

Challenge 2 (Third Challenge):

Read carefully. Ask questions for clarification where need be.

Find a link below to some data files to work with. This is a set of csv files from a daily scrapping of the corona dot help website tracking the numbers in the covid-19 pandemic.

There are two sets of files named covid* and worldwide* Both sets of files have two versions at two different times of the day: around 11pm and early in the morning.

The tasks:

1. Organise the files into 4 groups (show formula + number of files each). a. covid* 11pm+ b. world* 11pm+ c. covid* other time d. world* other time
2. Using 1b above, how many cases did the US and Germany each record in April?
3. Using 1b above, how many cases did the US and Germany each record in May?
4. Using 1c above, how many cases did Italy and Russia each record in April?
5. Using 1c above, how many cases did Italy and Russia each record in May?


















#Kicker Using 1b, how many cases were recorded by US, Italy, China, Russia, Germany in April + May?

Max plots: 4. I will be looking at the numbers.

This is an exercise in Wrangling so I am not going to bother about making any plots. I will dwell on the sorting and cleaning moves as well as the numbers associated with each question.

```
In [1]: 1 import IPython
        2 IPython.display.Image("files.png")
```

Out[1]:

		covid_help_data_28-Apr-2020 23-40-33	5/25/2020 11:59 AM	Microsoft Excel C...	11 KB
		Azure Fesh			
		Database for SQL			
		Dropbox			
		pythonanywhere_data_			
		Dropbox			
		OneDrive			
		This PC			
		covid_help_data_29-Apr-2020 10-02-29	5/25/2020 11:59 AM	Microsoft Excel C...	11 KB
		covid_help_data_29-Apr-2020 23-40-32	5/25/2020 11:59 AM	Microsoft Excel C...	11 KB
		covid_help_data_30-Apr-2020 10-02-28	5/25/2020 11:59 AM	Microsoft Excel C...	11 KB
		covid_help_data_30-Apr-2020 23-40-33	5/25/2020 11:59 AM	Microsoft Excel C...	11 KB
		worldwide_01-May-2020 10-11-07	5/25/2020 11:59 AM	Microsoft Excel C...	14 KB
		worldwide_01-May-2020 23-49-16	5/25/2020 11:59 AM	Microsoft Excel C...	14 KB
		worldwide_02-Apr-2020 23-56-11	4/13/2020 12:44 PM	Microsoft Excel C...	6 KB
		worldwide_02-May-2020 10-11-19	5/25/2020 11:59 AM	Microsoft Excel C...	14 KB
		worldwide_02-May-2020 23-49-17	5/25/2020 11:59 AM	Microsoft Excel C...	14 KB

```
In [2]: 1 #Import Needed Libraries
        2 import os
        3 import pandas as pd
```

```
In [3]: 1 #Assemble all csv files in the folder
2 #Using the os module, we find all files ending with 'csv'.
3 data = [line for line in os.listdir() if line.endswith('csv')]
4 print(f"There are {len(data)} csv files in the folder")
```

There are 181 csv files in the folder

```
In [4]: 1 #To sort the files holding the covid* and world* groups, we find files whose name h
2 cov = [line for line in data if line.startswith('cov')]
3 ww = [line for line in data if line.startswith('world')]
4
5 #For perspective and to ensure we are in line, the sum of both groups should be equal
6 len(cov), len(ww), len(cov) + len(ww), len(cov) + len(ww) == len(data)
```

Out[4]: (93, 88, 181, True)

```
In [5]: 1 #Then we find in each group, the files for the different months
2 #It's simple, going by the naming of the files, each one has at least the first 3 s
3 #Let's make two groups each of covid* and world* groups corresponding to april and m
4 cov_Apr = [line for line in cov if 'Apr' in line]
5 cov_May = [line for line in cov if 'May' in line]
6 ww_Apr = [line for line in ww if 'Apr' in line]
7 ww_May = [line for line in ww if 'May' in line]
```

```
In [6]: 1 #For this section, remember the variable names. We will use them down the line
2
3 #Find the 23 Hours(11pm) and Other times for the world* group. Again, the clues are
4 #For May
5 ww_May23 = [line for line in ww_May if '23' in line.split()[-1]]
6 ww_Mayother = [line for line in ww_May if not '23' in line.split()[-1]]
7
8 #For Apr
9 ww_Apr23 = [line for line in ww_Apr if '23' in line.split()[-1]]
10 ww_Aprother = [line for line in ww_Apr if not '23' in line.split()[-1]]
11
12
13 #Find the 23 Hours(11pm) and Other times for the covid* group
14 #For May
15 cov_May23 = [line for line in cov_May if '23' in line.split()[-1]]
16 cov_Mayother = [line for line in cov_May if not '23' in line.split()[-1]]
17
18 #For Apr
19 cov_Apr23 = [line for line in cov_Apr if '23' in line.split()[-1]]
20 cov_Aprother = [line for line in cov_Apr if not '23' in line.split()[-1]]
21
22
23 # Let's summarize the numbers we have for each group
24 print(f"Covid*\nApril 11pm > {len(cov_May23) + len(cov_Apr23)}.\nApril Other Time > {len(cov_Mayother) + len(cov_Aprother)}")
25 print(f"\nWorld*\nMay 11pm > {len(ww_May23) + len(ww_Apr23)}.\nMay Other Time > {len(ww_Mayother) + len(ww_Aprother)}")
```

```
Covid*
April 11pm > 46.
April Other Time > 47

World*
May 11pm > 43.
May Other Time > 45
```

Notice that so far, I have not read a single file. It is a choice. I wanted to sort and organise the files first.

- #Let's take a look at how a SAMPLE of the first csv file of each group looks. This will affect how we decide to work and what columns we wanna use.

```
In [7]: 1 #Apr 23 hours aka 11pm
2 df = pd.read_csv(ww_Apr23[0])
3 df.sample(5)
```

Out[7]:

	Country	Total confirmed cases	Total deaths	Confirmed recoveries	Cases confirmed today	Deaths today	Recoveries confirmed today
160	mongolia	14	0	2	0	0	0
56	egypt	865	58	201	86	6	22
111	brunei	133	1	56	2	0	4
27	malaysia	3,116	50	767	208	5	122
83	réunion	308	0	40	27	0	0

```
In [8]: 1 #Apr Other hours
2 df = pd.read_csv(ww_Aprother[0])
3 df.sample(5)
```

```
Out[8]:
```

	Country	Total confirmed cases	Total deaths	Confirmed recoveries	Cases confirmed today	Deaths today	Recoveries confirmed today
138	djibouti	40	0	0	7	0	0
9	turkey	18,135	356	415	2,456	79	82
165	libya	11	1	0	1	1	0
5	france	59,105	5,387	12,428	2,116	1,355	1,493
142	bermuda	35	0	11	3	0	1

```
In [9]: 1 #May 23 hours aka 11pm
2 df = pd.read_csv(ww_May23[0])
3 df.sample(5)
```

```
Out[9]:
```

	Country	Total confirmed	Total_confirmed_today	Total_deaths	Total_deaths_today	Total_Recoveries	Total_R
29	qatar	14,096	687	12	2	1,436	
32	united-arab-emirates	13,038	557	111	6	2,543	
86	hong-kong	1,040	2	4	0	859	
160	haiti	81	5	8	2	8	
6	turkey	122,392	2,188	3,258	84	53,808	

```
In [10]: 1 #May Other hours
2 df = pd.read_csv(ww_Mayother[0])
3 df.sample(5)
```

```
Out[10]:
```

	Country	Total confirmed	Total_confirmed_today	Total_deaths	Total_deaths_today	Total_Recoveries	Total_R
161	mozambique	76	0	0	0	12	
31	united-arab-emirates	13,038	557	111	6	2,543	
186	saint-lucia	17	0	0	0	15	
142	ethiopia	133	2	3	0	66	
66	croatia	2,076	14	69	2	1,348	

In [11]:

```
1 #May 23 hours aka 11pm
2 df = pd.read_csv(ww_May23[0])
3 df.sample(5)
```

Out[11]:

	Country	Total confirmed	Total_confirmed_today	Total_deaths	Total_deaths_today	Total_Recoveries	Total_F
201	seychelles	11	0	0	0	6	
105	guatemala	599	14	16	0	66	
133	congo-brazzaville	229	9	9	0	25	
125	montenegro	322	0	7	0	214	
7	ruissia	114,431	7,933	1,169	96	13,220	

In [12]:

```
1 #May Other hours
2 df = pd.read_csv(ww_Mayother[0])
3 df.sample(5)
```

Out[12]:

	Country	Total confirmed	Total_confirmed_today	Total_deaths	Total_deaths_today	Total_Recoveries	Total_F
46	colombia	6,507	300	293	15	1,439	
155	uganda	83	2	0	0	52	
196	gambia	12	1	1	0	8	
28	japan	14,088	352	430	36	2,460	
158	haiti	81	5	8	2	8	

In [13]:

```
1 #Some numbers come with commas in them. I could use the 'thousand' argument while re
2 #However, I want to throw in some measure of control by myself by defining how I wa
3 def number(df):
4     #April entries has two different names for total confirmed so I convert that co
5     try: #Apr format
6         #overwrite the rows of these columns with their respective values replacing
7         df['total_confirmed_cases'] = [int(line.replace(',', '')) for line in df['to
8         df['cases_confirmed_today'] = [int(line.replace(',', '')) for line in df['ca
9     except: #May format
10        df['total_confirmed'] = [int(line.replace(',', '')) for line in df['total_co
11        df['total_confirmed_today'] = [int(line.replace(',', '')) for line in df['to
12
13    #Change the date to datetime extracting the date
14    #Ended up not using this but I am leaving it around all the same
15    df['date'] = [pd.to_datetime(line).date() for line in df['date']]
16
17    return df
```

Question 2: Using 1b above (world* 11pm+), how many cases did the US and Germany each record in April?

The Flow

- Again, I refuse to read and merge all the files. Reason is simple: They have varying types and formats
- For each group, I read a single csv file, extract the country I'm interested in and keep piling the df in a list
- I ensure the dfs I am concatenating have same shape and same column names
- Since I ensured the last line, I can then concat them (merge them one under the other), reset the index and drop theirs.
- I am using the total confirmed for each country which is given for each day. This is the total from inception.
- With the above, I can take the first total on record and get the difference from the last on record. This gives me total for that period.

In [14]:

```
1 #Disable the copy warning on this one
2 pd.options.mode.chained_assignment = None
3
4 #Open two different lists to hold the dataframes from reading and extracting the new
5 us_a, germany_a = [], []
6
7 #Let's go with the world* apr taken at 23 hours. Rememebr that group from above?
8 for line in ww_Apr23:
9     #The US numbers
10    df = pd.read_csv(line) #For each item in the group, read the csv file
11    df.columns = [line.lower().replace(' ', '_') for line in df.columns] #convert col
12    df = df[df['country'] == 'united-states'] #Pick the country we are interested in
13    try:
14        #For the files where confirmed col names is 'total_confirmed_cases'. Let's g
15        df = df[['country', 'total_confirmed_cases', 'cases_confirmed_today']]
16        df['date'] = line.split()[0].split('_')[1] #add a date col stripping it from
17        df['time'] = line.split()[1][:-4] #ditto for time
18        us_a.append(df)
19    except:
20        #In case the total confirmed columns is name 'total_confirmed'. We grab that
21        df = df[['country', 'total_confirmed', 'total_confirmed_today']]
22        df['date'] = line.split()[0].split('_')[1]
23        df['time'] = line.split()[1][:-4]
24        #To ensure consistent naming, we rename the col to same format as the above
25        df.columns = ['country', 'total_confirmed_cases', 'cases_confirmed_today', '
26        us_a.append(df)
27
28    #We repeat the above for Germany
29    #I know. Should've been a function. I am leaving it like this for repetition and
30    df = pd.read_csv(line)
31    df.columns = [line.lower().replace(' ', '_') for line in df.columns]
32    df = df[df['country'] == 'germany']
33    try:
34        #For confirmed col names 'total_confirmed_cases'
35        df = df[['country', 'total_confirmed_cases', 'cases_confirmed_today']]
36        df['date'] = line.split()[0].split('_')[1]
37        df['time'] = line.split()[1][:-4]
38        germany_a.append(df)
39    except:
40        #For confirmed col names 'total_confirmed_cases'
41        df = df[['country', 'total_confirmed', 'total_confirmed_today']]
42        df['date'] = line.split()[0].split('_')[1]
43        df['time'] = line.split()[1][:-4]
44        #rename the cols to same format as the above dataframe for easy concat
45        df.columns = ['country', 'total_confirmed_cases', 'cases_confirmed_today',
46        germany_a.append(df)
47
48    #We can now concat the list of dataframes for each country and convert our numbers
49    us_data_a = pd.concat(us_a).reset_index(drop=True)
50    us_data_a = number(us_data_a)
51
52    germany_data_a = pd.concat(germany_a).reset_index(drop=True)
53    germany_data_a = number(germany_data_a)
```

In [15]:

```
1  #To the Last part of our Logic/flow. Get the total for the month by getting difference
2  #The total for the month of April will be last date's record minus the first date's
3  germany_april_start = germany_data_a.total_confirmed_cases.tolist()[0]
4  germany_april_end = germany_data_a.total_confirmed_cases.tolist()[-1]
5
6  #Option2: Going with cases confirmed per day column.
7  #Notice that will return a different number. It's expected as the end of day for each
8  germany_april = germany_data_a.cases_confirmed_today.sum()
9
10 print(f"Total cases for Germany in April is {germany_april_end - germany_april_start}")
11
12 #Since there are 13 items on the list, total for the month of April will be last date's record minus the first date's
13 us_april_start = us_data_a.total_confirmed_cases.tolist()[0]
14 us_april_end = us_data_a.total_confirmed_cases.tolist()[-1]
15
16 #Option2: Going with cases confirmed per day column.
17 us_april = us_data_a.cases_confirmed_today.sum()
18
19 print(f"Total cases for US in April is {us_april_end - us_april_start:,}")
20 print(f"\nPer 'Case Confirmed Today' col:\nGermany April confirmed> {germany_april_end - germany_april_start:,}")
```

Total cases for Germany in April is 78,215

Total cases for US in April is 851,822

Per 'Case Confirmed Today' col:

Germany April confirmed> 45,008.

US April confirmed> 563,709.

Question 3: Using 1b above (world* 11pm+), how many cases did the US and Germany each record in May?

- exact same flow as question 2

In [16]:

```
1 us_m, germany_m = [], []
2 for line in ww_May23:
3     #US
4     df = pd.read_csv(line)
5     df.columns = [line.lower().replace(' ', '_') for line in df.columns]
6     df = df[df['country'] == 'united-states']
7
8     df = df[['country', 'total_confirmed', 'total_confirmed_today']]
9     df['date'] = line.split()[0].split('_')[1]
10    df['time'] = line.split()[1][:4]
11    #rename the cols to same format as the above dataframe for easy concat
12    us_m.append(df)
13
14    #Germany
15    df = pd.read_csv(line)
16    df.columns = [line.lower().replace(' ', '_') for line in df.columns]
17    df = df[df['country'] == 'germany']
18
19    df = df[['country', 'total_confirmed', 'total_confirmed_today']]
20    df['date'] = line.split()[0].split('_')[1]
21    df['time'] = line.split()[1][:4]
22    #rename the cols to same format as the above dataframe for easy concat
23    germany_m.append(df)
24
25 us_data = pd.concat(us_m).reset_index(drop=True)
26 us_data = number(us_data)
27 germany_data = pd.concat(germany_m).reset_index(drop=True)
28 germany_data = number(germany_data)
```

In [17]:

```
1 #Total for the month of April will be last date's record minus the first date's
2 germany_may_start = germany_data.total_confirmed.tolist()[0]
3 germany_may_end = germany_data.total_confirmed.tolist()[-1]
4
5 #Style2 Going with cases confirmed per day
6 germany_may = germany_data.total_confirmed_today.sum()
7
8 print(f"Total cases for Germany in May is {germany_may_end - germany_may_start:,}")
9
10 #Since there are 13 items on the list, total for the month of April will be last date's
11 us_may_start = us_data.total_confirmed.tolist()[0]
12 us_may_end = us_data.total_confirmed.tolist()[-1]
13
14 #Style2 Going with cases confirmed per day
15 us_may = us_data.total_confirmed_today.sum()
16
17 print(f"Total cases for US in May is {us_may_end - us_may_start:,}")
18 print(f"\nPer 'Cases Confirmed Today' col:\nGermany May confirmed> {germany_may:,},\nUS April confirmed> {us_may:,}")
```

Total cases for Germany in May is 16,251

Total cases for US in May is 555,285

Per 'Cases Confirmed Today' col:

Germany May confirmed> 17,319.

US April confirmed> 585,611.

Question 4: Using 1c above (covid* other time), how many cases did Italy and Russia each record in April?

- Same flow as 2 and 3 minus converting the numbers. Our numbers are good here
- The column of interest here is 'infected'

In [18]:

```
1 italy, russia = [], []
2 for line in cov_Aprother:
3     #Italy
4     df = pd.read_csv(line)
5     df.columns = [line.lower().replace(' ', '_') for line in df.columns]
6     df = df[df['countries'] == 'Italy']
7
8     #Take cols we need
9     df = df[['countries', 'infected']]
10    df['date'] = line.split()[0].split('_')[1]
11    df['time'] = line.split()[1][:-4]
12    italy.append(df)
13
14    #Russia
15    df = pd.read_csv(line)
16    df.columns = [line.lower().replace(' ', '_') for line in df.columns]
17    df = df[df['countries'] == 'Russia']
18
19    #Take cols we need
20    df = df[['countries', 'infected']]
21    df['date'] = line.split()[0].split('_')[-1]
22    df['time'] = line.split()[1][:-4]
23    russia.append(df)
24
25    italy_data = pd.concat(italy).reset_index(drop=True)
26    russia_data = pd.concat(russia).reset_index(drop=True)
```

In [19]:

```
1 #total for the month of April will be last date's record minus the first date's
2 italy_april_start = italy_data.infected.tolist()[0]
3 italy_april_end = italy_data.infected.tolist()[-1]
4
5 print(f"Total cases for Italy in April is {italy_april_end - italy_april_start:,}")
6
7 #total for the month of April will be last date's record minus the first date's
8 russia_april_start = russia_data.infected.tolist()[0]
9 russia_april_end = russia_data.infected.tolist()[-1]
10
11 print(f"Total cases for Russia in April is {russia_april_end - russia_april_start:,}")
```

Total cases for Italy in April is 88,349

Total cases for Russia in April is 102,950

Question 5: Using 1c above (covid* other time), how many cases did Italy and Russia each record in May?

- Same logic/flow from question 4 and working with May

In [20]:

```
1 italy, russia = [], []
2 for line in cov_Mayother:
3     #Italy
4     df = pd.read_csv(line)
5     df.columns = [line.lower().replace(' ', '_') for line in df.columns]
6     df = df[df['countries'] == 'Italy']
7
8     #Take cols we need
9     df = df[['countries', 'infected']]
10    df['date'] = line.split()[0].split('_')[1]
11    df['time'] = line.split()[1][-4]
12    italy.append(df)
13
14    #Russia
15    df = pd.read_csv(line)
16    df.columns = [line.lower().replace(' ', '_') for line in df.columns]
17    df = df[df['countries'] == 'Russia']
18
19    #Take cols we need
20    df = df[['countries', 'infected']]
21    df['date'] = line.split()[0].split('_')[-1]
22    df['time'] = line.split()[1][-4]
23    russia.append(df)
24
25    italy_data = pd.concat(italy).reset_index(drop=True)
26    # italy_data = number_v2(italy_data)
27    russia_data = pd.concat(russia).reset_index(drop=True)
28    # russia_data = number_v2(russia_data)
```

In [21]:

```
1 #Since there are 13 items on the list, total for the month of April will be last day
2 italy_may_start = italy_data.infected.tolist()[0]
3 italy_may_end = italy_data.infected.tolist()[-1]
4
5 print(f"Total cases for Italy in May is {italy_may_end - italy_may_start:,}")
6
7 #Since there are 13 items on the list, total for the month of April will be last day
8 russia_may_start = russia_data.infected.tolist()[0]
9 russia_may_end = russia_data.infected.tolist()[-1]
10
11 print(f"Total cases for Russia in April is {russia_may_end - russia_may_start:,}")
```

Total cases for Italy in May is 24,395

Total cases for Russia in April is 238,996

Kicker

The Flow

- We make a list of the countries we are interested in (Notice I added some extra countries plus my country Nigeria)
- We define a kicker() taking 3: the list of countries we are looking for, the list of apr file and list of may files. Remember them from somewhere before.
- Same idea of reading each file to find the country we want, put the df from each file into a list of dfs, concat them and get the list. We do this for each of the months.

- Finally we add both months to get our total tally for the country.

In [22]:

```
1 def kicker(target_list, apr=ww_Apr23, may=ww_May23):
2     for target in target_list: #For each country on on List
3         #make an empty list for each month to hold the df of numbers from each file
4         data_apr = []
5         data_may = []
6
7         #For each file in the month of april
8         for day in apr:
9             df = pd.read_csv(day)
10            df.columns = [line.lower().replace(' ', '_') for line in df.columns]
11            df = df[df['country'] == target]
12            try:
13                #For confirmed col names
14                df = df[['country', 'total_confirmed_cases']]
15                data_apr.append(df)
16            except:
17                #For confirmed col names 'total_confirmed_cases'
18                df = df[['country', 'total_confirmed']]
19                #rename the cols to same format as the above dataframe for easy conc
20                df.columns = ['country', 'total_confirmed_cases']
21                data_apr.append(df)
22
23            #Concat the files in the list of dfs
24            c_data_apr = pd.concat(data_apr).reset_index(drop=True) #reset index and drop
25
26            #convert the col holding our numbers replacing commas and making them ints
27            c_data_apr['total_confirmed_cases']=[int(line.replace(',','')) for line in c_data_apr['total_confirmed_cases']]
28            c_data_apr_start = c_data_apr.total_confirmed_cases.tolist()[0] #get first
29            c_data_apr_end = c_data_apr.total_confirmed_cases.tolist()[-1] #get last
30            c_data_apr_total = c_data_apr_end - c_data_apr_start #Get their difference
31
32            #Repeat for the month of may
33            for day in may:
34                df = pd.read_csv(day)
35                df.columns = [line.lower().replace(' ', '_') for line in df.columns]
36                df = df[df['country'] == target]
37                try:
38                    #For confirmed col names
39                    df = df[['country', 'total_confirmed_cases']]
40                    data_may.append(df)
41                except:
42                    #For confirmed col names 'total_confirmed_cases'
43                    df = df[['country', 'total_confirmed']]
44                    #rename the cols to same format as the above dataframe for easy conc
45                    df.columns = ['country', 'total_confirmed_cases']
46                    data_may.append(df)
47
48            #Repeat the concat and number column conversion
49            c_data_may = pd.concat(data_may).reset_index(drop=True)
50            c_data_may['total_confirmed_cases']=[int(line.replace(',','')) for line in c_data_may['total_confirmed_cases']]
51            c_data_may_start = c_data_may.total_confirmed_cases.tolist()[0]
52            c_data_may_end = c_data_may.total_confirmed_cases.tolist()[-1]
53            c_data_may_total = c_data_may_end - c_data_may_start
54
55            #print the result. Remember we are in a loop of listed countries. The first
56            print(f"Total confirmed for {target} in April and May is {c_data_apr_total + c_data_may_total}")
```

In [23]:

```
1 us = []
2 germany = []
3 china = []
4 italy = []
5 russia = []
6 targets = ['united-states', 'germany', 'china', 'italy', 'russia', 'brazil', 'india']
7 #add more countires on the list as you like
```

In [24]:

```
1 #We call the kicker function passing as target list of countries 'targets'
2 kicker(targets)
```

Total confirmed for united-states in April and May is 1,407,107
Total confirmed for germany in April and May is 94,466
Total confirmed for china in April and May is 623
Total confirmed for italy in April and May is 112,651
Total confirmed for russia in April and May is 333,000
Total confirmed for brazil in April and May is 348,572
Total confirmed for india in April and May is 133,598
Total confirmed for nigeria in April and May is 7,417