

COVID Global Forecast: SIR model + ML regressions

In the context of the global COVID-19 pandemic, Kaggle has launched several challenges in order to provide useful insights that may answer some of the open scientific questions about the virus. This is the case of the [COVID19 Global Forecasting](#), in which participants are encouraged to fit worldwide data in order to predict the pandemic evolution, hopefully helping to determine which factors impact the transmission behavior of COVID-19.

TABLE OF CONTENTS

1. [Exploratory data analysis \(EDA\)](#)
 - 1.1. [COVID-19 global tendency excluding China](#)
 - 1.2. [COVID-19 tendency in China](#)
 - 1.3. [Italy, Spain, UK and Singapore](#)
 2. [SIR model](#)
 - 2.1. [Implementing the SIR model](#)
 - 2.2. [Fit SIR parameters to real data](#)
 3. [Data enrichment](#)
 - 3.1. [Join data, filter dates and clean missing](#)
 - 3.2. [Compute lags and trends](#)
 - 3.3. [Add country details](#)
 4. [Predictions for the early stages of the transmission](#)
 - 4.1. [Linear Regression for one country](#)
 - 4.2. [Linear Regression for all countries \(method 1\)](#)
 - 4.3. [Linear Regression for all countries \(method 2\)](#)
 - 4.4. [Linear regression with lags](#)
 5. [Predictions for the late stages of the transmission](#)
 - 5.1. [Logistic curve fit](#)
 - 5.2. [Logistic curve fit for all countries](#)
 - 5.3. [APIMs](#)
 6. [Statement of the author](#)

Disclaimer 1: this notebook is being updated frequently with the objective of improving predictions by using new models.

Disclaimer 2: the training dataset is also updated on a daily basis in order to include the most recent cases. In order to be up to date and prevent data leaking and other potential problems, daily updates on "filtered dates" will be applied.

Disclaimer 3: the COVID Global Forecasting competition is updated week by week (with a new competition). I'll move the notebook from previous weeks to the new one, so that it only appears in the most recent competition.

```
In [1]: import numpy as np
import pandas as pd
train = pd.read_csv("../input/covid19-global-forecasting-week-4/train.csv")
test = pd.read_csv("../input/covid19-global-forecasting-week-4/test.csv")

import seaborn as sns
from sklearn import preprocessing
from datetime import datetime
from scipy import integrate, optimize
import warnings
warnings.filterwarnings('ignore')

# ML libraries
import lightgbm as lgb
import xgboost as xgb
from xgboost import plot_importance, plot_tree
from sklearn.model_selection import RandomizedSearchCV, GridSearchCV
from sklearn import linear_model
from sklearn.metrics import mean_squared_error
```

1. Exploratory data analysis (EDA)

First of all, let's take a look on the data structure:

```
In [2]: submission_example = pd.read_csv("../input/covid19-global-forecasting-week-4/submission.csv")
train = pd.read_csv("../input/covid19-global-forecasting-week-4/train.csv")
train.Province_State.fillna("None", inplace=True)
display(train.describe())
print("Number of Country_Region: ", train['Country_Region'].nunique())
print("Dates go from day", max(train['Date']), "to day", min(train['Date']), ", a total of", train['Date'].nunique())
print("Countries with Province/State informed: ", train.loc[train['Province_State']!=None]['Country_Region'].nunique())
```

Id	Province_State	Country_Region	Date	ConfirmedCases	Fatalities
0	None	Afghanistan	2020-01-23	0.0	0.0
1	None	Afghanistan	2020-01-24	0.0	0.0
2	None	Afghanistan	2020-01-25	0.0	0.0
3	None	Afghanistan	2020-01-26	0.0	0.0
4	None	Afghanistan	2020-01-26	0.0	0.0

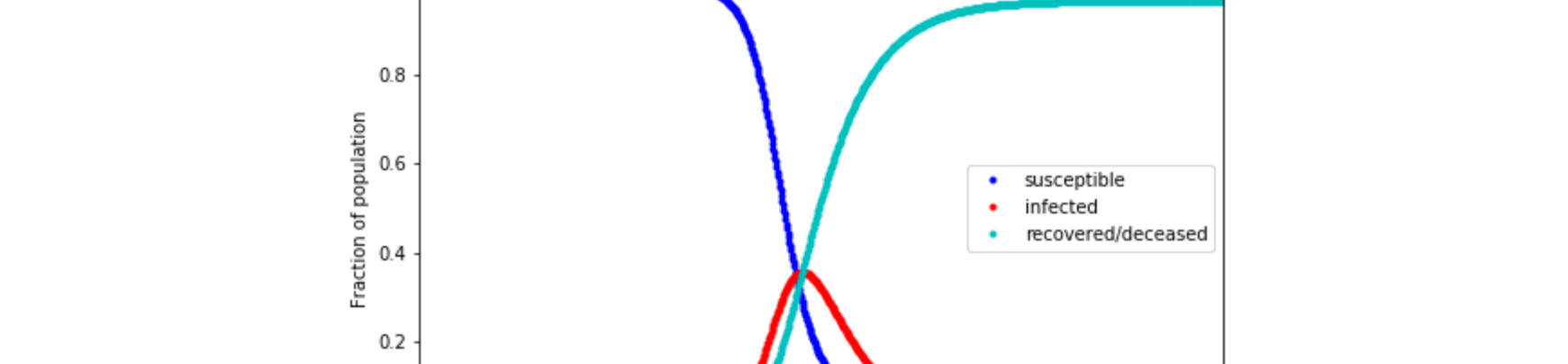
	Id	ConfirmedCases	Fatalities
count	26292.000000	26292.000000	26292.000000
mean	17826.500000	1186.183896	62.265594
std	10300.678012	8549.128727	605.049077
min	1.000000	0.000000	0.000000
50%	8913.750000	0.000000	0.000000
95%	17826.500000	1.000000	0.000000
75%	26739.250000	117.000000	1.000000
max	35552.000000	20320.000000	21067.000000

Number of Country_Region: 184
Dates go from day 2020-01-23 to day 2020-01-22, a total of 84 days
Countries with Province/State informed: 'Australia' 'Canada' 'China' 'Denmark' 'France' 'Netherlands' 'United Kingdom'

The dataset covers 163 countries and almost 2 full months from 2020, which is enough data to get some clues about the pandemic. Let's see a few plots of the worldwide tendency to see if we can extract some insights:

```
In [3]: confirmed_country = train.groupby(['Country_Region', 'Province_State']).agg({'ConfirmedCases': ['sum']})
#Fatalities_country_China = train[train['Country_Region']=='China'].groupby(['Country_Region', 'Province_State']).agg({'Fatalities': ['sum']})
confirmed_total_date = train.groupby(['Date']).agg({'ConfirmedCases': ['sum']})
confirmed_total_date_noChina = train[train['Country_Region']!='China'].groupby(['Date']).agg({'ConfirmedCases': ['sum']})
fatalities_total_date_China = train[train['Country_Region']=='China'].groupby(['Date']).agg({'Fatalities': ['sum']})
total_date_China = confirmed_total_date_China.join(fatalities_total_date_China)

fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(17,7))
total_date.plot(ax=ax1)
ax1.set_title("Global confirmed cases", size=13)
ax1.set_ylabel("Number of cases", size=13)
ax1.set_xlabel("Date", size=13)
fatalities_total_date.plot(ax=ax2, color='orange')
ax2.set_title("Global deceased cases", size=13)
ax2.set_ylabel("Number of cases", size=13)
ax2.set_xlabel("Date", size=13)
```



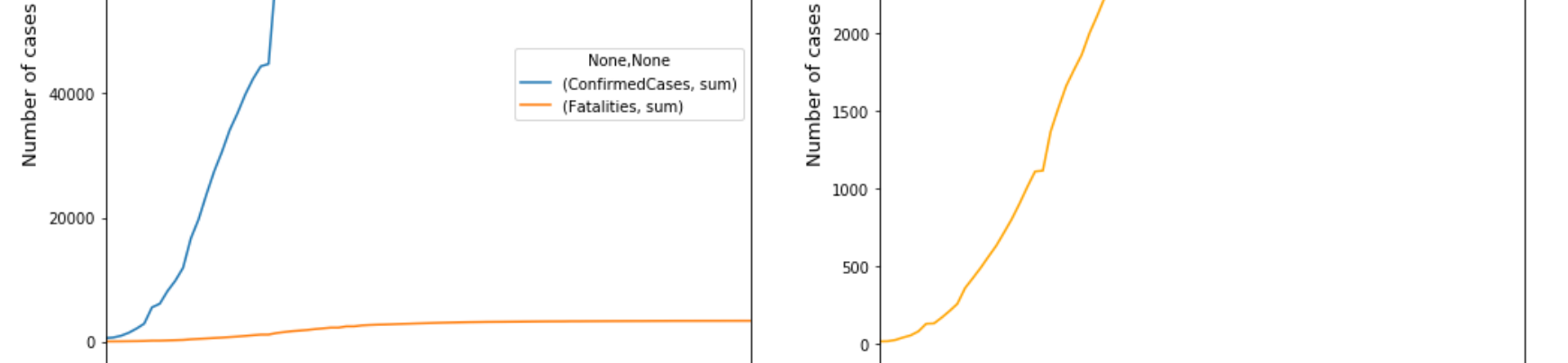
Observations: The global curve shows a rich fine structure, but these numbers are strongly affected by the vector zero country, China. Given that COVID-19 started there, during the initial expansion of the virus there was no reliable information about the real infected cases. In fact, the criteria to consider infection cases was modified around 2020-02-11, which strongly perturbed the curve as you can see from the figure.

1.1. COVID-19 global tendency excluding China

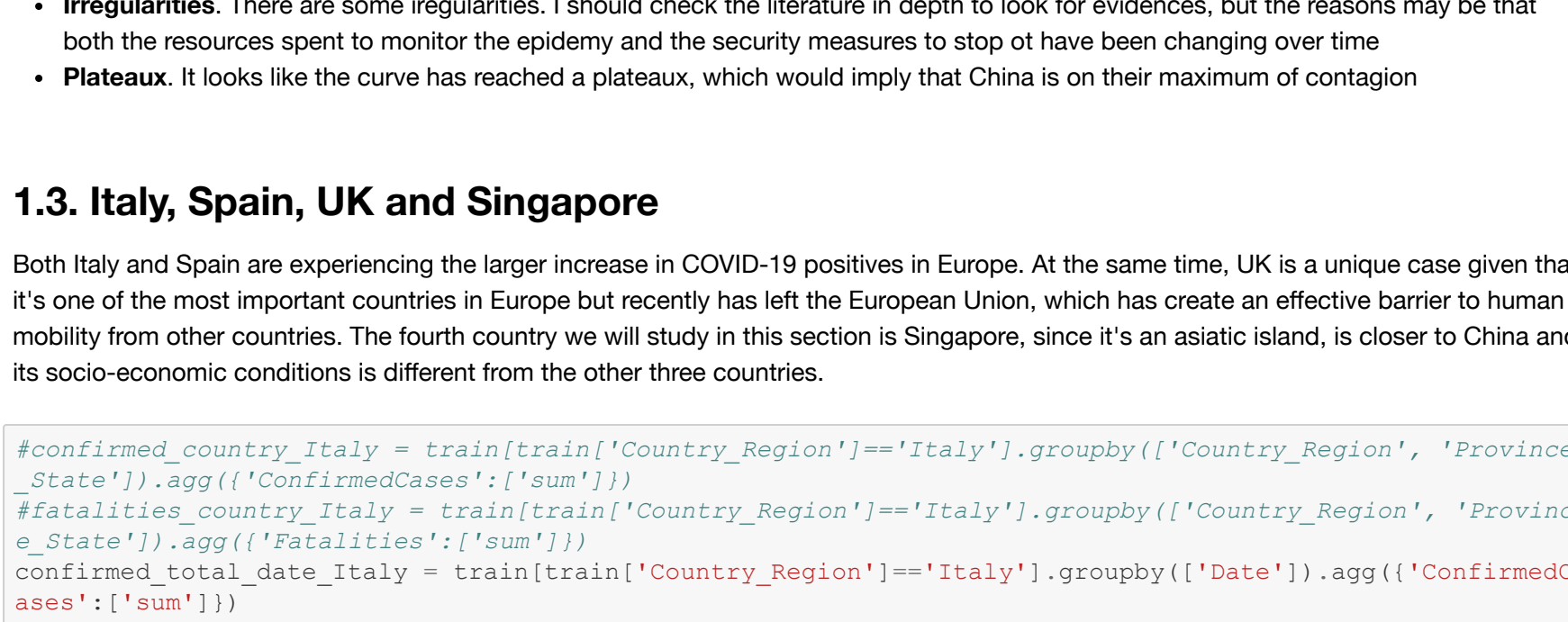
Since details of the initial breakthrough strongly interfere with the results, it's recommended to analyze China independently. Let's first see the results without China:

```
In [4]: confirmed_country_noChina = train[train['Country_Region']!='China'].groupby(['Country_Region', 'Province_State']).agg({'ConfirmedCases': ['sum']})
#Fatalities_country_noChina = train[train['Country_Region']!='China'].groupby(['Country_Region', 'Province_State']).agg({'Fatalities': ['sum']})
confirmed_total_date_noChina = train[train['Country_Region']!='China'].groupby(['Date']).agg({'ConfirmedCases': ['sum']})
fatalities_total_date_noChina = train[train['Country_Region']!='China'].groupby(['Date']).agg({'Fatalities': ['sum']})
total_date_noChina = confirmed_total_date_noChina.join(fatalities_total_date_noChina)

fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(17,7))
total_date_noChina.plot(ax=ax1)
ax1.set_title("Global confirmed cases excluding China", size=13)
ax1.set_ylabel("Number of cases", size=13)
ax1.set_xlabel("Date", size=13)
fatalities_total_date_noChina.plot(ax=ax2, color='orange')
ax2.set_title("Global deceased cases excluding China", size=13)
ax2.set_ylabel("Number of cases", size=13)
ax2.set_xlabel("Date", size=13)
```



Observations: In this case the general behavior looks cleaner, and in fact the curve resembles a typical epidemiology model like SIR. SIR present a large increasing number of infections that, once it reaches the maximum of the contagion, decreases with a lower slope. For comparison, a SIR simulation from section [2. SIR model](#):



1.2. COVID-19 tendency in China

Since China was the initial infected country, the COVID-19 behavior is different from the rest of the world. The medical system was not prepared for the pandemic and the number of infections that, once it reaches the maximum of the contagion, decreases with a lower slope. For comparison, a SIR simulation from section [2. SIR model](#):

```
In [5]: confirmed_country_China = train[train['Country_Region']=='China'].groupby(['Country_Region', 'Province_State']).agg({'ConfirmedCases': ['sum']})
#Fatalities_country_China = train[train['Country_Region']=='China'].groupby(['Country_Region', 'Province_State']).agg({'Fatalities': ['sum']})
confirmed_total_date_China = train[train['Country_Region']=='China'].groupby(['Date']).agg({'ConfirmedCases': ['sum']})
fatalities_total_date_China = train[train['Country_Region']=='China'].groupby(['Date']).agg({'Fatalities': ['sum']})
total_date_China = confirmed_total_date_China.join(fatalities_total_date_China)

fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(17,7))
total_date_China.plot(ax=ax1)
ax1.set_title("China confirmed cases", size=13)
ax1.set_ylabel("Number of cases", size=13)
ax1.set_xlabel("Date", size=13)
fatalities_total_date_China.plot(ax=ax2, color='orange')
ax2.set_title("China deceased cases", size=13)
ax2.set_ylabel("Number of cases", size=13)
ax2.set_xlabel("Date", size=13)
```



Observations:

- **Smoothness.** Both plots are less smooth than theoretical simulations or the curve from the rest of the world cumulative
- **Infected criteria.** The moment in which the criteria to consider an infected case was changed is directly spotted
- **Regularity.** There are some regularities I should check the literature in depth to look for evidences, but the reasons may be that
 - both the resources spent to monitor the epidemic and the security measures to stop it have been changing over time
 - Plateaux. It looks like the curve has reached a plateau, which would imply that China is on their maximum of contagion

1.3. Italy, Spain, UK and Singapore

Both the Italy and Spain are experiencing the larger increase in COVID-19 positives in Europe. At the same time, UK is a unique case given that it's one of the most important countries in Europe but recently has left the European Union, which has created an effective barrier to human mobility from other countries. The fourth country we will study in this section is Singapore, since it's an asiatic island, is closer to China and its socio-economic conditions are different from the other three countries.

```
In [6]: confirmed_country_Italy = train[train['Country_Region']=='Italy'].groupby(['Country_Region', 'Province_State']).agg({'ConfirmedCases': ['sum']})
#Fatalities_country_Italy = train[train['Country_Region']=='Italy'].groupby(['Country_Region', 'Province_State']).agg({'Fatalities': ['sum']})
confirmed_total_date_Italy = train[train['Country_Region']=='Italy'].groupby(['Date']).agg({'ConfirmedCases': ['sum']})
fatalities_total_date_Italy = train[train['Country_Region']=='Italy'].groupby(['Date']).agg({'Fatalities': ['sum']})
total_date_Italy = confirmed_total_date_Italy.join(fatalities_total_date_Italy)

confirmed_country_Spain = train[train['Country_Region']=='Spain'].groupby(['Country_Region', 'Province_State']).agg({'ConfirmedCases': ['sum']})
#Fatalities_country_Spain = train[train['Country_Region']=='Spain'].groupby(['Country_Region', 'Province_State']).agg({'Fatalities': ['sum']})
confirmed_total_date_Spain = train[train['Country_Region']=='Spain'].groupby(['Date']).agg({'ConfirmedCases': ['sum']})
fatalities_total_date_Spain = train[train['Country_Region']=='Spain'].groupby(['Date']).agg({'Fatalities': ['sum']})
total_date_Spain = confirmed_total_date_Spain.join(fatalities_total_date_Spain)

confirmed_country_UK = train[train['Country_Region']=='United Kingdom'].groupby(['Country_Region', 'Province_State']).agg({'ConfirmedCases': ['sum']})
#Fatalities_country_UK = train[train['Country_Region']=='United Kingdom'].groupby(['Country_Region', 'Province_State']).agg({'Fatalities': ['sum']})
confirmed_total_date_UK = train[train['Country_Region']=='United Kingdom'].groupby(['Date']).agg({'ConfirmedCases': ['sum']})
fatalities_total_date_UK = train[train['Country_Region']=='United Kingdom'].groupby(['Date']).agg({'Fatalities': ['sum']})
total_date_UK = confirmed_total_date_UK.join(fatalities_total_date_UK)

confirmed_country_Australia = train[train['Country_Region']=='Australia'].groupby(['Country_Region', 'Province_State']).agg({'ConfirmedCases': ['sum']})
#Fatalities_country_Australia = train[train['Country_Region']=='Australia'].groupby(['Country_Region', 'Province_State']).agg({'Fatalities': ['sum']})
confirmed_total_date_Australia = train[train['Country_Region']=='Australia'].groupby(['Date']).agg({'ConfirmedCases': ['sum']})
fatalities_total_date_Australia = train[train['Country_Region']=='Australia'].groupby(['Date']).agg({'Fatalities': ['sum']})
total_date_Australia = confirmed_total_date_Australia.join(fatalities_total_date_Australia)

confirmed_country_Singapore = train[train['Country_Region']=='Singapore'].groupby(['Country_Region', 'Province_State']).agg({'ConfirmedCases': ['sum']})
#Fatalities_country_Singapore = train[train['Country_Region']=='Singapore'].groupby(['Country_Region', 'Province_State']).agg({'Fatalities': ['sum']})
confirmed_total_date_Singapore = train[train['Country_Region']=='Singapore'].groupby(['Date']).agg({'ConfirmedCases': ['sum']})
fatalities_total_date_Singapore = train[train['Country_Region']=='Singapore'].groupby(['Date']).agg({'Fatalities': ['sum']})
total_date_Singapore = confirmed_total_date_Singapore.join(fatalities_total_date_Singapore)
```



As a fraction of the total population of each country:

```
In [7]: pop_italy = 60486683.
pop_spain = 46749596.
pop_UK = 67764927.
pop_singapore = 5837230.

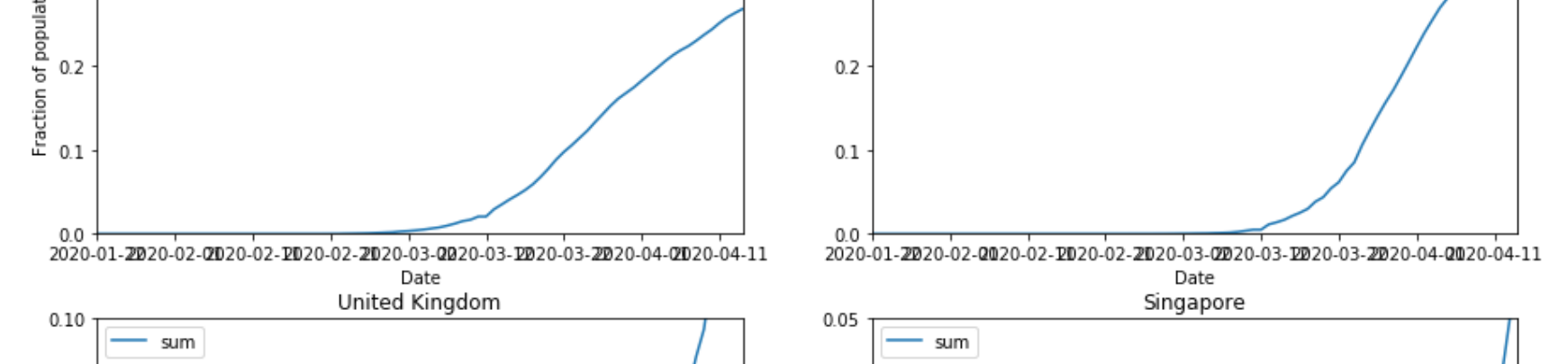
total_date_Italy.ConfirmedCases = total_date_Italy.ConfirmedCases/pop_italy*100.
total_date_Spain.ConfirmedCases = total_date_Spain.ConfirmedCases/pop_spain*100.
total_date_UK.ConfirmedCases = total_date_UK.ConfirmedCases/pop_UK*100.
total_date_Singapore.ConfirmedCases = total_date_Singapore.ConfirmedCases/pop_singapore*100.

plt.figure(figsize=(15,10))
plt.subplot(1, 2, 1)
total_date_Italy.ConfirmedCases.plot(ax=plt.gca(), title='Italy')
plt.ylabel('Fraction of population infected')
plt.ylim(0, 0.5)

plt.subplot(2, 2, 2)
total_date_Spain.ConfirmedCases.plot(ax=plt.gca(), title='Spain')
plt.ylim(0, 0.5)

plt.subplot(2, 2, 3)
total_date_UK.ConfirmedCases.plot(ax=plt.gca(), title='United Kingdom')
plt.ylabel('Fraction of population infected')
plt.ylim(0, 0.1)

plt.subplot(2, 2, 4)
total_date_Singapore.ConfirmedCases.plot(ax=plt.gca(), title='Singapore')
plt.ylim(0, 0.05)
```



In order to compare the 4 countries, it's also interesting to see the evolution of the infections from the first confirmed case:

```
In [8]: confirmed_country_Italy = train[train['Country_Region']=='Italy'] & train['ConfirmedCases']!=0
confirmed_country_Spain = train[train['Country_Region']=='Spain'] & train['ConfirmedCases']!=0
confirmed_country_UK = train[train['Country_Region']=='United Kingdom'] & train['ConfirmedCases']!=0
confirmed_country_Australia = train[train['Country_Region']=='Australia'] & train['ConfirmedCases']!=0
confirmed_country_Singapore = train[train['Country_Region']=='Singapore'] & train['ConfirmedCases']!=0

total_date_Italy = confirmed_total_date_Italy.join(fatalities_total_date_Italy)
total_date_Spain = confirmed_total_date_Spain.join(fatalities_total_date_Spain)
total_date_UK = confirmed_total_date_UK.join(fatalities_total_date_UK)
total_date_Australia = confirmed_total_date_Australia.join(fatalities_total_date_Australia)
total_date_Singapore = confirmed_total_date_Singapore.join(fatalities_total_date_Singapore)

# Plots
plt.figure(figsize=(12,6))
plt.plot(italy_30, label='Italy')
plt.plot(spain_30, label='Spain')
plt.plot(uk_30, label='UK')
plt.plot(singapore_30, label='Singapore')
plt.legend(['Italy', 'Spain', 'UK', 'Singapore'], loc='upper left')
plt.title('COVID-19 infections from the first confirmed case', size=15)
plt.xlabel('Days', size=13)
plt.ylabel('Infected cases', size=13)
plt.ylim(0, 130000)
plt.savefig('SIR_example.png')
plt.show()
```



Observations:

- **Italy.** With almost 120,000 confirmed cases, Italy shows one of the most alarming scenarios of COVID-19. The infections curve is very steep, and more than 2% of population has been infected
- **Spain.** Spain has the same number of cumulative infected cases than Italy, near 120,000. However, Spain's total population is lower (around 42 millions) and hence the percentage of population that has been infected rises up to 3%.
- **United Kingdom.** Despite not being very far from them, the UK shows less cases. This may be due to the number of tests performed, but it's soon to know for sure. The number of cases is around 40,000, this is, a 0.6 of the total population.
- **Singapore.** Singapore is relatively isolated given that it's an island, and the number of international travels is lower than for the other 3 countries. The number of cases is still very low (<1000), despite the general tendency is to increase. However, the infections started faster in the beginning, but the slope of the infections curve hasn't increased very much in the past weeks. A 0.2% of the population was infected

2. SIR model

We have seen some general behavior of the virus in aggregated data, for the country where the coronavirus was originated and for four other interesting countries. There's a lot of information to be extracted from this data: for example, we haven't analyzed the effects of longlat of countries. However, since our main purpose is to develop a predictive model in order to understand the key factors that impact the COVID-19 transmission, I'll move on to one of the most famous epidemiological models: SIR.

SIR is a simple model that considers a population that belongs to one of the following states:

1. **Susceptible (S).** The individual hasn't contracted the disease, but she can be infected due to transmission from infected people
2. **Infected (I).** The person has contracted the disease
3. **Recovered/Deceased (R).** The disease may lead to one of two destinies: either the person survives, hence developing immunity to the disease, or the person is deceased

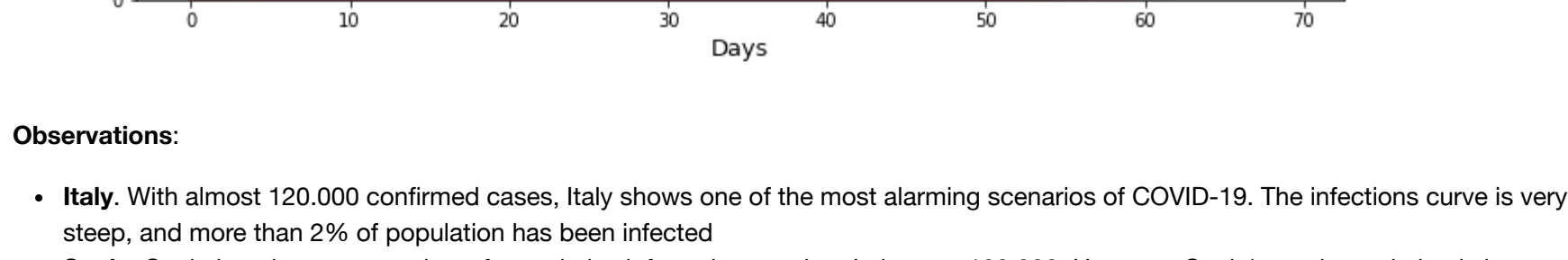


Image by Kai Sasaki from [jupyter.com](#)

There are many versions of this model, considering birth and death (SIRD with demography), with intermediate states, etc. However, since we are in the early stages of the COVID-19 expansion and our interest is focused in the short term, we will consider that people develop immunity (in the long term, immunity may be lost) and the COVID-19 may come back within a certain seasonality (like the common flu) and there is no transition from recovered to the remaining two states. With this, the differential equations that govern the system are:

$$\begin{aligned} \frac{dS}{dt} &= -\beta SI \\ \frac{dI}{dt} &= \beta SI - \gamma I \\ \frac{dR}{dt} &= \gamma I \end{aligned}$$

Where β is the contagion rate of the pathogen and γ is the recovery rate.

2.1. Implementing the SIR model

SIR model can be implemented in many ways: from the differential equations governing the system, within a mean field approximation or running the dynamics in a social network (graph). For the sake of simplicity, I've chosen the first option, and we will simply run a numerical method (Runge-Kutta) to solve the differential equations system.

The functions governing the diff. eqs. are:

```
In [9]: # Susceptible equation
def fS(N, a, b, c, fa, fb, fc, beta, gamma, hs):
    fa = -beta*a*b
    return fa

# Infected equation
def fI(N, a, b, beta, gamma):
    fb = beta*a*b - gamma*b
    return fb

# Recovered/deceased equation
def fR(N, b, gamma):
    fc = gamma*b
    return fc

In order to solve the differential equations system, we develop a 4th order Runge-Kutta method:
```

```
In [10]: # Runge-Kutta method of 4th order for 3 dimensions (susceptible a, infected b and recovered r)
def rk4(N, a, b, c, fa, fb, fc, beta, gamma, hs):
    a1 = fa(N, a, b, beta)
    b1 = fb(N, a, b, beta, gamma)
    c1 = fc(N, b, gamma)
    ak = a + a1*0.5
    bk = b + b1*0.5
    ck = c + c1*0.5
    a2 = fa(N, ak, bk, beta)
    b2 = fb(N, ak, bk, beta, gamma)
    c2 = fc(N, bk, gamma)
    ak = a + a2*0.5
    bk = b + b2*0.5
    ck = c + c2*0.5
    a3 = fa(N, ak, bk, beta)
    b3 = fb(N, ak, bk, beta, gamma)
    c3 = fc(N, bk, gamma)
    ak = a + a3
    bk = b + b3
    ck = c + c3
    a4 = fa(N, ak, bk, beta)
    b4 = fb(N, ak, bk, beta, gamma)
    c4 = fc(N, bk, gamma)
    b = b + (b1 + 2*(b2 + b3) + b4)/6
    c = c + (c1 + 2*(c2 + c3) + c4)/6
    return a, b, c
```

And finally, to obtain the evolution of the disease we simply define the initial conditions and call the rk4 method:

```
In [11]: def SIR(N, b0, beta, gamma, hs):
    """
    N = total number of population
    beta = transition rate S->I
    gamma = transition rate I->R
    k = denotes the constant degree distribution of the network (average value for networks in which the probability of finding a node with a different connectivity decays exponentially fast)
    hs = jump step of the numerical integration

    # Initial condition
    a = float(1)/N - b0
    b = float(1)/N * b0
    c = 0

    sus, inf, rec = [0,0,0]
    for i in range(1, N): # Run for a certain number of time-steps
        sus.append(a)
        inf.append(b)
        rec.append(c)
        a, b, c = rk4(N, a, b, c, fa, fb, fc, beta, gamma, hs)
    return sus, inf, rec
```

Results obtained for N-world population, only one initial infected case, $\beta = 0.3$, $\gamma = 0.5$ and a leap pass $h_t = 0.1$ are shown below:

```
In [12]: # Parameters of the model
N = 78000*(10**6)
b0 = 0.7
beta = 0.2
gamma = 0.1
hs = 0.1

sus, inf, rec = SIR(N, b0, beta, gamma, hs)

# Plots
plt.figure(figsize=(10,5))
plt.plot(sus, label='susceptible')
plt.plot(inf, label='infected')
plt.plot(rec, label='recovered/deceased')
plt.title('SIR model')
plt.xlabel('time', fontsize=10)
plt.ylabel('Fraction of population', fontsize=10)
plt.legend(['susceptible', 'infected', 'recovered/deceased'])
plt.ylim(0, 1.0)
plt.xlim(0, 1000)
plt.savefig('SIR_example.png')
plt.show()
```



Observations:

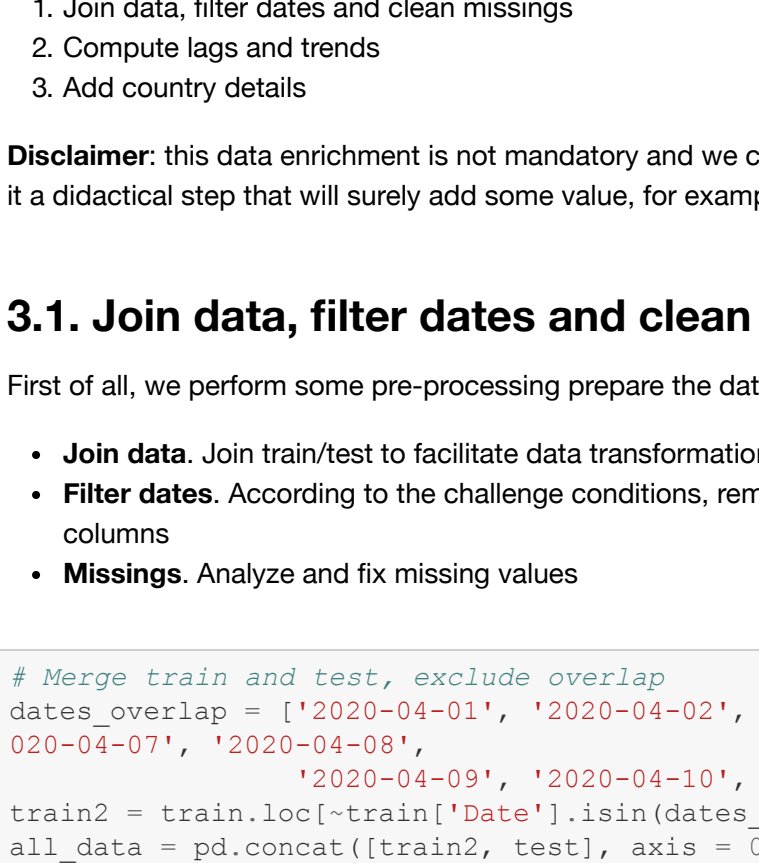
- The number of infected cases increases for a certain time period, and then eventually decreases given that individuals recover/decease from the disease
- The susceptible fraction of population decreases as the virus is transmitted, to eventually drop to the absorbent state 0
- The opposite happens for the recovered/deceased case

Notice that different initial conditions and parameter values will lead to other scenarios, feel free to play with these numbers to study the system.

2.2. Fit SIR parameters to real data

The SIR model is purely theoretical, and we are interested in to a real approximation of the COVID-19 expansion in order to extract insights and understand the transmission of the virus. Hence, we need to extract the β and γ parameters for each case if we hope to be able to predict the evolution of the system.

[13]: population = float(46750238)
country_df = pd.DataFrame(
country_df['ConfirmedCases'] = train.loc[train['Country_Region']=='Spain'].ConfirmedCases.diff().fillna(
country_df['day_count'] = list(range(1,len(country_df)+1))
country_df = country_df.iloc[1:]
country_df['day_count'] = list(range(1,len(country_df)+1))
ydata = [i for i in country_df.ConfirmedCases]
info = ydata[0:100]
xdata = country_df.day_count
ydata = np.array(ydata, dtype=float)
xdata = np.array(xdata, dtype=float)
N = population
info = ydata[0:100]
sus0 = N - info
rec0 = 0.0
def sir_model(y, x, beta, gamma):
sus = -beta * y[0] * y[1] / N
rec = gamma * y[1]
inf = -sus + rec
return sus, inf, rec
def fit_odeint(y, x, beta, gamma):
return integrate.odeint(sir_model, (sus0, info, rec0), x, args=(beta, gamma))[1:]
popt, pcov = optimize.curve_fit(fit_odeint, xdata, ydata)
fitted = fit_odeint(xdata, "popt")
plt.plot(xdata, ydata, 'o')
plt.plot(xdata, fitted)
plt.title("Fit of SIR model for Spain infected cases")
plt.ylabel("Population infected")
plt.xlabel("Days")
plt.show()
print("Optimal parameters: beta =", popt[0], " and gamma =", popt[1])



Optimal parameters: beta = 9.4492339303625 and gamma = 9.27220462762098

I'm not happy with the fit of parameters and I want to work more on this, since I'm not properly reproducing the curves. I'll keep working on this for curiosity, but in the meanwhile I'll develop a data-centric approach to the prediction.

3. Data enrichment

Analyzing SIR simulations was meant to understand a model that approximately resembles the transmission mechanism of many virus, including the COVID-19. However, there are alternative methods that may prove being equally useful both to predict and to understand the pandemic evolution. Many of these methods rely on having rich data to extract conclusions and allow algorithms to extrapolate patterns in data, and that is exactly what we are going to do.

Main workflow of this section:

1. Join data, filter dates and clean missings
2. Compute lags and trends
3. Add country details

Disclaimer: the data enrichment is not mandatory and we could end up not using all of the new features in our models. However I consider it a didactical step that will surely add some value, for example in an in-depth exploratory analysis.

3.1. Join data, filter dates and clean missings

First of all, we perform some pre-processing prepare the dataset, consisting on:

- **Join data.** Join train/test to facilitate data transformations
- **Filter dates.** According to the challenge conditions, remove ConfirmedCases and Fatalities post 2020-03-12. Create additional date columns
- **Missings.** Analyze and fix missing values

```
In [14]: # Merge train and test, exclude overlap  
dates_overlap = ['2020-04-01', '2020-04-02', '2020-04-03', '2020-04-04', '2020-04-05', '2020-04-06', '2020-04-07', '2020-04-08', '2020-04-09', '2020-04-10', '2020-04-11', '2020-04-12', '2020-04-13', '2020-04-14']  
train2 = train.loc[~train['Date'].isin(dates_overlap)]  
test_data = pd.concat([train2, test], axis = 0, sort=False)  
  
# Double check that there are no informed ConfirmedCases and Fatalities after 2020-03-11  
all_data.loc[all_data['Date'] >= '2020-04-01', 'ConfirmedCases'] = 0  
all_data.loc[all_data['Date'] >= '2020-04-01', 'Fatalities'] = 0  
all_data['Date'] = pd.to_datetime(all_data['Date'])  
  
# Create date column  
le = preprocessing.LabelEncoder()  
all_data['Day_num'] = le.fit_transform(all_data.Date)  
all_data['Month'] = le.fit_transform(all_data.Month)  
all_data['Year'] = all_data['Date'].dt.month  
all_data['Year'] = all_data['Date'].dt.year  
  
# Fill null values given that we merged train-test datasets  
all_data['Province_State'].fillna('None', inplace=True)  
all_data['ConfirmedCases'].fillna(0, inplace=True)  
all_data['Fatalities'].fillna(0, inplace=True)  
all_data['ForecastId'].fillna(-1, inplace=True)  
all_data['Day_num'].fillna(-1, inplace=True)  
all_data['Month'].fillna(-1, inplace=True)  
all_data['Year'].fillna(-1, inplace=True)  
display(all_data)  
display(all_data.loc[all_data['Date'] == '2020-04-01'])
```

	Id	Province_State	Country_Region	Date	ConfirmedCases	Fatalities	ForecastId	Day_num	Day	Month	Year
0	1.0	None	Afghanistan	2020-01-23	0.0	0.0	-1.0	0	22	1	2020
1	2.0	None	Afghanistan	2020-01-22	0.0	0.0	-1.0	1	23	1	2020
2	3.0	None	Afghanistan	2020-01-24	0.0	0.0	-1.0	2	24	1	2020
3	4.0	None	Afghanistan	2020-01-25	0.0	0.0	-1.0	3	25	1	2020
4	5.0	None	Afghanistan	2020-01-26	0.0	0.0	-1.0	4	26	1	2020
...
13454	-1.0	None	Zimbabwe	2020-05-10	0.0	0.0	13455.0	108	10	5	2020
13455	-1.0	None	Zimbabwe	2020-05-11	0.0	0.0	13456.0	109	11	5	2020
13456	-1.0	None	Zimbabwe	2020-05-12	0.0	0.0	13457.0	110	12	5	2020
13457	-1.0	None	Zimbabwe	2020-05-13	0.0	0.0	13458.0	111	13	5	2020
13458	-1.0	None	Zimbabwe	2020-05-14	0.0	0.0	13459.0	112	14	5	2020

35369 rows x 11 columns

Id	Province_State	Country_Region	Date	ConfirmedCases	Fatalities	ForecastId	Day_num	Day	Month	Year	
0	1.0	None	Afghanistan	2020-01-23	0.0	0.0	-1.0	0	22	1	2020
1	2.0	None	Afghanistan	2020-01-22	0.0	0.0	-1.0	1	23	1	2020
2	3.0	None	Afghanistan	2020-01-24	0.0	0.0	-1.0	2	24	1	2020
3	4.0	None	Afghanistan	2020-01-25	0.0	0.0	-1.0	3	25	1	2020
4	5.0	None	Afghanistan	2020-01-26	0.0	0.0	-1.0	4	26	1	2020
...	
13454	-1.0	None	Zimbabwe	2020-05-10	0.0	0.0	13455.0	108	10	5	2020
13455	-1.0	None	Zimbabwe	2020-05-11	0.0	0.0	13456.0	109	11	5	2020
13456	-1.0	None	Zimbabwe	2020-05-12	0.0	0.0	13457.0	110	12	5	2020
13457	-1.0	None	Zimbabwe	2020-05-13	0.0	0.0	13458.0	111	13	5	2020
13458	-1.0	None	Zimbabwe	2020-05-14	0.0	0.0	13459.0	112	14	5	2020

Observations:

- 'ConfirmedCases' and 'Fatalities' are now only informed for dates previous to 2020-03-12
- The dataset includes all countries and dates, which is required for the lag/trend step
- Missing values for 'ConfirmedCases' and 'Fatalities' have been replaced by 0, which may be dangerous if we do not remember it at the end of the process. However, since we will train only on dates previous to 2020-03-12, this won't impact our prediction algorithm
- A new column 'Day' has been created, as a day counter starting from the first date

Double-check that there are no remaining missing values:

```
In [15]: missings_count = (col:all_data[col].isnull().sum() for col in all_data.columns)  
missings = pd.Series.from_dict(missings_count, orient='index')  
print(missings.nlargest(30, 0))
```

	Id	Province_State	Country_Region	Date	ConfirmedCases	Fatalities	ForecastId	Day_num	Day	Month	Year
0	1.0	None	Afghanistan	2020-01-23	0.0	0.0	-1.0	0	22	1	2020
1	2.0	None	Afghanistan	2020-01-22	0.0	0.0	-1.0	1	23	1	2020
2	3.0	None	Afghanistan	2020-01-24	0.0	0.0	-1.0	2	24	1	2020
3	4.0	None	Afghanistan	2020-01-25	0.0	0.0	-1.0	3	25	1	2020
4	5.0	None	Afghanistan	2020-01-26	0.0	0.0	-1.0	4	26	1	2020
...
13454	-1.0	None	Zimbabwe	2020-05-10	0.0	0.0	13455.0	108	10	5	2020
13455	-1.0	None	Zimbabwe	2020-05-11	0.0	0.0	13456.0	109	11	5	2020
13456	-1.0	None	Zimbabwe	2020-05-12	0.0	0.0	13457.0	110	12	5	2020
13457	-1.0	None	Zimbabwe	2020-05-13	0.0	0.0	13458.0	111	13	5	2020
13458	-1.0	None	Zimbabwe	2020-05-14	0.0	0.0	13459.0	112	14	5	2020

Time spent: 6.311013221740723

3.2. Compute lags and trends

Enriching a dataset is key to obtain good results. In this case we will apply 2 different transformations:

Lag. Lags are a way to compute the previous value of a column, so that the lag 1 for ConfirmedCases would inform the this column from the previous day. The lag 3 of a feature X is simply:

$$X_{lag}(t) = X(t-3)$$

Trend. Transforming a column into its trend gives the natural tendency of this column, which is different from the raw value. The definition of trend I will apply is:

$$Trend_x = \frac{X(t) - X(t-1)}{X(t-1)}$$

The backlog of lags I'll apply is 14 days, while for trends is 7 days. For ConfirmedCases and Fatalities:

```
In [16]: def calculate_lag(df, lag_list, column):  
    for lag in lag_list:  
        column_lag = column + "_" + str(lag)  
        df[column_lag] = df.groupby(['Country_Region', 'Province_State'])[column].shift(lag, fill_value=0)  
    return df  
  
def calculate_trend(df, lag_list, column):  
    for lag in lag_list:  
        trend_column_lag = "trend_" + column + "_" + str(lag)  
        df[trend_column_lag] = (df[column + "_" + str(lag)] - df[column]).shift(lag, fill_value=0)  
    return df  
  
ts = time.time()  
all_data = calculate_lag(all_data.reset_index(), range(1,7), 'ConfirmedCases')  
all_data = calculate_trend(all_data.reset_index(), range(1,7), 'ConfirmedCases')  
all_data = calculate_lag(all_data.reset_index(), range(1,7), 'Fatalities')  
all_data = calculate_trend(all_data.reset_index(), range(1,7), 'Fatalities')  
all_data.fillna(0, inplace=True)  
all_data.fillna(0, inplace=True)  
print("Time spent: ", time.time()-ts)
```

Time spent: 6.311013221740723

As you see, the process is really fast. An example of some of the lag/trend columns for Spain:

```
In [17]: all_data[all_data['Country_Region']=='Spain'].iloc[40:50][['Id', 'Province_State', 'Country_Region', 'Date', 'ConfirmedCases', 'Fatalities', 'ForecastId', 'Day_num', 'ConfirmedCases_1', 'ConfirmedCases_2', 'ConfirmedCases_3', 'Fatalities_1', 'Fatalities_2', 'Fatalities_3']]
```

	Id	Province_State	Country_Region	Date	ConfirmedCases	Fatalities	ForecastId	Day_num	Day	Month	Year
15580	25348.0	None	Spain	2020-03-02	120.0	0.0	-1.0	40	40	3	2020
15581	25350.0	None	Spain	2020-03-03	165.0	1.0	-1.0	41	41	3	2020
15582	25351.0	None	Spain	2020-03-04	222.0	2.0	-1.0	42	42	3	2020
15583	25352.0	None	Spain	2020-03-05	259.0	3.0	-1.0	43	43	3	2020
15584	25353.0	None	Spain	2020-03-06	400.0	5.0	-1.0	44	44	3	2020
15585	25354.0	None	Spain	2020-03-07	500.0	10.0	-1.0	45	45	3	2020
15586	25355.0	None	Spain	2020-03-08	673.0	17.0	-1.0	46	46	3	2020
15587	25356.0	None	Spain	2020-03-09	1073.0	28.0	-1.0	47	47	3	2020
15588	25357.0	None	Spain	2020-03-10	1695.0	35.0	-1.0	48	48	3	2020
15589	25358.0	None	Spain	2020-03-11	2277.0	54.0	-1.0	49	49	3	2020

3.3. Add country details

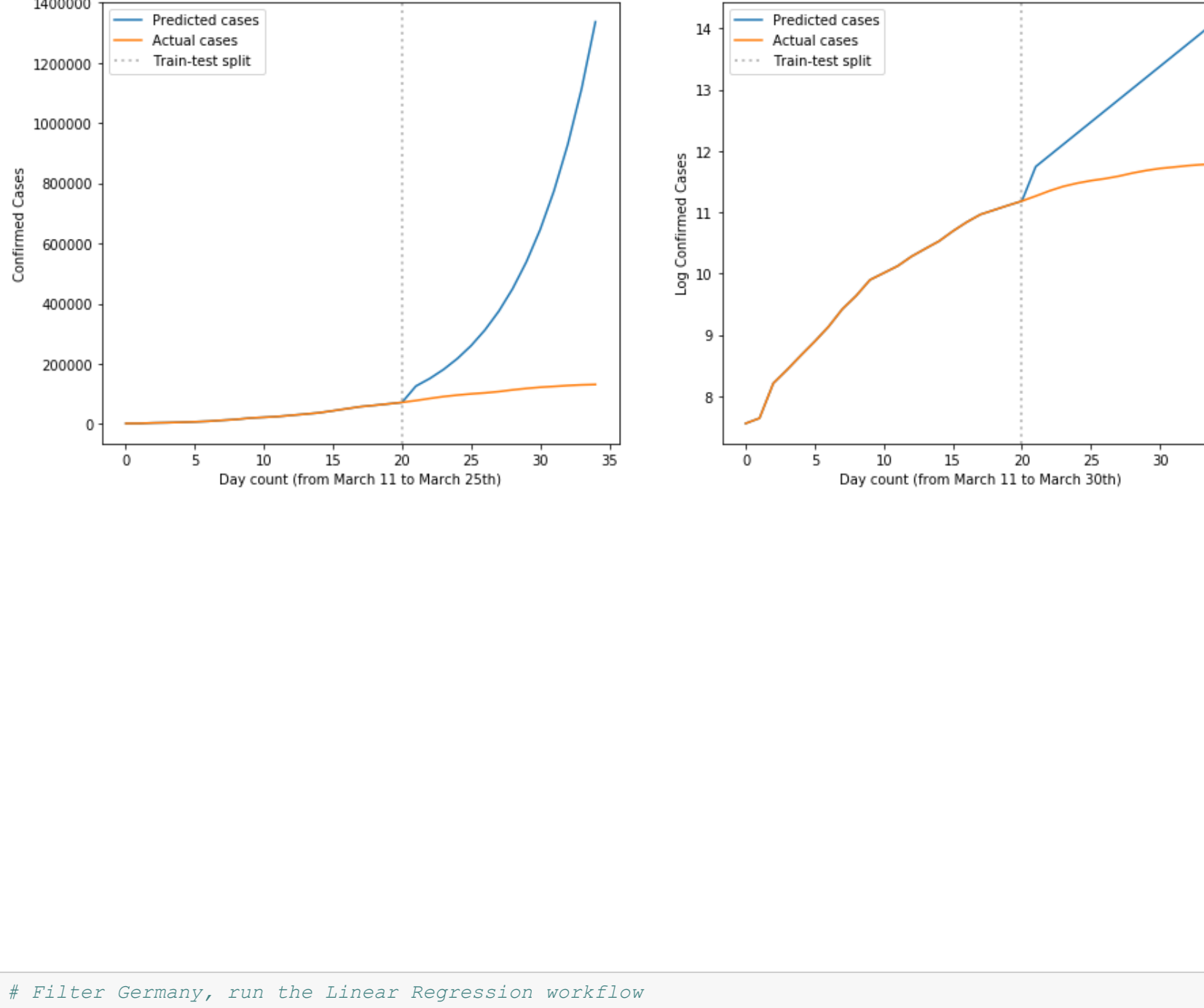
Variables like the total population of a country, the average age of citizens or the fraction of people living in cities may strongly impact on the COVID-19 transmission behavior. Hence, it's important to consider these factors. I'm using [Kaggle's dataset](#) based on Web Scrapping for this purpose.

```
In [18]: # Load countries data file  
world_population = pd.read_csv("../kaggle/input/population-by-country-2020/population_by_country_2020.csv")  
  
# Select desired columns and rename some of them  
world_population = world_population[['Country (or dependency)', 'Population (2020)', 'Density (P/Km²)', 'Land Area (Km²)', 'Med. Age', 'Urban Pop %']]  
world_population.columns = ['Country (or dependency)', 'Population (2020)', 'Density', 'Land Area', 'Med Age', 'Urban Pop']  
  
# Replace United States by US  
world_population.loc[world_population['Country (or dependency)']=='United States', 'Country (or dependency)'] = "US"  
  
# Remove the # character from Urban Pop values  
world_population['Urban Pop'] = world_population['Urban Pop'].str.rstrip("#")  
  
# Replace Urban Pop and Med Age "N.A." by their respective modes, then transform to int  
world_population.loc[world_population['Urban Pop']=='N.A.', 'Urban Pop'] = int(world_population.loc[world_population['Urban Pop']=='N.A.'].Urban Pop.mode()[0])  
world_population.loc[world_population['Med Age']=='N.A.', 'Med Age'] = int(world_population.loc[world_population['Med Age']=='N.A.'].Med Age.mode()[0])  
world_population['Med Age'] = world_population['Med Age'].astype('int16')  
  
print("Cleaned country details dataset")  
display(world_population)  
  
# Now join the dataset to our previous DataFrame and clean missings (not match in left join) - label encode cities  
all_data = all_data.merge(world_population, left_on='Country_Region', right_on='Country (or dependency)', how='left')  
all_data['Population (2020)'] = all_data['Population (2020)']  
all_data['Density'] = all_data['Density']  
all_data['Land Area'] = all_data['Land Area']  
all_data['Med Age'] = all_data['Med Age']  
all_data['Urban Pop'] = all_data['Urban Pop']  
all_data['ConfirmedCases'] = all_data['ConfirmedCases']  
all_data['Fatalities'] = all_data['Fatalities']  
all_data['ForecastId'] = all_data['ForecastId']  
all_data['Day_num'] = all_data['Day_num']  
all_data['Day'] = all_data['Day']  
all_data['Month'] = all_data['Month']  
all_data['Year'] = all_data['Year']  
all_data['ConfirmedCases_1'] = all_data['ConfirmedCases_1']  
all_data['ConfirmedCases_2'] = all_data['ConfirmedCases_2']  
all_data['ConfirmedCases_3'] = all_data['ConfirmedCases_3']  
all_data['Fatalities_1'] = all_data['Fatalities_1']  
all_data['Fatalities_2'] = all_data['Fatalities_2']  
all_data['Fatalities_3'] = all_data['Fatalities_3']  
all_data['ForecastId_1'] = all_data['ForecastId_1']  
all_data['ForecastId_2'] = all_data['ForecastId_2']  
all_data['ForecastId_3'] = all_data['ForecastId_3']  
all_data['Day_num_1'] = all_data['Day_num_1']  
all_data['Day_num_2'] = all_data['Day_num_2']  
all_data['Day_num_3'] = all_data['Day_num_3']  
all_data['Day_1'] = all_data['Day_1']  
all_data['Day_2'] = all_data['Day_2']  
all_data['Day_3'] = all_data['Day_3']  
all_data['Month_1'] = all_data['Month_1']  
all_data['Month_2'] = all_data['Month_2']  
all_data['Month_3'] = all_data['Month_3']  
all_data['Year_1'] = all_data['Year_1']  
all_data['Year_2'] = all_data['Year_2']  
all_data['Year_3'] = all_data['Year_3']  
all_data['ConfirmedCases_1'] = all_data['ConfirmedCases_1']  
all_data['ConfirmedCases_2'] = all_data['ConfirmedCases_2']  
all_data['ConfirmedCases_3'] = all_data['ConfirmedCases_3']  
all_data['Fatalities_1'] = all_data['Fatalities_1']  
all_data['Fatalities_2'] = all_data['Fatalities_2']  
all_data['Fatalities_3'] = all_data['Fatalities_3']  
all_data['ForecastId_1'] = all_data['ForecastId_1']  
all_data['ForecastId_2'] = all_data['ForecastId_2']  
all_data['ForecastId_3'] = all_data['ForecastId_3']  
all_data['Day_num_1'] = all_data['Day_num_1']  
all_data['Day_num_2'] = all_data['Day_num_2']  
all_data['Day_num_3'] = all_data['Day_num_3']  
all_data['Day_1'] = all_data['Day_1']  
all_data['Day_2'] = all_data['Day_2']  
all_data['Day_3'] = all_data['Day_3']  
all_data['Month_1'] = all_data['Month_1']  
all_data['Month_2'] = all_data['Month_2']  
all_data['Month_3'] = all_data['Month_3']  
all_data['Year_1'] = all_data['Year_1']  
all_data['Year_2'] = all_data['Year_2']  
all_data['Year_3'] = all_data['Year_3']  
all_data['ConfirmedCases_1'] = all_data['ConfirmedCases_1']  
all_data['ConfirmedCases_2'] = all_data['ConfirmedCases_2']  
all_data['ConfirmedCases_3'] = all_data['ConfirmedCases_3']  
all_data['Fatalities_1'] = all_data['Fatalities_1']  
all_data['Fatalities_2'] = all_data['Fatalities_2']  
all_data['Fatalities_3'] = all_data['Fatalities_3']  
all_data['ForecastId_1'] = all_data['ForecastId_1']  
all_data['ForecastId_2'] = all_data['ForecastId_2']  
all_data['ForecastId_3'] = all_data['ForecastId_3']  
all_data['Day_num_1'] = all_data['Day_num_1']  
all_data['Day_num_2'] = all_data['Day_num_2']  
all_data['Day_num_3'] = all_data['Day_num_3']  
all_data['Day_1'] = all_data['Day_1']  
all_data['Day_2'] = all_data['Day_2']  
all_data['Day_3'] = all_data['Day_3']  
all_data['Month_1'] = all_data['Month_1']  
all_data['Month_2'] = all_data['Month_2']  
all_data['Month_3'] = all_data['Month_3']  
all_data['Year_1'] = all_data['Year_1']  
all_data['Year_2'] = all_data['Year_2']  
all_data['Year_3'] = all_data['Year_3']  
all_data['ConfirmedCases_1'] = all_data['ConfirmedCases_1']  
all_data['ConfirmedCases_2'] = all_data['ConfirmedCases_2']  
all_data['ConfirmedCases_3'] = all_data['ConfirmedCases_3']  
all_data['Fatalities_1'] = all_data['Fatalities_1']  
all_data['Fatalities_2'] = all_data['Fatalities_2']  
all_data['Fatalities_3'] = all_data['Fatalities_3']  
all_data['ForecastId_1'] = all_data['ForecastId_1']  
all_data['ForecastId_2'] = all_data['ForecastId_2']  
all_data['ForecastId_3'] = all_data['ForecastId_3']  
all_data['Day_num_1'] = all_data['Day_num_1']  
all_data['Day_num_2'] = all_data['Day_num_2']  
all_data['Day_num_3'] = all_data['Day_num_3']  
all_data['Day_1'] = all_data['Day_1']  
all_data['Day_2'] = all_data['Day_2']  
all_data['Day_3'] = all_data['Day_3']  
all_data['Month_1'] = all_data['Month_1']  
all_data['Month_2'] = all_data['Month_2']  
all_data['Month_3'] = all_data['Month_3']  
all_data['Year_1'] = all_data['Year_1']  
all_data['Year_2'] = all_data['Year_2']  
all_data['Year_3'] = all_data['Year_3']  
all_data['ConfirmedCases_1'] = all_data['ConfirmedCases_1']  
all_data['ConfirmedCases_2'] = all_data['ConfirmedCases_2']  
all_data['ConfirmedCases_3'] = all_data['ConfirmedCases_3']  
all_data['Fatalities_1'] = all_data['Fatalities_1']  
all_data['Fatalities_2'] = all_data['Fatalities_2']  
all_data['Fatalities_3'] = all_data['Fatalities_3']  
all_data['ForecastId_1'] = all_data['ForecastId_1']  
all_data['ForecastId_2'] = all_data['ForecastId_2']  
all_data['ForecastId_3'] = all_data['ForecastId_3']  
all_data['Day_num_1'] = all_data['Day_num_1']  
all_data['Day_num_2'] = all_data['Day_num_2']  
all_data['Day_num_3'] = all_data['Day_num_3']  
all_data['Day_1'] = all_data['Day_1']  
all_data['Day_2'] = all_data['Day_2']  
all_data['Day_3'] = all_data['Day_3']  
all_data['Month_1'] = all_data['Month_1']  
all_data['Month_2'] = all_data['Month_2']  
all_data['Month_3'] = all_data['Month_3']  
all_data['Year_1'] = all_data['Year_1']  
all_data['Year_2'] = all_data['Year_2']  
all_data['Year_3'] = all_data['Year_3']  
all_data['ConfirmedCases_1'] = all_data['ConfirmedCases_1']  
all_data['ConfirmedCases_2'] = all_data['ConfirmedCases_2']  
all_data['ConfirmedCases_3'] = all_data['ConfirmedCases_3']  
all_data['Fatalities_1'] = all_data['Fatalities_1']  
all_data['Fatalities_2'] = all_data['Fatalities_2']  
all_data['Fatalities_3'] = all_data['Fatalities_3']  
all_data['ForecastId_1'] = all_data['ForecastId_1']  
all_data['ForecastId_2'] = all_data['ForecastId_2']  
all_data['ForecastId_3'] = all_data['ForecastId_3']  
all_data['Day_num_1'] = all_data['Day_num_1']  
all_data['Day_num_2'] = all_data['Day_num_2']  
all_data['Day_num_3'] = all_data['Day_num_3']  
all_data['Day_1'] = all_data['Day_1']  
all_data['Day_2'] = all_data['Day_2']  
all_data['Day_3'] = all_data['Day_3']  
all_data['Month_1'] = all_data['Month_1']  
all_data['Month_2'] = all_data['Month_2']  
all_data['Month_3'] = all_data['Month_3']  
all_data['Year_1'] = all_data['Year_1']  
all_data['Year_2'] = all_data['Year_2']  
all_data['Year_3'] = all_data['Year_3']  
all_data['ConfirmedCases_1'] = all_data['ConfirmedCases_1']  
all_data['ConfirmedCases_2'] = all_data['ConfirmedCases_2']  
all_data['ConfirmedCases_3'] = all_data['ConfirmedCases_3']  
all_data['Fatalities_1'] = all_data['Fatalities_1']  
all_data['Fatalities_2'] = all_data['Fatalities_2']  
all_data['Fatalities_3'] = all_data['Fatalities_3']  
all_data['ForecastId_1'] = all_data['ForecastId_1']  
all_data['ForecastId_2'] = all_data['ForecastId_2']  
all_data['ForecastId_3'] = all_data['ForecastId_3']  
all_data['Day_num_1'] = all_data['Day_num_1']  
all_data['Day_num_2'] = all_data['Day_num_2']  
all_data['Day_num_3'] = all_data['Day_num_3']  
all_data['Day_1'] = all_data['Day_1']  
all_data['Day_2'] = all_data['Day_2']  
all_data['Day_3'] = all_data['Day_3']  
all_data['Month_1'] = all_data['Month_1']  
all_data['Month_2'] = all_data['Month_2']  
all_data['Month_3'] = all_data['Month_3']  
all_data['Year_1'] = all_data['Year_1']  
all_data['Year_2'] = all_data['Year_2']  
all_data['Year_3'] = all_data['Year_3']  
all_data['ConfirmedCases_1'] = all_data['ConfirmedCases_1']  
all_data['ConfirmedCases_2'] = all_data['ConfirmedCases_2']  
all_data['ConfirmedCases_3'] = all_data['ConfirmedCases_3']  
all_data['Fatalities_1'] = all_data['Fatalities_1']  
all_data['Fatalities_2'] = all_data['Fatalities_2']  
all_data['Fatalities_3'] = all_data['Fatalities_3']  
all_data['ForecastId_1'] = all_data['ForecastId_1']  
all_data['ForecastId_2'] = all_data['ForecastId_2']  
all_data['ForecastId_3'] = all_data['ForecastId_3']  
all_data['Day_num_1'] = all_data['Day_num_1']  
all_data['Day_num_2'] = all_data['Day_num_2']  
all_data['Day_num_3'] = all_data['Day_num_3']  
all_data['Day_1'] = all_data['Day_1']  
all_data['Day_2'] = all_data['Day_2']  
all_data['Day_3'] = all_data['Day_3']  
all_data['Month_1'] = all_data['Month_1']  
all_data['Month_2'] = all_data['Month_2']  
all_data['Month_3'] = all_data['Month_3']  
all_data['Year_1'] = all_data['Year_1']  
all_data['Year_2'] = all_data['Year_2']  
all_data['Year_3'] = all_data['Year_3']  
all_data['ConfirmedCases_1'] = all_data['ConfirmedCases_1']  
all_data['ConfirmedCases_2'] = all_data['ConfirmedCases_2']  
all_data['ConfirmedCases_3'] = all_data['ConfirmedCases_3']  
all_data['Fatalities_1'] = all_data['Fatalities_1']  
all_data['Fatalities_2'] = all_data['Fatalities_2']  
all_data['Fatalities_3'] = all_data['Fatalities_3']  
all_data['ForecastId_1'] = all_data['ForecastId_1']  
all_data['ForecastId_2'] = all_data['ForecastId_2']  
all_data['ForecastId_3'] = all_data['ForecastId_3']  
all_data['Day_num_1'] = all_data['Day_num_1']  
all_data['Day_num_2'] = all_data['Day_num_2']  
all_data['Day_num_3'] = all_data['Day_num_3']  
all_data['Day_1'] = all_data['Day_1']  
all_data['Day_2'] = all_data['Day_2']  
all_data['Day_3'] = all_data['Day_3']  
all_data['Month_1'] = all_data['Month_1']  
all_data['Month_2'] = all_data['Month_2']  
all_data['Month_3'] = all_data['Month_3']  
all_data['Year_1'] = all_data['Year_1']  
all_data['Year_2'] = all_data['Year_2']  
all_data['Year_3'] = all_data['Year_3']  
all_data['ConfirmedCases_1'] = all_data['ConfirmedCases_1']  
all_data['ConfirmedCases_2'] = all_data['ConfirmedCases_2']  
all
```



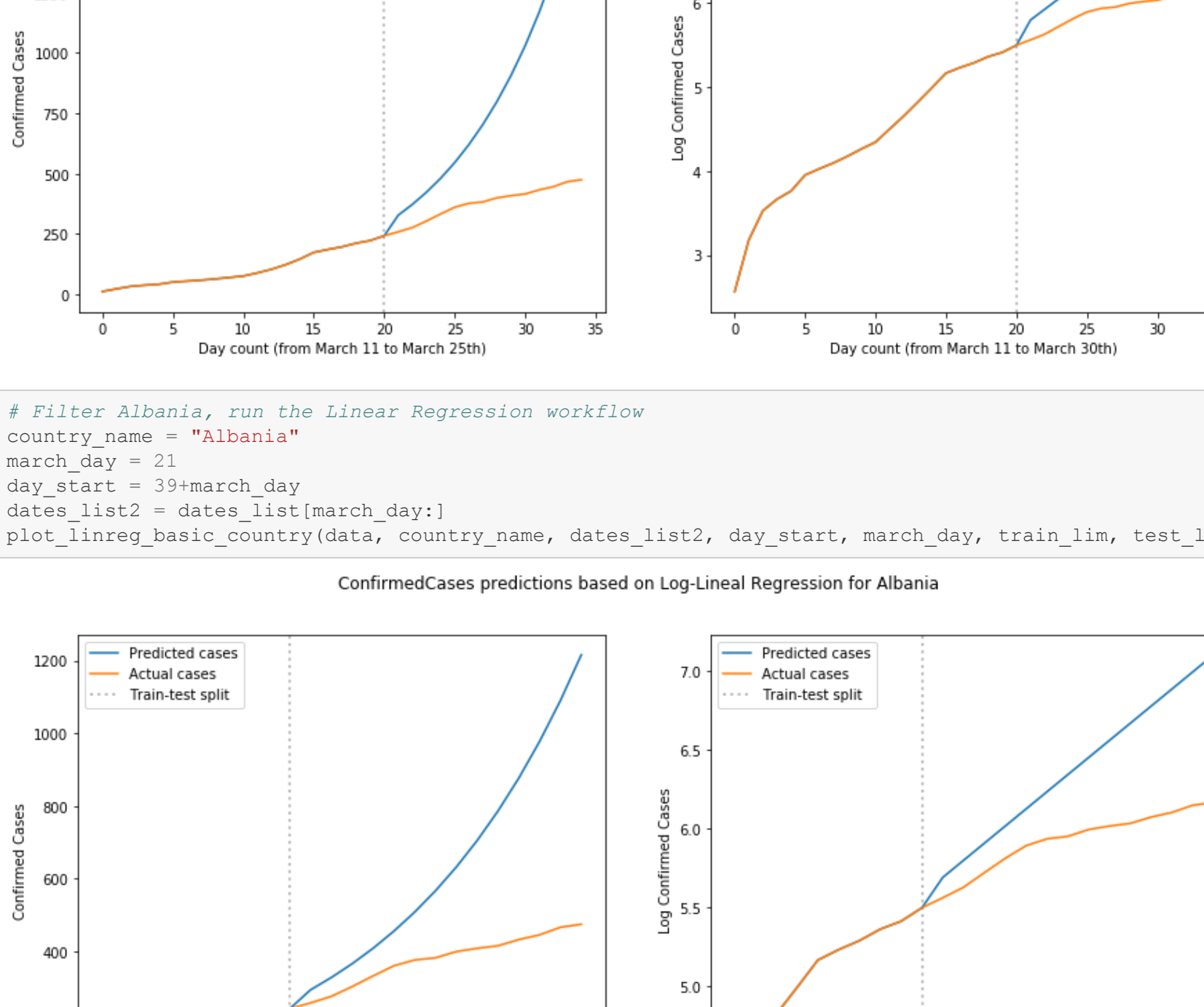
```
In [27]: # Filter Germany, run the Linear Regression workflow
country_name = "Germany"
march_day = 10
dates_list2 = dates_list[march_day:]
plot_linreg_basic_country(data, country_name, dates_list2, day_start, march_day, train_lim, test_lim)
```

ConfirmedCases predictions based on Log-Linear Regression for Germany



```
In [28]: # Filter Germany, run the Linear Regression workflow
country_name = "Germany"
march_day = 10
day_start = 39=march_day
dates_list2 = dates_list[march_day:]
plot_linreg_basic_country(data, country_name, dates_list2, day_start, march_day, train_lim, test_lim)
```

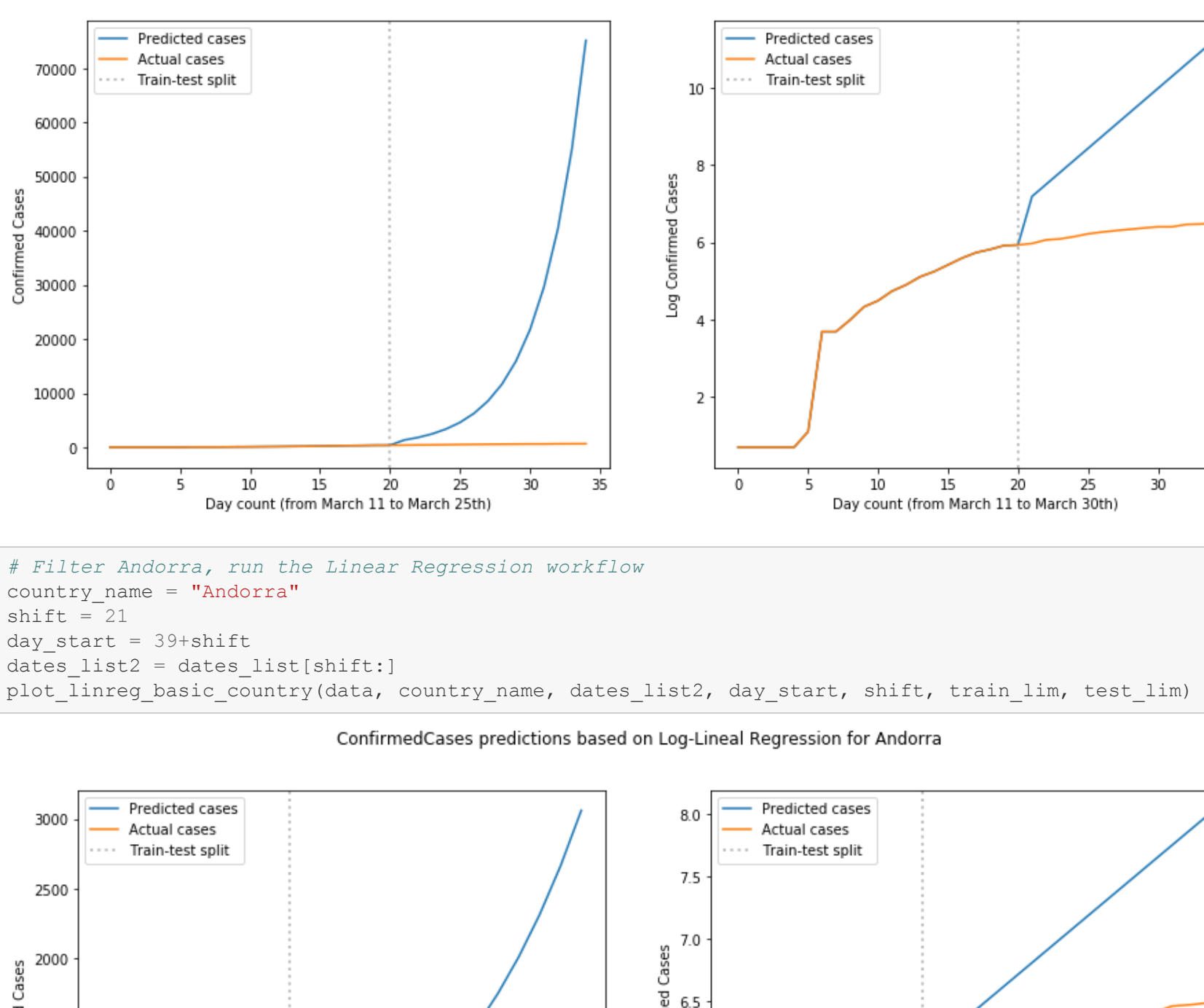
ConfirmedCases predictions based on Log-Linear Regression for Germany



• Albania

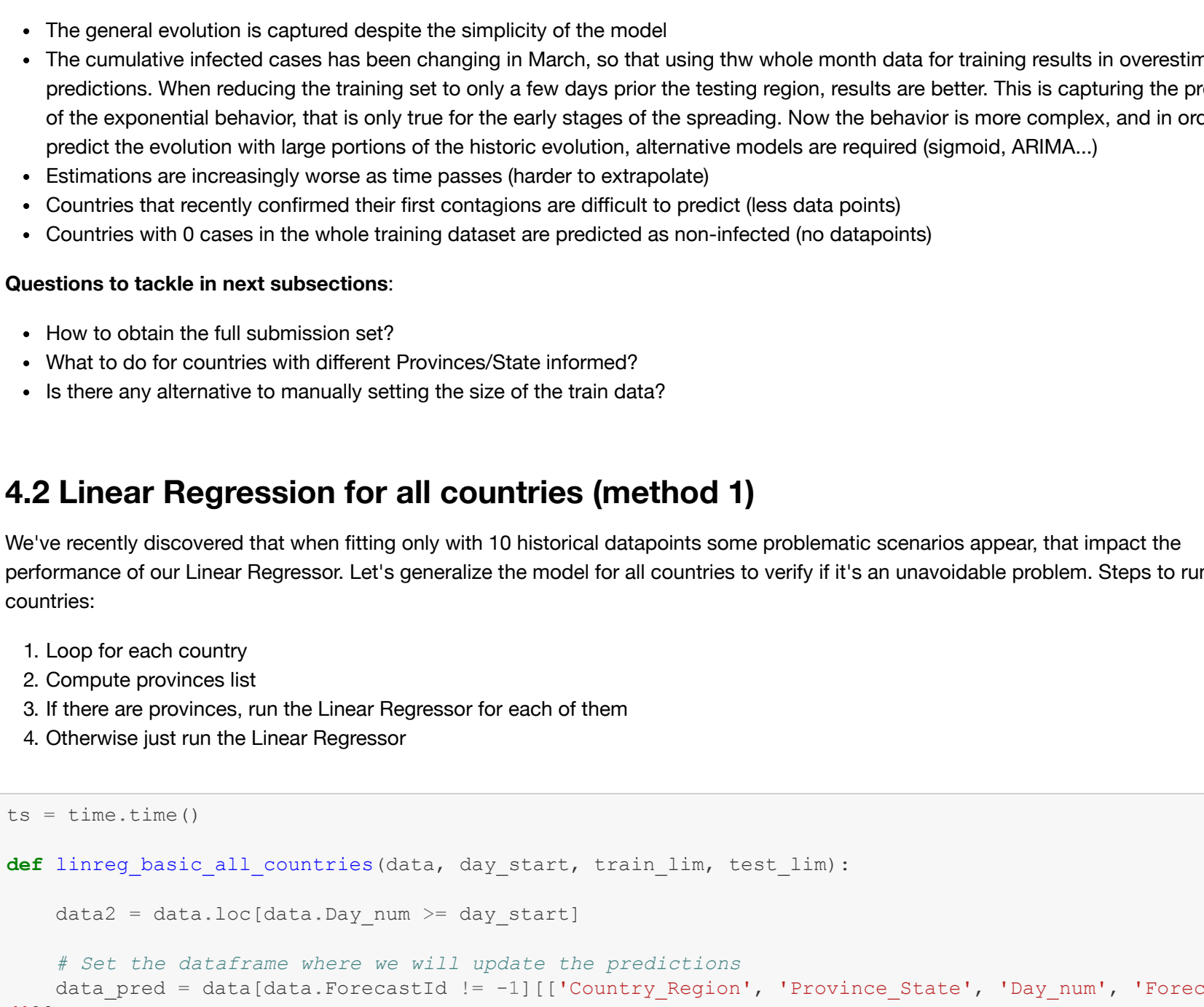
```
In [29]: # Filter Albania, run the Linear Regression workflow
country_name = "Albania"
march_day = 10
day_start = 39=march_day
dates_list2 = dates_list[march_day:]
plot_linreg_basic_country(data, country_name, dates_list2, day_start, march_day, train_lim, test_lim)
```

ConfirmedCases predictions based on Log-Linear Regression for Albania



```
In [30]: # Filter Albania, run the Linear Regression workflow
country_name = "Albania"
march_day = 10
day_start = 39=march_day
dates_list2 = dates_list[march_day:]
plot_linreg_basic_country(data, country_name, dates_list2, day_start, march_day, train_lim, test_lim)
```

ConfirmedCases predictions based on Log-Linear Regression for Albania



• Andorra

```
In [31]: # Filter Andorra, run the Linear Regression workflow
country_name = "Andorra"
shift = 10
day_start = 39=shift
dates_list2 = dates_list[shift:]
plot_linreg_basic_country(data, country_name, dates_list2, day_start, shift, train_lim, test_lim)
```

ConfirmedCases predictions based on Log-Linear Regression for Andorra



```
In [32]: # Filter Andorra, run the Linear Regression workflow
country_name = "Andorra"
shift = 10
day_start = 39=shift
dates_list2 = dates_list[shift:]
plot_linreg_basic_country(data, country_name, dates_list2, day_start, shift, train_lim, test_lim)
```

ConfirmedCases predictions based on Log-Linear Regression for Andorra



Observations:

- The general evolution is captured despite the simplicity of the model
- The cumulative infected cases has been changing in March, so that using thru whole month data for training results in overestimated predictions. When reducing the training set to only a few days prior the testing region, results are better. This is capturing the problem of the exponential behavior, that is only true for the early stages of the spreading. Now the behavior is more complex, and in order to predict the evolution with large portions of the historic evolution, alternative models are required (sigmoid, ARIMA,...)
- Estimations are increasingly worse as time passes (harder to extrapolate)
- Countries that recently confirmed their first contagions are difficult to predict (few data points)
- Countries with 0 cases in the whole training dataset are predicted as non-infected (no datapoints)

Questions to tackle in next subsections:

- How to obtain the full submission set?
- What to do for countries with different Provinces/State informed?
- Is there any alternative to manually setting the size of the train data?

4.2 Linear Regression for all countries (method 1)

We've recently discovered that when fitting only with 10 historical datapoints some problematic scenarios appear, that impact the performance of our Linear Regressor. Let's generalize the model for all countries to verify if it's an unavoidable problem. Steps to run for all countries:

1. Loop for each country
2. Compute provinces list
3. If there are provinces, run the Linear Regressor for each of them
4. Otherwise just run the Linear Regressor

```
In [33]: ts = time.time()

def linreg_basic_all_countries(data, day_start, train_lim, test_lim):
    data2 = data.loc[(data.Day_num >= day_start)]

    # Set the dataframe where we will update the predictions
    data_pred = data.loc[(data.ForecastId != -1)][['Country_Region', 'Province_State', 'Day_num', 'ForecastId']]
    data_pred['data_pred'] = data_pred['data_pred'].fillna(0)
    data_pred['Predicted_ConfirmedCases'] = [0]*len(data_pred)
    data_pred['Predicted_Fatalities'] = [0]*len(data_pred)
    print("Currently running Linear Regression for all countries")

    # Main loop for countries
    for c in data2['Country_Region'].unique():
        # List of provinces
        provinces_list = data2[data2['Country_Region']==c]['Province_State'].unique()

        # If the country has several Province/State informed
        if len(provinces_list)>1:
            for p in provinces_list:
                data_cp = data2[(data2['Country_Region']==c) & (data2['Province_State']==p)]
                X_train, Y_train_1, Y_train_2, X_test = split_data(data_cp, train_lim, test_lim)
                model_1, pred_1 = lin_reg(X_train, Y_train_1, X_test)
                model_2, pred_2 = lin_reg(X_train, Y_train_2, X_test)
                data_pred.loc[(data_pred['Country_Region']==c) & (data2['Province_State']==p)], ['Predicted_ConfirmedCases'] = pred_1
                data_pred.loc[(data_pred['Country_Region']==c) & (data2['Province_State']==p)], ['Predicted_Fatalities'] = pred_2
            else:
                # No Province/State informed
                data_c = data2[(data2['Country_Region']==c)]
                X_train, Y_train_1, Y_train_2, X_test = split_data(data_c, train_lim, test_lim)
                model_1, pred_1 = lin_reg(X_train, Y_train_1, X_test)
                model_2, pred_2 = lin_reg(X_train, Y_train_2, X_test)
                data_pred.loc[(data_pred['Country_Region']==c)], ['Predicted_ConfirmedCases'] = pred_1
                data_pred.loc[(data_pred['Country_Region']==c)], ['Predicted_Fatalities'] = pred_2

        # Apply exponential transf, and clean potential infinities due to final numerical precision
        data_pred[['Predicted_ConfirmedCases', 'Predicted_Fatalities']] = data_pred[['Predicted_ConfirmedCases', 'Predicted_Fatalities']].apply(lambda x: np.expmin(x))
        data_pred.replace([np.inf, -np.inf], 0, inplace=True)

    return data_pred

day_start = 65
#data_pred = linreg_basic_all_countries(data, day_start, train_lim, test_lim)
#get_submission(data_pred, 'ConfirmedCases', 'Predicted_ConfirmedCases', 'Predicted_Fatalities')
print("Process finished in ", round(time.time() - ts, 2), " seconds")

Process finished in 0.0 seconds

Final LMSE score for week 2, with training data prior to 2020-03-19 and measures on date 2020-04-01:

1.19651
```

4.3 Linear Regression for all countries (method 2)

An alternative method to setting the number of days for the training set is to simply keep all data for each country since the first case was confirmed. However, since there were certain countries where the initial outbreak was very smooth (i.e. in Spain there was only one confirmed case for 7 days in a row), predictions may be biased by these initial periods.

```
In [34]: ts = time.time()

# Set the dataframe where we will update the predictions
day_start = 65
data2 = data.loc[(data.Day_num >= day_start)]
data_pred3 = data2[(data2.ForecastId != -1)][['Country_Region', 'Province_State', 'Day_num', 'ForecastId']]
data_pred3['data_pred3'] = data_pred3['data_pred3'].fillna(0)
data_pred3['Predicted_ConfirmedCases'] = [0]*len(data_pred3)
data_pred3['Predicted_Fatalities'] = [0]*len(data_pred3)
how_many_days = test.Date.nunique()

print("Currently running Linear Regression for all countries")

# Main loop for countries
for c in data['Country_Region'].unique():
    # List of provinces
    provinces_list = data2[data2['Country_Region']==c]['Province_State'].unique()

    # If the country has several Province/State informed
    if len(provinces_list)>1:
        for p in provinces_list:
            # Only fit starting from the first confirmed case in the country
            train_countries_no0 = data.loc[(data['Country_Region']==c) & (data['Province_State']==p) & (data.ConfirmedCases!=0) & (data.ForecastId!=-1)]
            test_countries_no0 = data.loc[(data['Country_Region']==c) & (data['Province_State']==p) & (data.ForecastId!=-1)]
            data2 = pd.concat([train_countries_no0, test_countries_no0])

            # If there are no previous cases, predict 0
            if len(train_countries_no0) == 0:
                data_pred3.loc[(data_pred3['Country_Region']==c) & (data_pred3['Province_State']==p)], ['Predicted_ConfirmedCases'] = [0]*how_many_days
                data_pred3.loc[(data_pred3['Country_Region']==c) & (data_pred3['Province_State']==p)], ['Predicted_Fatalities'] = [0]*how_many_days
            else:
                # Else run linReg
                data_cp = data2[(data2['Country_Region']==c) & (data2['Province_State']==p)]
                X_train, Y_train_1, Y_train_2, X_test = split_data(data_cp, train_lim, test_lim)
                model_1, pred_1 = lin_reg(X_train, Y_train_1, X_test)
                model_2, pred_2 = lin_reg(X_train, Y_train_2, X_test)
                data_pred3.loc[(data_pred3['Country_Region']==c) & (data_pred3['Province_State']==p)], ['Predicted_ConfirmedCases'] = pred_1
                data_pred3.loc[(data_pred3['Country_Region']==c) & (data_pred3['Province_State']==p)], ['Predicted_Fatalities'] = pred_2

        # No Province/State informed
        else:
            # Only fit starting from the first confirmed case in the country
            train_countries_no0 = data.loc[(data['Country_Region']==c) & (data.ConfirmedCases!=0) & (data.ForecastId!=-1)]
            test_countries_no0 = data.loc[(data['Country_Region']==c) & (data.ForecastId!=-1)]
            data2 = pd.concat([train_countries_no0, test_countries_no0])

            # If there are no previous cases, predict 0
            if len(train_countries_no0) == 0:
                data_pred3.loc[(data_pred3['Country_Region']==c)], ['Predicted_ConfirmedCases'] = [0]*how_many_days
                data_pred3.loc[(data_pred3['Country_Region']==c)], ['Predicted_Fatalities'] = [0]*how_many_days
            else:
                # Else, run linReg
                data_c = data2[(data2['Country_Region']==c)]
                X_train, Y_train_1, Y_train_2, X_test = split_data(data_c, train_lim, test_lim)
                model_1, pred_1 = lin_reg(X_train, Y_train_1, X_test)
                model_2, pred_2 = lin_reg(X_train, Y_train_2, X_test)
                data_pred3.loc[(data_pred3['Country_Region']==c)], ['Predicted_ConfirmedCases'] = pred_1
                data_pred3.loc[(data_pred3['Country_Region']==c)], ['Predicted_Fatalities'] = pred_2

        # Apply exponential transf, and clean potential infinities due to final numerical precision
        data_pred3[['Predicted_ConfirmedCases', 'Predicted_Fatalities']] = data_pred3[['Predicted_ConfirmedCases', 'Predicted_Fatalities']].apply(lambda x: np.expmin(x))
        data_pred3.replace([np.inf, -np.inf], 0, inplace=True)

    #get_submission(data_pred3, 'ConfirmedCases', 'Predicted_ConfirmedCases', 'Predicted_Fatalities')

print("Process finished in ", round(time.time() - ts, 2), " seconds")

Currently running Linear Regression for all countries
Process finished in 6.01 seconds
```

From my experiments, this approach is not suitable for all countries linear regression model. In many cases there are strong transitional periods at the beginning, which frequently biases the regression. Hence, I will not apply this method on following sections, but you are welcome to use it for any other purposes.

Final LMSE score for week 2, with training data prior to 2020-03-19 and measures on date 2020-04-01:

1.62190

4.4. Linear regression with lags

With all the previous results in mind, I quite believe that Linear Regression is not a good approach for the early stages of the COVID-19 spread. Of course, this is only true for the initial outbreak where we are analyzing, and there's no way our model could predict when the number of new infections is going to decrease. But for short-term prediction purposes everything is fine, and we are in disposition to try to improve the results.

We remember those lagged variables we computed some sections before? Now it's time to use them, but first there's a problem to solve. If we use our dataset to predict the next following days of contagions, for the first day all the lags will be reported from the previous days, but what about the next days? Many of the lags will be unknown (flagged as 0), since the number of ConfirmedCases is only known for the train subset. The most simple approach to overcome this is:

1. Begin with the train dataset, with all cases and lags reported
2. Forecast only the following day, through the Linear Regression
3. Set the new prediction as a confirmed case
4. Recompute lags
5. Repeat from step 2 to step 4 for all remaining days

As usual, I'll start training on single countries in order to analyze the behavior of the model with these new features.

```
In [35]: # New split function, for one forecast day
def split_data_one_day(df, train_lim, test_lim):
    df_train = df[(df['Day_num'] <= train_lim, 'ForecastId' != -1)]
    df = df[(df['Day_num'] > test_lim)]

    #Train
    X_train = df[df['Day_num'] <= train_lim]
    y_train_1 = X_train.ConfirmedCases
    y_train_2 = X_train.Fatalities
    X_train.drop(['ConfirmedCases', 'Fatalities'], axis=1, inplace=True)

    #Test
    X_test = df[(df['Day_num'] >= test_lim)]
    X_test.drop(['ConfirmedCases', 'Fatalities'], axis=1, inplace=True)

    # Clean Id columns and keep ForecastId as index
    X_train.drop(['Id', 'inplace=True', errors='ignore'], axis=1)
    X_train.drop(['ForecastId', 'inplace=True', errors='ignore'], axis=1)
    X_test.drop(['Id', 'inplace=True', errors='ignore'], axis=1)
    X_test.drop(['ForecastId', 'inplace=True', errors='ignore'], axis=1)

    return X_train, y_train_1, y_train_2, X_test

def plot_real_vs_prediction_country(data, train, country_name, day_start, dates_list, march_day):
    # Select predictions from March list to March 25th
    predicted_data = data.loc[(data['Day_num'] <= march_day)]
    confirmed_data = data.loc[(data['Day_num'] > march_day)]
    real_data = train.loc[(train['Country_Region']==country_name) & (train['Date'].isin(dates_list)) & (train['Province_State']!=None)]
    dates_list_num = list(range(0, len(dates_list)))

    # Plot results
    fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(15,6))

    #Train
    ax1.plot(dates_list_num, np.expmin(predicted_data))
    ax1.plot(dates_list_num, np.expmin(real_data))
    ax1.axvline(30=march_day, linewidth=2, ls='-', color='grey', alpha=0.5)
    ax1.legend(['Predicted cases', 'Actual cases', 'Train-test split'], loc='upper left')
    ax1.set_xlabel('Day count (starting on March ' + str(march_day) + ')')
    ax1.set_ylabel('Confirmed Cases')

    ax2.plot(dates_list_num, predicted_data)
    ax2.plot(dates_list_num, np.expmin(real_data))
    ax2.axvline(30=march_day, linewidth=2, ls='-', color='grey', alpha=0.5)
    ax2.legend(['Predicted cases', 'Actual cases', 'Train-test split'], loc='upper left')
    ax2.set_xlabel('Day count (starting on March ' + str(march_day) + ')')
    ax2.set_ylabel('Log Confirmed Cases')

    plt.suptitle(("ConfirmedCases predictions based on Log-Linear Regression for "+country_name))

def plot_real_vs_prediction_country_fatalities(data, train, country_name, day_start, dates_list, march_day):
    # Select predictions from March list to March 25th
    predicted_data = data.loc[(data['Day_num'] <= march_day)]
    confirmed_data = data.loc[(data['Day_num'] > march_day)]
    real_data = train.loc[(train['Country_Region']==country_name) & (train['Date'].isin(dates_list)) & (train['Province_State']!=None)]
    dates_list_num = list(range(0, len(dates_list)))

    # Plot results
    fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(15,6))

    #Train
    ax1.plot(dates_list_num, np.expmin(predicted_data))
    ax1.plot(dates_list_num, np.expmin(real_data))
    ax1.axvline(30=march_day, linewidth=2, ls='-', color='grey', alpha=0.5)
    ax1.legend(['Predicted cases', 'Actual cases', 'Train-test split'], loc='upper left')
    ax1.set_xlabel('Day count (starting on March ' + str(march_day) + ')')
    ax1.set_ylabel('Fatalities Cases')

    ax2.plot(dates_list_num, predicted_data)
    ax2.plot(dates_list_num, np.expmin(real_data))
    ax2.axvline(30=march_day, linewidth=2, ls='-', color='grey', alpha=0.5)
    ax2.legend(['Predicted cases', 'Actual cases', 'Train-test split'], loc='upper left')
    ax2.set_xlabel('Day count (starting on March ' + str(march_day) + ')')
    ax2.set_ylabel('Log Fatalities Cases')

    plt.suptitle(("Fatalities predictions based on Log-Linear Regression for "+country_name))
```

• Spain

```
In [36]: # Function to compute the Linear Regression predictions with lags, for a certain Country/Region and Province
def lin_reg_with_lags_country(all_data, country_name, day_start, lag_size, country_dict, train_lim, test_lim):
    ts = time.time()

    # All country and features from all data (dataset without data leaking)
    data = all_data.copy()
    features = ['Id', 'Province_State', 'Country_Region', 'ConfirmedCases', 'Fatalities', 'ForecastId', 'Day_num']
    data = data[features]

    # Select country and data start (all days)
    data = data[(data['Country_Region']==country_name)]
    data = data.loc[(data['Day_num'] >= day_start)]

    # Lags
    data = calculate_lag(data, range(1, lag_size), 'ConfirmedCases')
    data = calculate_lag(data, range(1, lag_size), 'Fatalities')

    filter_col_confirmed = [col for col in data if col.startswith('Confirmed')]
    filter_col_fatalities = [col for col in data if col.startswith('Fataliti')]
    filter_col = np.append(filter_col_confirmed, filter_col_fatalities)
    data[filter_col] = data[filter_col].apply(lambda x: np.log1p(x))
    data.replace([np.inf, -np.inf], 0, inplace=True)
    data.fillna(0, inplace=True)

    # Start/end of forecast
    start_fcst = all_data[all_data['Id']==-1].Day_num.min()
    end_fcst = all_data[all_data['Id']==-1].Day_num.max()

    for d in list(range(start_fcst, end_fcst+1)):
        X_train, Y_train_1, Y_train_2, X_test = split_data_one_day(data, d, train_lim, test_lim)
        model_1, pred_1 = lin_reg(X_train, Y_train_1, X_test)
        data.loc[(data['Country_Region']==country_name)] = pred_1[0]
        model_2, pred_2 = lin_reg(X_train, Y_train_2, X_test)
        data.loc[(data['Country_Region']==country_name)] = pred_2[0]
        & (data['Day_num']==d), 'Fatalities'] = pred_2[0]

        # Recompute lags
        data = calculate_lag(data, range(1, lag_size), 'ConfirmedCases')
        data = calculate_lag(data, range(1, lag_size), 'Fatalities')
        data.replace([np.inf, -np.inf], 0, inplace=True)
        data.fillna(0, inplace=True)

    #print("Process for ", country_name, " finished in ", round(time.time() - ts, 2), " seconds")

    return data

# Function to compute the Linear Regression predictions with lags, for a certain Country/Region and Province
def lin_reg_with_lags_country_province(all_data, country_name, province_name, day_start, lag_size, country_dict, train_lim, test_lim):
    ts = time.time()

    # Filter country and features from all data (dataset without data leaking)
    data = all_data.copy()
    features = ['Id', 'Province_State', 'Country_Region', 'ConfirmedCases', 'Fatalities', 'ForecastId', 'Day_num']
    data = data[features]

    # Select country and data start (all days)
    data = data[(data['Country_Region']==country_name)]
    data = data.loc[(data['Day_num'] >= day_start)]

    # Lags
    data = calculate_lag(data, range(1, lag_size), 'ConfirmedCases')
    data = calculate_lag(data, range(1, lag_size), 'Fatalities')

    filter_col_confirmed = [col for col in data if col.startswith('Confirmed')]
    filter_col_fatalities = [col for col in data if col.startswith('Fataliti')]
    filter_col = np.append(filter_col_confirmed, filter_col_fatalities)
    data[filter_col] = data[filter_col].apply(lambda x: np.log1p(x))
    data.replace([np.inf, -np.inf], 0, inplace=True)
    data.fillna(0, inplace=True)

    # Start/end of forecast
    start_fcst = all_data[all_data['Id']==-1].Day_num.min()
    end_fcst = all_data[all_data['Id']==-1].Day_num.max()

    for d in list(range(start_fcst, end_fcst+1)):
        X_train, Y_train_1, Y_train_2, X_test = split_data_one_day(data, d, train_lim, test_lim)
        model_1, pred_1 = lin_reg(X_train, Y_train_1, X_test)
        data.loc[(data['Country_Region']==country_name)] = pred_1[0]
        model_2, pred_2 = lin_reg(X_train, Y_train_2, X_test)
        data.loc[(data['Country_Region']==country_name)] = pred_2[0]
        & (data['Day_num']==d), 'Fatalities'] = pred_2[0]

        # Recompute lags
        data = calculate_lag(data, range(1, lag_size), 'ConfirmedCases')
        data = calculate_lag(data, range(1, lag_size), 'Fatalities')
        data.replace([np.inf, -np.inf], 0, inplace=True)
        data.fillna(0, inplace=True)

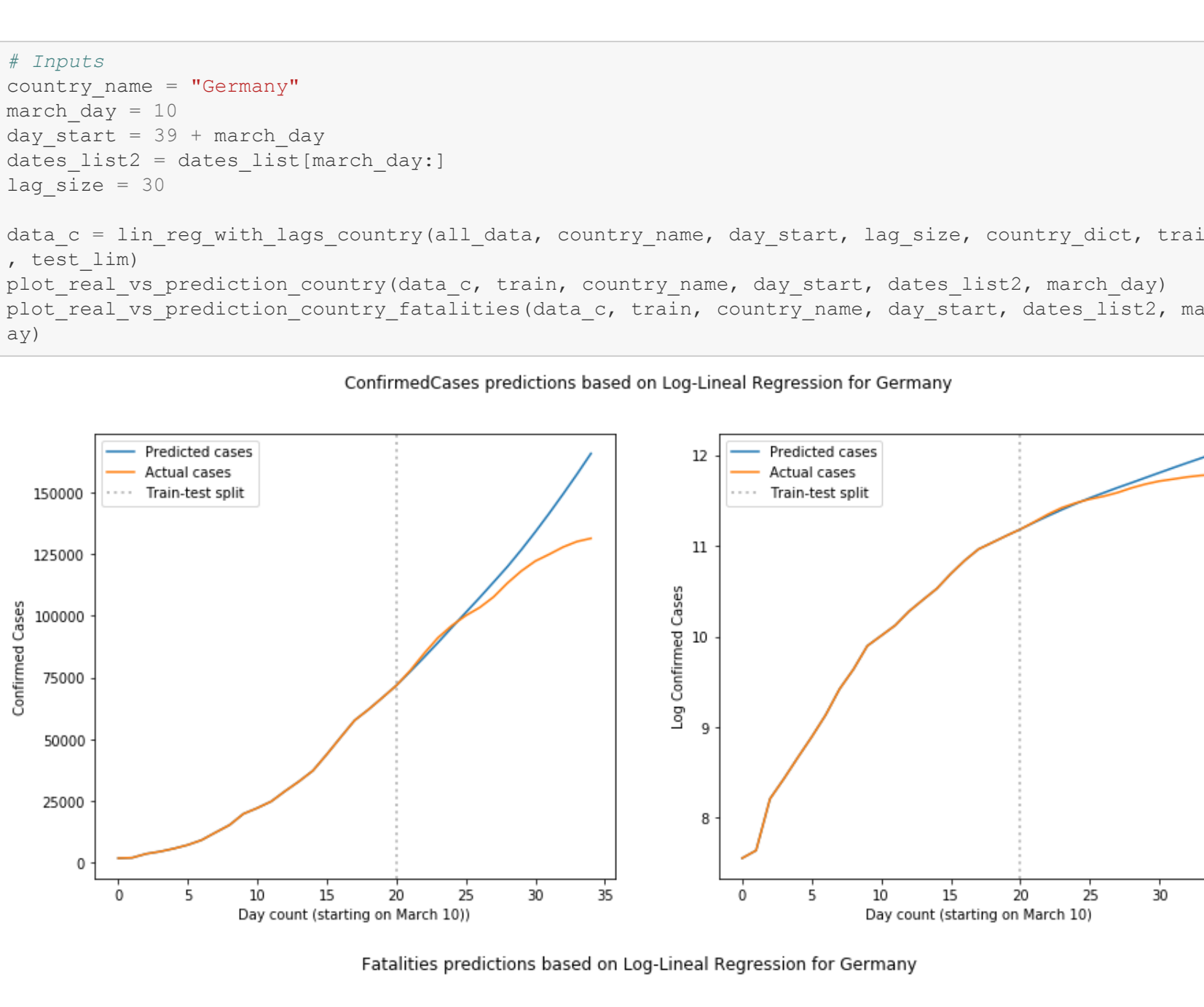
    #print("Process for ", country_name, " / ", province_name, " finished in ", round(time.time() - ts, 2), " seconds")

    return data

# Run the model for Spain
country_name = 'Spain'
march_day = 10
day_start = 39 = march_day
dates_list2 = dates_list[march_day:]
lag_size = 30

data_c = lin_reg_with_lags_country(all_data, country_name, day_start, lag_size, country_dict, train_lim, test_lim)
plot_real_vs_prediction_country(data_c, train, country_name, day_start, dates_list2, march_day)
plot_real_vs_prediction_country_fatalities(data_c, train, country_name, day_start, dates_list2, march_day)
```

ConfirmedCases predictions based on Log-Linear Regression for Spain



• Italy

```
In [37]: ts = time.time()

# Inputs
country_name = "Italy"
march_day = 10
day_start = 39 = march_day
dates_list2 = dates_list[march_day:]
lag_size = 30

data_c = lin_reg_with_lags_country(all_data, country_name, day_start, lag_size, country_dict, train_lim, test_lim)
plot_real_vs_prediction_country(data_c, train, country_name, day_start, dates_list2, march_day)
plot_real_vs_prediction_country_fatalities(data_c, train, country_name, day_start, dates_list2, march_day)
```

ConfirmedCases predictions based on Log-Linear Regression for Italy

• Germany

```
In [38]: # Inputs
country_name = "Germany"
march_day = 10
day_start = 39 = march_day
dates_list2 = dates_list[march_day:]
lag_size = 30

data_c = lin_reg_with_lags_country(all_data, country_name, day_start, lag_size, country_dict, train_lim, test_lim)
plot_real_vs_prediction_country(data_c, train, country_name, day_start, dates_list2, march_day)
plot_real_vs_prediction_country_fatalities(data_c, train, country_name, day_start, dates_list2, march_day)
```

ConfirmedCases predictions based on Log-Linear Regression for Germany

• Albania

