

Análisis de la segunda ola de la COVID-19

Objetivo

El objetivo principal de este proyecto es poder analizar la evolución de la pandemia COVID-19 en esta segunda ola, haciendo uso de algoritmos de Machine Learning, y así poder realizar predicciones de su evolución.

Prerrequisitos

Para poder realizar el proyecto debemos tener instalados en nuestro sistema una serie de librerías de Data Science y Machine Learning. A continuación quedan indicadas las mismas:

- Numpy
- Pandas
- Scikit-Learn
- Matplotlib
- Seaborn

Todas estas librerías estan ya instaladas en el entorno de trabajo de ANACONDA 3. Solo deberemos importar las mismas.

Configuración del entorno

```
In [1]: import warnings
import warnings as ignore
warnings.filterwarnings("ignore")
import datetime
import numpy as np
import pandas as pd
import sklearn
import matplotlib.pyplot as plt
import seaborn as sns
import matplotlib inline
```

Dataset

Para analizar este problema vamos a descargarlos un dataset que pone a disposición la Unión Europea en la siguiente url: <https://www.ecdc.europa.eu/en/publications-data/download-todays-data-geographic-distribution-covid-19-cases-worldwide>

```
In [2]: df = pd.read_excel("../Modulo 3/Dataset/COVID-19-geographic-distribution-worldwide-2020-11-01.xlsx")
```

Comprobamos los datos que hemos cargado en nuestro dataframe

```
In [3]: df.head(10)
```

	dateRep	day	month	year	cases	deaths	countriesAndTerritories	geoid	countryterritoryCode	popData2019	continentExp	Cumulative_n
0	2020-11-03	3	11	2020	95	3	Afghanistan	AF	AFG	38041757.0	Asia	
1	2020-11-02	2	11	2020	132	5	Afghanistan	AF	AFG	38041757.0	Asia	
2	2020-11-01	1	11	2020	76	0	Afghanistan	AF	AFG	38041757.0	Asia	
3	2020-10-31	31	10	2020	157	4	Afghanistan	AF	AFG	38041757.0	Asia	
4	2020-10-30	30	10	2020	123	3	Afghanistan	AF	AFG	38041757.0	Asia	
5	2020-10-29	29	10	2020	0	0	Afghanistan	AF	AFG	38041757.0	Asia	
6	2020-10-28	28	10	2020	113	7	Afghanistan	AF	AFG	38041757.0	Asia	
7	2020-10-27	27	10	2020	199	8	Afghanistan	AF	AFG	38041757.0	Asia	
8	2020-10-26	26	10	2020	65	3	Afghanistan	AF	AFG	38041757.0	Asia	
9	2020-10-15	25	10	2020	81	4	Afghanistan	AF	AFG	38041757.0	Asia	

```
In [4]: df.tail(10)
```

	dateRep	day	month	year	cases	deaths	countriesAndTerritories	geoid	countryterritoryCode	popData2019	continentExp	Cumulative_n
53156	2020-09-30	30	3	2020	0	0	Zimbabwe	ZW	ZWE	14645473.0	Africa	
53157	2020-09-29	29	3	2020	2	0	Zimbabwe	ZW	ZWE	14645473.0	Africa	
53158	2020-03-28	28	3	2020	2	0	Zimbabwe	ZW	ZWE	14645473.0	Africa	
53159	2020-03-27	27	3	2020	0	0	Zimbabwe	ZW	ZWE	14645473.0	Africa	
53160	2020-03-26	26	3	2020	1	0	Zimbabwe	ZW	ZWE	14645473.0	Africa	
53161	2020-03-25	25	3	2020	0	0	Zimbabwe	ZW	ZWE	14645473.0	Africa	
53162	2020-03-24	24	3	2020	0	1	Zimbabwe	ZW	ZWE	14645473.0	Africa	
53163	2020-03-23	23	3	2020	0	0	Zimbabwe	ZW	ZWE	14645473.0	Africa	
53164	2020-03-22	22	3	2020	1	0	Zimbabwe	ZW	ZWE	14645473.0	Africa	
53165	2020-03-21	21	3	2020	1	0	Zimbabwe	ZW	ZWE	14645473.0	Africa	

Comprobamos la info del DataFrame con el objetivo de determinar las columnas del df y el tipo de datos que contiene cada una de ellas. De esta forma, podremos realizar cualquier cambio de forma correcta según la tipología de los datos.

También se comprobamos los datos estadísticos principales.

```
In [5]: df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 53166 entries, 0 to 53165
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype
---  --
0   dateRep               non-null       datetime64[ns]
1   day                  non-null       int64
2   month                non-null       int64
3   year                 non-null       int64
4   cases                non-null       int64
5   deaths               non-null       int64
6   countriesAndTerritories non-null       object
7   geoid                non-null       object
8   countryterritoryCode non-null       object
9   popData2019          non-null       float64
10  continentExp          non-null       object
11  Cumulative number for 14 days of COVID-19 cases per 100000 float64
dtypes: datetime64[ns](1), float64(2), int64(5), object(4)
memory usage: 4.9+ MB
```

```
In [6]: df.describe()
```

	day	month	year	cases	deaths	popData2019	Cumulative number, for 14 days of COVID-19 cases, per 100000
count	53166.000000	53166.000000	53166.000000	53166.000000	53166.000000	5.308400e+04	50548.000000
mean	15.936708	6.366512	2019.998740	885.777038	22.707934	4.175643e+07	45.824057
std	8.918328	2.575420	0.035477	4977.075371	124.073255	1.550772e+08	112.109316
min	1.000000	1.000000	2019.000000	-8261.000000	-1918.000000	1.324800e+02	-147.419587
5th	8.000000	4.000000	2020.000000	0.000000	0.000000	1.324800e+02	0.555858
50th	16.000000	6.000000	2020.000000	12.000000	0.000000	7.813207e+06	5.652028
max	34.000000	9.000000	2020.000000	202.000000	3.000000	2.860872e+07	39.292371
max	34.000000	12.000000	2020.000000	101273.000000	49328.000000	1.433784e+09	1900.836210

El primer paso será cambiar el nombre de la columna 'Cumulative number for days of COVID-19 cases, per 100000'

```
In [7]: df = df.rename(
df.columns = df.columns.str.replace('Cumulative number for 14 days of COVID-19 cases, per 100000', 'Cum14_cases_100k'), axis = 1)
df
```

	dateRep	day	month	year	cases	deaths	countriesAndTerritories	geoid	countryterritoryCode	popData2019	continentExp	Cum14_cases_100k
0	2020-11-03	3	11	2020	95	3	Afghanistan	AF	AFG	38041757.0	Asia	
1	2020-11-02	2	11	2020	132	5	Afghanistan	AF	AFG	38041757.0	Asia	
2	2020-11-01	1	11	2020	76	0	Afghanistan	AF	AFG	38041757.0	Asia	
3	2020-10-31	31	10	2020	157	4	Afghanistan	AF	AFG	38041757.0	Asia	
4	2020-10-30	30	10	2020	123	3	Afghanistan	AF	AFG	38041757.0	Asia	

53166 rows x 12 columns

También se utilizará la columna dateRep como índice del DataFrame.

```
In [8]: df.set_index('dateRep', inplace=True)
```

Y ordenamos el DataFrame por el nuevo índice marcado. De esta forma tendremos ordenado el df mediante la fecha de los datos.

```
In [9]: df.sort_index(inplace=True)
df
```

	day	month	year	cases	deaths	countriesAndTerritories	geoid	countryterritoryCode	popData2019	continentExp	Cum14_cases_100k
2019-12-31	31	12	2019	0	0	Nigeria	NG	NGA	200963603.0	Africa	Na
2019-12-31	31	12	2019	0	0	Finland	FI	FIN	5517919.0	Europe	Na
2019-12-31	31	12	2019	0	0	Singapore	SG	SGP	5804343.0	Asia	Na
2019-12-31	31	12	2019	0	0	Luxemburg	LU	LUX	618940.0	Europe	Na
2019-12-31	31	12	2019	0	0	Netherlands	NL	NLD	17282163.0	Europe	Na
...
2020-11-03	3	11	2020	0	0	Saint Lucia	LC	LCA	182795.0	America	26.25890
2020-11-03	3	11	2020	28	0	Mali	ML	MLI	19658023.0	Africa	0.84461
2020-11-03	3	11	2020	0	0	Haiti	HT	HTI	11263079.0	America	0.69252
2020-11-03	3	11	2020	193	13	Honduras	HN	HND	9746115.0	America	92.59070
2020-11-03	3	11	2020	95	3	Afghanistan	AF	AFG	38041757.0	Asia	3.78791

53166 rows x 11 columns

Análisis de los datos

Qué países son los que tienen mayor incidencia de COVID con los datos disponibles?

Generamos un dataframe por país:

```
In [10]: df.countries = df[['countriesAndTerritories', 'Cum14_cases_100k']].\
groupby('countriesAndTerritories').last()
df.countries
```

Cum14_cases_100k

Columbia	31670	5.02944e+07	1099126	25.020827	4291289	2171	768170
Argentina	31670	4.47806e+07	1811818	40306431	70617515	2642	028062
Russia	28473	4.55272e+08	665508	164.36934	195.91319	1134	580351
South Africa	19465	5.855827e+07	727395	38.151744	33.240396	282	514571
Chile	14302	1.895204e+07	513188	104.912217	75.464192	2707	825307
Indonesia	14044	2.70625e+08	415402	18.355573	5.189458	53	496953
Ecuador	12692	1.73736e+07	605932	92.88951	73.053129	975	971841

213 rows x 1 columns

Generamos un dataframe con los datos del indicador de Casos por 100.000 habitantes en los últimos 14 días

```
In [11]: df.countries.sort_values(by=['Cum14_cases_100k'], ascending=False, inplace=True)
df.countries
```

	Cum14_cases_100k
countriesAndTerritories	
Belgium	1753.181152
Andorra	1660.606220
Czechia	1586.264531
Luxembourg	1337.853115
French Polynesia	1240.668135
Slovenia	1082.460157
Switzerland	1081.522710
Armenia	991.162192
Liechtenstein	990.150607
France	829.924001
Netherlands	762.410353
Croatia	654.695521
Georgia	598.384392
Slovakia	571.221929
Spain	567.244732
Guam	560.088486
Poland	558.905145
Montenegro	554.500130
Austria	537.866347
Bosnia and Herzegovina	533.293265
Italy	510.292109
San Marino	470.205788
United Kingdom	469.115601
North Macedonia	443.351698
Portugal	437.760792

Analizamos la situación actual por país de la pandemia del COVID

```
In [12]: #Generamos un nuevo dataframe con los datos agrupados
df_agrup = df.groupby('countriesAndTerritories').agg(
muerres_COVID = ('deaths', 'sum'),
# Población total
poblacion = ('popData2019', 'last'),
# Total enfermos COVID
enfermos_COVID = ('cases', 'sum'),
# Situación actual casos por 100k
cum14_cases_100k_ultimo_dia = ('Cum14_cases_100k', 'last')
df_agrup
```

muerres_COVID población enfermos_COVID cum14_cases_100k_ultimo_dia

Italy	510.292109	12.12.050203	64.710560	5.338934
France	829.920401	21.88.265109	55.862393	2.552793
Spain	567.244732	26.43.320651	77.245997	2.922309
Iran	113.550331	758.355006	43.102547	5.683705
Peru	109.715451	2788.471600	106.381140	3.815034
Colombia	253.028227	2171.768170	62.913982	2.896851
Argentina	403.006431	2642.028062	76.175715	2.672853

213 rows x 4 columns

Se crean nuevas variables para poder entender mejor el impacto de la pandemia. Las nuevas variables son las siguientes:

- Muertos por cada 100k habitantes
- Casos por cada 100k habitantes
- Índice mortalidad Covid

```
In [13]: df_agrup['muerres_100k'] = df_agrup['muerres_COVID']/df_agrup['poblacion']*100000
df_agrup['casos_100k'] = df_agrup['enfermos_COVID']/df_agrup['poblacion']*100000
df_agrup['indice_mortalidad_COVID'] = df_agrup['muerres_COVID']/df_agrup['enfermos_COVID']*100
df_agrup
```

muerres_COVID población enfermos_COVID cum14_cases_100k_ultimo_dia muerres_100k casos_100k indice_mortalidad_COVID

Guam	560.088466	2805.24304	47.221973	1.683358
Poland	558.905145	1041.481995	15.471596	1.485537
Montenegro	554.500130	3087.521015	15.413896	1.665799
Austria	537.986347	1289.907465	12.818155	0.993262
Bosnia and Herzegovina	533.293265	1583.430223	38.776152	2.448870

Calculamos la correlación entre las variables

```
df.clustering_corr()
```

213 rows x 7 columns

Ordenamos el df por el número de muertes provocadas por la COVID-19 y así vemos los países con mas muertes y vemos también que relación tienen con el índice de mortalidad.

```
In [14]: df_agrup = df_agrup.sort_values(by = ['muerres_COVID'], ascending=False)
df_agrup.head(20)
```

muerres_COVID población enfermos_COVID cum14_cases_100k_ultimo_dia muerres_100k casos_100k indice_mortalidad_COVID

```

Primero se comprobará que no haya datos nulos

df_clustering.isnull().sum()

sort_values(ascending=False)

cum14_casos_100k_ultimo_dia    3
muertes_100k                  2
casos_100k                     2
indice_mortalidad_COVID       0
dtype: int64

Borraremos los datos nulos para evitar interferencias con los análisis

df_clustering = df_clustering.dropna()

Normalizaremos los datos para que las diferentes unidades de medida que tienen nuestras variables estén a la misma escala y no interfieran con el modelo

from sklearn.preprocessing import StandardScaler, normalize
scaler = StandardScaler()
scaled_df = scaler.fit_transform(df_clustering)

normalized_df = normalize(scaled_df)

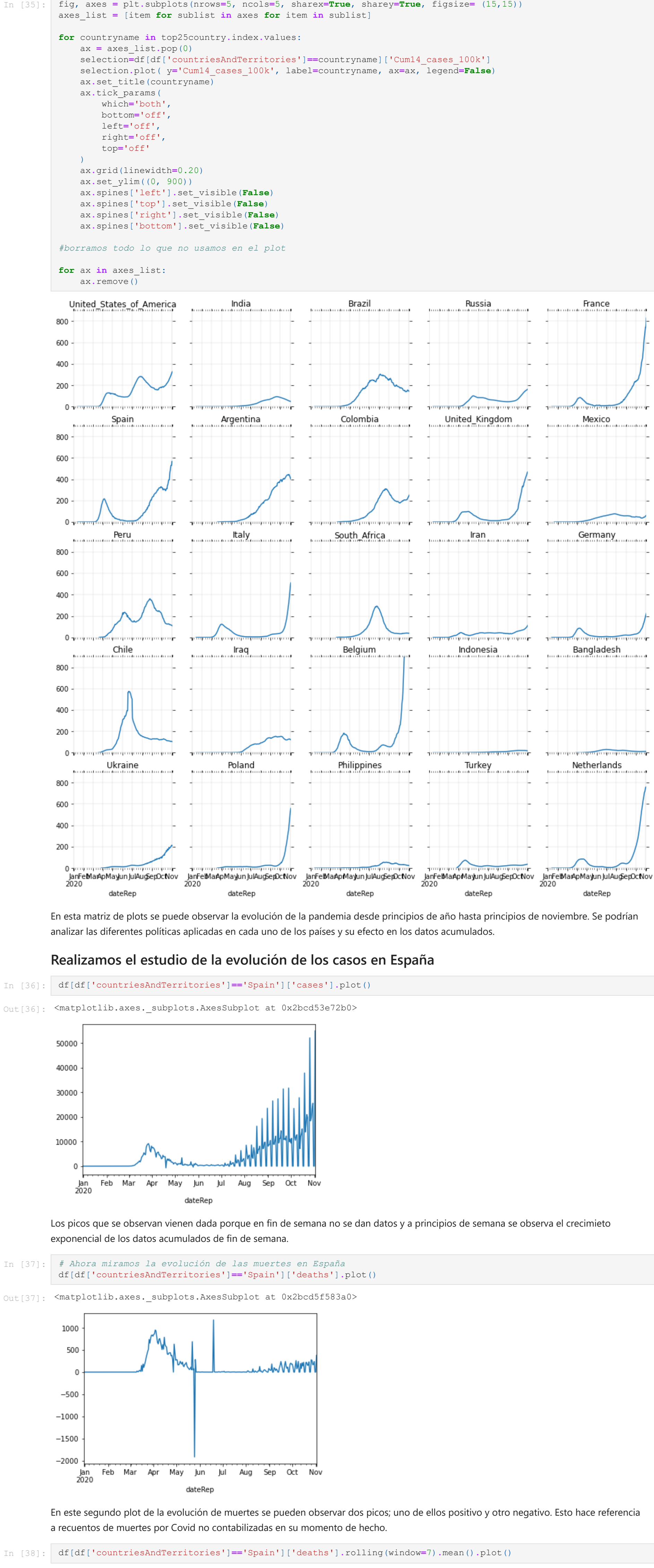
# Convertimos la matriz normalizada en un DataFrame
normalized_df = pd.DataFrame(normalized_df, columns = df_clustering.columns)

normalized_df.head(10)

```

	cum14_casos_100k_ultimo_dia	casos_100k	muertes_100k	indice_mortalidad_COVID
0	0.747726	0.392394	0.534849	0.029490
1	0.626096	0.637425	0.448146	-0.029527
2	0.892077	0.418168	0.150817	-0.081171
3	0.875089	0.459521	0.099628	-0.114604
4	0.899371	0.397668	-0.058030	-0.172116

Se realizará una matriz de gráficos para ver la evolución de la tasa de casos de los últimos 14 días por 100k habitantes.



En esta matriz de plots se puede observar la evolución de la pandemia desde principios de año hasta principios de noviembre. Se podrían analizar las diferentes políticas aplicadas en cada uno de los países y su efecto en los datos acumulados.

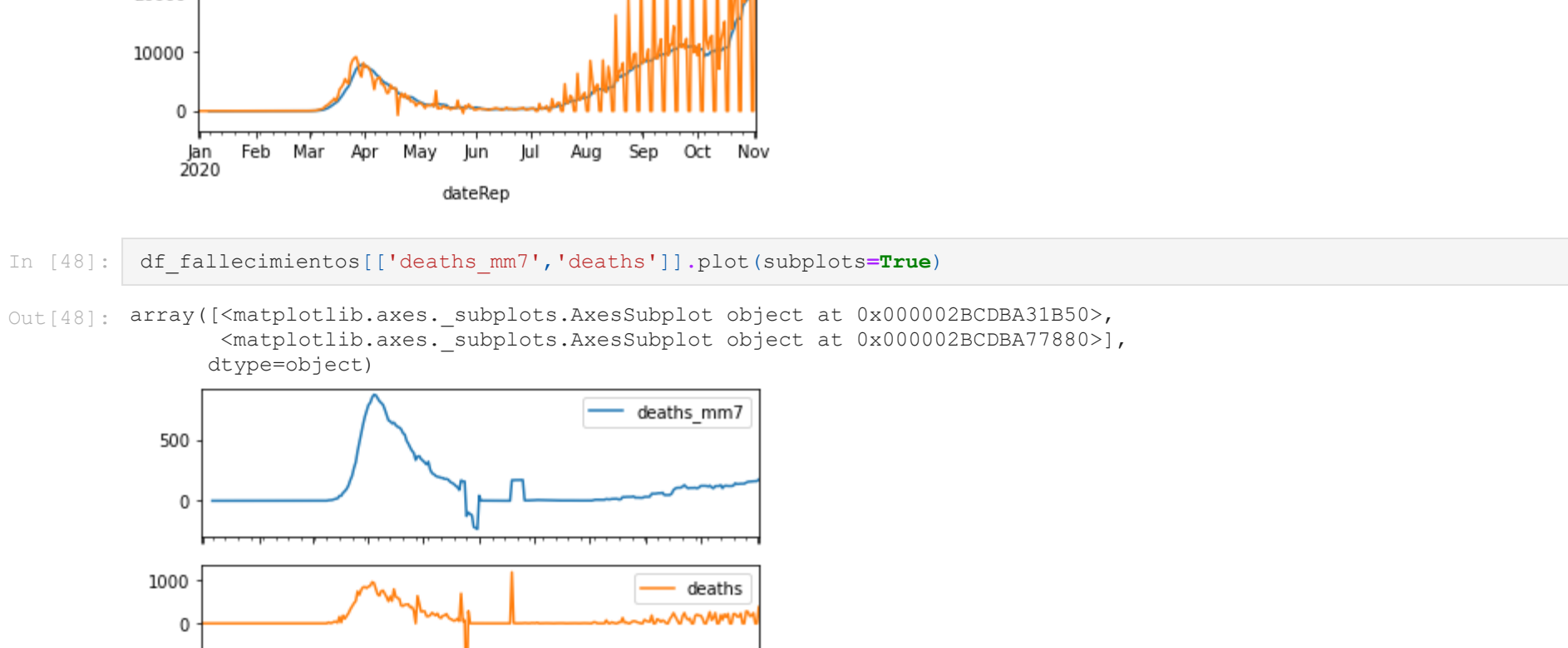
Realizamos el estudio de la evolución de los casos en España



Los picos que se observan vienen dada porque en fin de semana no se dan datos y a principios de semana se observa el crecimiento exponencial de los datos acumulados de fin de semana.



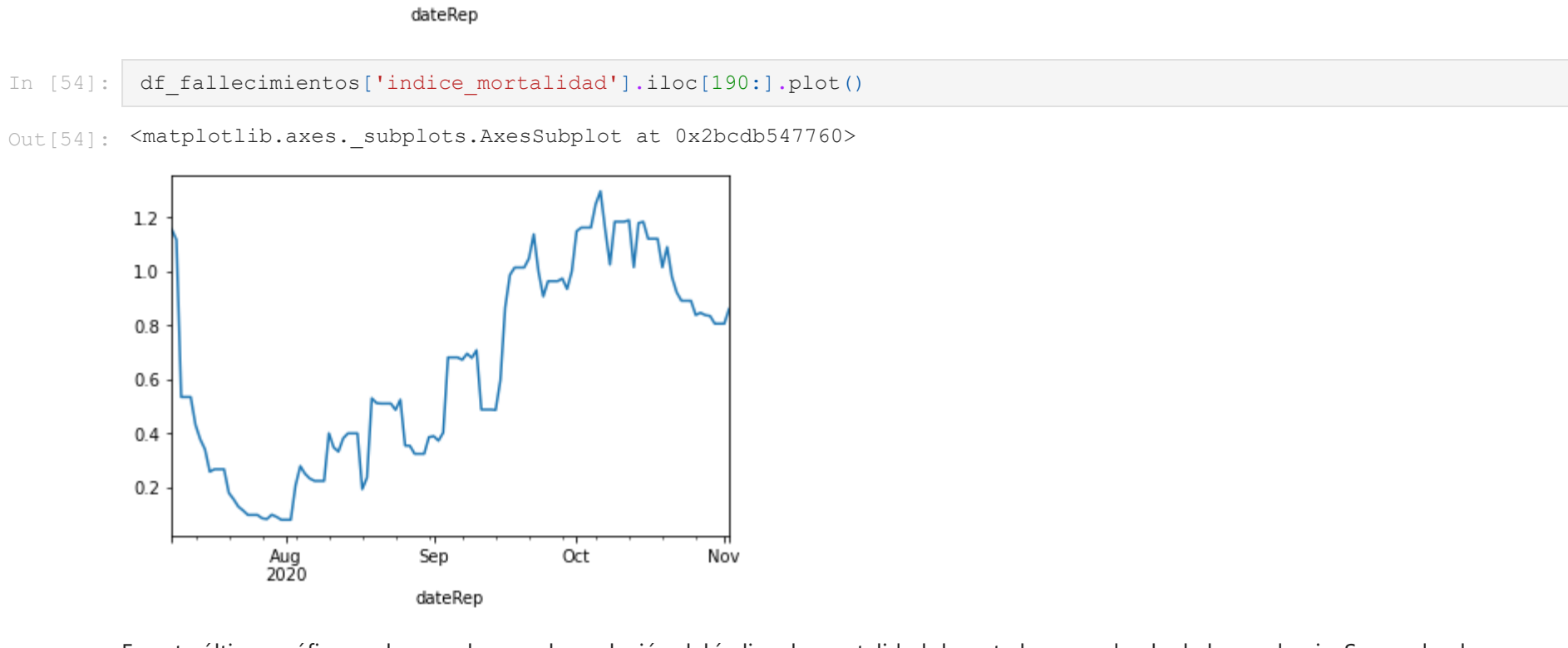
En este segundo plot de la evolución de muertes se pueden observar dos picos; uno de ellos positivo y otro negativo. Esto hace referencia a recuentos de muertes por Covid no contabilizadas en su momento de hecho.



Como ha evolucionado el índice de mortalidad en España



Debido a que España es un país que no ha comunicado los datos de muertes y casos de forma diaria, ya que durante el fin de semana no ha comunicado ningún caso o fallecimiento, calculamos la media móvil para interpretar mejor los gráficos.

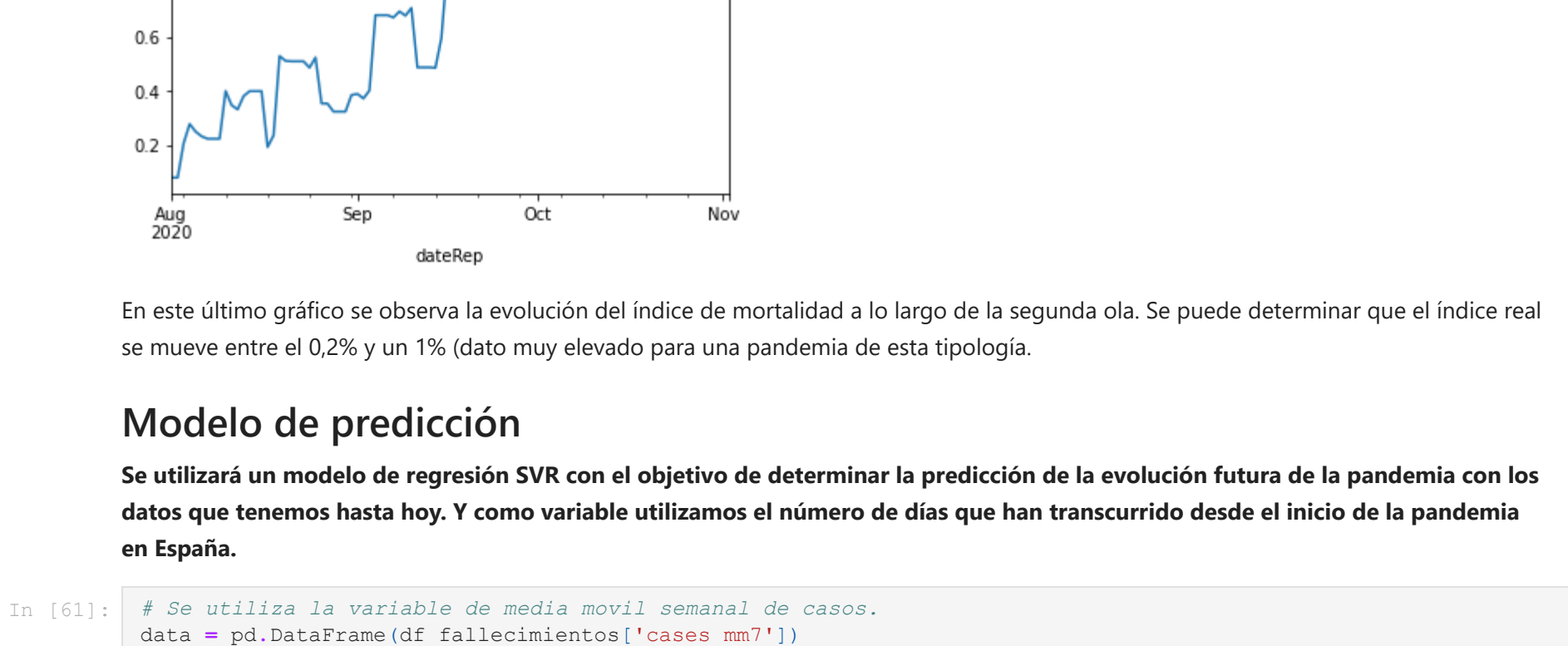


Esta diferencia en la gráfica es debida a que los datos estadísticos de la pandemia no se han tratado de la forma correcta. Otro factor que influye en los datos es la evolución de los protocolos seguidos para la detección de casos, ya que han ido variando a lo largo del tiempo.

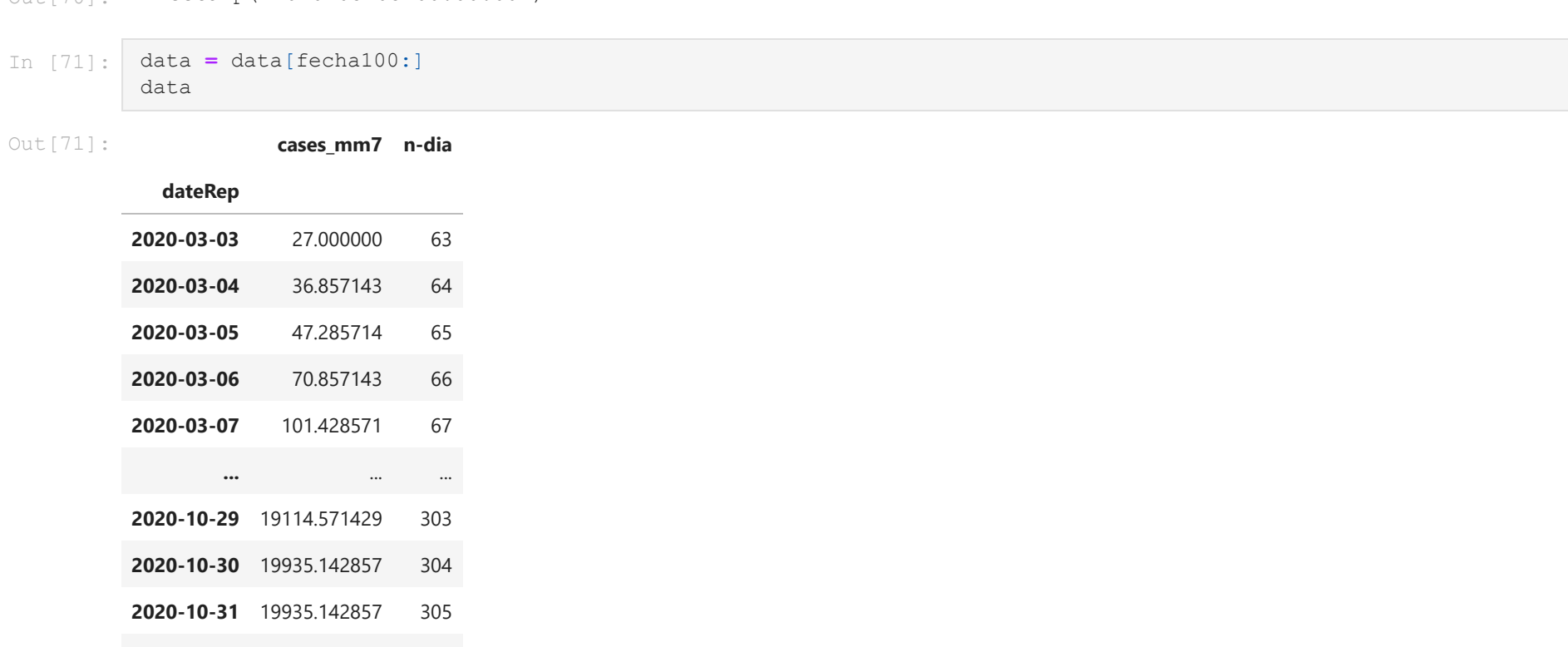
Se procede a dividir los datos en las dos oleadas de la pandemia, y así eliminar la parte donde España realizó el ajuste en las muertes



En este último gráfico podemos observar la evolución del índice de mortalidad durante la segunda ola de la pandemia. Se puede observar que un índice de mortalidad acorde a la realidad, ya que no existen ajustes en los datos, como si ocurre durante la primera ola, donde no se apuntaban todos los casos reales, ya que no se hacían pruebas PCR.



Aquí podemos observar, como después de los ajustes en los datos de casos y muertos, el índice de mortalidad se va ajustando a datos reales. También coincidir con el cambio en los protocolos de detección de casos mediante las pruebas PCR.



En este último gráfico se observa la evolución del índice de mortalidad a lo largo de la segunda ola. Se puede determinar que el índice real se mueve entre el 0,2% y un 1% (dato muy elevado para una pandemia de esta tipología).

Modelo de predicción

Se utilizará un modelo de regresión SVR con el objetivo de determinar la predicción de la evolución futura de la pandemia con los datos que tenemos hasta hoy. Y como variable utilizamos el número de días que han transcurrido desde el inicio de la pandemia en España.

