# Coronavirus World Data Analysis

KATE expects your code to define variables with specific names that correspond to certain things we are interested in.

KATE will run your notebook from top to bottom and check the latest value of those variables, so make sure you don't overwrite them.

- Remember to uncomment the line assigning the variable to your answer and don't change the variable or function names.
- Use copies of the original or previous DataFrames to make sure you do not overwrite them by mistake.

You will find instructions below about how to define each variable.

Once you're happy with your code, upload your notebook to KATE to check your feedback.

First of all, run the following cell to:

- import `pandas` with an alias of `pd`
- read a CSV containing the data to work with
- convert the `date` column to the `datetime` format
- create a DataFrame `df` containing the data for only **1st July 2020**
- take a look at the first few rows of the DataFrame

```
In [1]:    import pandas as pd

           data = pd.read_csv('data/owid-covid-data.csv')
           data['date'] = pd.to_datetime(data['date'])
           df = data[data['date'] == '2020-07-01']

           df.head()
```

Out[1]:

| | iso_code | continent | location | date | total_cases | new_cases | total_deaths | new_deaths | total_cases |
|---|---|---|---|---|---|---|---|---|---|
| 173 | AFG | Asia | Afghanistan | 2020-07-01 | 31517.0 | 279.0 | 746.0 | 13.0 | |
| 300 | ALB | Europe | Albania | 2020-07-01 | 2535.0 | 69.0 | 62.0 | 4.0 | |
| 491 | DZA | Africa | Algeria | 2020-07-01 | 13907.0 | 336.0 | 912.0 | 7.0 | |
| 613 | AND | Europe | Andorra | 2020-07-01 | 855.0 | 0.0 | 52.0 | 0.0 | |
| 727 | AGO | Africa | Angola | 2020-07-01 | 284.0 | 8.0 | 13.0 | 2.0 | |

5 rows × 34 columns

- `df` DataFrame now has one row of data for each country with data present for **July 1st 2020**
- however, it also has a row with a `location` of `World` which contains aggregated values for all countries
- `df.tail()`, `df.info()` and `df.shape` will allow for further exploration of the structure of the DataFrame

In [2]: ▶| `#df.tail()`
`df.tail()`

Out[2]:

| | iso_code | continent | location | date | total_cases | new_cases | total_deaths | new_deaths | total_ca |
|---|---|---|---|---|---|---|---|---|---|
| **29411** | ESH | Africa | Western Sahara | 2020-07-01 | 380.0 | 172.0 | 1.0 | 0.0 | |
| **29506** | YEM | Asia | Yemen | 2020-07-01 | 1158.0 | 30.0 | 312.0 | 8.0 | |
| **29623** | ZMB | Africa | Zambia | 2020-07-01 | 1594.0 | 26.0 | 24.0 | 2.0 | |
| **29738** | ZWE | Africa | Zimbabwe | 2020-07-01 | 591.0 | 17.0 | 7.0 | 0.0 | |
| **29934** | OWID_WRL | NaN | World | 2020-07-01 | 10465987.0 | 192563.0 | 511041.0 | 5732.0 | |

5 rows × 34 columns

In [3]: ▶| `#df.info()`
`df.info()`

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 211 entries, 173 to 29934
Data columns (total 34 columns):
 #   Column                        Non-Null Count  Dtype
---  ------                        --------------  -----
 0   iso_code                      211 non-null    object
 1   continent                     210 non-null    object
 2   location                      211 non-null    object
 3   date                          211 non-null    datetime64[ns]
 4   total_cases                   210 non-null    float64
 5   new_cases                     210 non-null    float64
 6   total_deaths                  210 non-null    float64
 7   new_deaths                    210 non-null    float64
 8   total_cases_per_million       210 non-null    float64
 9   new_cases_per_million         210 non-null    float64
 10  total_deaths_per_million      210 non-null    float64
 11  new_deaths_per_million        210 non-null    float64
 12  total_tests                   73 non-null     float64
 13  new_tests                     73 non-null     float64
 14  total_tests_per_thousand      73 non-null     float64
 15  new_tests_per_thousand        73 non-null     float64
 16  new_tests_smoothed            83 non-null     float64
 17  new_tests_smoothed_per_thousand  83 non-null  float64
 18  tests_units                   85 non-null     object
 19  stringency_index              155 non-null    float64
 20  population                    211 non-null    float64
 21  population_density            200 non-null    float64
 22  median_age                    187 non-null    float64
 23  aged_65_older                 184 non-null    float64
 24  aged_70_older                 186 non-null    float64
 25  gdp_per_capita                184 non-null    float64
 26  extreme_poverty               122 non-null    float64
 27  cvd_death_rate                186 non-null    float64
 28  diabetes_prevalence           194 non-null    float64
 29  female_smokers                141 non-null    float64
 30  male_smokers                  139 non-null    float64
 31  handwashing_facilities        92 non-null     float64
 32  hospital_beds_per_thousand    165 non-null    float64
 33  life_expectancy               208 non-null    float64
dtypes: datetime64[ns](1), float64(29), object(4)
memory usage: 57.7+ KB
```

```
In [5]:  ▶ #df.shape
           df.shape
```

Out[5]: (211, 34)

**Q1. Create a new DataFrame called `countries` which is the same as `df` but with the `World` row removed.**

- Use the `.copy()` method to ensure you have a distinct DataFrame in memory
- Assign this new DataFrame to the variable `countries`; do not modify `df`

See below code syntax for some guidance:

```
countries['location'] != 'World'
```

```
In [17]:  ▶ #add your code below
            #countries = df.copy()
            #countries = ...

            countries = df.copy()
            countries = df[df['location'] != 'World'].copy()
            countries.head()
```

Out[17]:

| | iso_code | continent | location | date | total_cases | new_cases | total_deaths | new_deaths | total_cases |
|---|---|---|---|---|---|---|---|---|---|
| **173** | AFG | Asia | Afghanistan | 2020-07-01 | 31517.0 | 279.0 | 746.0 | 13.0 | |
| **300** | ALB | Europe | Albania | 2020-07-01 | 2535.0 | 69.0 | 62.0 | 4.0 | |
| **491** | DZA | Africa | Algeria | 2020-07-01 | 13907.0 | 336.0 | 912.0 | 7.0 | |
| **613** | AND | Europe | Andorra | 2020-07-01 | 855.0 | 0.0 | 52.0 | 0.0 | |
| **727** | AGO | Africa | Angola | 2020-07-01 | 284.0 | 8.0 | 13.0 | 2.0 | |

5 rows × 34 columns

◀ ▬▬▬▬▬▬▬▬▬▬▬▬▬▬ ▶

**Q2. Check the shape of your DataFrame to confirm that `countries` has one row fewer than `df`:**

Please note you have been provided with the code for this question to carry out the necessary analysis. Simply uncomment the line of code and run the code cell to produce the desired results.

```
In [12]:  ▶ #print(df.shape, countries.shape)

            print(df.shape, countries.shape)
```

(211, 34) (210, 34)

**Q3. Define a DataFrame based on the `countries` DataFrame, but which only contains the columns in `cols` (defined below) and assign this to a variable called `countries_dr`**

- Order this DataFrame by `'total_deaths_per_million'`, with the highest numbers at the top.

See below code syntax for some guidance:

```
DataFrame_name[column_names].sort_values(by=..., ascending=False)
```

```
cols = ['continent', 'location', 'total_deaths_per_million']

#add your code below
#countries_dr

countries_dr = countries[cols].sort_values(by="total_deaths_per_million", ascending=Fa
print(countries_dr)
```

```
           continent                             location  \
23306        Europe                           San Marino
2917         Europe                              Belgium
613          Europe                              Andorra
28347        Europe                       United Kingdom
25362        Europe                                Spain
...             ...                                  ...
23111  North America  Saint Vincent and the Grenadines
23926        Africa                           Seychelles
15734        Africa                              Lesotho
10808        Europe                            Gibraltar
12195          Asia                            Hong Kong

       total_deaths_per_million
23306                  1237.551
2917                    841.615
613                     673.008
28347                   644.168
25362                   606.633
...                          ...
23111                     0.000
23926                     0.000
15734                     0.000
10808                     0.000
12195                       NaN

[210 rows x 3 columns]
```

**Q4.** Using the `countries` DataFrame we created earlier, find the sum of `total_tests` for countries in `Africa`, assigning the result, *as an integer*, to `africa_tests`.

- Use `.sum()` method calculate the sum for `total_tests` column
- Use `.astype(int)` method or `int()` function to convert results to an integer

See below code syntax for some guidance:

```
countries['continent'] == 'Africa'
```

```
#add your code below
#africa_tests

africa_tests = countries[countries['continent'] == 'Africa']["total_tests"].sum().asty
africa_tests
```

3445134

**Q5.** How many countries in Africa have no value recorded for the number of `total_tests` column? Assign the result to `africa_missing_test_data`.

- You may find the pandas `.isna()` method and python `len()` function useful

See below code syntax for some guidance:

```
len(DataFrame_name[column_name].isna())
```

```
In [116]:  ▶ africa = countries[countries['continent'] == 'Africa']

            africa_tests = africa['total_tests'].sum()

            africa_missing_test_data = africa['total_tests'].isna().sum()
            africa_missing_test_data
```

Out[116]: 45

**Q6. How many countries have a higher value for `total_tests` than the `United Kingdom`? Assign your answer to a variable called `countries_more_tests`.**

Remember to work from the `countries` DataFrame rather than `df`. You should avoid modifying any existing DataFrames.

```
In [34]:  ▶ #add your code below
            #countries_more_tests

            uk_tests = countries[countries['location'] == 'United Kingdom']['total_tests'].values[
            countries_more_tests = len(countries[countries['total_tests'] > uk_tests])
            countries_more_tests
```

Out[34]: 3

**Q7. Create a DataFrame called `beds_dr` which is based on the `countries` DataFrame, but contains only the columns `hospital_beds_per_thousand` and `total_deaths_per_million`.**

- Your answer should only include rows where there are values present in both of these columns
- You may find the `.dropna()` method useful

See below code syntax for some guidance:

```
DataFrame_name.dropna()
```

```
In [39]:  ▶ #add your code below
            #beds_dr

            beds_dr = countries[["hospital_beds_per_thousand", "total_deaths_per_million"]].dropna
            beds_dr
```

Out[39]:

| | hospital_beds_per_thousand | total_deaths_per_million |
|---|---|---|
| 173 | 0.50 | 19.163 |
| 300 | 2.89 | 21.544 |
| 491 | 1.90 | 20.798 |
| 952 | 3.80 | 30.635 |
| 1081 | 5.00 | 28.919 |
| ... | ... | ... |
| 29136 | 0.80 | 1.794 |
| 29332 | 2.60 | 0.000 |
| 29506 | 0.70 | 10.461 |
| 29623 | 2.00 | 1.305 |
| 29738 | 1.70 | 0.471 |

164 rows × 2 columns

**Q8. Refer to the `beds_dr` DataFrame. What is the average `total_deaths_per_million` for entries in `beds_dr` where `hospital_beds_per_thousand` is greater than the mean?**

- Save the results to a new variable called `dr_high_bed_ratio`

See below code syntax for some guidance:

```
beds_dr['hospital_beds_per_thousand'] > beds_dr['hospital_beds_per_thousand'].mean()
```

In [49]: ▶|
```
#add your code below
#dr_high_bed_ratio

x = beds_dr['hospital_beds_per_thousand'] > beds_dr['hospital_beds_per_thousand'].mean
dr_high_bed_ratio =beds_dr[x]['total_deaths_per_million'].mean()
dr_high_bed_ratio
```

Out[49]: 98.18423728813559

**Q9. Refer to the `beds_dr` DataFrame. What is the average `total_deaths_per_million` for entries in `beds_dr` where `hospital_beds_per_thousand` is less than the mean?**

- Save the results to a new variable called `dr_low_bed_ratio`

See below code syntax for some guidance:

```
beds_dr['hospital_beds_per_thousand'] < beds_dr['hospital_beds_per_thousand'].mean()
```

In [50]: ▶|
```
#add your code below
#dr_low_bed_ratio
y = beds_dr['hospital_beds_per_thousand'] < beds_dr['hospital_beds_per_thousand'].mean
dr_low_bed_ratio = beds_dr[y]['total_deaths_per_million'].mean()
dr_low_bed_ratio
```

Out[50]: 56.294057142857135

**Q10. Refer to the `countries` DataFrame. Create a new DataFrame called `no_new_cases` which contains only rows from `countries` with zero `new_cases`.**

Please note you have been provided with the code for this question to carry out the necessary analysis. Simply uncomment the lines of code and run the code cell to produce the desired results.

In [51]: ▶| `countries.head()`

Out[51]:

| | iso_code | continent | location | date | total_cases | new_cases | total_deaths | new_deaths | total_case: |
|---|---|---|---|---|---|---|---|---|---|
| **173** | AFG | Asia | Afghanistan | 2020-07-01 | 31517.0 | 279.0 | 746.0 | 13.0 | |
| **300** | ALB | Europe | Albania | 2020-07-01 | 2535.0 | 69.0 | 62.0 | 4.0 | |
| **491** | DZA | Africa | Algeria | 2020-07-01 | 13907.0 | 336.0 | 912.0 | 7.0 | |
| **613** | AND | Europe | Andorra | 2020-07-01 | 855.0 | 0.0 | 52.0 | 0.0 | |
| **727** | AGO | Africa | Angola | 2020-07-01 | 284.0 | 8.0 | 13.0 | 2.0 | |

5 rows × 34 columns

```
In [52]:    #add your code below
            #no_new_cases = countries[countries['new_cases'] == 0]
            #no_new_cases.head()

            no_new_cases = countries[countries['new_cases'] == 0]
            no_new_cases.head()
```

Out[52]:

| | iso_code | continent | location | date | total_cases | new_cases | total_deaths | new_deaths | total_cases_ |
|---|---|---|---|---|---|---|---|---|---|
| 613 | AND | Europe | Andorra | 2020-07-01 | 855.0 | 0.0 | 52.0 | 0.0 | |
| 836 | AIA | North America | Anguilla | 2020-07-01 | 3.0 | 0.0 | 0.0 | 0.0 | |
| 952 | ATG | North America | Antigua and Barbuda | 2020-07-01 | 66.0 | 0.0 | 3.0 | 0.0 | |
| 1381 | ABW | North America | Aruba | 2020-07-01 | 103.0 | 0.0 | 3.0 | 0.0 | |
| 2080 | BHS | North America | Bahamas | 2020-07-01 | 104.0 | 0.0 | 11.0 | 0.0 | |

5 rows × 34 columns

**Q11. Refer to the `no_new_cases` DataFrame. Which country in `no_new_cases` DataFrame has had the highest number of `total_cases` ?**

- Save the results to a new variable called `highest_no_new`

See below code syntax for some guidance:

```
no_new_cases['total_cases'] == no_new_cases['total_cases'].max()
```

```
In [117]:   highest_no_new = no_new_cases.loc[no_new_cases['total_cases'].idxmax(), 'location']
            print(highest_no_new)

            Cameroon
```

**Q12. Refer to the `countries` DataFrame. What is the sum of the `population` of all countries which have had zero `total_deaths` ?**

- Assign your answer to `sum_populations_no_deaths` variable
- Your answer should be in millions, rounded to the nearest whole number, and converted to an integer

```
In [63]:    #add your code below
            #sum_populations_no_deaths

            sum_populations_no_deaths = countries[countries['total_deaths'] == 0]['population'].su
            sum_populations_no_deaths = round(sum_populations_no_deaths)
            sum_populations_no_deaths
```

Out[63]:  192

**Q13. Create a function called `country_metric` which accepts the following three parameters:**

- a DataFrame (which can be assumed to be of a similar format to `countries` )
- a location (i.e. a string which will be found in the `location` column of the DataFrame)
- a metric (i.e. a string which will be found in any column (other than `location` ) in the DataFrame)

The function should return only the value from the first row for a given `location` and `metric` . *You may find `.iloc[]` useful.*

See below code syntax for some guidance:

```python
def country_metric(df, location, metric):

    return df[df['location'] == location].iloc[0][metric]
```

In [65]: ▶
```python
#add your code below
#def country_metric(df, location, metric):

def country_metric(df, location, metric):
    return countries[countries["location"] == location].iloc[0][metric]
```

**Q.14 Use your function to collect the value for** `Vietnam` **for the metric** `aged_70_older` **, assigning the result to** `vietnam_older_70` **.**

Please note you have been provided with the code for this question to carry out the necessary analysis. Simply uncomment the lines of code and run the code cell to produce the desired results.

In [67]: ▶
```python
#add your code below
#vietnam_older_70 = country_metric(countries, 'Vietnam', 'aged_70_older')
#vietnam_older_70

vietnam_older_70 = country_metric(countries, 'Vietnam', 'aged_70_older')
vietnam_older_70
```

Out[67]: 4.718

**Q.15 Create another function called** `countries_average` **, which accepts the following three parameters:**

- a DataFrame "df" (which can be assumed to be such as `countries` )
- a list of countries "countries" (which can be assumed to all be found in the `location` column of the DataFrame)
- a string "metric" (which can be assumed to be a column (other than `location` ) which will be found in the DataFrame) . For instance, this string value can be `life_expectancy` .

Note that for the test on KATE for this question to pass, you need to make sure the function accepts the three parameters in the following order: `countries_average(df, countries, metric)` . (You can call your parameters however you like as long as the type of these parameters are what was described above).

The function should return the average value for the given metric for the given list of countries.

You may find `.isin()` method useful while filtering for list of countries.

In [106]: ▶
```python
#add your code below
#def countries_average(df, countries, metric):

def countries_average(df, countries, metric):
    new_df = df[df['location'].isin(countries)]

    average_value = new_df[metric].mean()

    return average_value
```

**Q16. Use your** `countries_average` **function to find out the average** `life_expectancy` **of countries in the** `g7` **list defined below. Assign the result to the variable** `g7_avg_life_expectancy` **.**

Please note you have been provided with the code for this question to carry out the necessary analysis. Simply uncomment the lines of code and run the code cell to produce the desired results.

```
In [107]:  ▶|  g7 = ['United States', 'Italy', 'Canada', 'Japan', 'United Kingdom', 'Germany', 'Franc
```

```
In [108]:  ▶|  countries.head()
```

Out[108]:

| | iso_code | continent | location | date | total_cases | new_cases | total_deaths | new_deaths | total_case |
|---|---|---|---|---|---|---|---|---|---|
| 173 | AFG | Asia | Afghanistan | 2020-07-01 | 31517.0 | 279.0 | 746.0 | 13.0 | |
| 300 | ALB | Europe | Albania | 2020-07-01 | 2535.0 | 69.0 | 62.0 | 4.0 | |
| 491 | DZA | Africa | Algeria | 2020-07-01 | 13907.0 | 336.0 | 912.0 | 7.0 | |
| 613 | AND | Europe | Andorra | 2020-07-01 | 855.0 | 0.0 | 52.0 | 0.0 | |
| 727 | AGO | Africa | Angola | 2020-07-01 | 284.0 | 8.0 | 13.0 | 2.0 | |

5 rows × 34 columns

```
In [109]:  ▶|  #add your code below
              g7 = ['United States', 'Italy', 'Canada', 'Japan', 'United Kingdom', 'Germany', 'Franc
              g7_avg_life_expectancy = countries_average(df, g7, 'life_expectancy')
              g7_avg_life_expectancy
```

Out[109]:  82.10571428571428

**Q.17 Refer to the `countries` DataFrame. Find the country with lowest value for `life_expectancy` in the `countries` DataFrame, and create a string which is formatted as follows:**

'{country} has a life expectancy of {diff} years lower than the G7 average.'

Assign your string to the variable `headline` and ensure it is formatted exactly as above, with:

- use `f-strings` to format the string
- {country} being replaced by the value in the `location` column of the DataFrame
- {diff} being replaced by a float **rounded to one decimal place**, of the value from the `life_expectancy` column subtracted from `g7_avg_life_expectancy`. Please note that {diff} should be a positive value

```
diff = <G7 countries average life expectancy> - <value of the lowest life expe
ctancy country>
```

See below code syntax for some guidance:

```
lowest = countries[countries['life_expectancy'] == countries['life_expectancy'].m
in()].iloc[0]
country = lowest['location']
life_exp = lowest['life_expectancy']
```

```
In [110]:  ▶  #add your code below
              #headline = f'{country} has a life expectancy of {diff} years lower than the G7 average

              lowest = countries[countries['life_expectancy'] == countries['life_expectancy'].min()]
              country = lowest['location']
              life_exp = lowest['life_expectancy']

              diff = g7_avg_life_expectancy - life_exp


              headline = f"{country} has a life expectancy of {diff:.1f} years lower than the G7 ave

              print(headline)
```

Central African Republic has a life expectancy of 28.8 years lower than the G7 averag
e.

```
In [ ]:  ▶
```