

Initial Data Exploration for Covid19 pandemic in Malaysia

Data Dictionary

	Field	Description
In [1]:	Date	Date of incident
	Province/State	If there is a breakdown in states
Out [1]:	Lat	Latitude
	Lon	Longitude
In [2]:	Confirmed	Confirmed cases
	Recovered	Recovered cases
Out [2]:	Deaths	Death cases
	Active	Active cases

Analysis Summary

- From worldwide perspective, the confirmed cases are the highest range and many outliers found in boxplot.
- This analysis and machine learning model will predict COVID-19 deaths in Malaysia.
- From time plots, the biggest jump is from October 2020 onwards.
- By extracting individual cases by day, deaths peaked in April 2020 and October 2020 onwards.
- Separate charts by Month and Day are plotted to see the distribution.

Feature Engineering

- The dataset is transformed using diff method to extract per day for Confirmed, Recovered, Deaths and Active Cases.
- Due to diff method, I imputed a NaN value as 0.0 as starting point.
- Created separate month and day features from Date feature.

Import Libraries

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import statsmodels.api as sm
import datetime
import scipy.stats

#matplotlib inline
%matplotlib inline
#set the default autosave frequency in seconds
%autosave 60
sns.set_style('dark')
sns.set(font_scale=1.2)

import warnings
warnings.filterwarnings('ignore')

#Import feature_engine missing data imputers as mdi
#from feature_engine.outlier_removers import Winsorizer
#from feature_engine import categorical_encoders as ce

pd.set_option('display.max_columns',None)
#pd.set_option('display.max_rows',None)
pd.set_option('display.width', 1000)

np.random.seed(0)
np.set_printoptions(suppress=True)

Autosaving every 60 seconds
```

```
In [2]: df = pd.read_csv('covid19_data_cleaned.csv', parse_dates=['Date'], index_col=['Date'])
```

```
Out [3]:
```

	Province/State	Country	Lat	Long	Confirmed	Recovered	Deaths	Active
Date								
2020-01-22	NaN	Afghanistan	33.93911	67.709953	0	0	0	0
2020-01-23	NaN	Afghanistan	33.93911	67.709953	0	0	0	0
2020-01-24	NaN	Afghanistan	33.93911	67.709953	0	0	0	0
2020-01-25	NaN	Afghanistan	33.93911	67.709953	0	0	0	0
2020-01-26	NaN	Afghanistan	33.93911	67.709953	0	0	0	0
...
2020-11-21	NaN	Timor-Leste	-8.87420	125.727500	0	30	0	-30
2020-11-22	NaN	Timor-Leste	-8.87420	125.727500	0	30	0	-30
2020-11-23	NaN	Timor-Leste	-8.87420	125.727500	0	30	0	-30
2020-11-24	NaN	Timor-Leste	-8.87420	125.727500	0	30	0	-30
2020-11-25	NaN	Timor-Leste	-8.87420	125.727500	0	30	0	-30

85593 rows x 8 columns

Exploratory Data Analysis

```
In [4]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 85593 entries, 2020-01-22 to 2020-11-25
Data columns (total 8 columns):
 #   Column      Non-Null Count  Dtype
---  ---
 0   Province/State  26265 non-null    object
 1   Country        85593 non-null    object
 2   Lat            85593 non-null    float64
 3   Long           85593 non-null    float64
 4   Confirmed      85593 non-null    int64
 5   Recovered      85593 non-null    int64
 6   Deaths         85593 non-null    int64
 7   Active         85593 non-null    int64
dtypes: float64(2), int64(4), object(2)
memory usage: 5.3+ MB
```

```
In [5]: df.describe()
```

```
Out [5]:
```

	Lat	Long	Confirmed	Recovered	Deaths	Active
count	85593.000000	85593.000000	8.559300e+04	8.559300e+04	8.559300e+04	8.559300e+04
mean	20.916828	24.585736	5.867697e+04	3.708964e+04	1.926112e+06	1.966121e+04
std	25.065486	71.663375	4.359217e+05	2.917598e+05	1.198106235e	1.936343e+05
min	-51.796300	-135.000000	0.000000e+00	0.000000e+00	0.000000	-2.818830e+05
25%	6.423800	-15.180400	1.600000e+01	0.000000e+00	0.000000	1.000000e+00
50%	22.300000	21.745300	5.130000e+02	2.900000e+02	6.000000	6.800000e+01
75%	41.129000	85.240100	6.442000e+03	2.191000e+03	118.000000	1.667000e+03
max	71.706900	178.050000	1.272265e+07	8.679138e+06	26.2222.000000	7.674475e+06

```
In [6]: df.columns
```

```
Out [6]: Index(['Province/State', 'Country', 'Lat', 'Long', 'Confirmed', 'Recovered', 'Deaths', 'Active'], dtype=object)
```

```
In [7]: df["Country"].value_counts()
```

```
Out [7]:
```

Country	Count
China	10815
Canada	5253
France	3399
United Kingdom	3399
Australia	2472

```
In [8]: df["Country"].unique()
```

```
Out [8]: array(['Afghanistan', 'Albania', 'Algeria', 'Andorra', 'Angola', 'Antigua and Barbuda', 'Argentina', 'Armenia', 'Australia', 'Austria', 'Azerbaijan', 'Bahamas', 'Bahrain', 'Bangladesh', 'Barbados', 'Belarus', 'Belgium', 'Belize', 'Benin', 'Bhutan', 'Bolivia', 'Bosnia and Herzegovina', 'Botswana', 'Brazil', 'Brunei', 'Bulgaria', 'Burkina Faso', 'Burma', 'Burundi', 'Cabo Verde', 'Cambodia', 'Cameroon', 'Canada', 'Central African Republic', 'Chad', 'Chile', 'China', 'Colombia', 'Comoros', 'Congo (Brazzaville)', 'Congo (Kinshasa)', 'Costa Rica', 'Cote d'Ivoire', 'Croatia', 'Cuba', 'Cyprus', 'Czechia', 'Denmark', 'Diamond Princess', 'Djibouti', 'Dominica', 'Dominican Republic', 'Ecuador', 'Egypt', 'El Salvador', 'Equatorial Guinea', 'Eritrea', 'Estonia', 'Eswatini', 'Ethiopia', 'Fiji', 'Finland', 'France', 'Gabon', 'Gambia', 'Georgia', 'Germany', 'Ghana', 'Greece', 'Grenada', 'Guatemala', 'Guinea', 'Guinea-Bissau', 'Guyana', 'Honduras', 'Hungary', 'Iceland', 'India', 'Indonesia', 'Iran', 'Iraq', 'Ireland', 'Israel', 'Italy', 'Jamaica', 'Japan', 'Jordan', 'Kazakhstan', 'Kenya', 'Korea, South', 'Kosovo', 'Kuwait', 'Kyrgyzstan', 'Laos', 'Latvia', 'Lebanon', 'Lesotho', 'Liberia', 'Libya', 'Liechtenstein', 'Lithuania', 'Luxembourg', 'Madagascar', 'Malawi', 'Malaysia', 'Maldives', 'Mali', 'Malta', 'Marshall Islands', 'Mauritania', 'Mauritius', 'Mexico', 'Moldova', 'Monaco', 'Mongolia', 'Montenegro', 'Morocco', 'Mozambique', 'Namibia', 'Nepal', 'Netherlands', 'New Zealand', 'Nicaragua', 'Niger', 'Nigeria', 'North Macedonia', 'Norway', 'Oman', 'Pakistan', 'Panama', 'Papua New Guinea', 'Paraguay', 'Peru', 'Philippines', 'Poland', 'Portugal', 'Qatar', 'Romania', 'Russia', 'Rwanda', 'Saint Kitts and Nevis', 'Saint Lucia', 'Saint Vincent and the Grenadines', 'San Marino', 'Saudi Arabia', 'Senegal', 'Serbia', 'Seychelles', 'Sierra Leone', 'Singapore', 'Slovakia', 'Slovenia', 'Solomon Islands', 'Somalia', 'South Africa', 'South Sudan', 'Spain', 'Sri Lanka', 'Sudan', 'Suriname', 'Sweden', 'Switzerland', 'Syria', 'Taiwan', 'Tajikistan', 'Tanzania', 'Thailand', 'Timor-Leste', 'Togo', 'Trinidad and Tobago', 'Tunisia', 'Turkey', 'Uganda', 'Ukraine', 'United Arab Emirates', 'United Kingdom', 'Uruguay', 'Uzbekistan', 'Vanuatu', 'Venezuela', 'Vietnam', 'West Bank and Gaza', 'Western Sahara', 'Yemen', 'Zambia', 'Zimbabwe'], dtype=object)
```

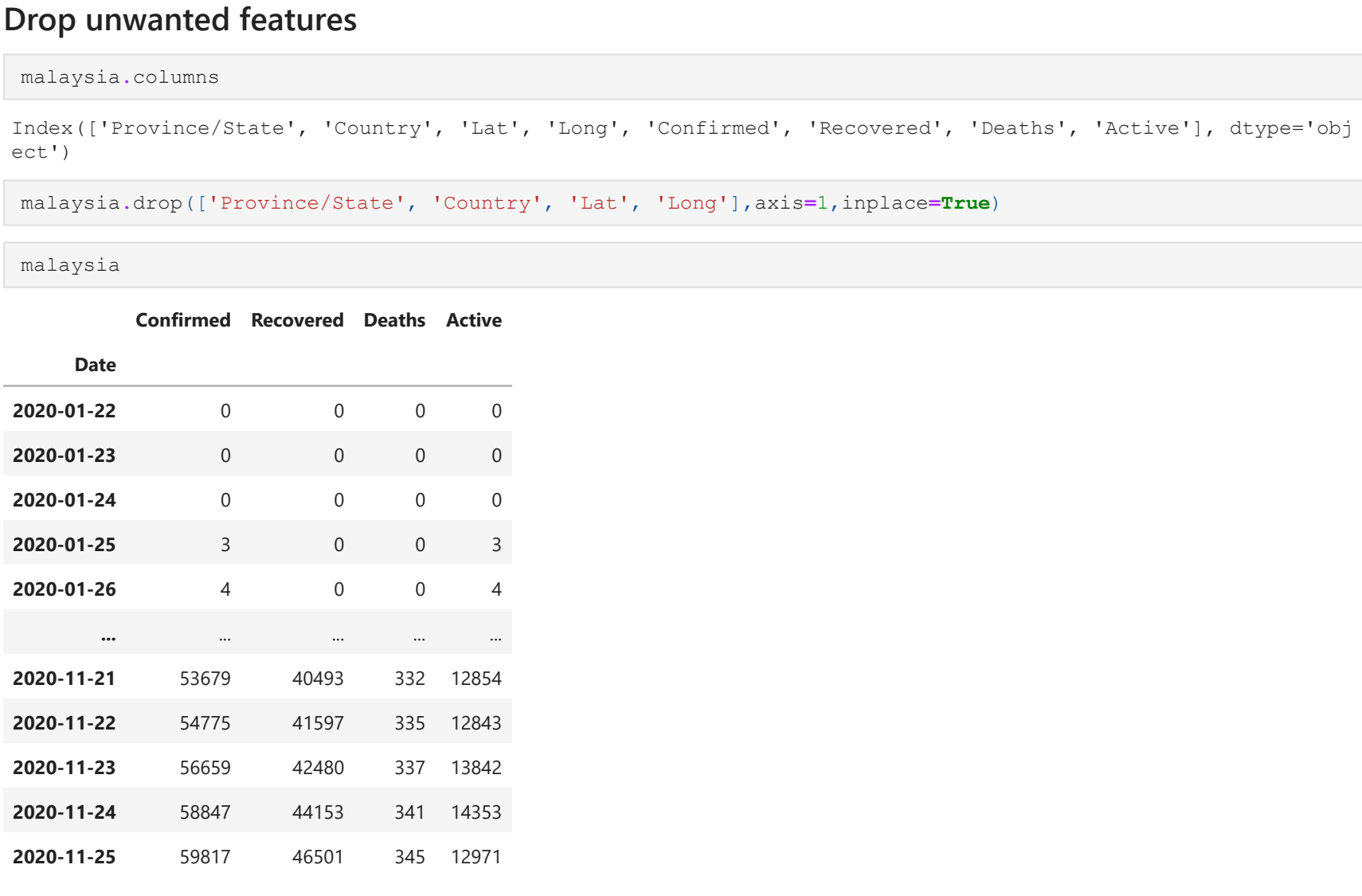
```
In [9]: df["Country"].nunique()
```

```
Out [9]: 191
```

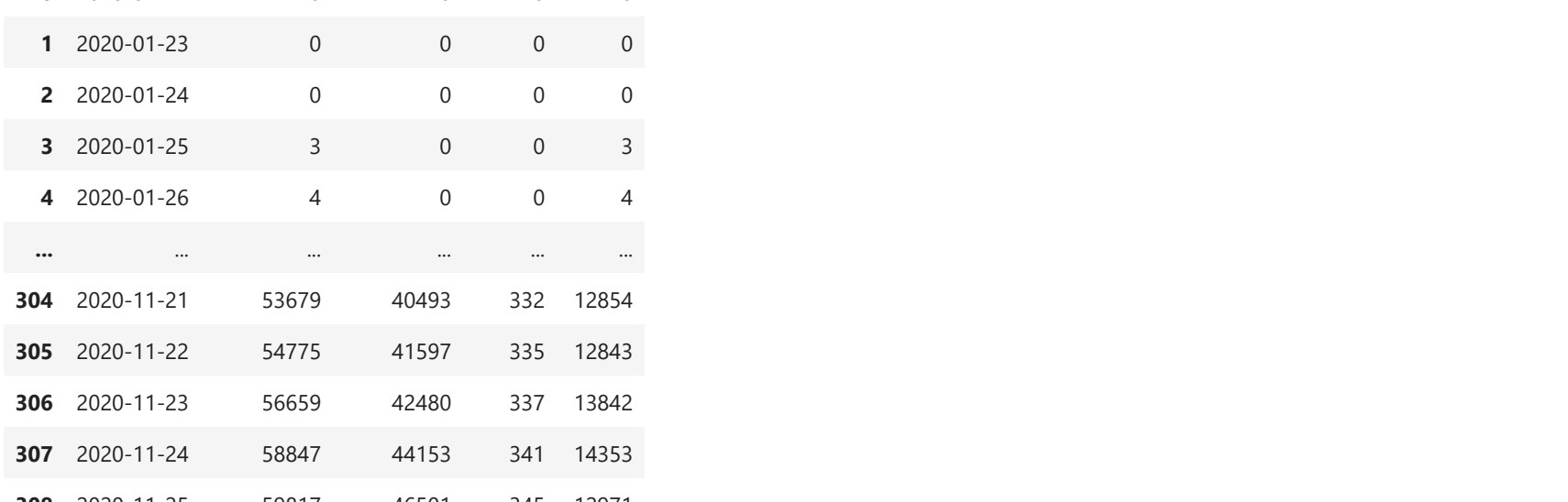
Data Visualization (Part 1)

Univariate Data Exploration

```
In [10]: df.hist(bins=50, figsize=(20,10))
plt.suptitle('Feature Distribution', x=0.5, y=1.02, ha='center', fontsize='large')
plt.tight_layout()
plt.show()
```



```
In [11]: df.boxplot(figsize=(20,10))
plt.suptitle('BoxPlot', x=0.5, y=1.02, ha='center', fontsize='large')
plt.tight_layout()
plt.show()
```



Extract country Malaysia from dataset for analysis and machine learning

```
In [12]: malaysia = df[df["Country"] == "Malaysia"]
```

```
Out [12]:
```

	Province/State	Country	Lat	Long	Confirmed	Recovered	Deaths	Active
Date								
2020-01-22	NaN	Malaysia	4.210484	101.975766	0	0	0	0
2020-01-23	NaN	Malaysia	4.210484	101.975766	0	0	0	0
2020-01-24	NaN	Malaysia	4.210484	101.975766	0	0	0	0
2020-01-25	NaN	Malaysia	4.210484	101.975766	3	0	0	3
2020-01-26	NaN	Malaysia	4.210484	101.975766	4	0	0	4
...
2020-11-21	NaN	Malaysia	4.210484	101.975766	53679	40493	332	12854
2020-11-22	NaN	Malaysia	4.210484	101.975766	54775	41597	335	12843
2020-11-23	NaN	Malaysia	4.210484	101.975766	56659	42480	337	13842
2020-11-24	NaN	Malaysia	4.210484	101.975766	58847	44153	341	14353
2020-11-25	NaN	Malaysia	4.210484	101.975766	59817	46501	345	12971

309 rows x 8 columns

Drop unwanted features

```
In [14]: malaysia.columns
```

```
Out [14]: Index(['Province/State', 'Country', 'Lat', 'Long', 'Confirmed', 'Recovered', 'Deaths', 'Active'], dtype=object)
```

```
In [15]: malaysia.drop(['Province/State', 'Country', 'Lat', 'Long'], axis=1, inplace=True)
```

```
Out [15]:
```

	Confirmed	Recovered	Deaths	Active
Date				
2020-01-22	0	0	0	0
2020-01-23	0	0	0	0
2020-01-24	0	0	0	0
2020-01-25	3	0	0	3
2020-01-26	4	0	0	4
...
2020-11-21	53679	40493	332	12854
2020-11-22	54775	41597	335	12843
2020-11-23	56659	42480	337	13842
2020-11-24	58847	44153	341	14353
2020-11-25	59817	46501	345	12971

309 rows x 4 columns

```
In [17]: malaysia.reset_index(inplace=True)
```

```
In [18]: #Save to csv
#malaysia.to_csv('malaysiacovid.csv', index=False)
```

Time-Series Analysis

```
In [19]: df = pd.read_csv('malaysiacovid.csv', parse_dates=["Date"])
```

```
Out [19]:
```

	Date	Confirmed	Recovered	Deaths	Active
0	2020-01-22	0	0	0	0
1	2020-01-23	0	0	0	0
2	2020-01-24	0	0	0	0
3	2020-01-25	3	0	0	3
4	2020-01-26	4	0	0	4
...
304	2020-11-21	53679	40493	332	12854
305	2020-11-22	54775	41597	335	12843
306	2020-11-23	56659	42480	337	13842
307	2020-11-24	58847	44153	341	14353
308	2020-11-25	59817	46501	345	12971

309 rows x 5 columns

```
In [21]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 309 entries, 0 to 308
Data columns (total 5 columns):
 #   Column      Non-Null Count  Dtype
---  ---
 0   Date        309 non-null    datetime64[ns]
 1   Confirmed   309 non-null    int64
 2   Recovered   309 non-null    int64
 3   Deaths     309 non-null    int64
 4   Active      309 non-null    int64
dtypes: datetime64(1), int64(4)
memory usage: 12.2 KB
```

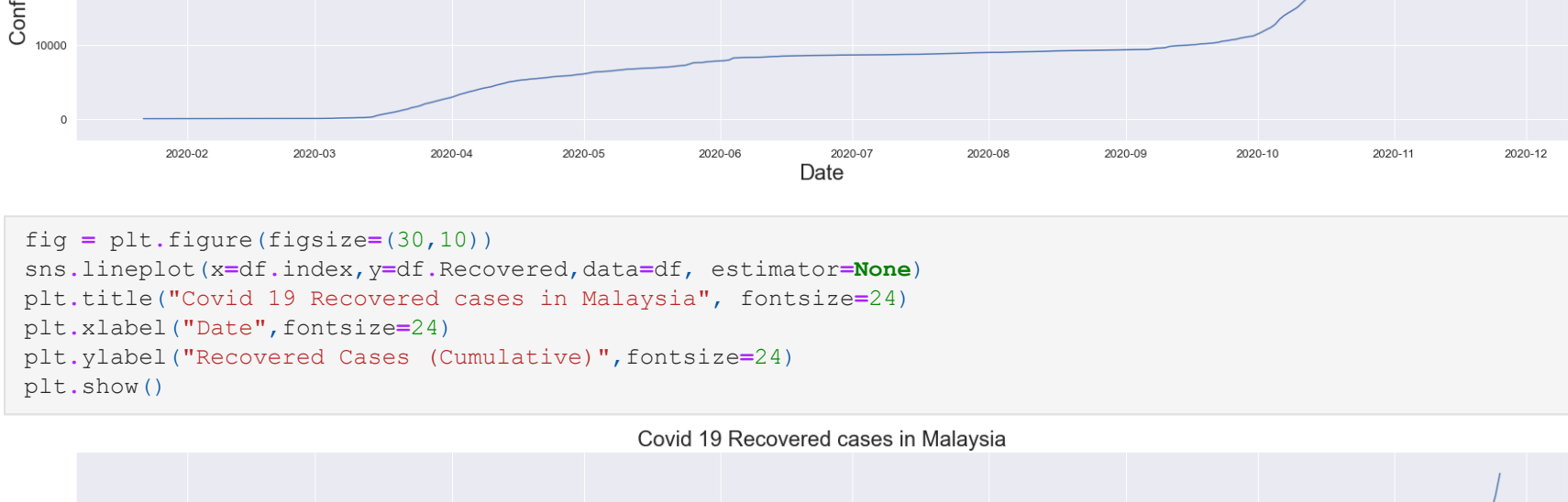
```
In [22]: df.set_index("Date", inplace=True)
```

```
Out [22]:
```

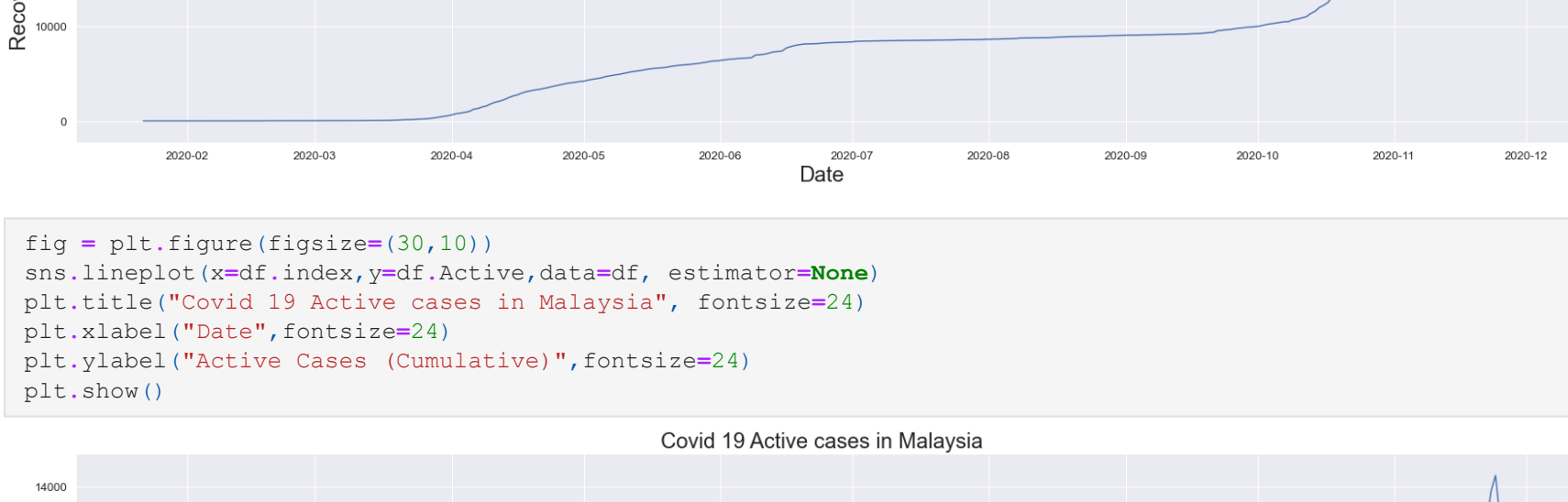
	Confirmed	Recovered	Deaths	Active
Date				
2020-01-22	0	0	0	0
2020-01-23	0	0	0	0
2020-01-24	0	0	0	0
2020-01-25	3	0	0	3
2020-01-26	4	0	0	4
...
2020-11-21	53679	40493	332	12854
2020-11-22	54775	41597	335	12843
2020-11-23	56659	42480	337	13842
2020-11-24	58847	44153	341	14353
2020-11-25	59817	46501	345	12971

309 rows x 4 columns

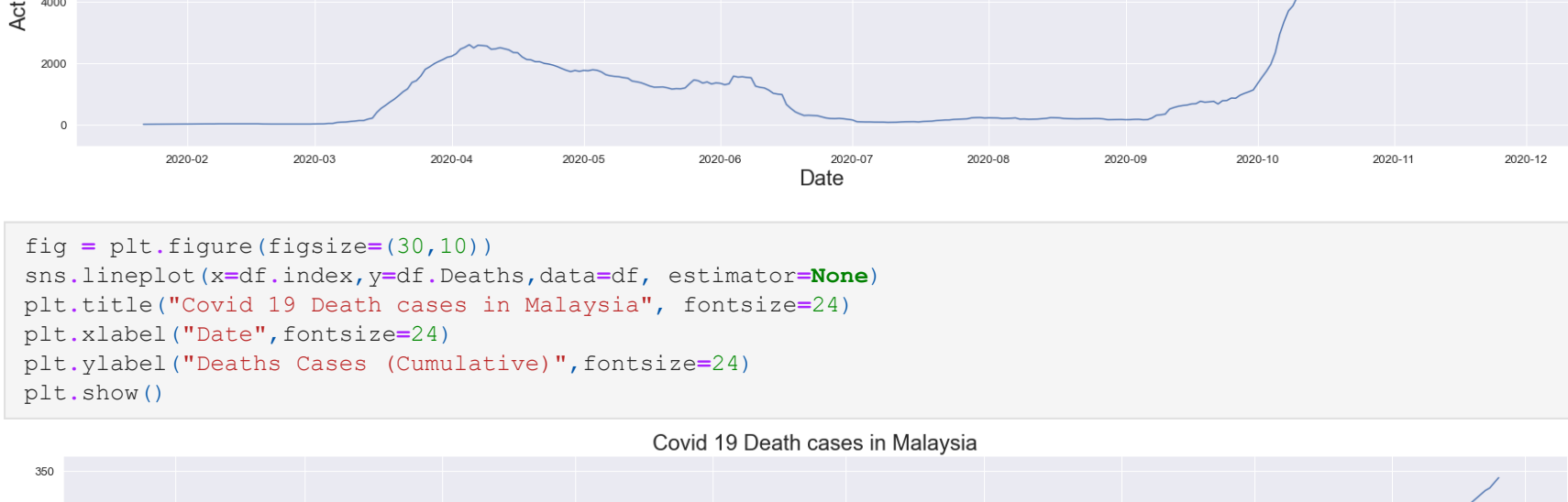
```
In [24]: fig = plt.figure(figsize=(30,10))
sns.lineplot(x=df.index, y=df.Confirmed, data=df, estimator=None)
plt.title("Covid 19 confirmed cases in Malaysia", fontsize=24)
plt.xlabel("Date", fontsize=24)
plt.ylabel("Confirmed Cases (Cumulative)", fontsize=24)
plt.show()
```



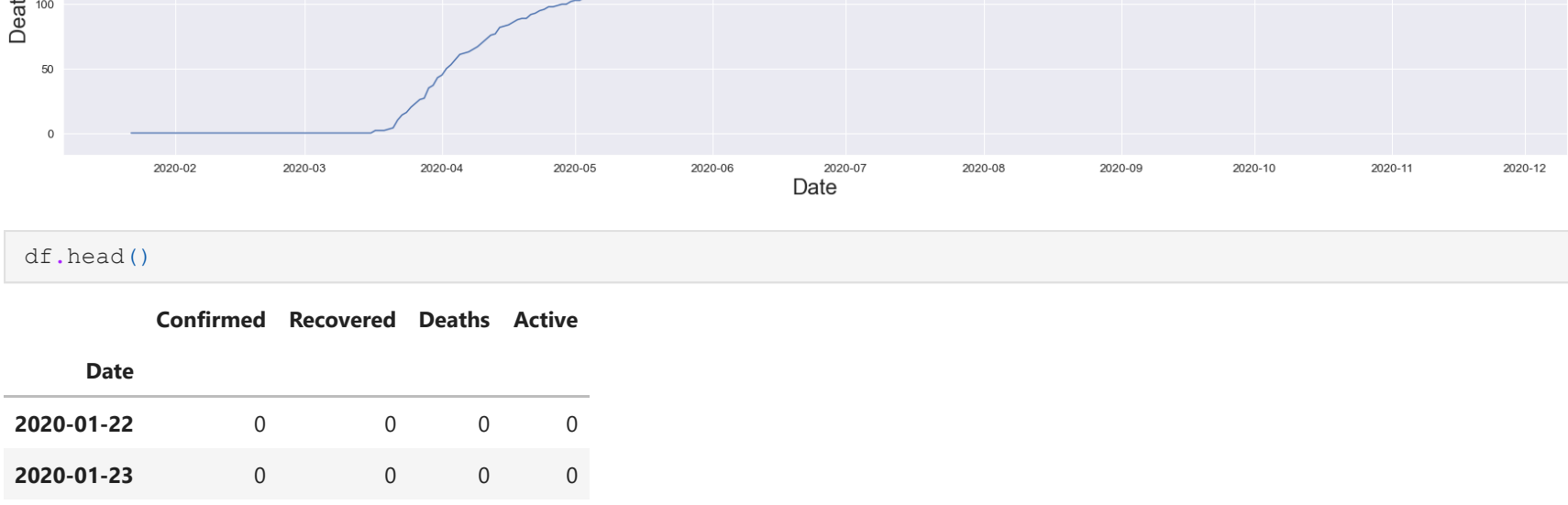
```
In [25]: fig = plt.figure(figsize=(30,10))
sns.lineplot(x=df.index, y=df.Recovered, data=df, estimator=None)
plt.title("Covid 19 Recovered cases in Malaysia", fontsize=24)
plt.xlabel("Date", fontsize=24)
plt.ylabel("Recovered Cases (Cumulative)", fontsize=24)
plt.show()
```



```
In [26]: fig = plt.figure(figsize=(30,10))
sns.lineplot(x=df.index, y=df.Active, data=df, estimator=None)
plt.title("Covid 19 Active cases in Malaysia", fontsize=24)
plt.xlabel("Date", fontsize=24)
plt.ylabel("Active Cases (Cumulative)", fontsize=24)
plt.show()
```



```
In [27]: fig = plt.figure(figsize=(30,10))
sns.lineplot(x=df.index, y=df.Deaths, data=df, estimator=None)
plt.title("Covid 19 Death cases in Malaysia", fontsize=24)
plt.xlabel("Date", fontsize=24)
plt.ylabel("Deaths Cases (Cumulative)", fontsize=24)
plt.show()
```



```
Out [27]:
```

	Confirmed	Recovered	Deaths	Active
Date				
2020-01-22	0	0	0	0
2020-01-23	0	0	0	0
2020-01-24	0	0	0	0
2020-01-25	3	0	0	3
2020-01-26	4	0	0	4
...
2020-11-21	53679	40493	332	12854
2020-11-22	54775	41597	335	12843
2020-11-23	56659	42480	337	13842
2020-11-24	58847	44153	341	14353
2020-11-25	59817	46501	345	12971

309 rows x 4 columns

```
In [28]: df.head()
```

```
Out [28]:
```

	Confirmed	Recovered	Deaths	Active
Date				
2020-01-22	0	0	0	0
2020-01-23	0	0	0	0
2020-01-24	0	0	0	0
2020-01-25	3	0	0	3
2020-01-26	4	0	0	4
...
2020-11-21	53679	40493	332	12854
2020-11-22	54775	41597	335	12843
2020-11-23	56659	42480	337	13842
2020-11-24	58847	44153	341	14353
2020-11-25	59817	46501	345	12971

309 rows x 4 columns

```
In [29]: df["ConfirmedDiff"] = df["Confirmed"].diff()
```

```
In [30]: df["DeathsDiff"] = df["Deaths"].diff()
```

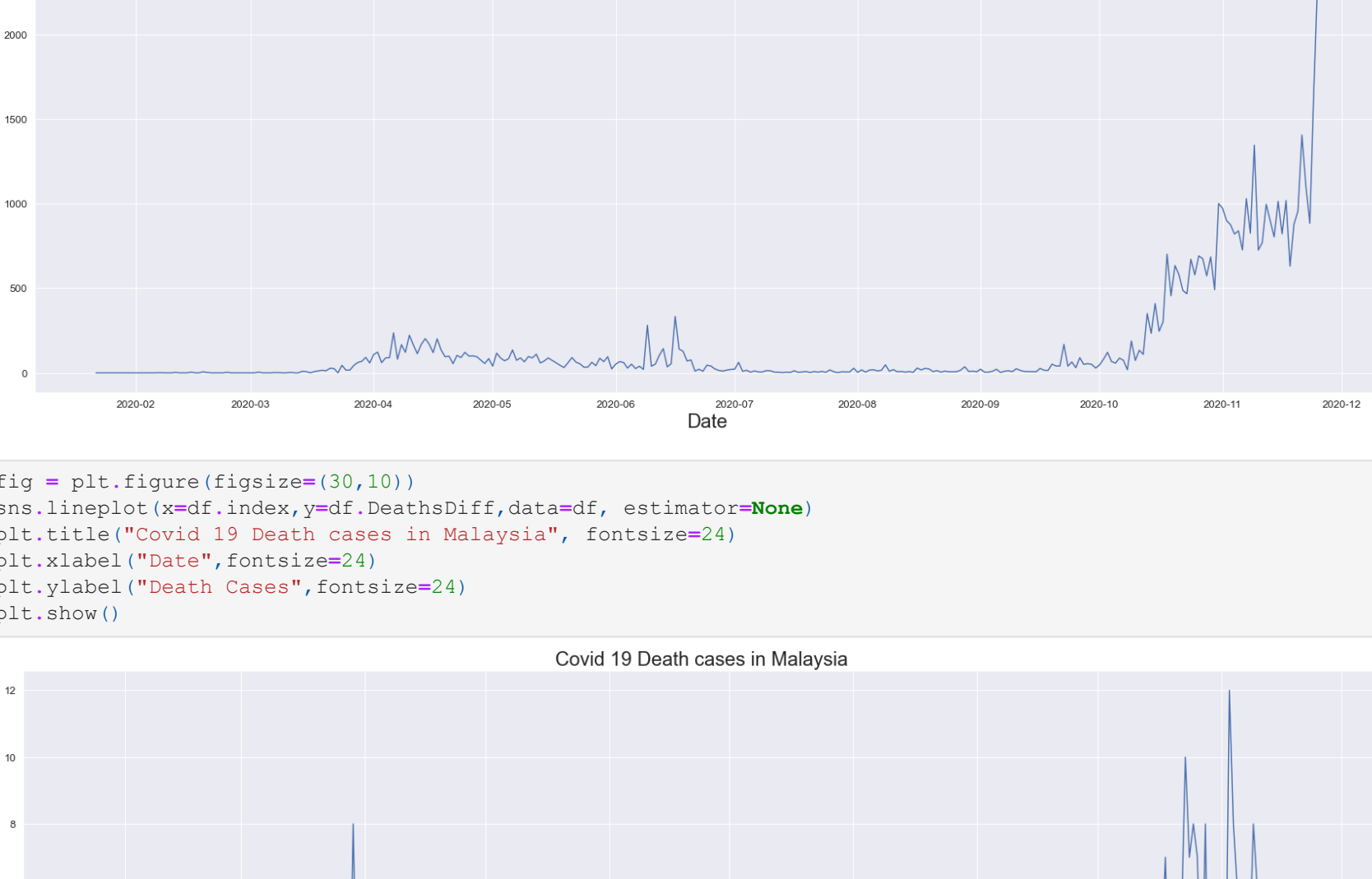
```
In [31]: df["RecoveredDiff"] = df["Recovered"].diff()
```

```
In [32]: df["ActiveDiff"] = df["Active"].diff()
```

```
Out [32]:
```

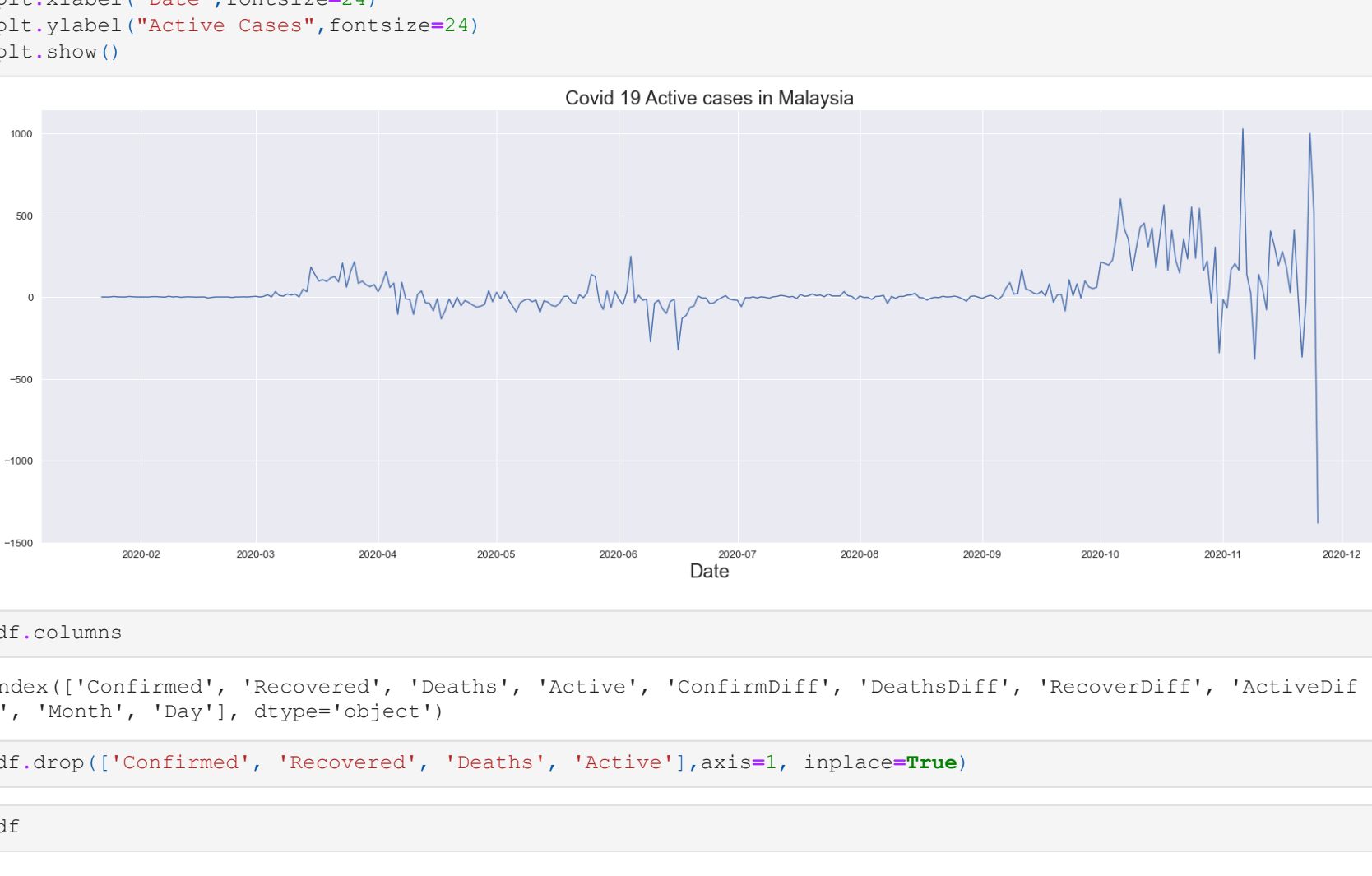
	Confirmed	Recovered	Deaths	Active	ConfirmedDiff	DeathsDiff	RecoveredDiff	ActiveDiff
Date								
2020-01-22	0	0	0	0	NaN	NaN	NaN	NaN
2020-01-23	0	0	0	0	0.0	0.0	0.0	0.0
2020-01-24	0	0	0	0	0.0	0.0	0.0	0.0
2020-01-25	3	0	0	3	3.0	0.0	0.0	3.0
2020-01-26	4	0	0	4	1.0	0.0	0.0	1.0
...
2020-11-21	53679	40493	332	12854	10410	3.0	1405.0	-367.0
2020-11-22	54775	41597	335	12843	1096.0	3.0	1104.0	-11.0
2020-11-23	56659	42480	337	13842	1884.0	2.0	883.0	999.0
2020-11-24	58847	44153	341	14353	2188.0	4.0	1673.0	511.0
2020								

Covid 19 Recovered cases in Malaysia



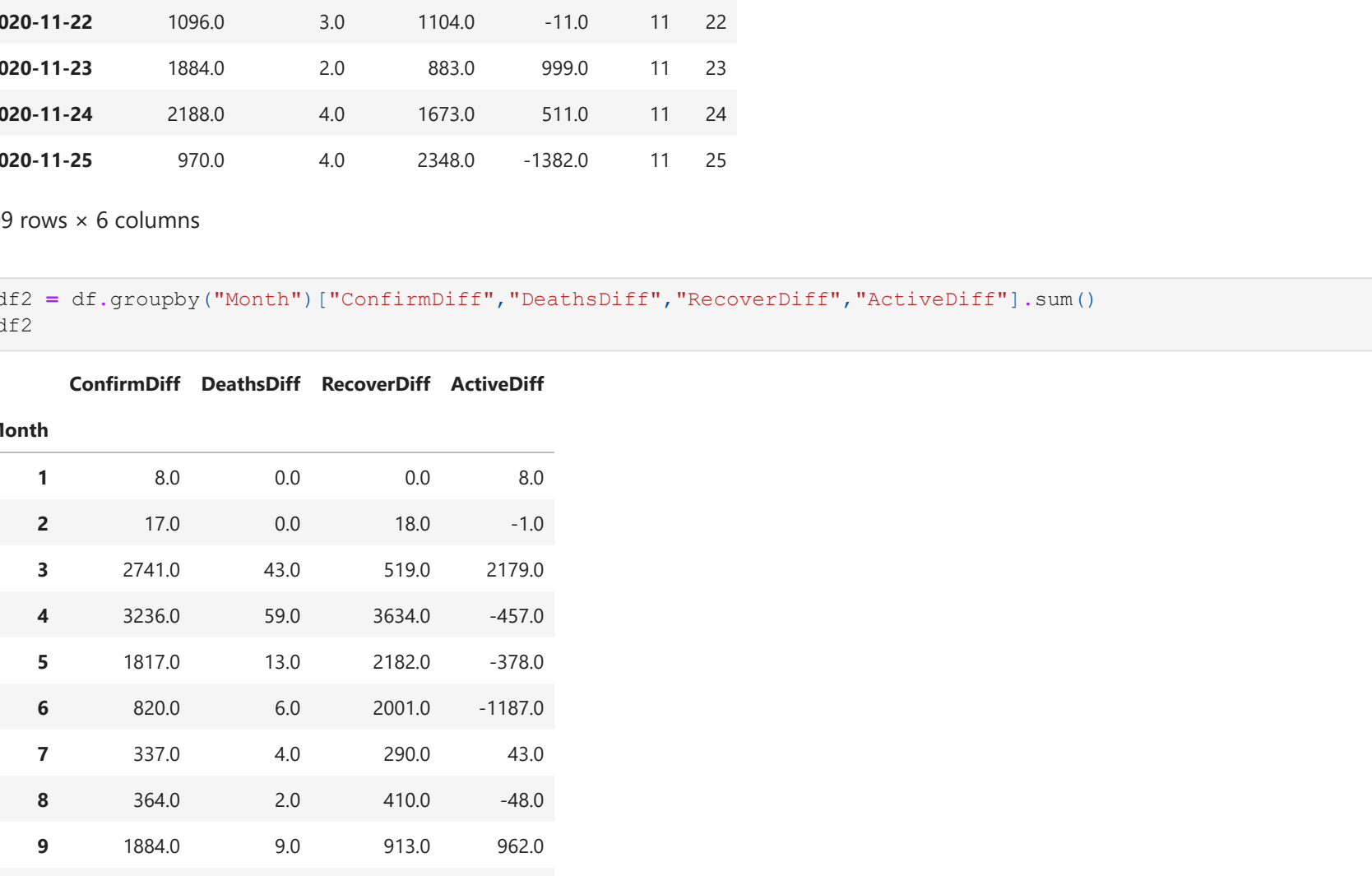
```
In [50]: fig=plt.figure(figsize=(30,10))
sns.lineplot(x=df2.index,y=df2.DeathsDiff,data=df, estimator=None)
plt.title("Covid 19 Death cases in Malaysia", fontsize=24)
plt.xlabel("Date",fontsize=24)
plt.ylabel("Death Cases",fontsize=24)
plt.show()
```

Covid 19 Death cases in Malaysia



```
In [51]: fig=plt.figure(figsize=(30,10))
sns.lineplot(x=df2.index,y=df2.ActiveDiff,data=df, estimator=None)
plt.title("Covid 19 Active cases in Malaysia", fontsize=24)
plt.xlabel("Date",fontsize=24)
plt.ylabel("Active Cases",fontsize=24)
plt.show()
```

Covid 19 Active cases in Malaysia



```
In [52]: df.columns
```

```
Out[52]: Index(['Confirmed', 'Recovered', 'Deaths', 'Active', 'ConfirmDiff', 'DeathsDiff', 'RecoverDiff', 'ActiveDiff', 'Month', 'Day'], dtype='object')
```

```
In [53]: df.drop(['Confirmed', 'Recovered', 'Deaths', 'Active'],axis=1, inplace=True)
```

```
In [54]: df
```

```
Out[54]:
```

	ConfirmDiff	DeathsDiff	RecoverDiff	ActiveDiff	Month	Day
Date						
2020-01-22	0.0	0.0	0.0	0.0	1	22
2020-01-23	0.0	0.0	0.0	0.0	1	23
2020-01-24	0.0	0.0	0.0	0.0	1	24
2020-01-25	3.0	0.0	0.0	3.0	1	25
2020-01-26	1.0	0.0	0.0	1.0	1	26
...
2020-11-21	1041.0	3.0	1405.0	-367.0	11	21
2020-11-22	1096.0	3.0	1104.0	-11.0	11	22
2020-11-23	1884.0	2.0	883.0	999.0	11	23
2020-11-24	2188.0	4.0	1673.0	511.0	11	24
2020-11-25	970.0	4.0	2349.0	-1382.0	11	25

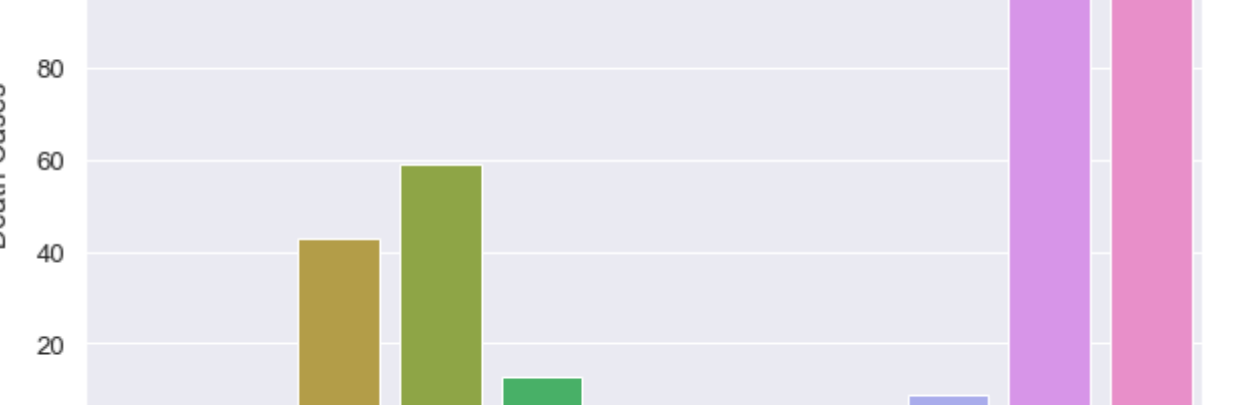
309 rows x 6 columns

```
In [55]: df2=df.groupby("Month")["ConfirmDiff","DeathsDiff","RecoverDiff","ActiveDiff"].sum()
```

```
Out[55]:
```

Month	ConfirmDiff	DeathsDiff	RecoverDiff	ActiveDiff
1	8.0	0.0	0.0	8.0
2	17.0	0.0	18.0	-1.0
3	2741.0	43.0	519.0	2179.0
4	3236.0	59.0	3634.0	-457.0
5	1817.0	13.0	2182.0	-378.0
6	820.0	6.0	2001.0	-1187.0
7	337.0	4.0	290.0	43.0
8	364.0	2.0	410.0	-48.0
9	1884.0	9.0	913.0	962.0
10	20324.0	113.0	11281.0	8930.0
11	28269.0	96.0	25253.0	2920.0

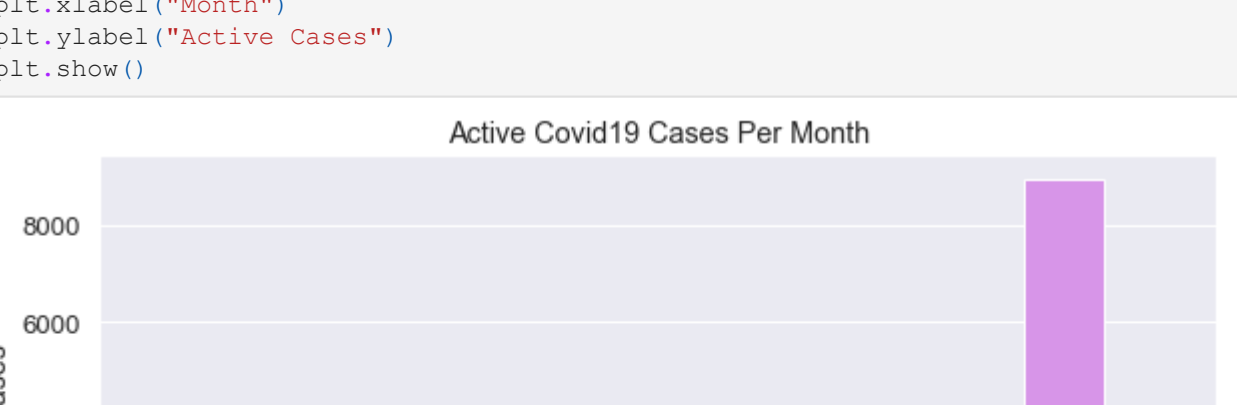
```
In [56]: plt.figure(figsize=(10,5))
sns.barplot(x=df2.index,y=df2.ConfirmDiff, data=df2, ci=None, estimator=sum)
plt.title("Confirmed Covid19 Cases Per Month")
plt.xlabel("Month")
plt.ylabel("Confirmed Cases")
plt.show()
```



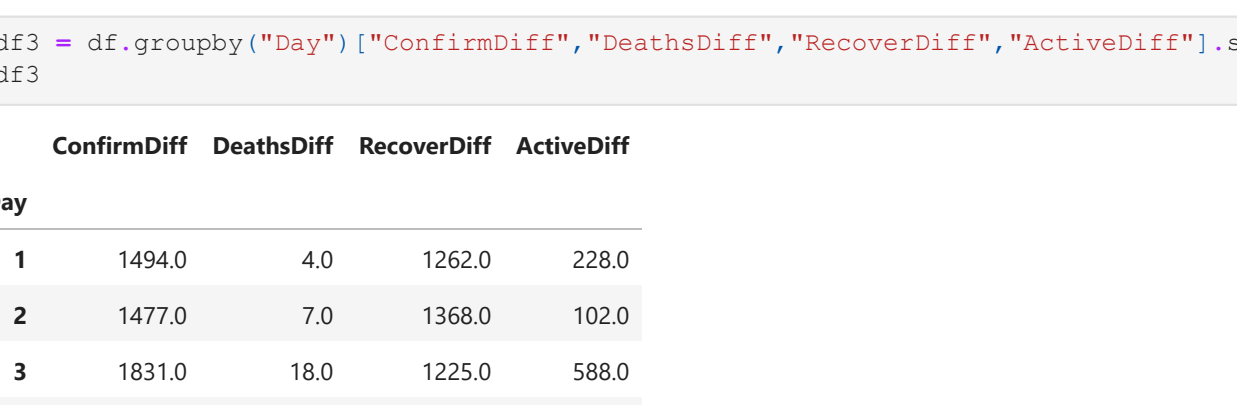
```
In [57]: plt.figure(figsize=(10,5))
sns.barplot(x=df2.index,y=df2.DeathsDiff, data=df2, ci=None, estimator=sum)
plt.title("Deaths Covid19 Cases Per Month")
plt.xlabel("Month")
plt.ylabel("Death Cases")
plt.show()
```



```
In [58]: plt.figure(figsize=(10,5))
sns.barplot(x=df2.index,y=df2.RecoverDiff, data=df2, ci=None, estimator=sum)
plt.title("Recovered Covid19 Cases Per Month")
plt.xlabel("Month")
plt.ylabel("Recovered Cases")
plt.show()
```



```
In [59]: plt.figure(figsize=(10,5))
sns.barplot(x=df3.index,y=df3.ActiveDiff, data=df3, ci=None, estimator=sum)
plt.title("Active Covid19 Cases Per Month")
plt.xlabel("Month")
plt.ylabel("Active Cases")
plt.show()
```

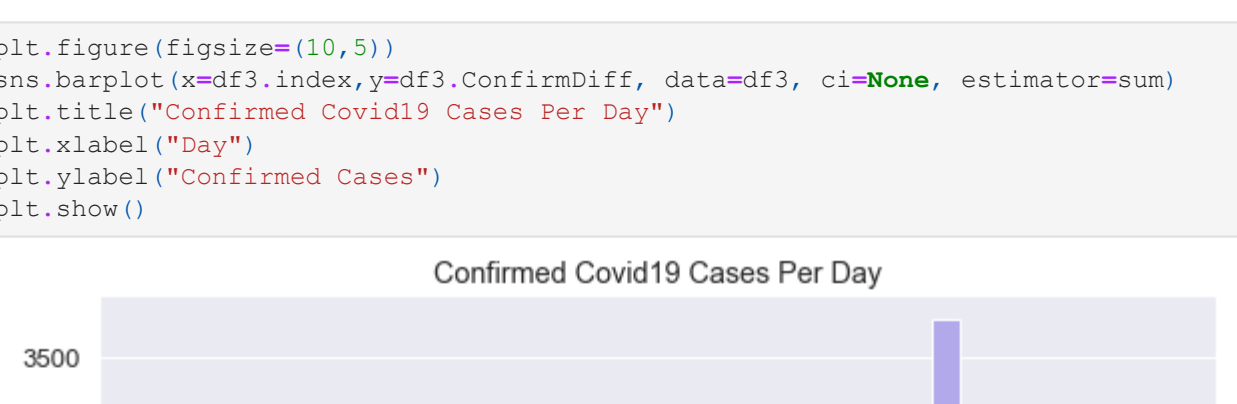


```
In [60]: df3=df.groupby("Day")["ConfirmDiff","DeathsDiff","RecoverDiff","ActiveDiff"].sum()
```

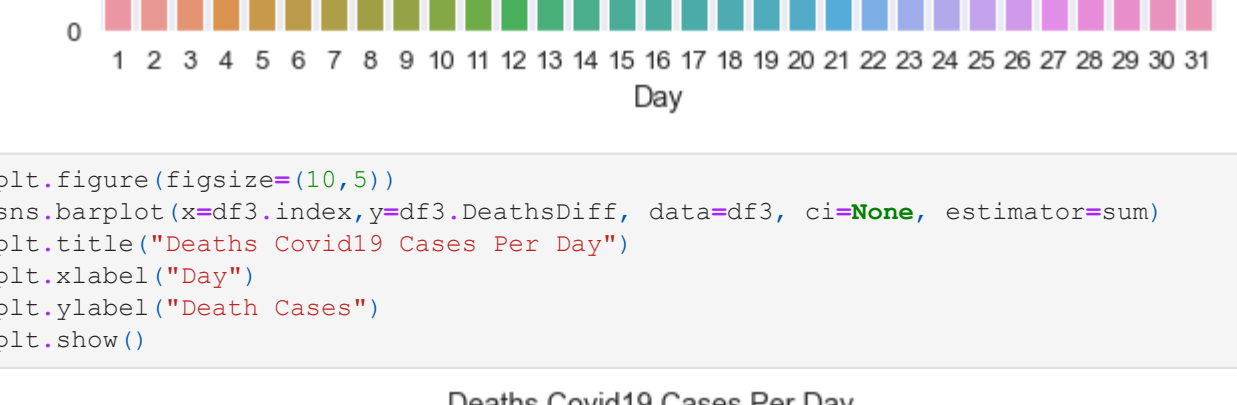
```
Out[60]:
```

Day	ConfirmDiff	DeathsDiff	RecoverDiff	ActiveDiff
1	1494.0	4.0	1262.0	228.0
2	1477.0	7.0	1368.0	102.0
3	1831.0	18.0	1225.0	588.0
4	1845.0	12.0	1114.0	719.0
5	1703.0	12.0	1163.0	528.0
6	2718.0	9.0	1233.0	1476.0
7	1988.0	4.0	1327.0	657.0
8	1578.0	11.0	1182.0	385.0
9	1557.0	17.0	2029.0	-489.0
10	1513.0	13.0	1210.0	290.0
11	1887.0	9.0	1234.0	644.0
12	1767.0	7.0	1451.0	309.0
13	2296.0	9.0	1631.0	656.0
14	2088.0	13.0	1370.0	705.0
15	2206.0	7.0	1657.0	444.0
16	2098.0	12.0	1757.0	429.0
17	2338.0	13.0	1762.0	563.0
18	1874.0	14.0	1699.0	161.0
19	2443.0	11.0	1627.0	805.0
20	2116.0	7.0	1885.0	234.0
21	2150.0	13.0	2226.0	-109.0
22	2315.0	16.0	1978.0	321.0
23	3094.0	21.0	1559.0	1514.0
24	3775.0	14.0	2663.0	1098.0
25	2340.0	18.0	3169.0	-847.0
26	1808.0	11.0	994.0	803.0
27	1193.0	7.0	936.0	250.0
28	1183.0	10.0	901.0	272.0
29	1129.0	8.0	953.0	168.0
30	1159.0	8.0	824.0	325.0
31	874.0	8.0	1114.0	-248.0

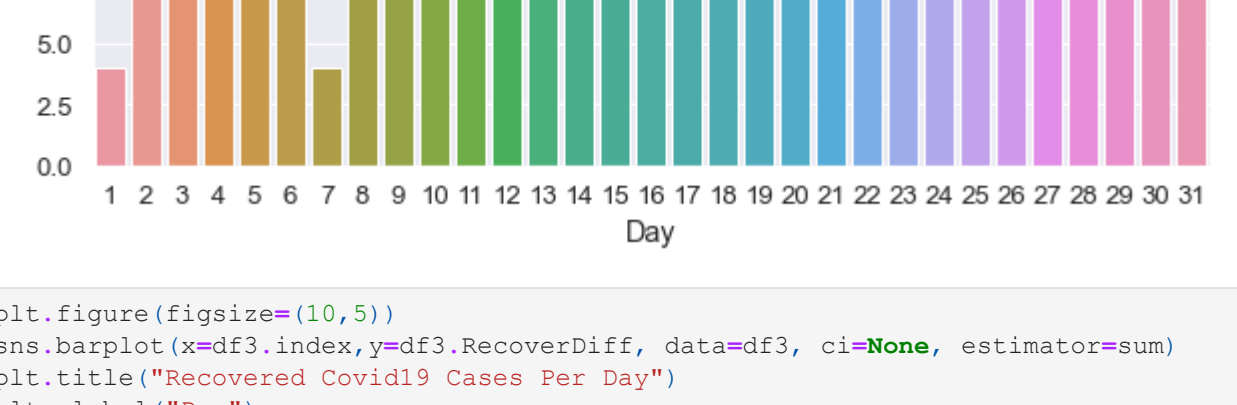
```
In [61]: plt.figure(figsize=(10,5))
sns.barplot(x=df3.index,y=df3.ConfirmDiff, data=df3, ci=None, estimator=sum)
plt.title("Confirmed Covid19 Cases Per Day")
plt.xlabel("Day")
plt.ylabel("Confirmed Cases")
plt.show()
```



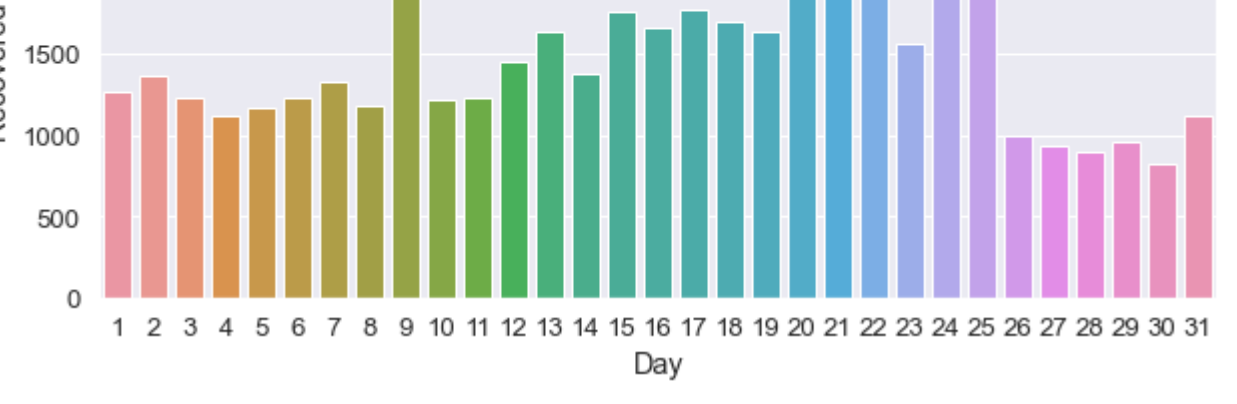
```
In [62]: plt.figure(figsize=(10,5))
sns.barplot(x=df3.index,y=df3.DeathsDiff, data=df3, ci=None, estimator=sum)
plt.title("Deaths Covid19 Cases Per Day")
plt.xlabel("Day")
plt.ylabel("Death Cases")
plt.show()
```



```
In [63]: plt.figure(figsize=(10,5))
sns.barplot(x=df3.index,y=df3.RecoverDiff, data=df3, ci=None, estimator=sum)
plt.title("Recovered Covid19 Cases Per Day")
plt.xlabel("Day")
plt.ylabel("Recovered Cases")
plt.show()
```



```
In [64]: plt.figure(figsize=(10,5))
sns.barplot(x=df3.index,y=df3.ActiveDiff, data=df3, ci=None, estimator=sum)
plt.title("Active Covid19 Cases Per Day")
plt.xlabel("Day")
plt.ylabel("Active Cases")
plt.show()
```



```
In [65]: df.head()
```

```
Out[65]:
```

	ConfirmDiff	DeathsDiff	RecoverDiff	ActiveDiff	Month	Day
Date						
2020-01-22	0.0	0.0	0.0	0.0	1	22
2020-01-23	0.0	0.0	0.0	0.0	1	23
2020-01-24	0.0	0.0	0.0	0.0	1	24
2020-01-25	3.0	0.0	0.0	3.0	1	25
2020-01-26	1.0	0.0	0.0	1.0	1	26

Correlation

```
In [66]: df.corr()
```

```
Out[66]:
```

	ConfirmDiff	DeathsDiff	RecoverDiff	ActiveDiff	Month	Day
ConfirmDiff	1.000000	0.641523	0.875959	0.517088	0.604269	0.002108
DeathsDiff	0.641523	1.000000	0.561145	0.326380	0.341420	0.035592
RecoverDiff	0.875959	0.561145	1.000000	0.040080	0.545945	0.027499
ActiveDiff	0.517088	0.326380	0.040080	1.000000	0.283484	-0.044819
Month	0.604269	0.341420	0.545945	0.283484	1.000000	-0.098535
Day	0.002108	0.035592	0.027499	-0.044819	-0.098535	1.000000

```
In [67]: plt.figure(figsize=(16,9))
sns.heatmap(df.corr(),cmap="coolwarm",annot=True,fmt='.2f',linewidths=2)
plt.show()
```

