique dates sorted
31 date = list(set(data['date']))
32 date = sorted(date, key=lambda x: datetime.datetime.strptime(x, '%d-%m-%Y'←
       ))
33
34 # Extracting dates and converting it to date time object
35 x_dates = []
36 eachProviceCases = {}
37
38 # append Unique dates in x_dates
39 for d in date:
40     splitedDate = d.split('-')
41     temp = datetime.date(int(splitedDate[2]),int(splitedDate[1]),int(←
           splitedDate[0]))
42     x_dates.append(temp)
43
```

```
44  ![alt text](Q1_plot.png)
```

Question 1 is basically related to keep data updated. So every time at the start of program this function will check that either there is any existing file with the name "covid-19" in our system, if is exist it will be deleted automatically and then function will send a request for new file on internet. After downloading file, it will be open in writing mode. Needed data extracted from the file on the bases of keys. Dates are formatted as well as converted into the numbers. One important thing is that province names and dates extracted uniquely.

## Total No of cases in each province over time

Listing 2: Question- 2 .

```
1   for prov in provences:
2
3       # Exclude Canada
4       if prov != 'Canada':
5           # add province name as key and empty list as value
6           eachProviceCases[prov] = []
7
8           # append list of zeros equal to no of dates against province
9           for d in date:
10              eachProviceCases[prov].append(0)
11
12          for index,d in enumerate(date): # iterate on sorted date list
13              for idx,da in enumerate(data["date"]): # match it within date ↩
                      column of data
14                  if da == d and prov == data['prname'][idx]: # if date ↩
                          match and province name is equal
15                      eachProviceCases[prov][index] = int(data['numtotal'][↩
                            idx]) # add data against province for given date
16                      break
17
18  # plot graph Total No of cases in each provence over time
19  plotChartFunction(x_dates,eachProviceCases,'Total No of cases in each ↩
        provence over time')
```

A unique name of provinces is created, and the names are extracted from those created province name. The Canada is excluded here, province named has been passed to an empty key value pair creating an empty list as a value. The number of zeros is being appended in the above empty list, which is equal to the number of dates. Sorted date list are being enumerated and each date is iterated through date column. The province name and the date are checked if they are equal. Then the data is converted into integer and is saved in the particular province name.

A common function called plotchartFunction is used, which takes x-dates as the parameter, representing the X-axis, each province cases representing the Y-axis and the last parameter is the title of the chart.
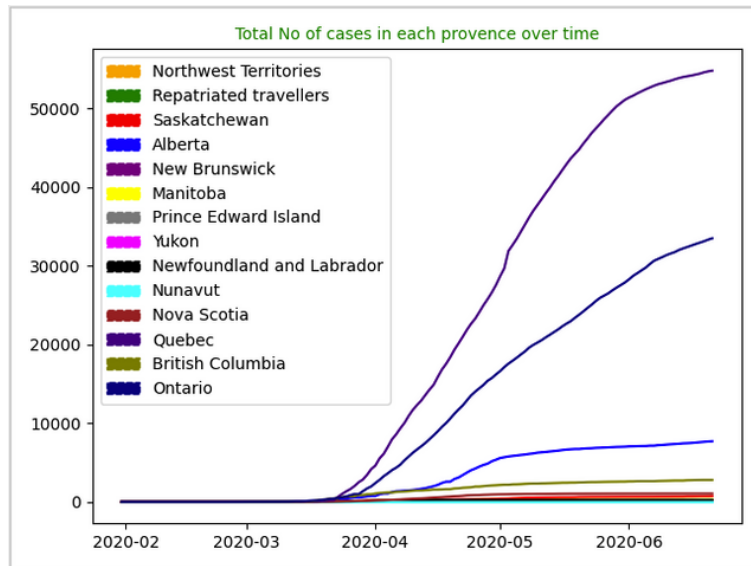
## Question 2 Graph



Figure 1: Total No of cases in each provence over time

Passing those data, we came up with this plot, you can see these purple lines indicates the number of cases in Ontario, which is more than 50,000 cases and so on. These datas are mainly rising from the month of march.

## Total Number of individuals tested in each province overtime

Listing 3: Question- 3 .

```
1  eachProviceTestedCases = {}
2
3  # iterate through each provence
4  for prov in provences:
5
6      # Exclude Canada
7      if prov != 'Canada':
8
9          # add province name as key and empty list as value
10         eachProviceTestedCases[prov] = []
11
12         # append list of zeros equal to no of dates against province
13         for d in date:
```

```
14                    eachProviceTestedCases[prov].append(0)
15
16         for index,d in enumerate(date):  # iterate on sorted date list
17             for idx,da in enumerate(data["date"]):  # match it within date↩
                   column of data
18                 if da == d and prov == data['prname'][idx]: # if date ↩
                      match and province name is equal
19                     if data['numtested'][idx] != '':
20                         eachProviceTestedCases[prov][index] = int(data['↩
                              numtested'][idx])  # add data against province↩
                              for given date
21                     else:
22                         eachProviceTestedCases[prov][index] = 0
23                     break
24 # plot graph Total No of cases tested each provence over time
25 plotChartFunction(x_dates,eachProviceTestedCases,'Total No of cases tested↩
      each provence over time')
```

The source used to solve question number 2 is very similar to the one that has been used to solve question 3. It, basically follows the same basic principle. The only major difference was numtotal column name was replaced with numtested. Where numtested column represents the total number of individuals tested in each province. Again, the plotchartFunction is used in order to plot and name the title of the chart below.
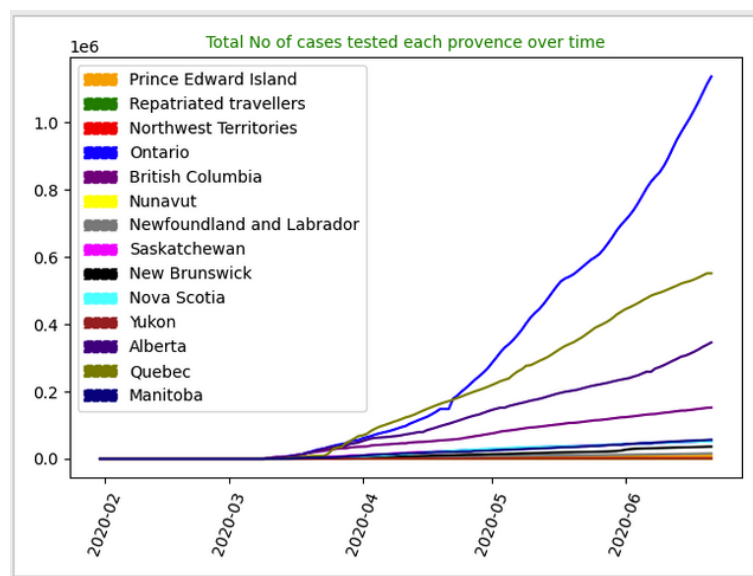
## Question 3 Graph



Figure 2: Total Number of individuals tested in each province overtime

From the above graph, we can clearly depict as Ontario leading the plots with the most number of individual tests. Quebec takes the second lead on the plot, followed by Alberta.

# Total No of new cases per day in each province over time

Listing 4: Question- 4 .

```python
1  newCasesPerDay = {}
2  dates = set()
3  # iterate through each provence
4  for prov in provences:
5
6      # Exclude Canada
7      if prov != 'Canada':
8
9          # add province name as key and empty list as value
10         newCasesPerDay[prov] = []
11
12         for index,d in enumerate(date): # iterate on sorted date list
13
14             proceed = False
15             # check if all provinces have zero data for a specific date
16             for idx2,checkDate in enumerate(data["date"]):
17                 if checkDate == d and int(data['numtoday'][idx2]) > 50:
18                     proceed = True
19                     break
20
21             if proceed:
22                 dates.add(d)
23                 ignore = False
24                 for idx,da in enumerate(data["date"]):  # match it within ↩
                        date column of data
25                     if da == d and prov == data['prname'][idx]: # if date ↩
                            match and province name is equal
26                         if data['numtoday'][idx] != '':
27                             newCasesPerDay[prov].append(int(data['numtoday↩
                                    '][idx])) # add data against province for ↩
                                    given date
28                         else:
29                             newCasesPerDay[prov].append(0)
30                         ignore=True
31                         break
32                 if not ignore:
33                     newCasesPerDay[prov].append(0)
34
35 date = list(dates)
36 date = sorted(date, key=lambda x: datetime.datetime.strptime(x, '%d-%m-%Y'↩
       ))
37 x_dates = []
```

```
38  for d in date:
39      splitedDate = d.split('-')
40      temp = datetime.date(int(splitedDate[2]),int(splitedDate[1]),int(↩
            splitedDate[0]))
41      x_dates.append(temp)
42
43
44  # plot graph Total No of new cases per day each provence over time
45  plotChartFunction(x_dates,newCasesPerDay,'No of new cases per day each ↩
        provence over time')
```

Again, the key value pairs is created and canada is being exculded from the list. For this question, a minimum cut-off score was determined in order to figure out to extract the value form the particular date. The data are being verified if they are equal, going through the data sets. If any of the province for a particular date is going to get a value greater than 50, the proceed check turns true. If the proceed is true, for that particular date the data is stored for the plot. Rest of the source code is similar to 2 and 3, but here equal number of dates is given excluding some dates above. Not including all the dates was our main concern here. Finally, the same plot chart Function is used to plot a graph against new cases per day over time.
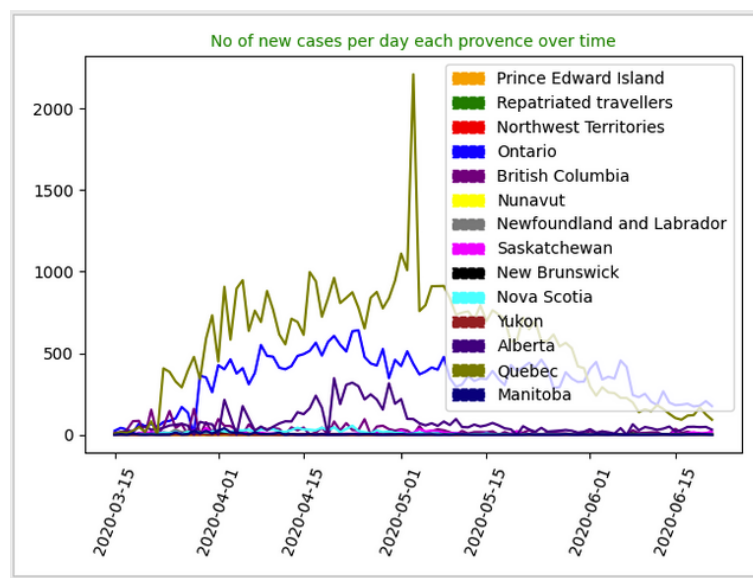
## Question 4 Graph



Figure 3: Total No of new cases per day in each province over time

Here we are excluding all the zeros starting from 50. Quebec are more cases compared to any other followed by Ontario.

## Doubling Rate of the Number of Cases

```
1  try:
2      os.remove("output.txt")
3  except:
4      print("File not exist")
5
6  os.system("bash doubling.sh  Canada 'number of cases' 01-04-2020")
7  os.system("bash doubling.sh  Alberta 'number of cases' 01-04-2020")
8  os.system("bash doubling.sh  Ontario 'number of cases' 01-04-2020")
9
10 f = open("output.txt", "r")
11 prvinceName = []
12 NoOfDays=[]
13 legends=[]
14 colors = ['r','g','b']
15 index=0
16 for x in f:
17     temp = x.split(" ")
18     a=mpatches.Patch(color=colors[index],linestyle='--',label=temp[0]+'('+↵
           temp[2]+') --> '+'('+temp[4]+')')
19     legends.append(a)
20     prvinceName.append(temp[0]+'('+temp[1]+')')
21     NoOfDays.append(int(temp[5].replace('\n', '')))
22     index=index+1
23 print(prvinceName)
24 print(NoOfDays)
25
26 barlist=plt.bar(prvinceName,NoOfDays)
27 barlist[0].set_color('r')
28 barlist[1].set_color('g')
29 barlist[2].set_color('b')
30
31 plt.legend(handles=legends)
32 plt.title('Doubling Rate From Given Date')
33 plt.xlabel('Provinces From Date')
34 plt.ylabel('No Of Days')
35
36 plt.show()
37 \begin{lstlisting}[label={list:fifth},caption= Question- 5 .]
```
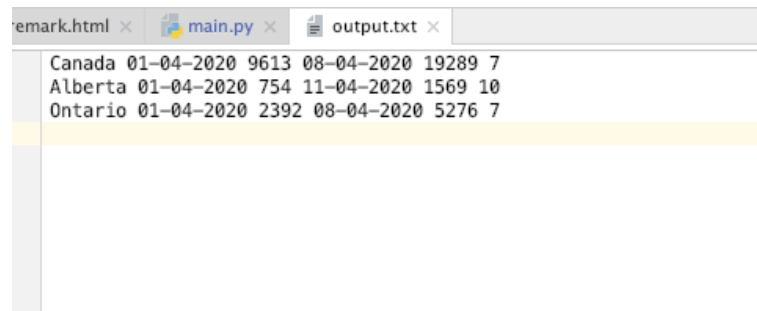
A shell script called 'doubling.sh' is used to compute the doubling rate of the cases. To compute the doubling rate through the shell script, three arguments are passed. For instance, the three arguments are province name, number of cases or number of deaths and date required to complete the doubling rate respectively. example

bash doubling.sh Canada'number of cases'01-04-2020

The output of doubling.sh is stored at output.txt. These data from the output.txt file are further used to plot the graph representing the doubling rate for minimum three provinces. To clarify, a screenshot is added below where the number of cases in Alberta was 754 and it took 10 days for it to exactly double the value and got close to 1569.

## Question 5 Graph



```
emark.html ×   main.py ×   output.txt ×
    Canada 01-04-2020 9613 08-04-2020 19289 7
    Alberta 01-04-2020 754 11-04-2020 1569 10
    Ontario 01-04-2020 2392 08-04-2020 5276 7
```

Figure 4: Doubling Rate of the Number of Cases



| A | B | C | D | E | F | G | H | I | J |
|---|---|---|---|---|---|---|---|---|---|
| 48 | Alberta | Alberta | 27/03/2020 | 542 | 0 | 2 | 542 | 38215 | 27 |
| 48 | Alberta | Alberta | 28/03/2020 | 542 | 0 | 2 | 542 | 38215 | 53 |
| 48 | Alberta | Alberta | 29/03/2020 | 621 | 0 | 2 | 621 | 44097 | 73 |
| 48 | Alberta | Alberta | 30/03/2020 | 690 | 0 | 8 | 690 | 46057 | 94 |
| 48 | Alberta | Alberta | 31/03/2020 | 754 | 0 | 9 | 754 | 48692 | 120 |
| 48 | Alberta | Alberta | 01/04/2020 | 754 | 0 | 9 | 754 | 53141 | 142 |
| 48 | Alberta | Alberta | 02/04/2020 | 968 | 0 | 13 | 968 | 57096 | 174 |
| 48 | Alberta | Alberta | 03/04/2020 | 1075 | 0 | 18 | 1075 | 60508 | |
| 48 | Alberta | Alberta | 04/04/2020 | 1075 | 0 | 18 | 1075 | 62520 | |
| 48 | Alberta | Alberta | 05/04/2020 | 1250 | 0 | 23 | 1250 | 63315 | 240 |
| 48 | Alberta | Alberta | 06/04/2020 | 1348 | 0 | 24 | 1348 | 64183 | 361 |
| 48 | Alberta | Alberta | 07/04/2020 | 1373 | 0 | 26 | 1373 | 65265 | |
| 48 | Alberta | Alberta | 08/04/2020 | 1423 | 0 | 29 | 1423 | 66783 | |
| 48 | Alberta | Alberta | 09/04/2020 | 1451 | 0 | 32 | 1451 | 68116 | |
| 48 | Alberta | Alberta | 10/04/2020 | 1500 | 0 | 39 | 1500 | 70080 | |
| 48 | Alberta | Alberta | 11/04/2020 | 1569 | 0 | 40 | 1569 | 72779 | |
| 48 | Alberta | Alberta | 12/04/2020 | 1651 | 0 | 44 | 1651 | 74709 | |
| 48 | Alberta | Alberta | 13/04/2020 | 1732 | 0 | 46 | 1732 | 77007 | |
| 48 | Alberta | Alberta | 14/04/2020 | 1870 | 0 | 48 | 1870 | 79695 | |
| 48 | Alberta | Alberta | 15/04/2020 | 1996 | 0 | 48 | 1996 | 79695 | |
| 48 | Alberta | Alberta | 16/04/2020 | 2158 | 0 | 50 | 2158 | 85502 | |
| 48 | Alberta | Alberta | 17/04/2020 | 2397 | 0 | 50 | 2397 | 89144 | |

Figure 5: Doubling Rate of the Number of Cases

## Doubling Rate (.sh scripts)

Listing 6: .sh scripts .

```
1  numTotalindex=0;
2  dateindex=0;
3  provinceindex=0;
4  numDeathsindex=0;
5
6  var=$(cat covid19.csv | head -n 1 | tr ',' ' ')
7  echo $var
8  for colNames in $var
```
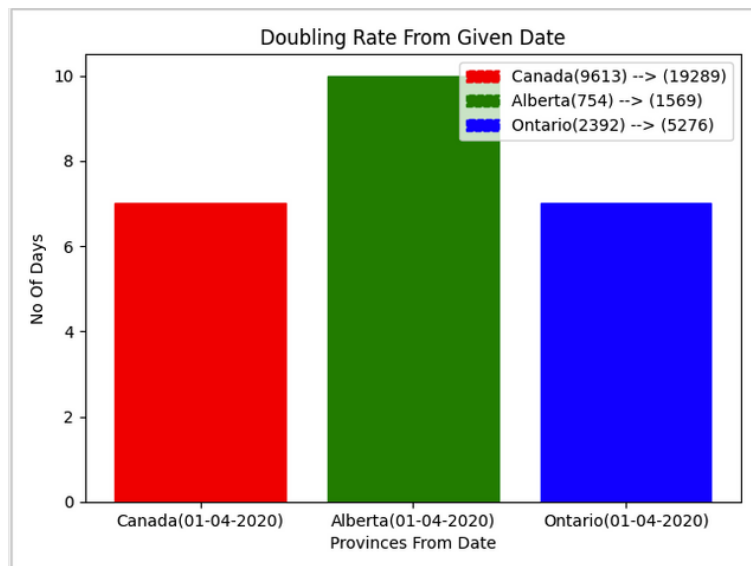
Figure 6: Doubling Rate Ploting

```
 9  do
10   if [ "numtotal" = $colNames ]
11   then
12     echo $colNames
13     echo $numTotalindex
14     break
15   fi
16   numTotalindex=$((numTotalindex+1))
17  done
18
19  for colNames in $var
20  do
21   if [ "date" = $colNames ]
22   then
23     echo $colNames
24     echo $dateindex
25     break
26   fi
27   dateindex=$((dateindex+1))
28  done
29
30  for colNames in $var
31  do
32   if [ "prname" = $colNames ]
33   then
34     echo $colNames
35     echo $provinceindex
36     break
37   fi
38   provinceindex=$((provinceindex+1))
```

```
39  done
40
41  for colNames in $var
42  do
43   if [ "numdeaths" = $colNames ]
44    then
45      echo $colNames
46      echo $numDeathsindex
47      break
48   fi
49   numDeathsindex=$((numDeathsindex+1))
50  done
51
52  convertDate()
53  {
54    var=$(echo $1 | tr "-" " " )
55    arr=($var)
56    d1=${arr[2]}"-"${arr[1]}"-"${arr[0]}
57    echo "$d1"
58  }
59
60  # Find Doubling Rate For Number Of Cases
61  if [ "$2" = "number of cases" ]
62  then
63      doublingRate=0
64      noOfdays=0
65      arr=''
66      echo "number of cases"
67      numberOfCsvRows=$(cat covid19.csv  | wc -l)
68      for((i=2;i<=numberOfCsvRows;i++))
69      do
70        var=$(cat covid19.csv | head -n "$i" | tail -n 1 | tr ',' ' ')
71        arr=($var)
72        d1=$(convertDate ${arr[dateindex]})
73        d2=$(convertDate "$3")
74        if [ "$1" = ${arr[provinceindex]} ] && [ "$d1" ">" "$d2" ]
75        then
76          #echo $d1
77          #echo $d2
78          echo $var
79          if [ ${arr[numTotalindex]} -gt "$doublingRate" ]
80          then
81            noOfdays=$((noOfdays+1))
82             break
83          else
84             noOfdays=$((noOfdays+1))
85           fi
```

```bash
 86      elif [ "$1" = ${arr[provinceindex]} ] && [ "$d1" "=" "$d2" ]
 87      then
 88       doublingRate=$((${arr[numTotalindex]} * 2))
 89       echo $doublingRate
 90      fi
 91    done
 92
 93    echo 'No Of Days = '$noOfdays
 94    echo $1 $3 $(($doublingRate / 2)) ${arr[dateindex]} ${arr[numTotalindex↩
         ]} $noOfdays  >> output.txt
 95
 96  else
 97    doublingRate=0
 98    noOfdays=0
 99    arr=''
100    echo "number of deaths"
101    numberOfCsvRows=$(cat covid19.csv  | wc -l)
102    for((i=2;i<=numberOfCsvRows;i++))
103    do
104      var=$(cat covid19.csv | head -n "$i" | tail -n 1 | tr ',' ' ')
105      arr=($var)
106      d1=$(convertDate ${arr[dateindex]})
107      d2=$(convertDate "$3")
108      if [ "$1" = ${arr[provinceindex]} ] && [ "$d1" ">" "$d2" ]
109      then
110        #echo $d1
111        #echo $d2
112        echo $var
113        if [ ${arr[numDeathsindex]} -gt "$doublingRate" ]
114        then
115          noOfdays=$((noOfdays+1))
116          break
117        else
118          noOfdays=$((noOfdays+1))
119        fi
120      elif [ "$1" = ${arr[provinceindex]} ] && [ "$d1" "=" "$d2" ]
121      then
122       doublingRate=$((${arr[numDeathsindex]} * 2))
123       echo $doublingRate
124      fi
125    done
126
127    echo 'No Of Days = '$noOfdays
128    echo $1 $3 $(($doublingRate / 2)) ${arr[dateindex]} ${arr[↩
         numDeathsindex]} $noOfdays  >> output.txt
129  fi
```

Doubling script is a shell script. It comprises of shell commands and with the help of pipelines and filters, doubling rate is being computed. Column index for required values like prname, numdeaths etc are identified to make the shell script dynamic. A loop is executed for number of rows and for each row, the code verifies whether the date is equal to the desired input date. The number of cases is computed on the basis of input date and it looks for the index where the value is being doubled. The number of days in parallel is counted and stored securely on Output.txt.

## Predicted data for next 7 days

Listing 7: FREE CHOICE .

```
1  eachProviceCases={}
2  for prov in provences:
3
4      # Exclude Canada
5      if prov != 'Canada':
6          # add province name as key and empty list as value
7          eachProviceCases[prov] = []
8
9          # append list of zeros equal to no of dates against province
10         for d in date:
11             eachProviceCases[prov].append(0)
12
13         for index,d in enumerate(date): # iterate on sorted date list
14             for idx,da in enumerate(data["date"]): # match it within date ↩
                   column of data
15                 if da == d and prov == data['prname'][idx]: # if date ↩
                       match and province name is equal
16                     eachProviceCases[prov][index] = int(data['numtotal'][↩
                           idx]) # add data against province for given date
17                     break
18
19 x_dates=[]
20
21 for i in range((0),14):
22     d=datetime.datetime.now() + datetime.timedelta(days=i)
23     x_dates.append(d.date())
24
25 for prov in provences:
26     # Exclude Canada
27     if prov != 'Canada':
28         eachProviceCasesUpdated={}
29         Stats = eachProviceCases[prov]
30         last7Values=[]
```

```
31          nextTwoWeekData=[]
32
33          for i in range((0), len(Stats)):
34              if i >= (len(Stats) - 7):
35                  last7Values.append(Stats[i])
36          print(prov)
37          print(last7Values)
38
39          #generate next 2 week data
40          for i in range((0),14):
41              count=0
42              sum=0
43              for j in range((i+1),len(last7Values)):
44                  sum = sum + (last7Values[j] - last7Values[j-1])
45                  count=count+1
46              last7Values.append(last7Values[len(last7Values)-1] + (sum/↩
                   count))
47              nextTwoWeekData.append(last7Values[len(last7Values)-1] + (sum/↩
                   count))
48              print(last7Values[len(last7Values)-1] + (sum/count))
49
50          eachProviceCasesUpdated[prov] = nextTwoWeekData
51          print(x_dates)
52          print(eachProviceCasesUpdated)
53          plotChartFunction(x_dates,eachProviceCasesUpdated,'No of cases in ↩
                   '+prov+' forecast data')
```

## Question Free Choice Graph

For the free choice task, the team decided to forecast the number of cases, data prediction for the upcoming next two weeks. A dictionary with sorted data for each province is ready and last seven days of data is excluded from the excel file. That data is used to take the average of difference in last 7 days values, adding it to the last value to predict the data for the next day. Moving are window across after adding the next value and repeating again until for next 14 days.
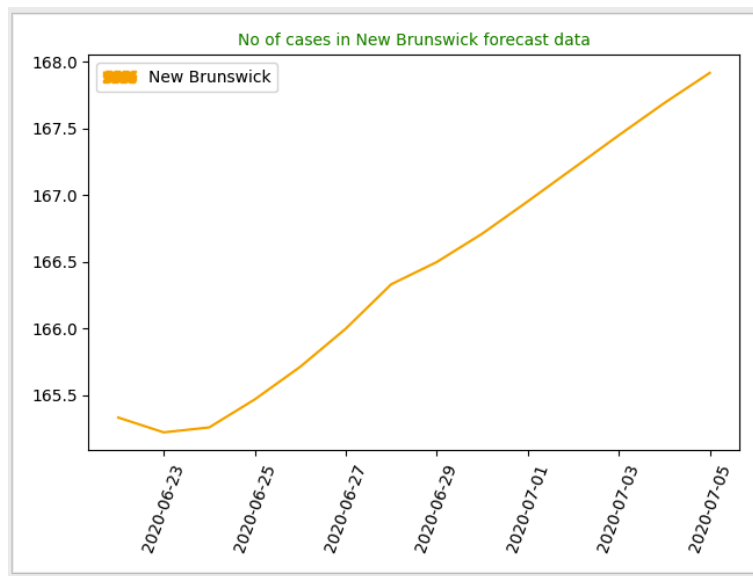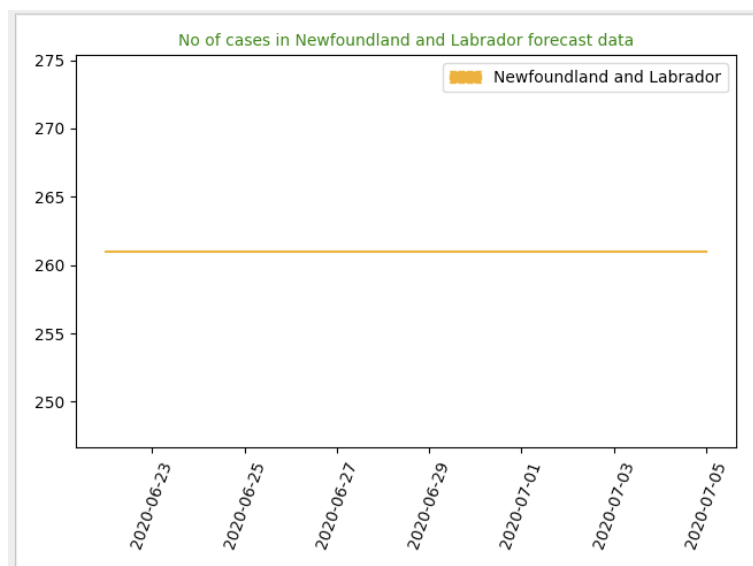
Figure 7: Forcast for burnswick



Figure 8: Forcast for NewFoundLand