

FORECASTING CASE COUNTS USING REGRESSION MODELS AND AN ANALYSIS ON COVID-19 DATASET

Venkateshwaran Thamilselvan

University of Luxembourg

113-6, Avenue du swing

L-4367, Belvaux, Luxembourg

+352 661364956

venkateshwaran.thamilselvan.001@student.uni.lu

ABSTRACT

As the COVID-19 pandemic tends to grow in number worldwide, we record huge data related to COVID-19. With the huge computing power which we have today, we can try to analyze the data collected either to find some useful facts or create some machine learning models with which we can predict the growth of pandemic in different parts of the world. In this use-case study, I have taken the John Hopkins datasets [1] to perform data analysis and to create machine learning model (mainly based on Linear regression and Random forest regression) to predict and forecast the growth of the pandemic in a particular locality or a country. Here we mainly predict and forecast the confirmed case count and the recovered case count.

PROBLEM STATEMENT

In this case study, we are going to do a detailed analysis on the following intriguing concepts-

1. With this huge data in our hand we should be able to rank the countries based on -
 - a. Current death count
 - b. Current confirmed case count
 - c. Current active case count
 - d. Current recovered case count
 - e. Case-Fatality-Ratio
 - f. Incidence rate

By performing the above analysis, we can understand different dynamics in different countries. For e.g. There may be countries who ranks higher in confirmed cases but lower fatality rate and death count. We can identify the nature of the pandemics in different countries.

2. Finding the correlation between the confirmed case counts, recovered case counts and death counts worldwide and for specific counties. With this analysis, we can understand how the different records correlate with each other.
3. Creating machine learning models to predict the confirmed case counts. These types of machine learning models will help us understand how pandemic is spreading in different regions around the world and help us forecast the confirmed case counts with

some accuracy for some future date. This task includes-

- a. Creating random forest regression model with the whole data and predict the confirmed cases for a particular location (using latitude and longitude) at a particular date. We should measure the accuracy by some evaluation metrics such as RMSE or MSE.
 - b. Creating the random forest regression model for the US specific data (using the US specific records in John Hopkins university data) and to predict the confirmed cases on a particular location in a state of the country. This makes sense because the US specific data has many features recorded such as hospitalization rate, testing rate etc.
4. Visualization of the trends-
 - a. Plot the worldwide trend of confirmed, recovered and death counts as a line chart.
 - b. Plot the trends of confirmed, recovered and death counts as a line chart separately for each country.
 5. Forecasting the confirmed case count and recovery case count using the linear regression model for Luxembourg for the next 10 days and plotting it as a line chart.
 6. Using the Mobility dataset [2] to find the correlation between the confirmed case counts and the following-
 - a. Retail and recreation percent change from baseline
 - b. Grocery and pharmacy percent change from baseline
 - c. Parks percent change from baseline
 - d. Transit stations percent change from baseline
 - e. Workplaces percent change from baseline
 - f. Residential percent change from baseline

APPROACH

Since we are handling huge data such as Mobility data (136MB – Single file), COVID-19 data provided by John Hopkins University (Totally 612MB), we should use some distributed systems to process and analyze the

data. Here in this case study, I have used the HPC cluster provided by University of Luxembourg [3] for distributed analysis.

Project Details-

- i. In this project, I have used Scala as a frontend language for Spark.
- ii. ML is used for creating regression models.
- iii. MLlib is used for getting the statistics of the data.
- iv. Scalaplot library [4] is used to plot the line charts.
- v. Spark SQL is used to perform queries to filter our data.
- vi. Spark SQL-Dataframe and RDD is used to store our data in the distributed manner.

Now we will see how to solve the above problem statement section-wise.

1. We have taken the `csse_covid_19_data` from JHU for this part as it contains all the required data for this section. In `csse_covid_19_daily_reports`, the date format is different for different period of time. Moreover, the date formats within a single file is same. So, we read each file separately using a loop and parse date with different date formats used and store it as timestamp. This while csv files are read as dataframes and the dataframe corresponding to each date are then merged with our original dataframe. This final dataframe contains all the data with proper timestamp. The unfilled entries in Confirmed, death and recovery are actually zeros, so we update nulls with 0 for those columns. The active is computed with Eq. 1, where *act* is active case count, *conf* is confirmed case count and *deaths* are the death count.

$$act = conf - recov - deaths \quad (1)$$

Since the all the data are accumulated, meaning that the count for today is computed by adding the previous count with the new cases today, the latest date has the total counts for all regions. So, to rank the countries, we take the maximum date in our data and select the rows with this maximum date. After that we group the data by countries and aggregate it to find the sum of confirmed, recovered, death and active case counts.

Since the Case-Fatality-Ratio (CFR) is computed with the Eq. 2, we cannot just sum up all the ratios. Since we have both the

death and confirmed case counts for a country, we compute CFR with the Eq. 2.

$$CFR = deaths/confirmed \quad (2)$$

$$IR = cases / 100,000 \text{ persons} \quad (3)$$

The Incidence Rate (IR) takes the population of the locality into account. As this is normalized, to compute IR for a country, we just take the average of IR of all the state/province of that country.

2. To find the correlation of a country, we first filter the data of that country from the whole data. Then, we group them by date and take aggregate to find the sum of confirmed, recovered, death and active case counts. In this aggregated count, we have the count of whole country for the corresponding date. Then we sort the date in descending order and compute the spearman correlation coefficient. To compute the correlation coefficient worldwide, we do exactly the above steps, except that at first, we do not select the data for a particular country.
3. a. To perform random forest regression, we first remove the string data types such as `Combined_Key`, `Admin2` and some data which have lots of null rows such as `Case-Fatality_Ratio` and `Incidence_Rate`. Then, we group the data frame by date, province/state, country, latitude and longitude. After that we aggregate to find the sum of confirmed, recovered, deaths and active cases. Finally, we will get the aggregated data for all localities for all dates.

Some entries in our data do not have the latitude and longitude. There is a csv file named `UID_ISO_FIPS_LookUp_Table.csv` which has the latitude and longitude for each locality. We lookup that csv file with the province/state and country as key and take its latitude and longitude and put it in our data. This drastically reduces the null occurrence and helps us a lot for our regression.

We also have the population of each locality in the lookup table. So, we perform left join on those two data based on the latitude and longitude as the key. By doing that we get the population column for each locality in our dataframe. This may drastically increase the performance of our regression as the population plays a key role in the spread of pandemic.

Finally, we select the data with dates less than 2020-08-15 as train data and greater than that date as the test data. We also create param grid for cross validation. We use the the combination in table 1 for param grid creation. Then we compute the root mean square error for our test data.

Hyperparameter	values
Impurity	"variance"
Max Depth	10, 20
Max Bins	10, 100, 300

Table 1

b. Similarly, we perform the regression for US specific data. Before that, for preprocessing US data, we round the latitude and longitude to 5 decimals so as they match with the lookup table, and also update null occurrences of recovery with 0 and null occurrences of active with the newly computed value from Eq. 1

4. **NOTE:** For plotting the line charts, we should have opened the spark-shell as below (with the Scalaplot library's package) -

```
spark-shell --packages =
"org.sameersingh.scalaplot:scalaplot:0.0.4"
```

For plotting the worldwide data, we group our dataframe by date and aggregate it to find the sum of Confirmed, recovered and deaths. Since we already have such a dataframe, we reuse it for this purpose. Thus, we will get the unique entries for each date.

Since only the double data format is supported in the x axis of the line chart, we sort the dataframe in ascending order of date. Then the first date is encoded as 1 and incremented by 1 for the following dates. By this way we encode the date to doubles for the ease of plotting in Scalaplot library.

To plot chart separately for each country, we first get the list of distinct country names from our dataframe. Then we iterate through each country name in the list and then select the data for that country and plot as above.

5. For forecasting, we use the timeseries data provided by JHU. The timeseries for confirmed and recovered case counts are given in separate files and each date is in separate columns. For our convenience we transpose the table such that the date columns come to row with its title as Date

and other columns such as Province/State, Country/Region, Latitude and Longitude lies as it is. By doing so, we can easily lookup the count for a location on a specific date. Then we parse the date column to timestamp data type. Since the string datatypes are not supported, we drop the Country and province/state columns. We can then lookup the country and province/state with the help of the latitude and longitudes.

For forecasting, we use the linear regression model. For cross validation we use the hyperparameter values from table 2.

Hyperparameter	values
Max Iteration	10, 20, 30
Tolerance	0.00000001, 0.0000001, 0.000001
Optimization solver	"l-bfgs", "normal", "auto"
Elastic net param	0.5, 0.6, 0.8
Regularization param	0.1, 0.3, 0.5

Table 2

The elastic net param, tolerance and regularization params should be chosen properly so that the overfitting of the data does not occur. The values in the above table are the perfect values for creating combination for training our model. We have filtered and taken the data only for Luxembourg for training our model. We put all of the data as the training data and we read the lux_test.csv. After creating the model for confirmed and recovered cases separately, we predict the confirmed and recovered cases for our test data. Finally we plot the actual confirmed and recovered case counts along with the forecasted confirmed and recovered case counts.

6. For mobility data analysis, we read the Global_Mobility_Report.csv file provided by google and then preprocess it by converting the date string to timestamp format. After that we group the dataframe by country and date so that we can analyze based on the country. Then we select only the Luxembourg's data and sort it based on the date. Then we take the required columns as rdd to compute the spearman correlation.

EVALUATION AND RESULT

1. Ranking

a. Ranking based on Death count

```
scaled > deathList.sort(desc("Sum_Deaths")).show(20, false)
```

Country_Region	Min_Deaths	Max_Deaths	Sum_Deaths
US	0	123602	116023
Brazil	1482	125369	104201
Mexico	1228	19598	54666
India	0	118650	47033
United Kingdom	0	142072	46791
Italy	123	116833	35225
France	0	130247	38375
Spain	0	18464	28579
Peru	0	110183	121713
Iran	118988	118988	118988
Russia	0	14611	115231
Colombia	0	13998	113837
South Africa	111010	111010	111010
Chile	11	17992	110205
Belgium	19900	19900	19900
Germany	0	12627	19213
Canada	0	15709	19052
Netherlands	0	11542	16182
Pakistan	158	12297	16139
Ecuador	15984	15984	15984

only showing top 20 rows

Global Deaths
751,887
166,971 deaths
US
104,201 deaths
Brazil
54,666 deaths
Mexico
47,033 deaths
India
46,791 deaths
United Kingdom
35,231 deaths
Italy
30,376 deaths
France
28,605 deaths
Spain

The actual deaths which is obtained from the website [5] source provided by John Hopkins University is on the right (as of 15th Aug 2020). Our output is shown in the left image.

- By comparing these two data, we can see that we have computed correctly (Sum_Deaths is similar to the website's data) and so we can proceed to study the dataset further.
- The Min_Deaths is the minimum number of deaths marked by at least one of the province/states of the country.
- The Max_Deaths is the maximum number of deaths marked by at least one of the province/states of the country.

b. Ranking based on confirmed cases

```
scaled > confirmedList.sort(desc("Sum_Confirmed")).show(20, false)
```

Country_Region	Min_Confirmed	Max_Confirmed	Sum_Confirmed
US	0	1220167	15248702
Brazil	122242	1674455	13274876
India	0	1560126	12461190
Russia	1162	1250303	1905762
South Africa	1572865	1572865	1572865
Mexico	12524	183683	1505751
Peru	0	1250708	1498555
Colombia	16	1419944	1433805
Chile	164	1264222	1380034
Spain	0	190890	1337334
Iran	1336324	1336324	1336324
United Kingdom	0	1270971	1315680
Saudi Arabia	1294519	1294519	1294519
Pakistan	12164	1124929	1286674
Argentina	1276072	1276072	1276072
Bangladesh	1269115	1269115	1269115
Italy	1478	197128	1252235
Turkey	1245635	1245635	1245635
France	14	1230778	1244096
Germany	1959	153123	1222281

only showing top 20 rows

Global Cases
5,254,878
Cases by Country/Region/State/province

- The actual confirmed cases which obtained from the website [5] source provided by John Hopkins University is on right (as of 15th Aug 2020). The Left image is our output.
- By comparing our output and their website output, we can conclude that our result is correct. But it may not be precisely equal, and the websites data will be mostly slightly larger than our data. This is because website data is updated more frequently than the git repository.

c. Ranking based on Active cases

```
scaled > activeList.sort(desc("Sum_Active")).show(20, false) //Top 20 rank list
```

Country_Region	Min_Active	Max_Active	Sum_Active
US	1-1774648	1209312	13306948
India	0	150105	1661595
Brazil	1-255	1190373	1598313
United Kingdom	0	1228899	1267330
Russia	10	152208	1175475
Colombia	14	164991	1169166
Spain	0	158978	1158353
Peru	1-341938	1240525	1134904
France	11	1128785	1130092
South Africa	1123978	1123978	1123978
Bangladesh	1110687	1110687	1110687
Argentina	178276	178276	178276
Sweden	1284	121548	178076
Philippines	174713	174713	174713
Bolivia	158855	158855	158855
Netherlands	11	114698	159966
Belgium	148362	148362	148362

Comparing active and confirmed case ranking, we can see that India is gaining in the active count compared to Brazil and thus in near future the total deaths in India is likely to increase than Brazil. But we cannot say anything with certainty, because the death rate is linked to-

- The percentage of old aged people affected
- percentage of COVID-19 patients with other background illnesses
- the health facility availability and affordability.

d. Ranking based on Recovered cases

```
scaled > recoveredList.sort(desc("Sum_Recovered")).show(20, false) //Top 20 rank list
```

Country_Region	Min_Recovered	Max_Recovered	Sum_Recovered
Brazil	19345	1457758	12616981
India	0	1401442	11808936
US	0	11796326	11796326
Russia	198	1194823	1721473
South Africa	1461734	1461734	1461734
Mexico	11681	168437	1410479
Chile	162	1251316	1355037
Peru	0	1354232	1354232
Iran	1293811	1293811	1293811
Pakistan	11981	1118924	1265215
Saudi Arabia	1262959	1262959	1262959
Colombia	12	184153	1261293
Turkey	1228980	1228980	1228980
Italy	1390	175012	1203326
Germany	0	148400	1200440
Argentina	1199005	1199005	1199005
Bangladesh	1156623	1156623	1156623
Spain	0	140736	1150376
Iraq	1120129	1120129	1120129
Qatar	1111258	1111258	1111258

only showing top 20 rows

e. Ranking based on CFR

```
scaled > cfrList.sort(desc("CFR")).show(20, false) //Top 20 rank list
```

Country_Region	Sum_Confirmed	Sum_Deaths	CFR
Yemen	11730	1494	128.55491329479769
MS Zaandam	19	12	122.22222222222222
United Kingdom	1305572	146278	115.144712211851871
Belgium	169402	19845	114.185470159361401
Italy	1247832	135146	114.181380935472417
France	1225198	130268	113.440616701746908
Hungary	14526	1597	113.190455148033584
Netherlands	155021	16167	111.208447683611713
Mexico	1434193	147472	110.933386765793092
Western Sahara	10	11	110.0
Spain	1288522	128445	109.85866914827985
Chad	1936	175	109.012820512820513
Canada	1118510	18886	107.582482490929036
Sweden	180422	15743	107.141080788786402
Ireland	126109	11763	106.752460837259183
Ecuador	186232	15736	106.651822989145561
Sudan	111738	1752	106.406542852746663
Liberia	11189	175	106.307821698906644
Niger	11136	169	106.073943661917831
San Marino	1699	142	106.00858369087124

only showing top 20 rows

In above figure the aggregated CFR is calculated by summing up deaths and confirmed cases in that region and dividing the death by confirmed cases. It is then multiplied

with 100 to represent as percentage. There can be many reasons for a very large CFR. They are-

- The percentage of people with underlying health conditions is larger.
- The preventive and precautionary measures are not properly undertaken.
- People of that country generally has bad hygiene.
- The health systems are not up to the par.
- The testing is not properly done, therefore actual active cases in that region is very much higher than the reported value.

f. Ranking based on IR

```
scala> incidenceList.sort(desc("Avg_Incidence_Rate")).show(20, false) //Top 20 rank list
+-----+-----+
|Country_Region|Avg_Incidence_Rate|
+-----+-----+
|Qatar|13975.342408696799|
|Bahrain|12706.421020896424|
|Brazil|12195.5207804316337|
|San Marino|12059.6381637102954|
|Panama|11840.2379919338431|
|Kuwait|11760.540706225385|
|Oman|11620.30790608742931|
|Holy See|11483.3127317676144|
|Armenia|11393.713760033801|
|Chile|11393.3450966070952|
|Peru|11374.5281504994193|
|Andorra|11280.0103539765742|
|Luxembourg|11182.9527010620216|
|US|11176.9962538803211|
|Israel|1052.273913323269|
|Maldives|1030.8172167935147|
|South Africa|976.4842217894288|
|Singapore|950.0297674854277|
|Saudi Arabia|849.954416152621|
|Bolivia|839.1138238412669|
+-----+-----+
only showing top 20 rows
```

This considers the population of that country in account which is much reasonable.

Incidence rate is already there in the dataset, but it is there only for separate province/state. We can compute the incidence rate of a country by grouping and taking average of all its provinces/states. This way of computing is correct because the incidence rate for all region is in per 100,000 persons.

- The global incidence rate is around 335. This means that every 335 persons out of 100,000 in the world has COID-19. This count can be used to compare the various pandemics.
- In figure 8, we can see that Luxembourg has 3.5 times the global average incidence rate which is an alarming count.

2. Correlation of Confirmed-Recovered-Death

For Luxembourg:

```
scala> println(s"Confirmed-Recovered correlation is: $confirmed_recovered_correlation")
Confirmed-Recovered correlation is: 0.9993150450452986

scala> println(s"Confirmed-Death correlation is: $confirmed_death_correlation")
Confirmed-Death correlation is: 0.986776276372286

scala> println(s"Death-Recovered correlation is: $death_recovered_correlation")
Death-Recovered correlation is: 0.9861071562346669
```

Worldwide:

```
scala> println(s"Confirmed-Recovered correlation is: $confirmed_recovered_correlation")
Confirmed-Recovered correlation is: 0.998808537278034

scala> println(s"Confirmed-Death correlation is: $confirmed_death_correlation")
Confirmed-Death correlation is: 0.999923173126207

scala> println(s"Death-Recovered correlation is: $death_recovered_correlation")
Death-Recovered correlation is: 0.9986347486690788
```

- From the above results we can see that Confirmed case and Death count has the highest correlation.
- When we increase the sample count, the correlation increases. This means that all these data tend to correlate between each other.

3. Predicting confirmed cases

a. Worldwide data

```
scala> predictions.drop("featureVector").show(5)
+-----+-----+-----+-----+-----+-----+-----+-----+
|Latitude|Longitude|Last_Update|Deaths|Recovered|Confirmed|Active|UID(code)|Population|prediction|
+-----+-----+-----+-----+-----+-----+-----+-----+
|38.05721|107.67412020-08-10 04:27:42|61|5771|5831|0|150801|150|31030000|580.9532040784281|
|132.00254883|85.35132240-2020-08-10 04:27:42|47|0|2760|273184001811|840|1045421247|172520119980|
|137.18049267|-89.32918818-2020-08-10 04:27:42|0|0|39|391640170031|840|5761139.12560809990216061|
|43.35073842|-92.11702714-2020-08-10 04:27:42|0|0|501|501844190091|840|9158154.94680547636381|
|38.44185191|-84.23604741-2020-08-10 04:27:42|0|0|1231|123184021007|840|38800125.112228265102361|
+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 5 rows

scala> val rmse = evaluator.evaluate(predictions)
rmse: Double = 7399.476193416496

scala> println(s"Root Mean Squared Error (RMSE) on test data = $rmse")
Root Mean Squared Error (RMSE) on test data = 7399.476193416496
```

b. US specific data

```
scala> predictions.drop("featureVector").show(5)
+-----+-----+-----+-----+-----+-----+-----+-----+
|Latitude|Longitude|Last_Update|Deaths|Recovered|Confirmed|Active|UID(code)|Population|prediction|
+-----+-----+-----+-----+-----+-----+-----+-----+
|40.32509351|-122.2370171-2020-08-10 04:27:56|11|0|318184001001|840|630841305.1470105017746|
|132.71301833|-85.30861208-2020-08-10 04:27:56|0|0|130184013079|840|124041124.5090027500799|
|132.74976524|-81.08779441-2020-08-10 04:27:56|91|0|240184013251|840|139061238.709493620652841|
|42.47204507|-92.3664032-2020-08-10 04:27:56|11|0|120184019051|840|27021128.4516432256420|
|137.19113093|-95.20840679-2020-08-10 04:27:56|11|0|159184020091|840|130181146.763962961866981|
+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 5 rows

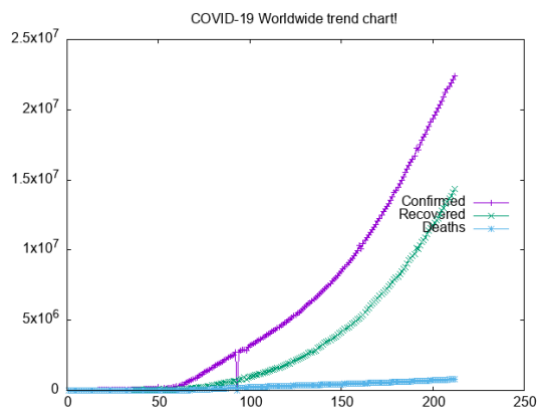
scala> val rmse = evaluator.evaluate(predictions)
rmse: Double = 7780.5982607258855

scala> println(s"Root Mean Squared Error (RMSE) on test data = $rmse")
Root Mean Squared Error (RMSE) on test data = 7780.5982607258855
```

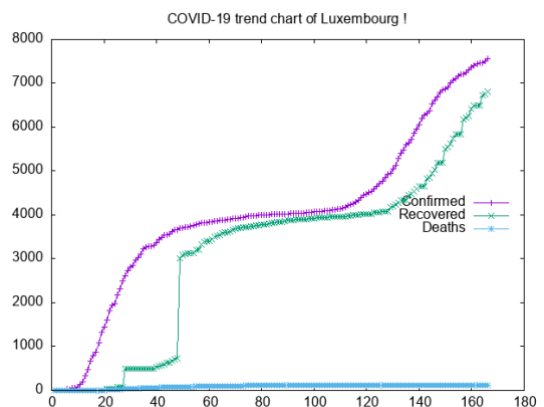
The above model predicts the confirmed cases for the past dates with some accuracy.

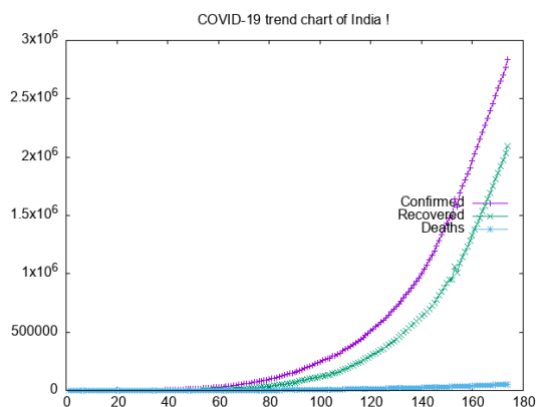
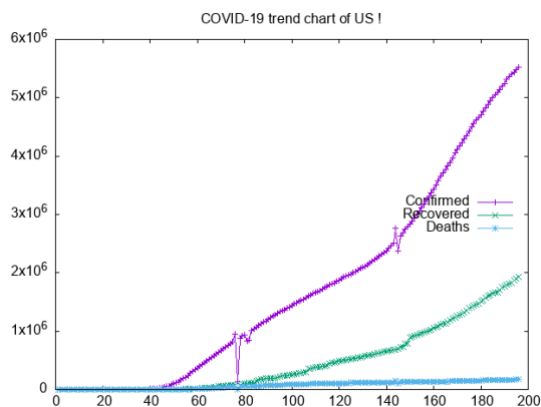
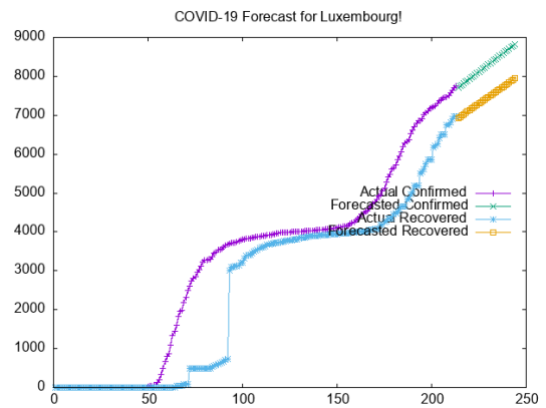
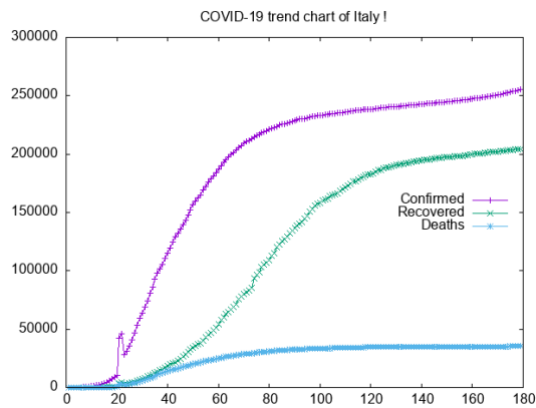
4. Trend charts

a. Worldwide trend chart



b. Country wise trend chart





The trend chart for all the countries is in [output/charts/countrywise](#).

5. Forecast for Luxembourg using linear regression

6. Correlation of Mobility and Confirmed

```
scala> println(s"Confirmed - Retail and Recreation correlation is: $confirmed_retail_correlation")
Confirmed - Retail and Recreation correlation is: 0.12892183742536742

scala> println(s"Confirmed - Grocery and Pharmacy correlation is: $confirmed_grocery_correlation")
Confirmed - Grocery and Pharmacy correlation is: 0.83291868988765576

scala> println(s"Confirmed - Park visit correlation is: $confirmed_parks_correlation")
Confirmed - Park visit correlation is: 0.74788845930992

scala> println(s"Confirmed - Transit station correlation is: $confirmed_transit_correlation")
Confirmed - Transit station correlation is: 0.1672224085827479

scala> println(s"Confirmed - Workplace visit correlation is: $confirmed_workplace_correlation")
Confirmed - Workplace visit correlation is: 0.12275871185083565

scala> println(s"Confirmed - Residential correlation is: $confirmed_residential_correlation")
Confirmed - Residential correlation is: -0.1310433252188186
```

Bibliography

- [1] J. H. University, "COVID-19 Data Repository by the Center for Systems Science and Engineering (CSSE) at Johns Hopkins University," [Online]. Available: <https://github.com/CSSEGISandData/COVID-19>.
- [2] Google, "COVID-19 Community Mobility Reports," [Online]. Available: <https://www.google.com/covid19/mobility/>.
- [3] U. o. Luxembourg, "High Performance Computing in Luxembourg," [Online]. Available: <https://hpc.uni.lu/>.
- [4] Sameersingh. [Online]. Available: <https://github.com/sameersingh/scalaplot>.
- [5] J. H. University, "John Hopkins university website for COVID-19 trends," [Online]. Available: <https://www.arcgis.com/apps/opsdashboard/index.html#/bda7594740fd40299423467b48e9ecf6>.