

Python初級數據分析員證書

(六) 數據分析及可視化專案

13. 數據分析專案

Demo5 - Covid-19 Statistics

Visualization

Review

- Statistics
- Hypothesis testing
- Algebra
- Linear regression
- Propositional logic
- Python
- R
- SQL
- Pandas, NumPy, SciPy
- Data Visualization, Matplotlib, Seaborn, Plotly
- Dashboard Visualization, Business Intelligence
- Storytelling



13. 數據分析專案 Data Analysis Project – Demo5

Chapter Summary

- Covid-19 Data Visualization
- Data Import
- Color Palettes
- Create plot function
- px.choropleth with time slider
- pd.melt()
- px.treemap

Data import

In this chapter, we are going to visualise practical Covid-19 statistics.

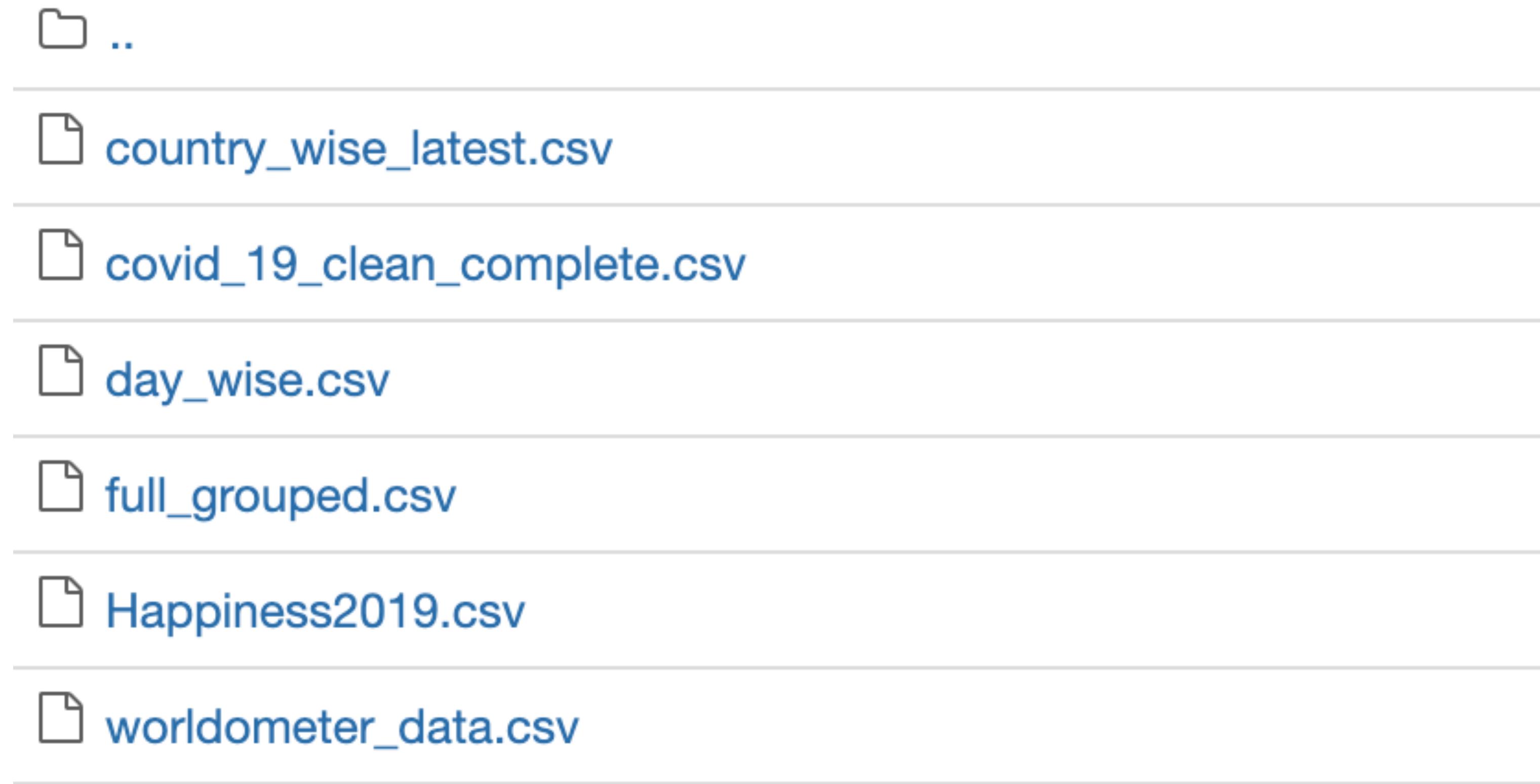
Source of data: <https://github.com/CSSEGISandData/COVID-19>

The data was collected from WHO and each countries health department and finalised with students in Johns Hopkins University and enthusiasts on GitHub.



Data Import

There are 6 csv datafile in Covid-19 folder. Firstly we import them and have a glance.



Data Import

```

1 import pandas as pd
2 import numpy as np
3
4 import plotly.express as px

```

Data Import

```

1 full_table = pd.read_csv('../covid-19/covid_19_clean_complete.csv')
2 full_table

```

	Province/State	Country/Region	Lat	Long	Date	Confirmed	Deaths	Recovered	Active	WHO Region
0	NaN	Afghanistan	33.939110	67.709953	2020-01-22	0	0	0	0	Eastern Mediterranean
1	NaN	Albania	41.153300	20.168300	2020-01-22	0	0	0	0	Europe
2	NaN	Algeria	28.033900	1.659600	2020-01-22	0	0	0	0	Africa
3	NaN	Andorra	42.506300	1.521800	2020-01-22	0	0	0	0	Europe
4	NaN	Angola	-11.202700	17.873900	2020-01-22	0	0	0	0	Africa

Data Import

Pd.to_datetime is usually a must for time series data, otherwise Pandas will treat it as string.

```
1 full_grouped = pd.read_csv('../covid-19/full_grouped.csv')
2 full_grouped['Date'] = pd.to_datetime(full_grouped['Date'])
3 full_grouped
```

	Date	Country/Region	Confirmed	Deaths	Recovered	Active	New cases	New deaths	New recovered	WHO Region
0	2020-01-22	Afghanistan	0	0	0	0	0	0	0	Eastern Mediterranean
1	2020-01-22	Albania	0	0	0	0	0	0	0	Europe
2	2020-01-22	Algeria	0	0	0	0	0	0	0	Africa
3	2020-01-22	Andorra	0	0	0	0	0	0	0	Europe
4	2020-01-22	Angola	0	0	0	0	0	0	0	Africa

Data Import

```

1 day_wise = pd.read_csv('../covid-19/day_wise.csv')
2 day_wise['Date'] = pd.to_datetime(day_wise['Date'])
3 day_wise

```

	Date	Confirmed	Deaths	Recovered	Active	New cases	New deaths	New recovered	Deaths / 100 Cases	Recovered / 100 Cases	Deaths / 100 Recovered	No. of countries
0	2020-01-22	555	17	28	510	0	0	0	3.06	5.05	60.71	6
1	2020-01-23	654	18	30	606	99	1	2	2.75	4.59	60.00	8
2	2020-01-24	941	26	36	879	287	8	6	2.76	3.83	72.22	9
3	2020-01-25	1434	42	39	1353	493	16	3	2.93	2.72	107.69	11
4	2020-01-26	2118	56	52	2010	684	14	13	2.64	2.46	107.69	13

Data Import

```

1 country_wise = pd.read_csv('../covid-19/country_wise_latest.csv')
2 country_wise = country_wise.replace('', np.nan).fillna(0)
3 country_wise

```

	Country/Region	Confirmed	Deaths	Recovered	Active	New cases	New deaths	New recovered	Deaths / 100 Cases	Recovered / 100 Cases	Deaths / 100 Recovered	Confirmed last week	1 week change	1 week % increase	WHO Re
0	Afghanistan	36263	1269	25198	9796	106	10	18	3.50	69.49	5.04	35526	737	2.07	Ea Mediterra
1	Albania	4880	144	2745	1991	117	6	63	2.95	56.25	5.25	4171	709	17.00	Eu
2	Algeria	27973	1163	18837	7973	616	8	749	4.16	67.34	6.17	23691	4282	18.07	A
3	Andorra	907	52	803	52	10	0	0	5.73	88.53	6.48	884	23	2.60	Eu
4	Angola	950	41	242	667	18	1	0	4.32	25.47	16.94	749	201	26.84	A

Data Import

```

1 worldometer_data = pd.read_csv('../covid-19/worldometer_data.csv')
2 worldometer_data = worldometer_data.replace('', np.nan).fillna(0)
3 worldometer_data

```

	Country/Region	Continent	Population	TotalCases	NewCases	TotalDeaths	NewDeaths	TotalRecovered	NewRecovered	ActiveCases	Serious,Critical	CasesPerCapita
0	USA	North America	3.311981e+08	5032179	0.0	162804.0	0.0	2576668.0	0.0	2292707.0	18296.0	18296.0
1	Brazil	South America	2.127107e+08	2917562	0.0	98644.0	0.0	2047660.0	0.0	771258.0	8318.0	8318.0
2	India	Asia	1.381345e+09	2025409	0.0	41638.0	0.0	1377384.0	0.0	606387.0	8944.0	8944.0
3	Russia	Europe	1.459409e+08	871894	0.0	14606.0	0.0	676357.0	0.0	180931.0	2300.0	2300.0
4	South Africa	Africa	5.938157e+07	538184	0.0	9604.0	0.0	387316.0	0.0	141264.0	539.0	539.0

Color Palettes

Coloring on graph is much important than many people think. Especially when we preparing visual aids for presentation and public report. In general you may follow this simple theory.

Red Excitement Strength Love Energy	Orange Confidence Success Bravery Sociability	Yellow Creativity Happiness Warmth Cheer	Green Nature Healing Freshness Quality	Blue Trust Peace Loyalty Competence
Pink Compassion Sincerity Sophstication Sweet	Purple Royalty Luxury Spirituality Ambition	Brown Dependable Rugged Trustworthy Simple	Black Formality Dramatic Sophistication Security	White Clean Simplicity Innocence Honest

Color Palettes

The screenshot shows the Color Hunt website interface. At the top, there's a navigation bar with a logo, a search bar containing 'Vintage', and an 'Add tag' button. Below the navigation, there's a sidebar with filters: 'New', 'Popular', 'Random', 'Collection', 'Pastel', 'Vintage' (which is highlighted), 'Retro', 'Neon', 'Gold', 'Light', 'Dark', 'Warm', 'Cold', 'Summer', 'Fall', and 'Winter'. The main area displays two color palettes. The first palette, under the 'Vintage' tab, has a total of 450 likes and was posted 2 days ago. It consists of five horizontal stripes: orange, red, dark red, teal, and light blue. The second palette, under the 'Dark' tab, has 497 likes and was posted 6 days ago. It consists of five horizontal stripes: pink (#FFEBEB), teal, blue, grey, and orange.

There are many color palettes (with
rgb#) online resources for you to
pick series of theme color.

<https://colorhunt.co/>

Colour Palettes

In our Covid-19 visualization, it is relatively simple. We fix these four colour on our decision, in case of picking by plotly default.

```
1 cnf, dth, rec, act = '#393e46', '#ff2e63', '#fe9801', '#21bf73'  
  
1 temp = day_wise[['Date', 'Deaths', 'Recovered', 'Active']].tail(1)  
2 temp = temp.melt(id_vars="Date", value_vars=['Active', 'Deaths', 'Recovered'])  
3 fig = px.treemap(temp, path=["variable"], values="value", height=225,  
4                   color_discrete_sequence=[act, rec, dth])  
5 fig.data[0].textinfo = 'label+text+value'  
6 fig.show()
```



Create a plot function

DRY, in programming norm, Don't Repeat Yourself. If you need to plot same kind of plots more than once, you should create a function.

```
1 def plot_map(df, col, pal):  
2     df = df[df[col]>0]  
3     fig = px.choropleth(df, locations="Country/Region", locationmode='country names',  
4                           color=col, hover_name="Country/Region",  
5                           title=col, hover_data=[col], color_continuous_scale=pal)  
6     #     fig.update_layout(coloraxis_showscale=False)  
7     fig.show()
```

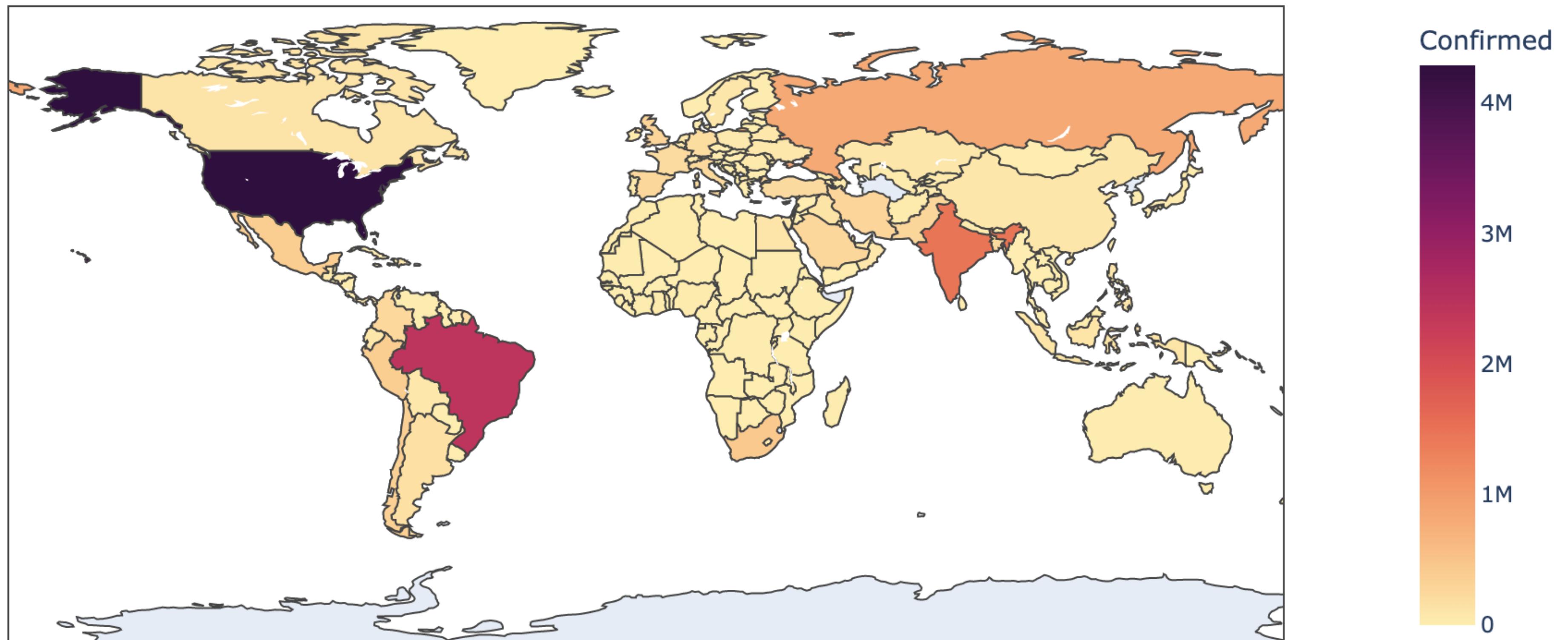
col is column in df; pal is color palettes series

For more choice on continuous color series, go to:
<https://plotly.com/python/builtin-colorscales/>

Plotting confirmed case

```
1 plot_map(country_wise, 'Confirmed', 'matter')
```

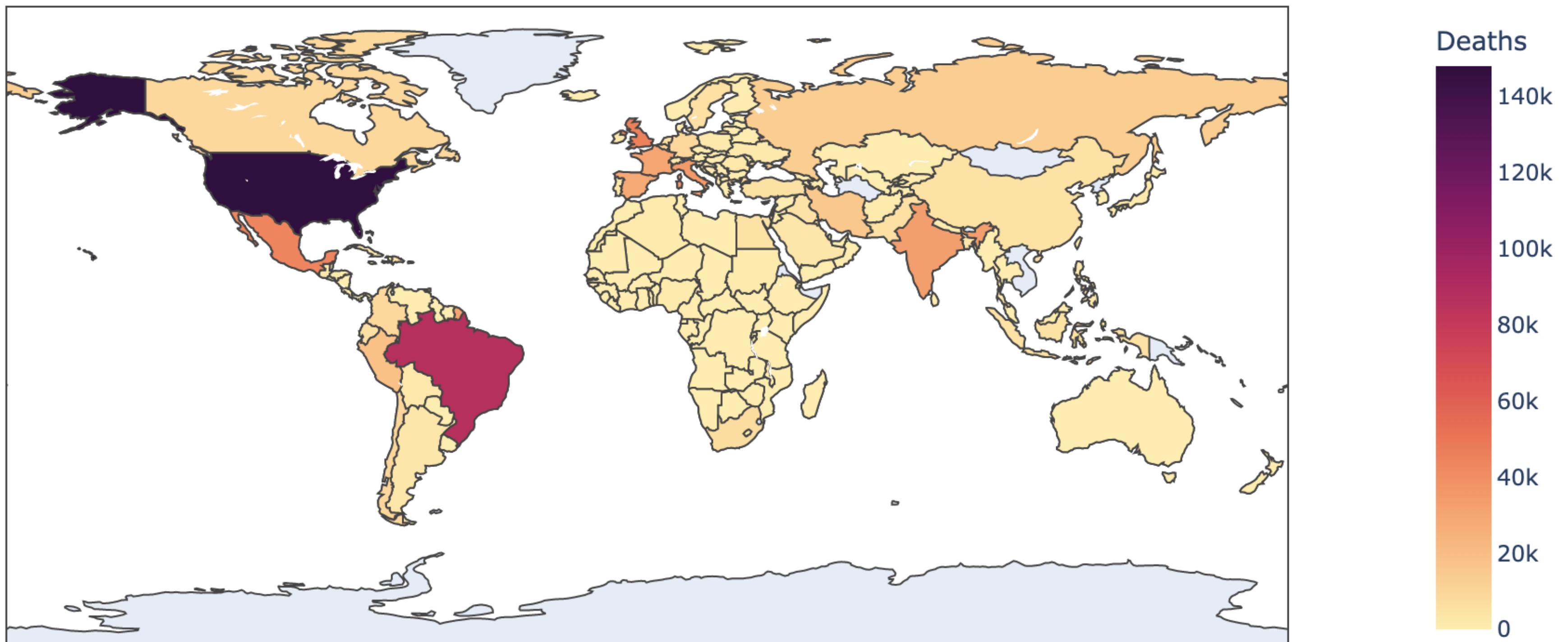
Confirmed



Plotting death cases

```
1 plot_map(country_wise, 'Deaths', 'matter')
```

Deaths



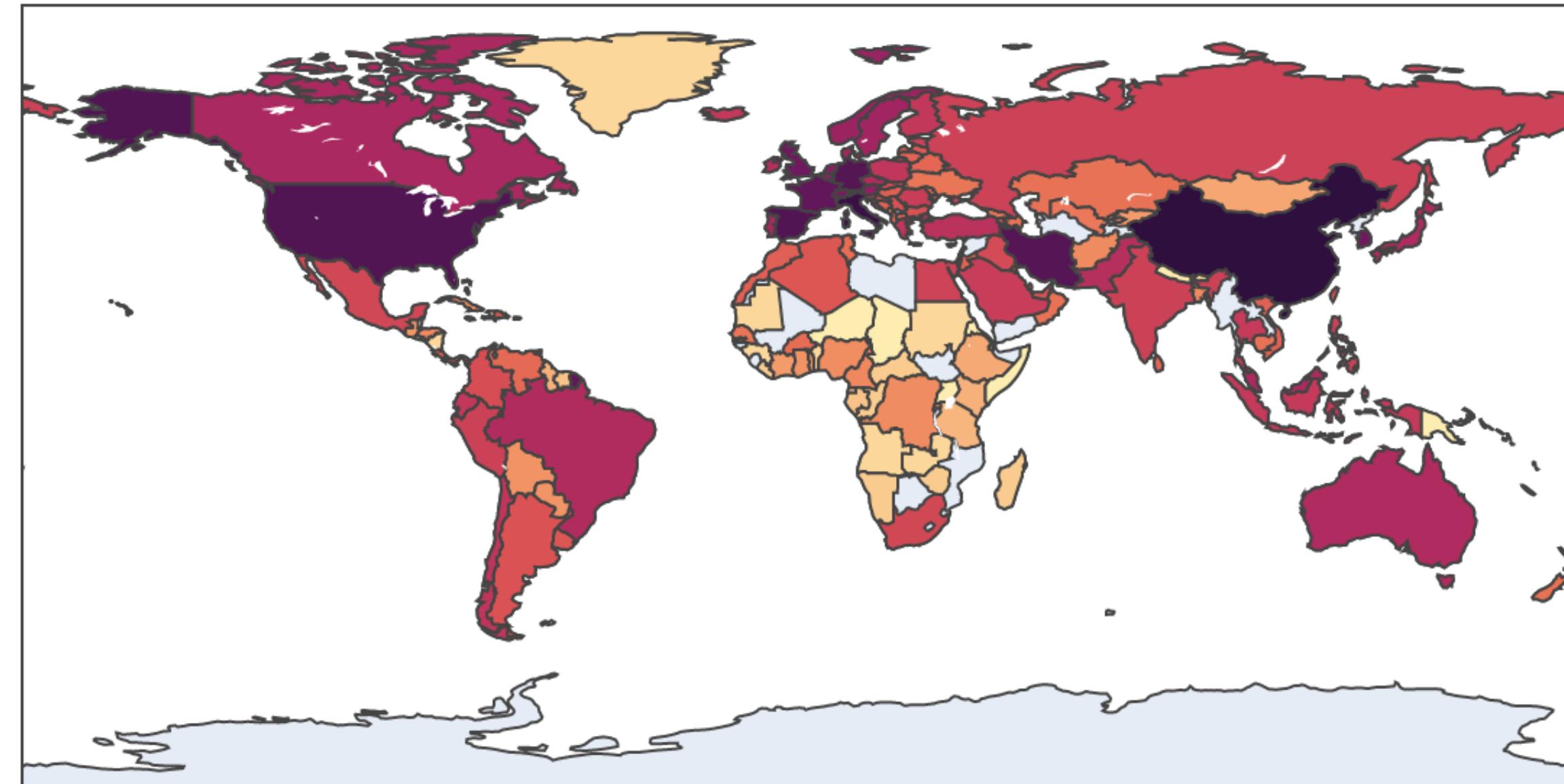
Choropleth with time slider

Covid-19 outbreak spread rapidly, we may plot an interactive graph to review how it did.

```
1 fig = px.choropleth(full_grouped, locations="Country/Region",
2                     color=np.log(full_grouped["Confirmed"]),
3                     locationmode='country names', hover_name="Country/Region",
4                     animation_frame=full_grouped["Date"].dt.strftime('%Y-%m-%d'),
5                     title='Cases over time', color_continuous_scale=px.colors.sequential.matter)
6 fig.update(layout_coloraxis_showscale=False)
7 fig.show()
```

Choropleth with time slider

Cases over time



Case over time

```
1 temp = full_grouped.groupby('Date')[['Recovered', 'Deaths', 'Active']].sum().reset_index()  
2 temp = temp.melt(id_vars="Date", value_vars=['Recovered', 'Deaths', 'Active'],  
3 var_name='Case', value_name='Count')  
4 temp.head()
```

	Date	Case	Count
0	2020-01-22	Recovered	28
1	2020-01-23	Recovered	30
2	2020-01-24	Recovered	36
3	2020-01-25	Recovered	39
4	2020-01-26	Recovered	52

`pd.melt(frame, id_vars=None, value_vars=None, var_name=None, value_name='value', col_level=None, ignore_index=True)`

Unpivot a DataFrame from wide to long format,
optionally leaving identifiers set.

Case over time

```
1 fig = px.area(temp, x="Date", y="Count", color='Case', height=600, width=700,  
2                     title='Cases over time', color_discrete_sequence = [rec, dth, act])  
3 fig.update_layout(xaxis_rangeslider_visible=True)  
4 fig.show()
```



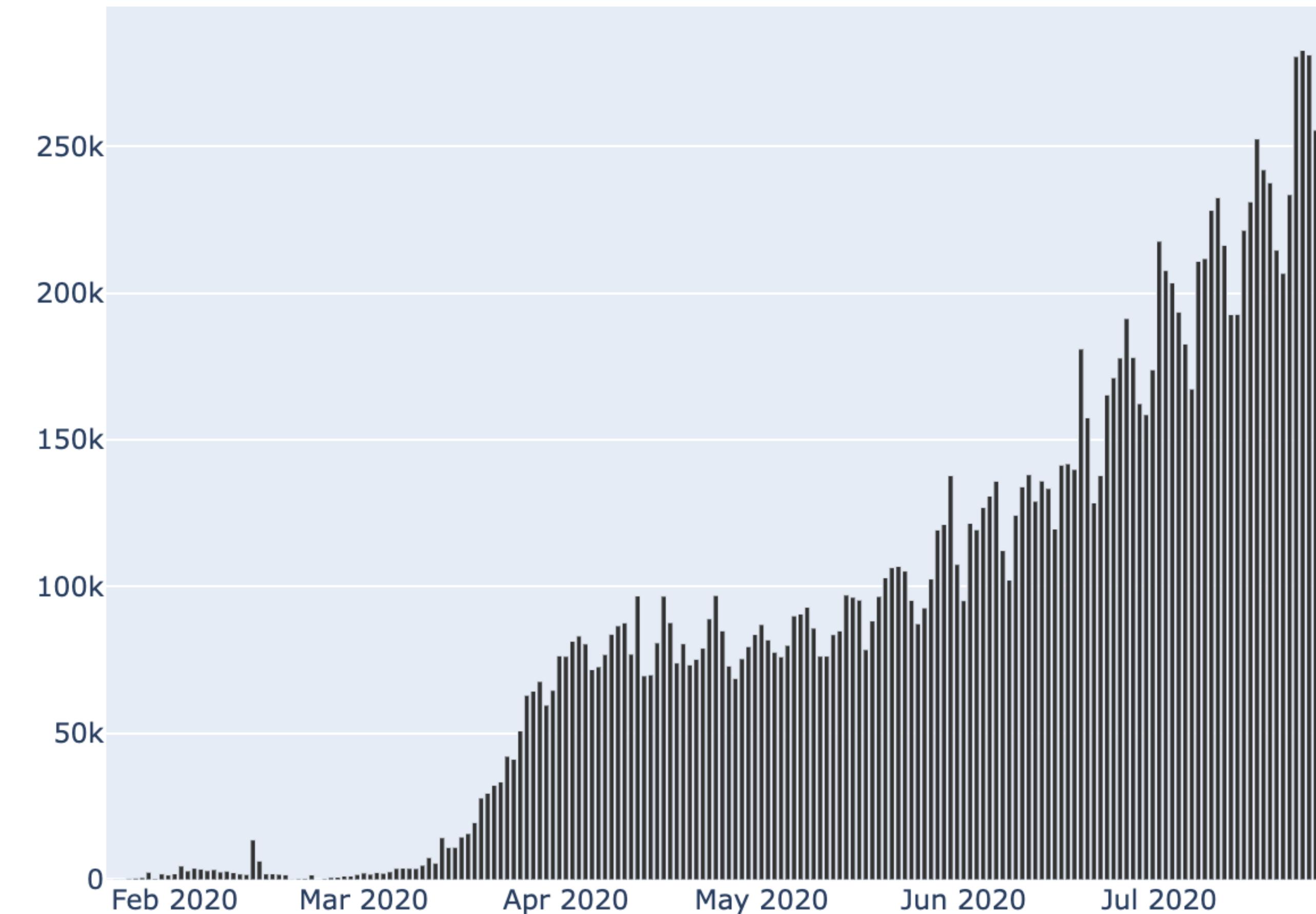
Create function to plot bar and line chart

```
1 def plot_daywise(col, hue):
2     fig = px.bar(day_wise, x="Date", y=col, width=700, color_discrete_sequence=[hue])
3     fig.update_layout(title=col, xaxis_title="", yaxis_title="")
4     fig.show()
```

```
1 def plot_daywise_line(col, hue):
2     fig = px.line(day_wise, x="Date", y=col, width=700, color_discrete_sequence=[hue])
3     fig.update_layout(title=col, xaxis_title="", yaxis_title="")
4     fig.show()
```

```
1 plot_daywise('New cases', '#333333')
```

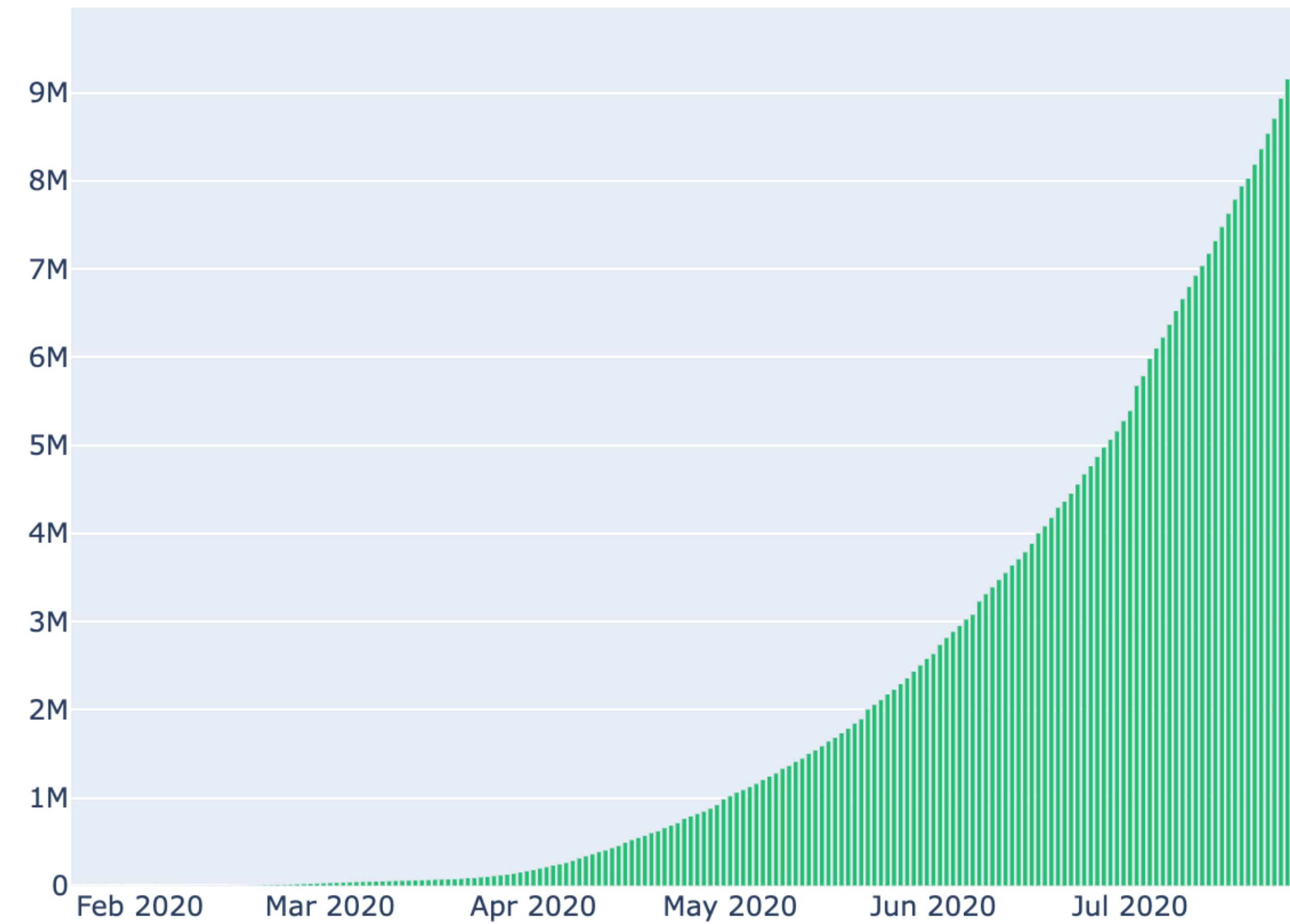
New cases



13. 數據分析專案 Data Analysis Project – Demo5

```
1 plot_daywise('Recovered', rec)
```

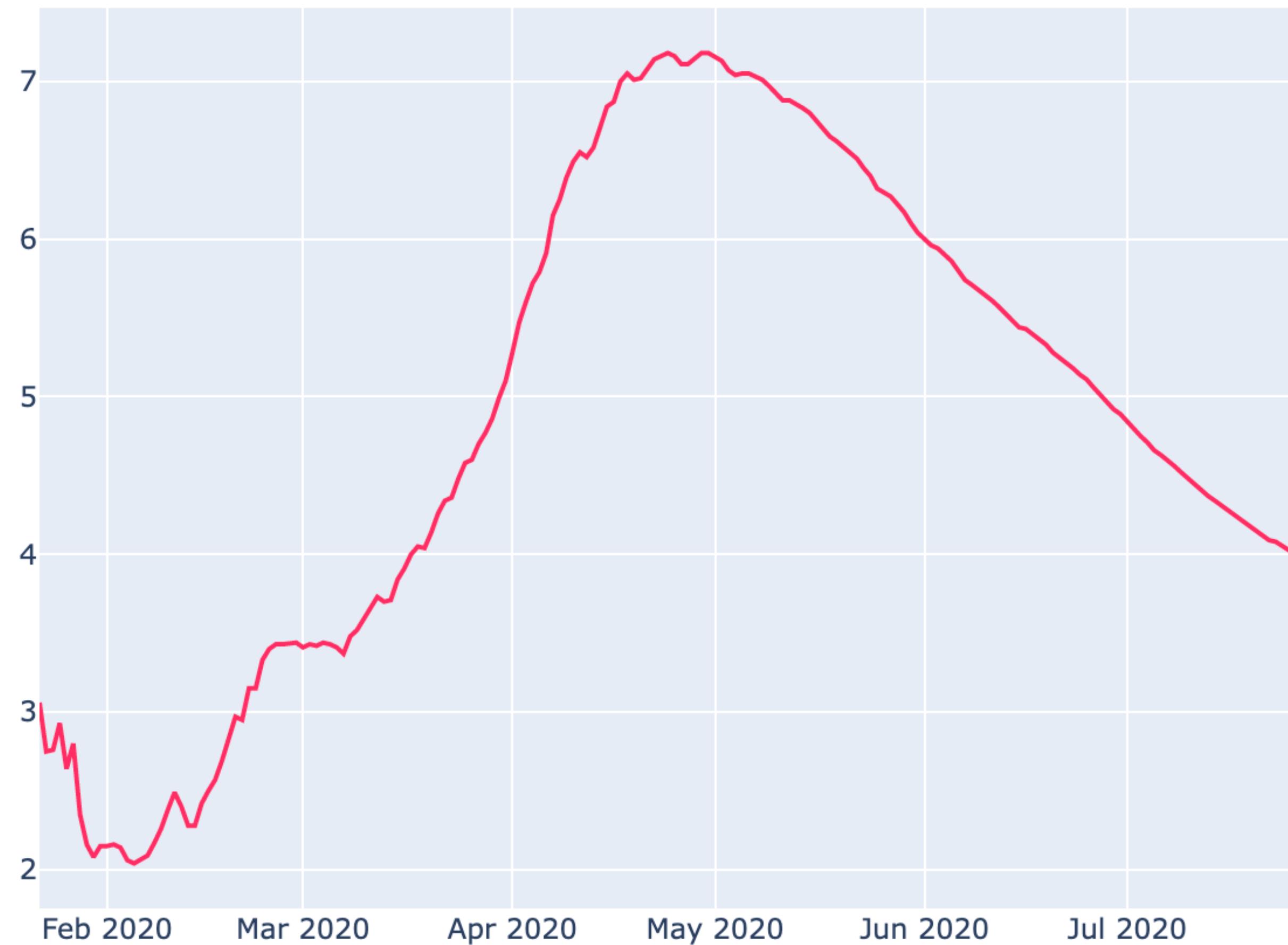
Recovered



13. 數據分析專案 Data Analysis Project – Demo5

```
1 | plot_daywise_line('Deaths / 100 Cases', dth)
```

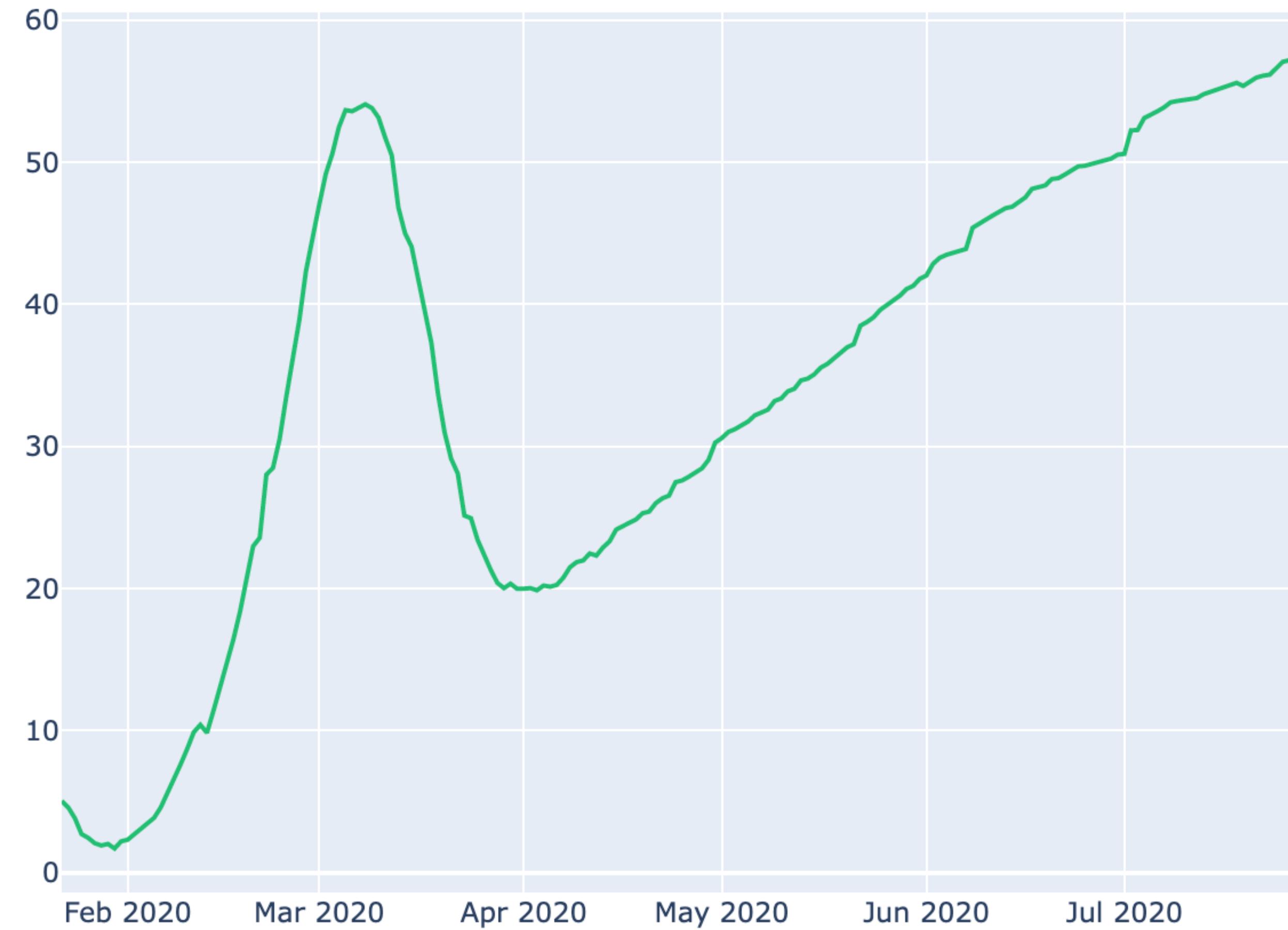
Deaths / 100 Cases



13. 數據分析專案 Data Analysis Project – Demo5

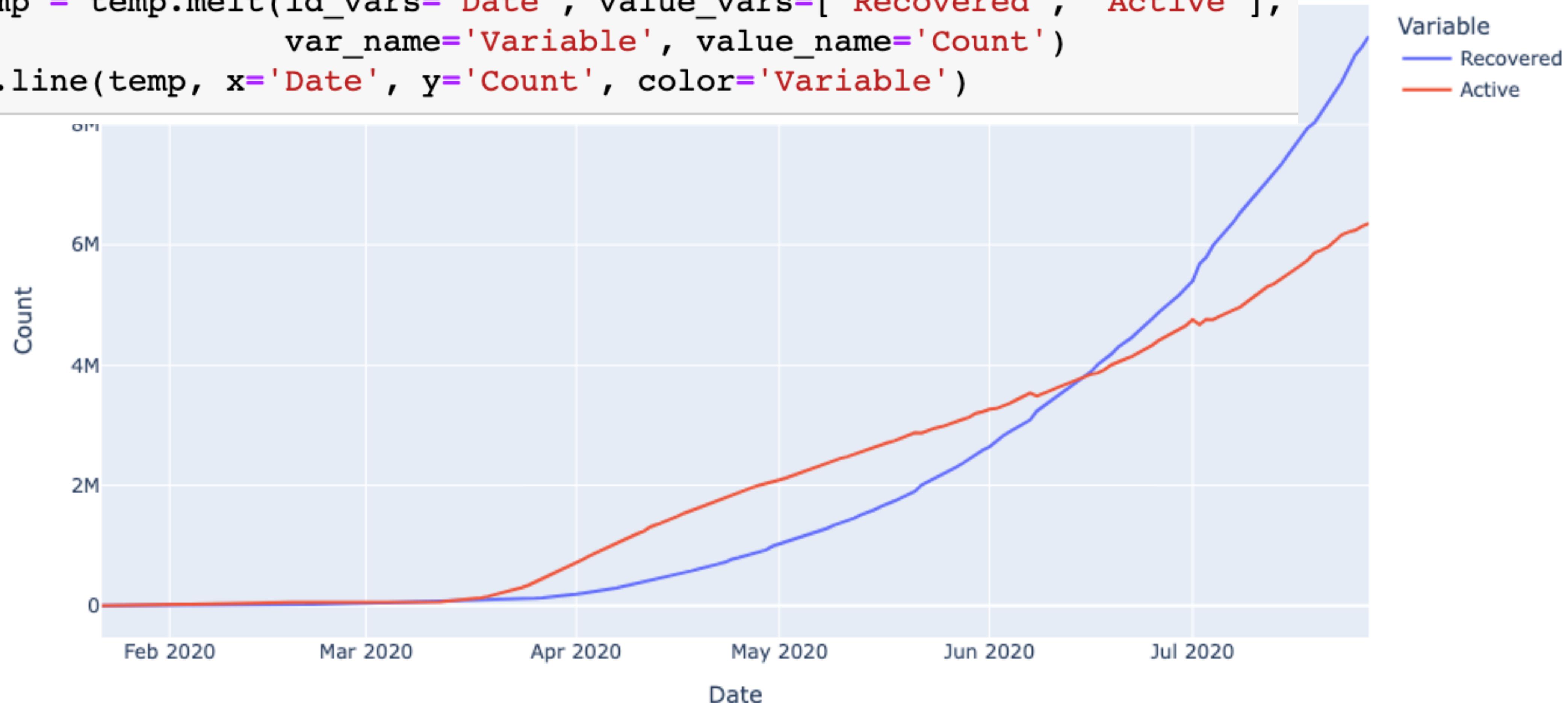
```
1 plot_daywise_line('Recovered / 100 Cases', rec)
```

Recovered / 100 Cases



Plotting recover and active case

```
1 temp = day_wise[['Date', 'Recovered', 'Active']]  
2 temp = temp.melt(id_vars='Date', value_vars=['Recovered', 'Active'],  
3                   var_name='Variable', value_name='Count')  
4 px.line(temp, x='Date', y='Count', color='Variable')
```



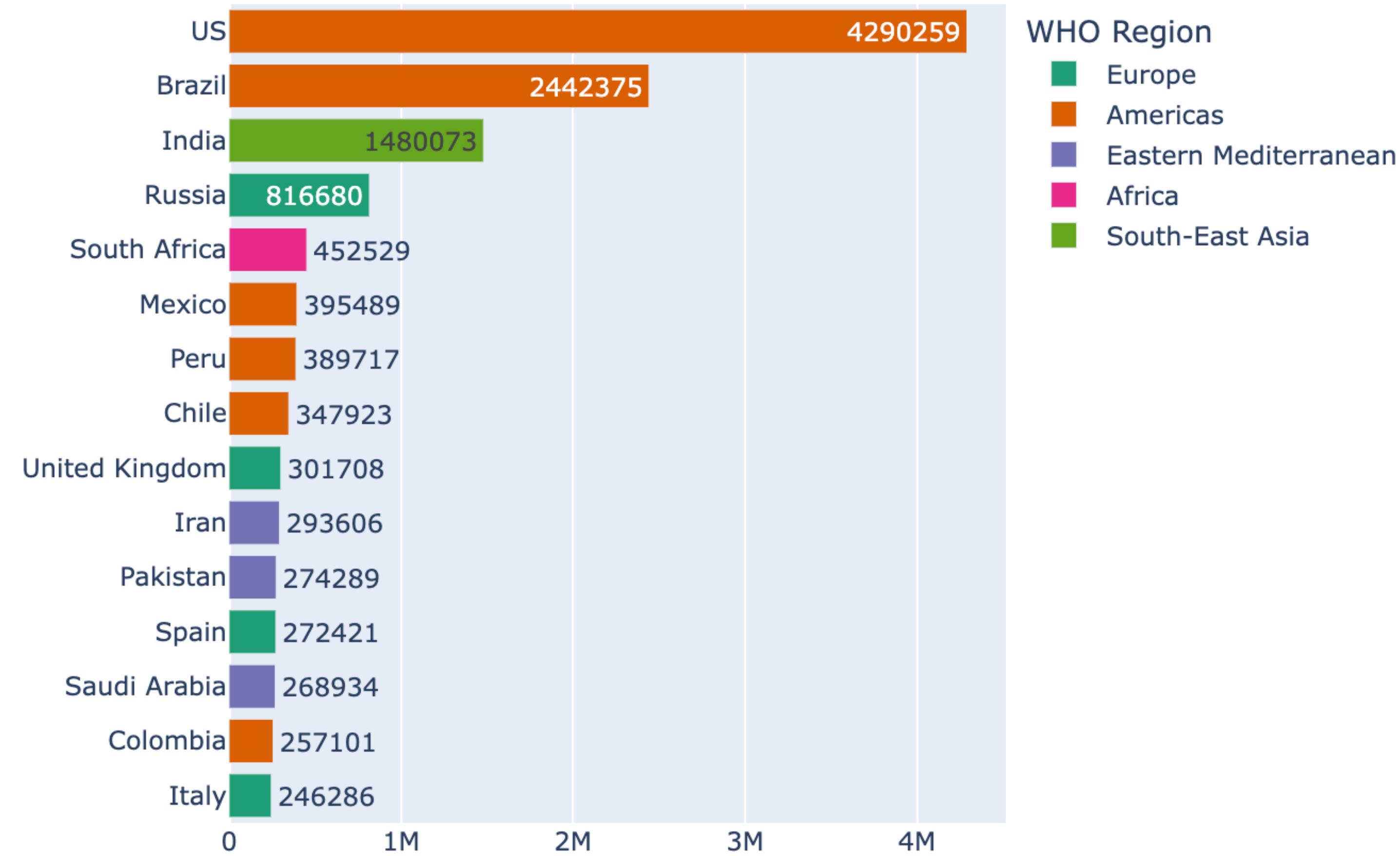
Horizontal Bar

```
1 def plot_hbar(df, col, n, hover_data=[]):
2     fig = px.bar(df.sort_values(col).tail(n),
3                   x=col, y="Country/Region", color='WHO Region',
4                   text=col, orientation='h', width=700, hover_data=hover_data,
5                   color_discrete_sequence = px.colors.qualitative.Dark2)
6     fig.update_layout(title=col, xaxis_title="", yaxis_title="",
7                       yaxis_categoryorder = 'total ascending',
8                       uniformtext_minsize=8, uniformtext_mode='hide')
9     fig.show()
```

Horizontal Bar

```
1 plot_hbar(country_wise, 'Confirmed', 15)
```

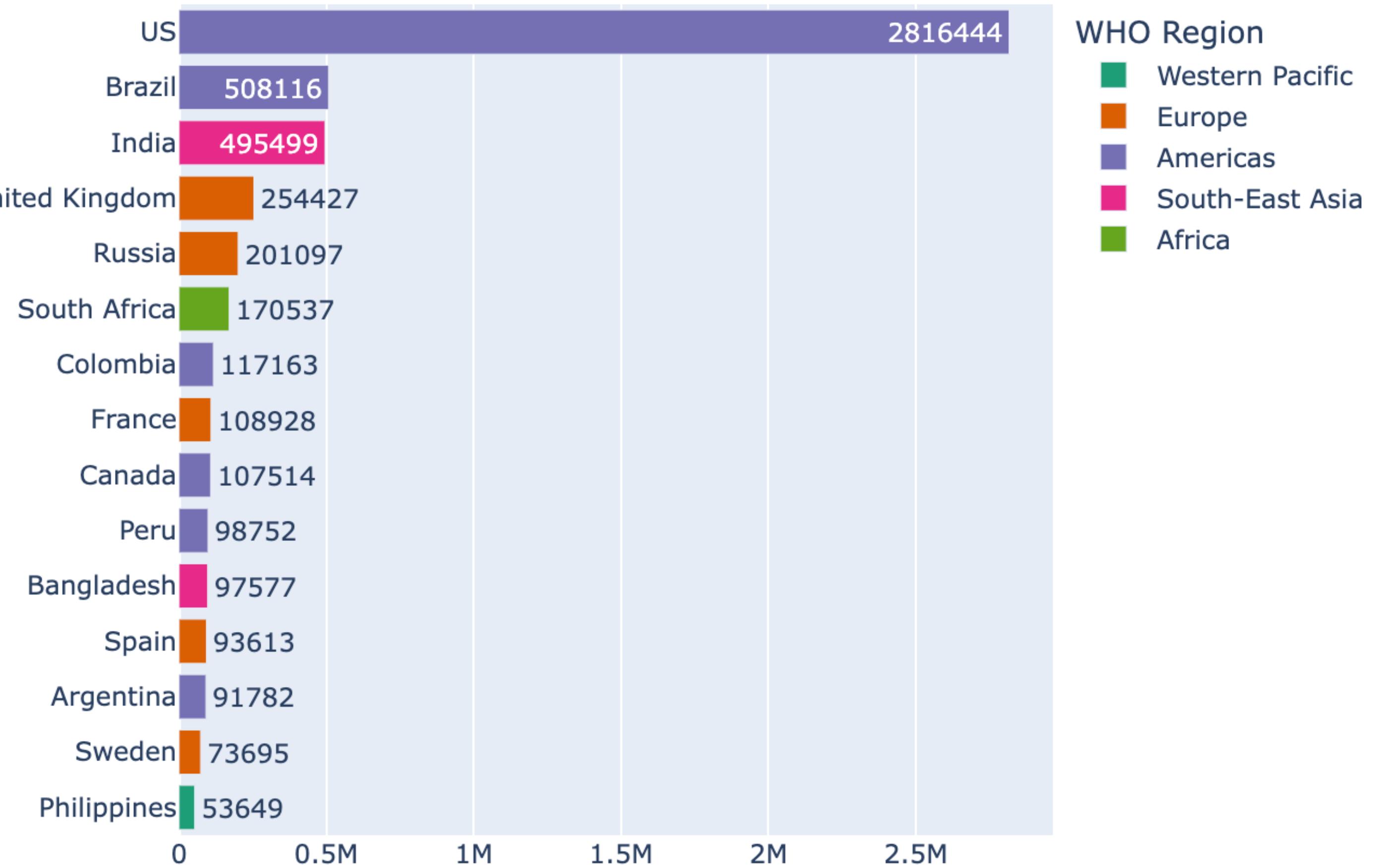
Confirmed



Horizontal Bar

```
1 plot_hbar(country_wise, 'Active', 15)
```

Active



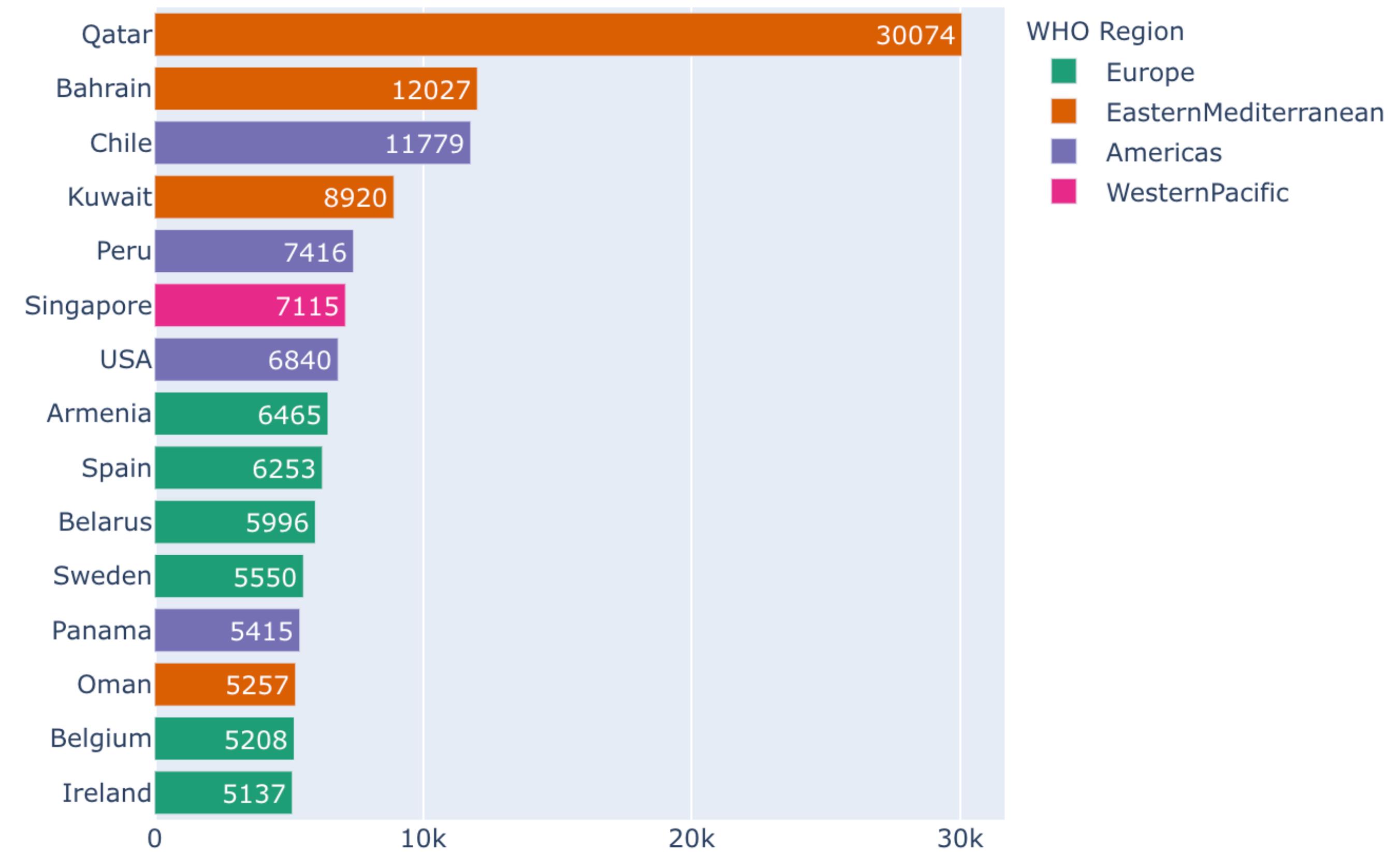
Horizontal Bar

```
1 def plot_hbar_wm(col, n, min_pop=1000000, sort='descending'):
2     df = worldometer_data[worldometer_data['Population']>min_pop]
3     df = df.sort_values(col, ascending=True).tail(n)
4     fig = px.bar(df,
5                   x=col, y="Country/Region", color='WHO Region',
6                   text=col, orientation='h', width=700,
7                   color_discrete_sequence = px.colors.qualitative.Dark2)
8     fig.update_layout(title=col+' (Only countries with > 1M Pop)',
9                       xaxis_title="", yaxis_title="",
10                      yaxis_categoryorder = 'total ascending',
11                      uniformtext_minsize=8, uniformtext_mode='hide')
12     fig.show()
```

Horizontal Bar

```
plot_hbar_wm('Tot Cases/1M pop', 15, 1000000)
```

Tot Cases/1M pop (Only countries with > 1M Pop)

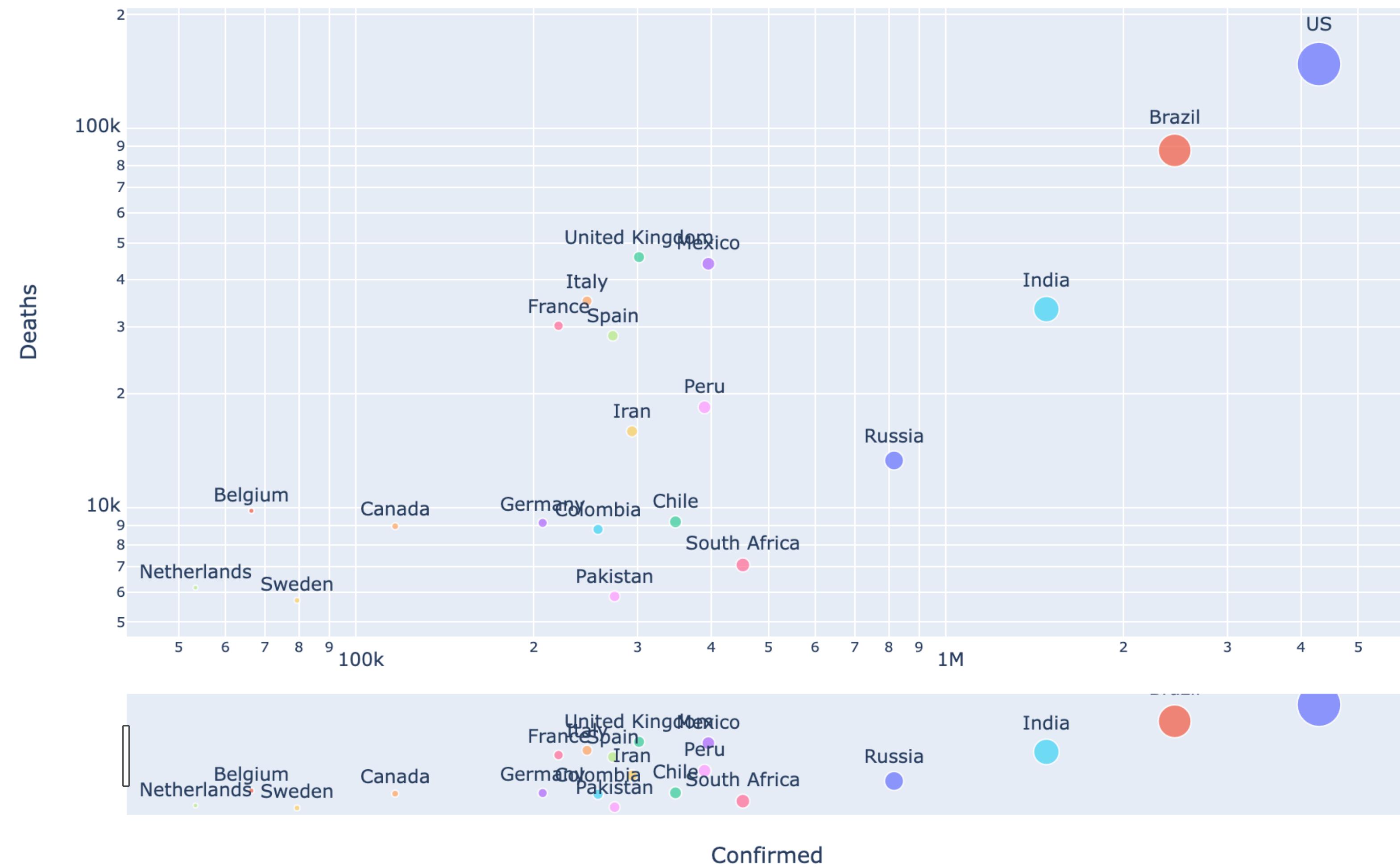


Scatter plot on Confirmed vs Deaths

```
1 fig = px.scatter(country_wise.sort_values('Deaths', ascending=False).iloc[:20, :],
2                   x='Confirmed', y='Deaths', color='Country/Region', size='Confirmed',
3                   height=700, text='Country/Region', log_x=True, log_y=True,
4                   title='Deaths vs Confirmed (Scale is in log10)')
5 fig.update_traces(textposition='top center')
6 fig.update_layout(showlegend=False)
7 fig.update_layout(xaxis_rangeslider_visible=True)
8 fig.show()
```

Scatter plot on Confirmed vs Deaths

Deaths vs Confirmed (Scale is in log10)



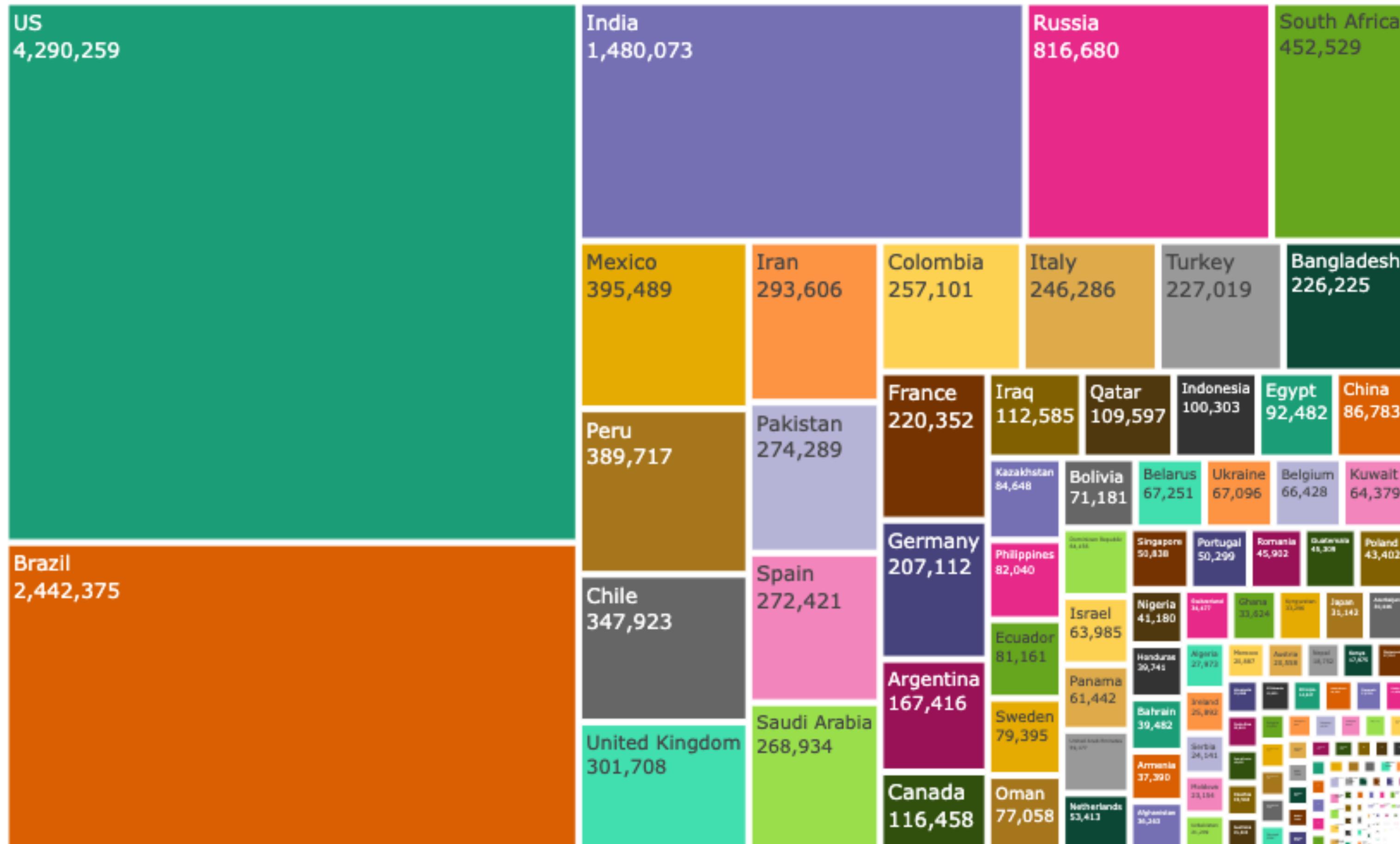
Composition with px.treemap

```
1 def plot_treemap(col):
2     fig = px.treemap(country_wise, path=["Country/Region"], values=col, height=700,
3                       title=col, color_discrete_sequence = px.colors.qualitative.Dark2)
4     fig.data[0].textinfo = 'label+text+value'
5     fig.show()
```

Composition with px.treemap

Confirmed

```
1 plot_treemap('Confirmed')
```



Bubble plot a long list

```
1 def plot_bubble(col, pal):
2     temp = full_grouped[full_grouped[col]>0].sort_values('Country/Region', ascending=False)
3     fig = px.scatter(temp, x='Date', y='Country/Region', size=col, color=col, height=3000,
4                       color_continuous_scale=pal)
5     fig.update_layout(yaxis = dict(dtick = 1))
6     fig.update(layout_coloraxis_showscale=False)
7     fig.show()
```

```
1 plot_bubble('New cases', 'Viridis')
```



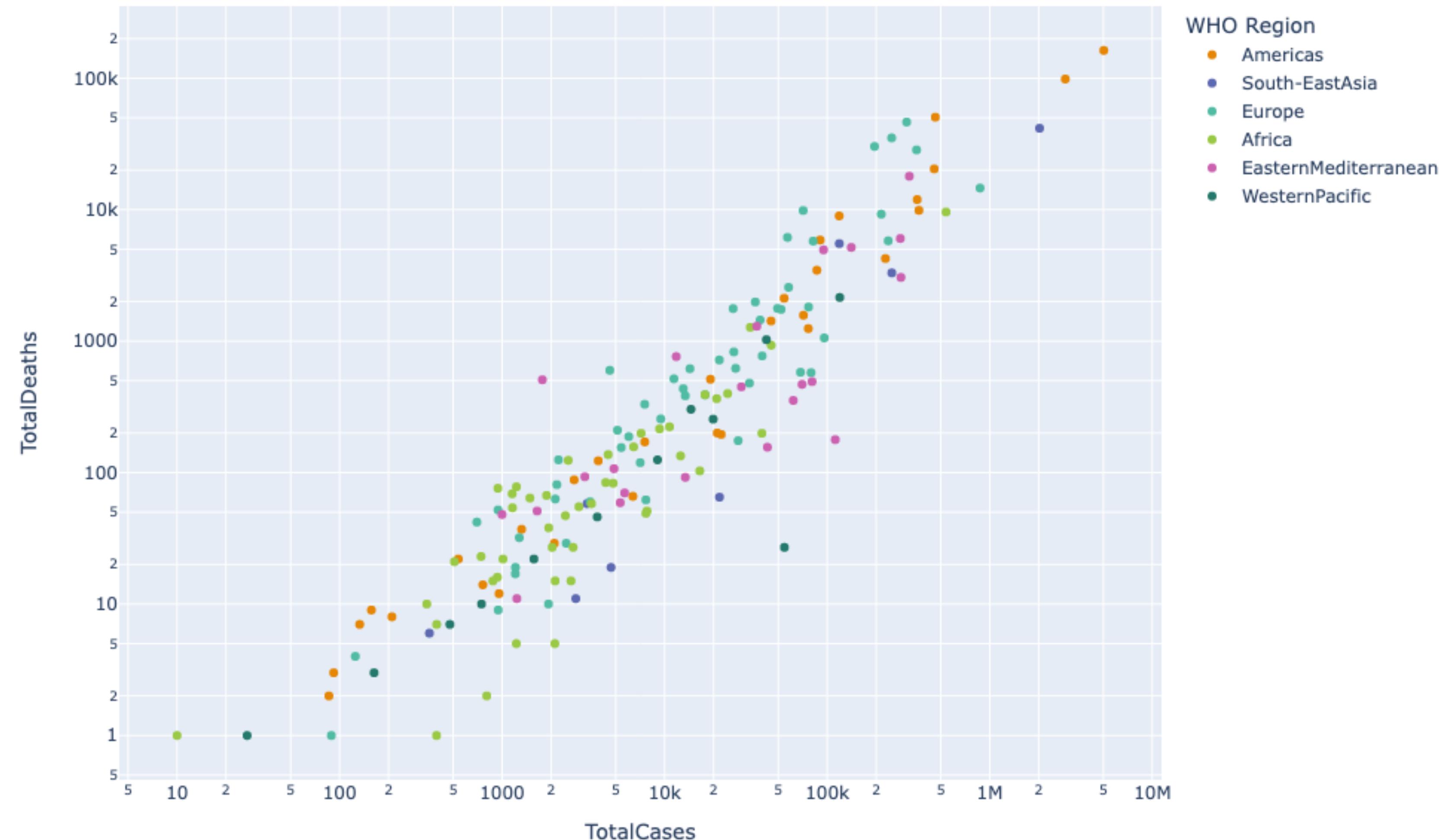
Confirmed vs Deaths

```
1 temp = worldometer_data[worldometer_data['WHO Region']!=0]

1 fig = px.scatter(temp, x='TotalCases', y='TotalDeaths', color='WHO Region',
2                   height=700, hover_name='Country/Region', log_x=True, log_y=True,
3                   title='Confirmed vs Deaths',
4                   color_discrete_sequence=px.colors.qualitative.Vivid)
5 fig.update_traces(textposition='top center')
6 # fig.update_layout(showlegend=False)
7 # fig.update_layout(xaxis_rangeslider_visible=True)
8 fig.show()
```

Confirmed vs Deaths

Confirmed vs Deaths

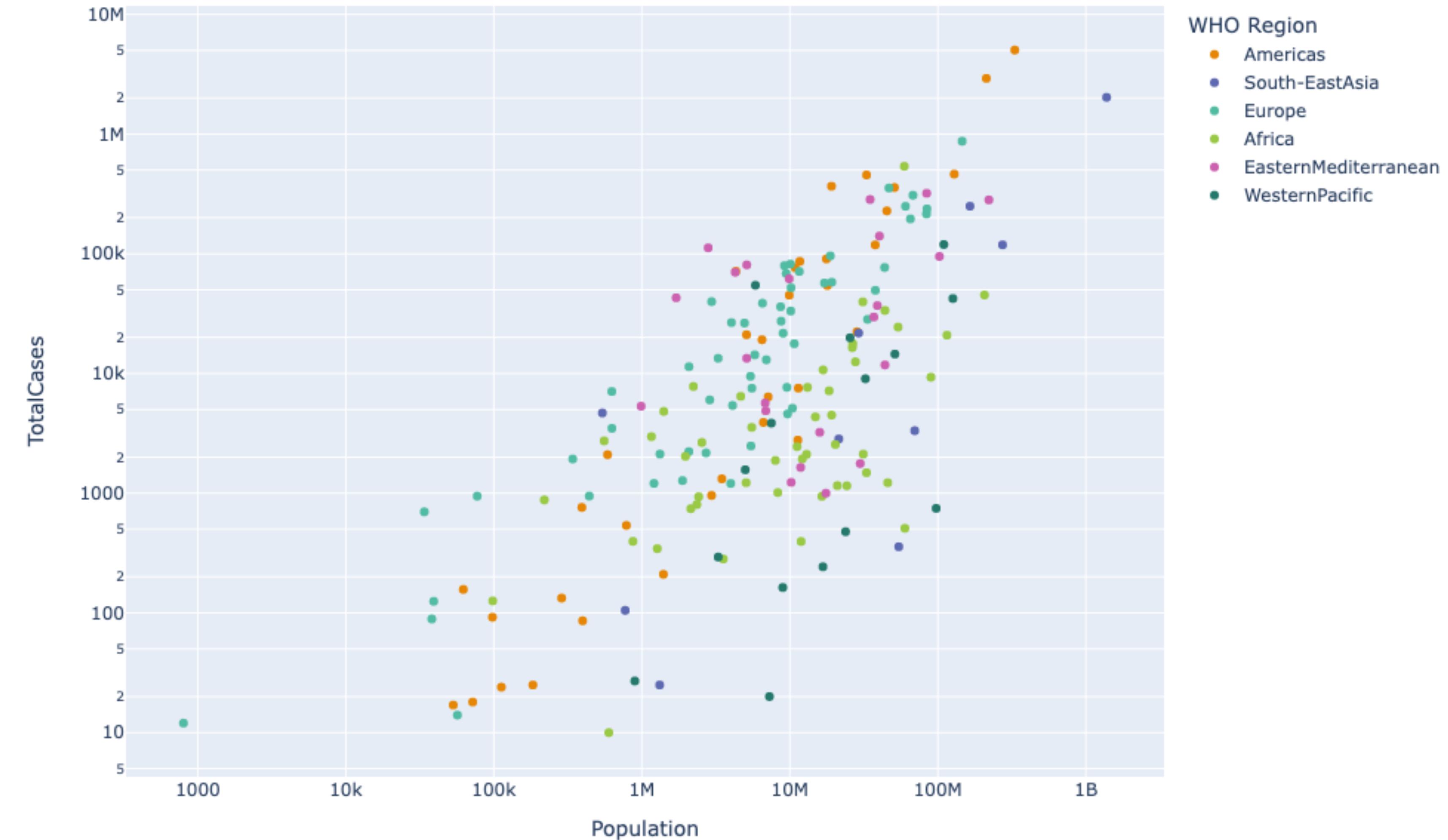


Populations vs Confirmed

```
1 fig = px.scatter(temp, x='Population', y='TotalCases', color='WHO Region',
2                   height=700, hover_name='Country/Region', log_x=True, log_y=True,
3                   title='Population vs Confirmed',
4                   color_discrete_sequence=px.colors.qualitative.Vivid)
5 fig.update_traces(textposition='top center')
6 # fig.update_layout(showlegend=False)
7 # fig.update_layout(xaxis_rangeslider_visible=True)
8 fig.show()
```

Populations vs Confirmed

Population vs Confirmed

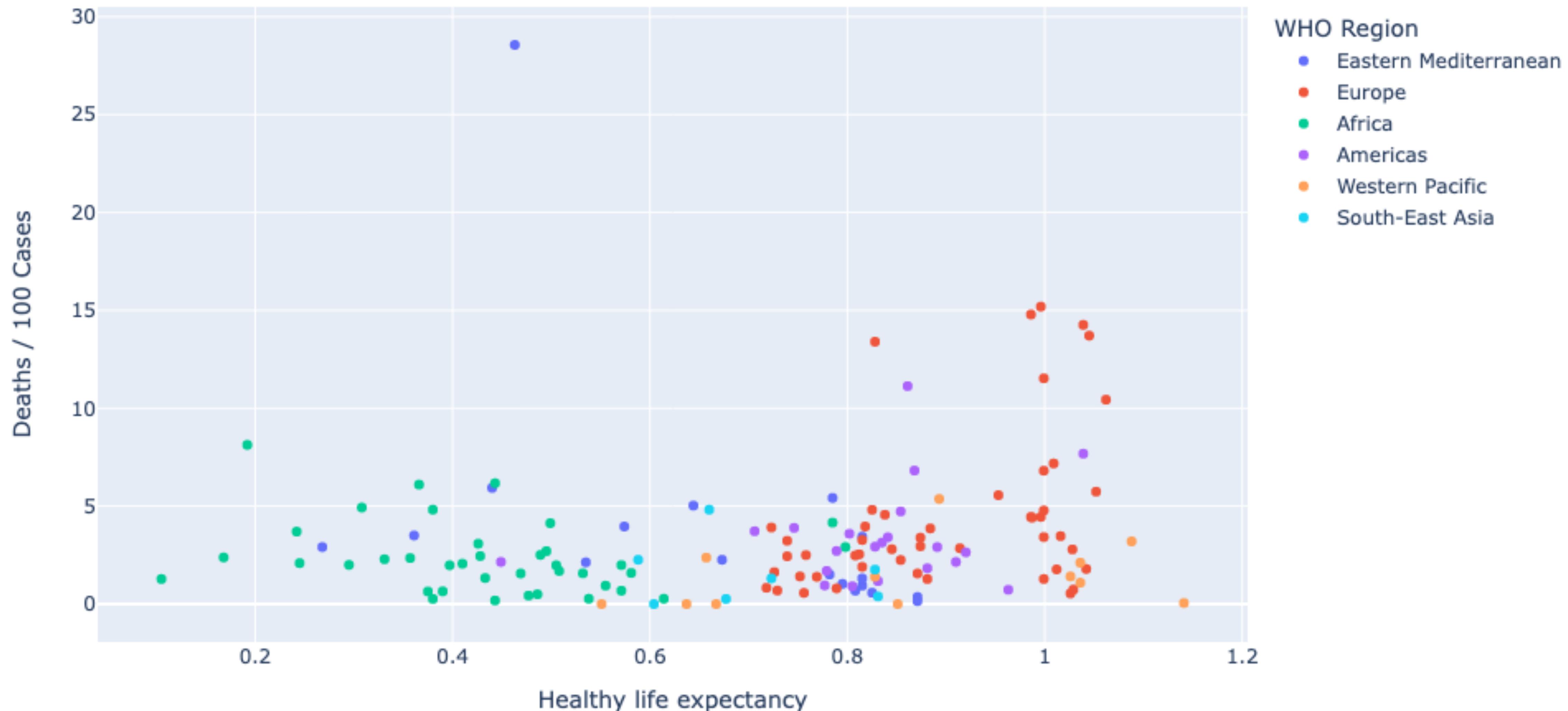


Deaths vs Health Life Expectancy

Pd.merge() to combine 2 DFs on key.

```
1 happiness_report = pd.read_csv('../Covid-19/Happiness2019.csv')
2 happiness_report = happiness_report[['Country or region', 'Healthy life expectancy']]
3
4 temp = country_wise.merge(happiness_report, left_on='Country/Region', right_on='Country or region')
5 px.scatter(temp, y='Deaths / 100 Cases', x='Healthy life expectancy', color='WHO Region',
6             hover_data=['Country/Region'])
```

Deaths vs Health Life Expectancy



Chapter Wrap Up

PX could be that easy and helpful. In this chapter, we demonstrate practical usage of data visualization of Covid-19 data.

Plotting function may be needed for plotting a few same kind of graph.

Once you clear the coding, everything is at your finger tip.

Reference & Resources

Official Website:

<https://plotly.com/python/>



Plotly Graph Objects:

<https://plotly.com/python/graph-objects/>

WorldBank API:

- <https://blogs.worldbank.org/opendata/introducing-wbgapi-new-python-package-accessing-world-bank-data>
- <https://nbviewer.org/github/tgherzog/wbgapi/blob/master/examples/wbgapi-cookbook.ipynb>
- <https://pypi.org/project/wbgapi/>

GitHub Open Source Code:

<https://github.com/plotly/plotly.py>

