









[illegible]



```
[33]: badly_affected_countries = ["France","Italy","United Kingdom","Spain","Iran","Germany"]
for country in badly_affected_countries:
    get_RMSLE_per_region(country, groundtruth_df, display_only=True)

France
RMSLE on cases: 0.6055889849865056
RMSLE on fatalities: 0.917603893958617

Italy
RMSLE on cases: 0.21993188939467181
RMSLE on fatalities: 0.27908283478945

United Kingdom
RMSLE on cases: 0.702326332485703
RMSLE on fatalities: 1.0985807570598574

Spain
RMSLE on cases: 0.3452272671599796
RMSLE on fatalities: 0.4824718447195897

Iran
RMSLE on cases: 0.2918973895938459
RMSLE on fatalities: 0.24915402316605503

Germany
RMSLE on cases: 0.36357028547960063
RMSLE on fatalities: 0.8511859973565415

In [34]: healthy_countries = ["Taiwan","Singapore","Kenya","Slovenia","Portugal", "Israel"]
for country in healthy_countries:
    get_RMSLE_per_region(country, groundtruth_df, display_only=True)

Taiwan*
RMSLE on cases: 0.4725040374958486
RMSLE on fatalities: 0.4312974401612734

Singapore
RMSLE on cases: 0.42418201374459735
RMSLE on fatalities: 0.3024619370927203

Kenya
RMSLE on cases: 0.453351452584227
RMSLE on fatalities: 1.074759898089988

Slovenia
RMSLE on cases: 0.12575939557457347
RMSLE on fatalities: 0.1639770020679584

Portugal
RMSLE on cases: 0.463214022221402
RMSLE on fatalities: 0.2443228659880897

Israel
RMSLE on cases: 0.4725838196781893
RMSLE on fatalities: 0.2281550922784558
```

## 8. Outputs: Observing the curves

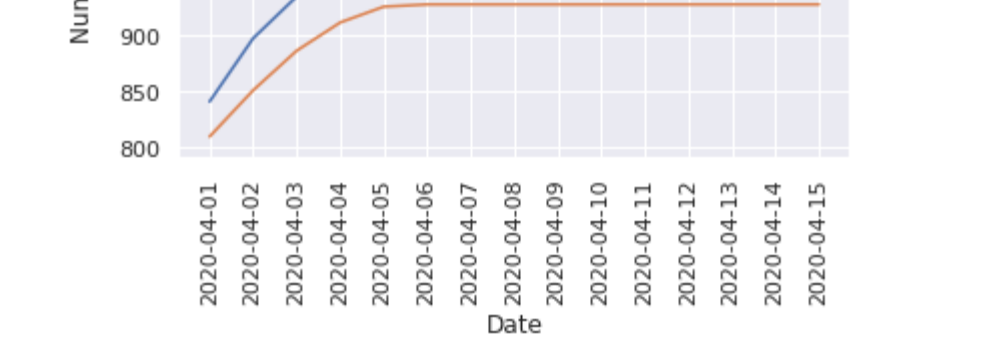
I am using the original training CSV file instead of the enriched dataset as it is up-to-date with the latest stats from this week.

```
In [35]: def display_comparison(region,groundtruth_df):
    groundtruth = groundtruth_df.query("Country_Region=="+"region+" and Date>='2020-04-01' and Date<
    = '2020-04-15'")
    prediction = copy_df.query("Country_Region=="+"region+" and Date>='2020-04-01' and Date<='2020-04-
    15'")

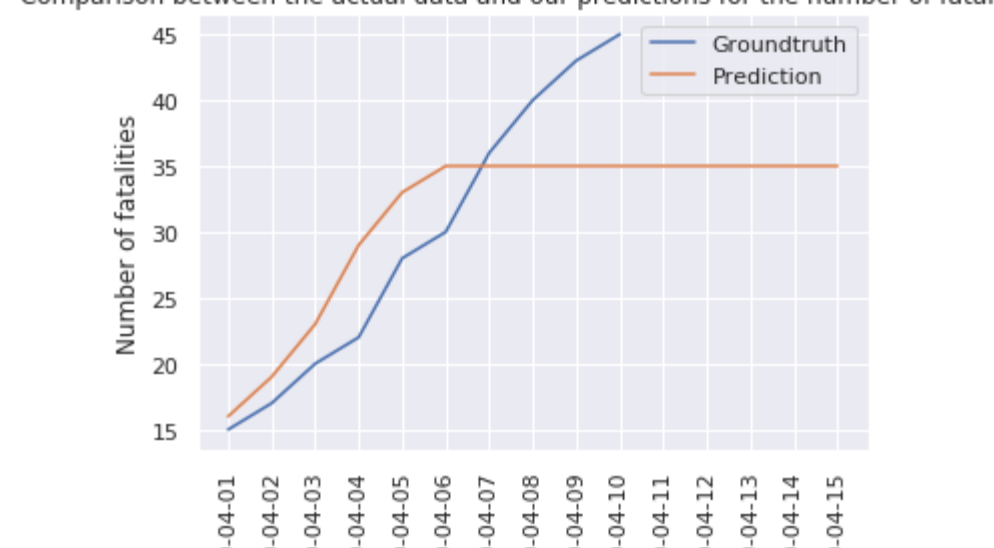
    plt.plot(groundtruth.ConfirmedCases.values)
    plt.plot(prediction.ConfirmedCases.values)
    plt.title("Comparison between the actual data and our predictions for the number of cases")
    plt.ylabel("Number of cases")
    plt.xlabel("Date")
    plt.xticks(range(len(prediction.Date.values)),prediction.Date.values,rotation='vertical')
    plt.legend(['Groundtruth', 'Prediction'], loc='best')
    plt.show()

    plt.plot(groundtruth.Fatalities.values)
    plt.plot(prediction.Fatalities.values)
    plt.title("Comparison between the actual data and our predictions for the number of fatalities")
    plt.ylabel("Number of fatalities")
    plt.xlabel("Date")
    plt.xticks(range(len(prediction.Date.values)),prediction.Date.values,rotation='vertical')
    plt.legend(['Groundtruth', 'Prediction'], loc='best')
    plt.show()

In [36]: display_comparison("Canada,Newfoundland and Labrador", groundtruth_df)
```

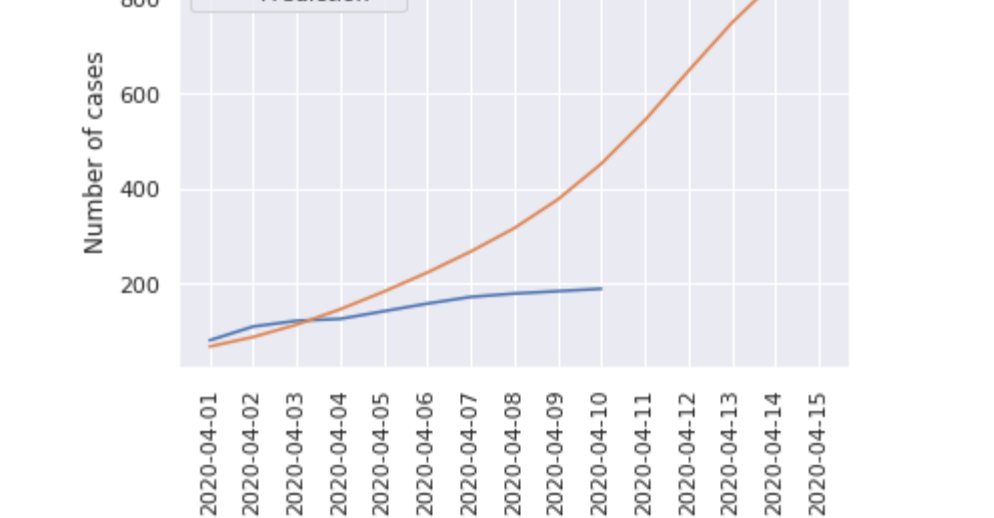


Comparison between the actual data and our predictions for the number of fatalities

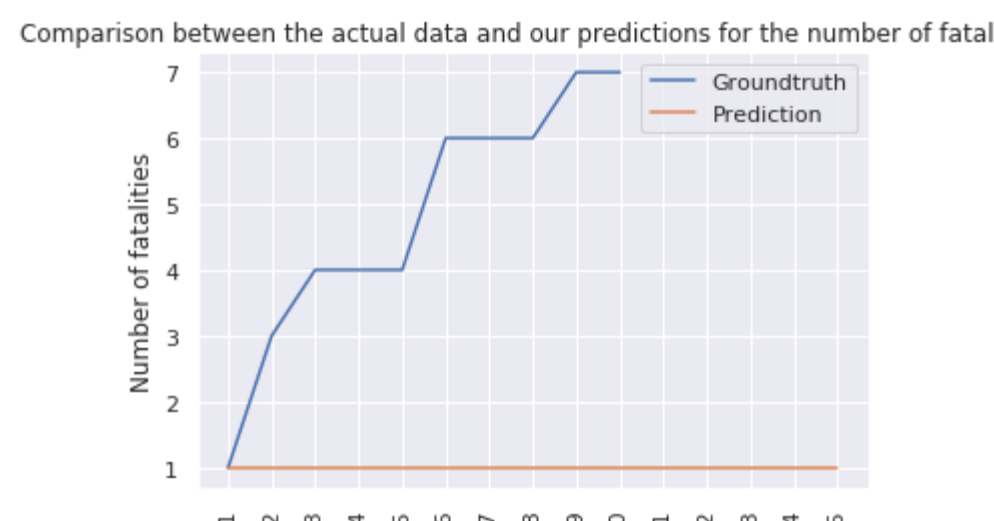


```
In [37]: display_comparison("Slovenia", groundtruth_df)
```

Comparison between the actual data and our predictions for the number of cases

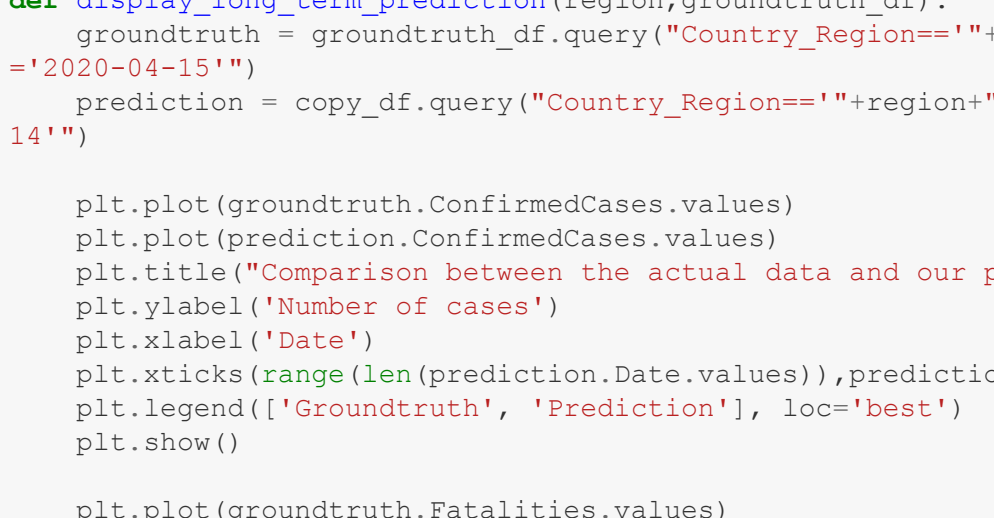


Comparison between the actual data and our predictions for the number of fatalities

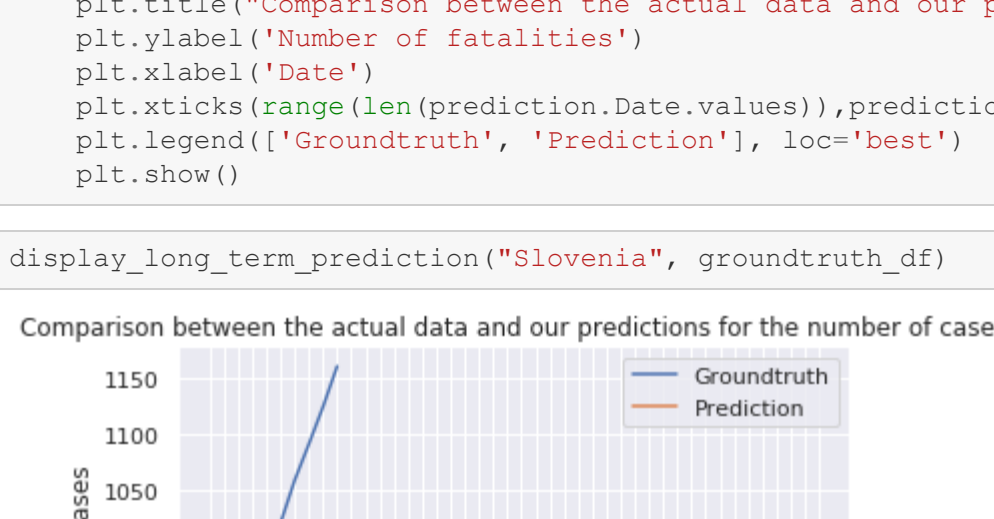


```
In [38]: display_comparison("Kenya", groundtruth_df)
```

Comparison between the actual data and our predictions for the number of cases



Comparison between the actual data and our predictions for the number of fatalities



Now let's see how well the model perform on a much longer period.

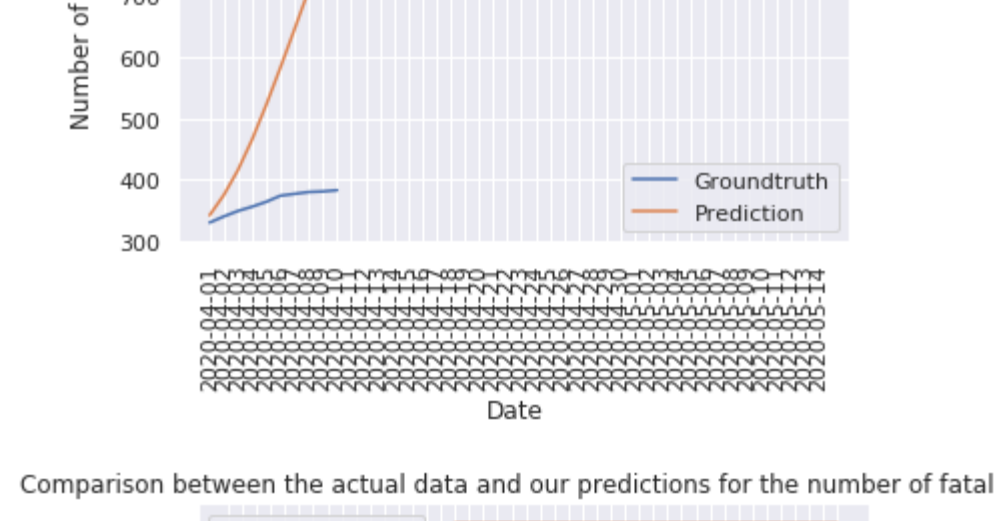
```
In [39]: def display_long_term_prediction(region,groundtruth_df):
    groundtruth = groundtruth_df.query("Country_Region=="+"region+" and Date>='2020-04-01' and Date<
    = '2020-04-15'")
    prediction = copy_df.query("Country_Region=="+"region+" and Date>='2020-04-01' and Date<='2020-05-
    14'")

    plt.plot(groundtruth.ConfirmedCases.values)
    plt.plot(prediction.ConfirmedCases.values)
    plt.title("Comparison between the actual data and our predictions for the number of cases")
    plt.ylabel("Number of cases")
    plt.xlabel("Date")
    plt.xticks(range(len(prediction.Date.values)),prediction.Date.values,rotation='vertical')
    plt.legend(['Groundtruth', 'Prediction'], loc='best')
    plt.show()

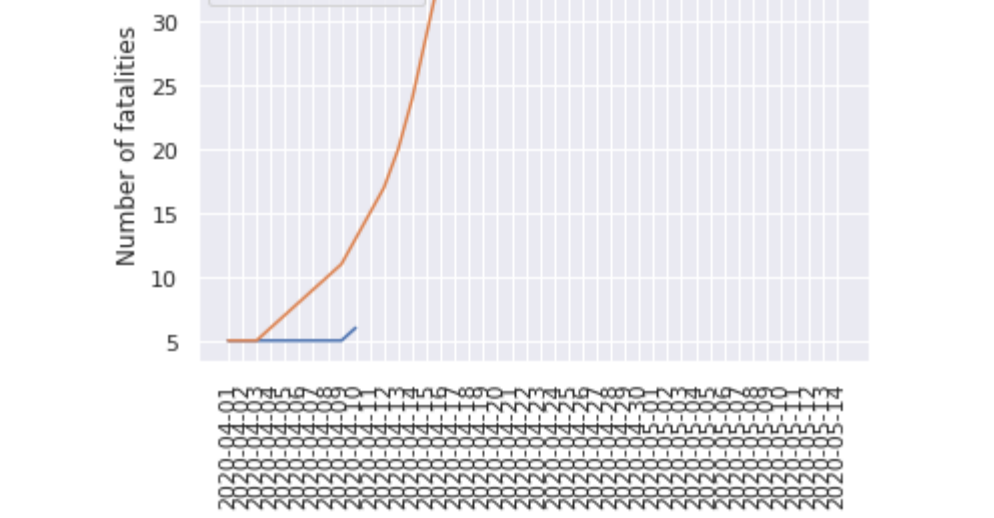
    plt.plot(groundtruth.Fatalities.values)
    plt.plot(prediction.Fatalities.values)
    plt.title("Comparison between the actual data and our predictions for the number of fatalities")
    plt.ylabel("Number of fatalities")
    plt.xlabel("Date")
    plt.xticks(range(len(prediction.Date.values)),prediction.Date.values,rotation='vertical')
    plt.legend(['Groundtruth', 'Prediction'], loc='best')
    plt.show()

In [40]: display_long_term_prediction("Slovenia", groundtruth_df)
```

Comparison between the actual data and our predictions for the number of cases

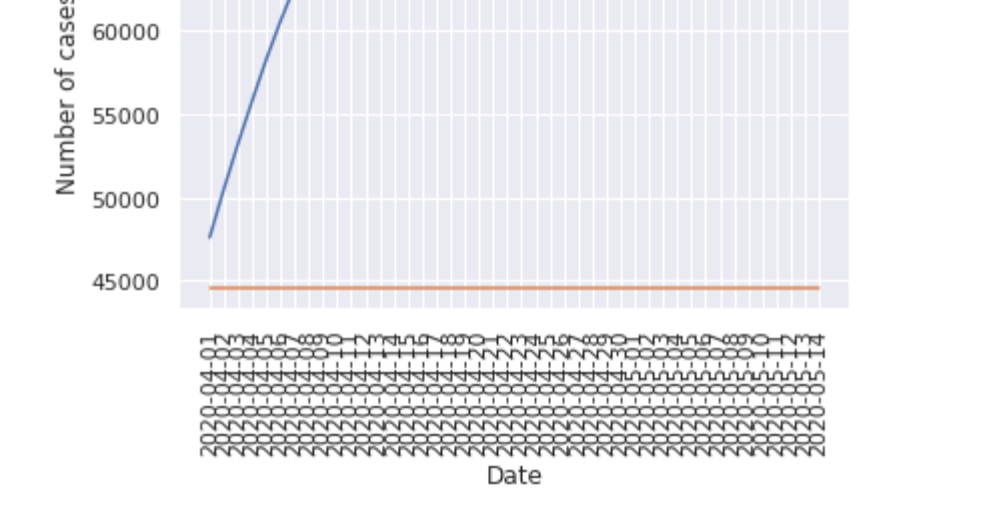


Comparison between the actual data and our predictions for the number of fatalities

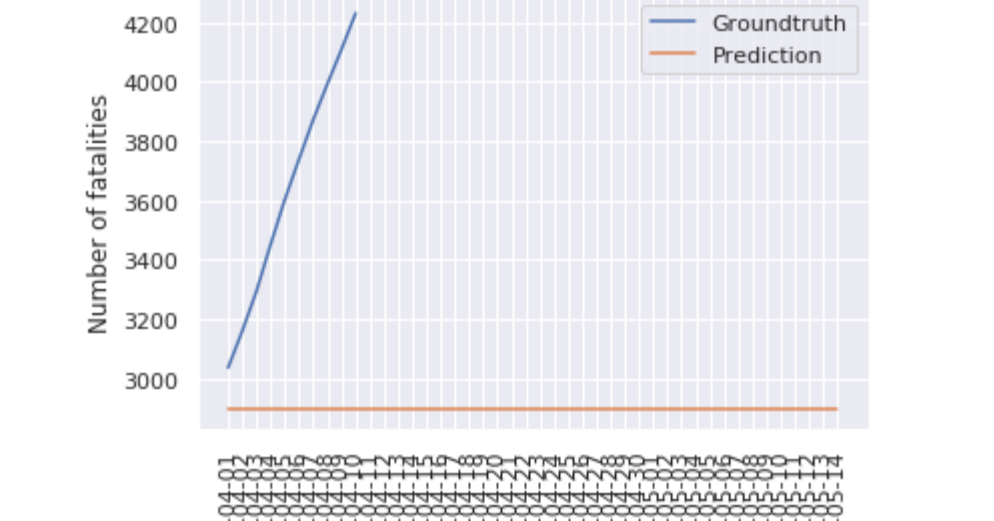


```
In [41]: display_long_term_prediction("Taiwan", groundtruth_df)
```

Comparison between the actual data and our predictions for the number of cases

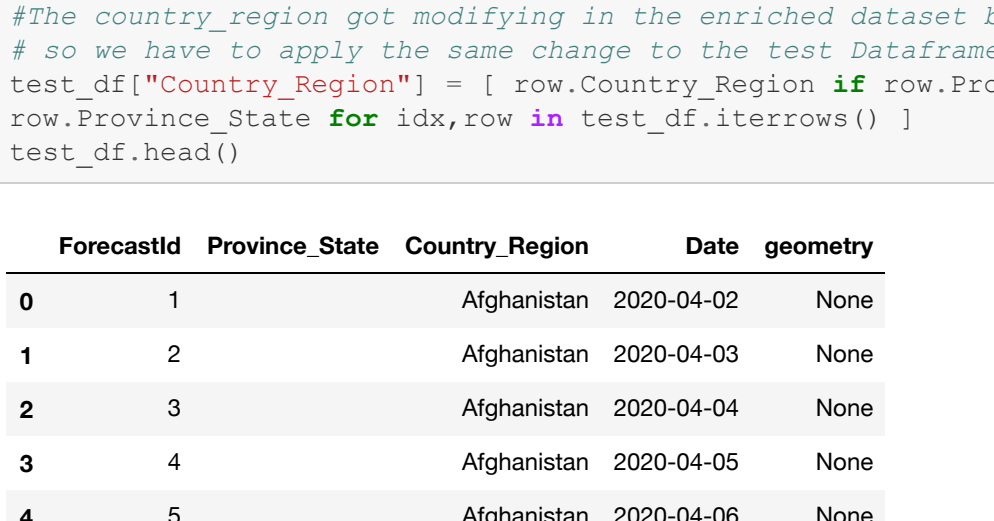


Comparison between the actual data and our predictions for the number of fatalities

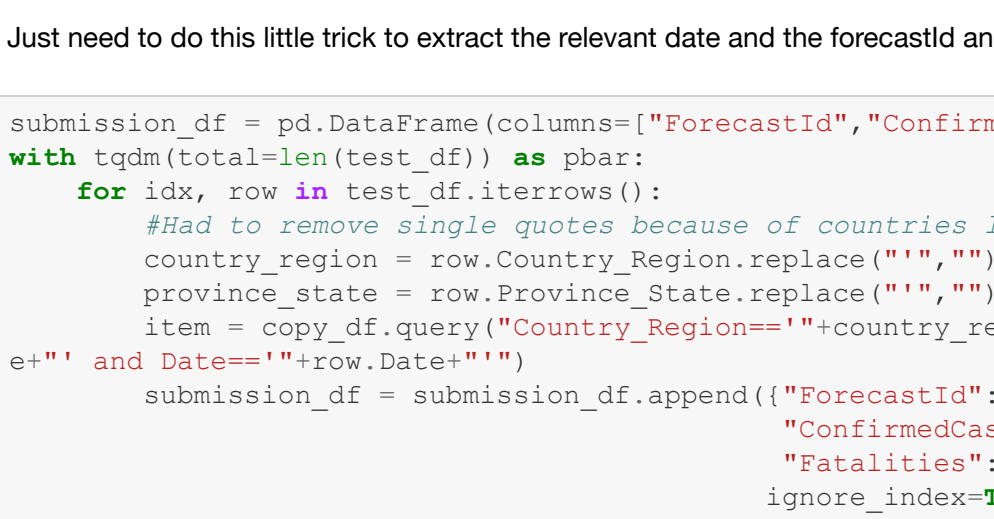


```
In [42]: display_long_term_prediction("Iran", groundtruth_df)
```

Comparison between the actual data and our predictions for the number of cases



Comparison between the actual data and our predictions for the number of fatalities



## 9. Generating the submission file

```
In [43]: test_df = gpd.read_file("../kaggle/input/covid19-global-forecasting-week-4/test.csv")
#The country_region got modifying in the enriched dataset by @Optimo,
# so we have to apply the same change to the test DataFrame.
test_df["Country_Region"] = [ row.Country_Region if row.Province_State=="else row.Country_Region" "+
row.Province_State for idx,row in test_df.iterrows() ]
test_df.head()
```

```
Out[43]:
```

	ForecastId	Province_State	Country_Region	Date	geomet
0	1		Afghanistan	2020-04-02	None
1	2		Afghanistan	2020-04-03	None
2	3		Afghanistan	2020-04-04	None
3	4		Afghanistan	2020-04-05	None
4	5		Afghanistan	2020-04-06	None

Just need to do this little trick to extract the relevant date and the forecastId and add that to the submission file.

```
In [44]: submission_df = pd.DataFrame(columns=["ForecastId","ConfirmedCases","Fatalities"])
for idx, row in test_df.iterrows():
    #to remove single quotes because of countries like Cote D'Ivoire for example
    country_region = row.Country_Region.replace("'", "").strip(" ")
    province_state = row.Province_State.replace("'", "").strip(" ")
    item = copy_df.query("Country_Region=="+"country_region" and Province_State==" "+
    province_state and Date==" "+row.Date+"")
    submission_df = submission_df.append({"ForecastId":row.ForecastId,
    "ConfirmedCases":int(item.ConfirmedCases.values[0]),
    "Fatalities":int(item.Fatalities.values[0])},
    ignore_index=True)

    pbar.update(1)

100%|██████████| 13459/13459 [03:34<00:00, 62.84it/s]
```

```
In [45]: submission_df.sample(20)
```

```
Out[45]:
```

	ForecastId	ConfirmedCases	Fatalities
10528	10529	954	44
9529	9530	0	0
5626	5627	981	0
10196	10197	987	47
8367	8368	1938	329
4995	4996	567	13
3881	3882	978	27
3323	3324	861	10
8092	8093	5	1
6418	6419	930	74
6159	6160	914	36
13274	13275	866	0
10595	10596	6741	135
1223	1224	933	41
2790	2791	943	35
11621	11622	75833	1550
7652	7653	949	37
10362	10363	960	42
1867	1868	970	1
8928	8929	983	0

```
In [46]: submission_df.to_csv("submission.csv",index=False)
```

## 10. Conclusion

As there is still "little" data available for such kind of model, it was clearly an ambitious approach but I was curious to see how far we could go with such solution. It appears that the model was reasonably good at predicting the earlier stages of the spread when the numbers are still in the few hundreds, but then it seems to plateau far too early. It also appears to be performing awfully in countries when the outbreaks are already at an advanced stage. When looking back at the version notes, it seems that adding information about the regions improved the model's performance in a minor with significant manner. Information about the restrictions on populations over time showed more impact than demographic information such as the number of hospital beds or density of population. While many other models have demonstrated far better results, I think this "naïve" model could be interesting on the long run as more data become available, mostly due to how easy it is to integrate more new inputs to the model.

If you found this notebook helpful, please give it an upvote. It will be greatly appreciated!