

WEEK8 Assignment-1:Data Wrangling with Pandas

Data Wrangling with Pandas Assignment problems

Problem-1

In [ ]: The European Centre for Disease Prevention and Control (ECDC) provides an open dataset on COVID-19 cases called, daily number of new reported cases of COVID-19 by country worldwide. This dataset is updated daily, but we will use a snapshot that contains data from January 1, 2020 through September 18, 2020. Clean and pivot the data so that it is in wide format:  
(Get covid19\_cases.csv file using this link:  
https://raw.githubusercontent.com/svkarthik86/Advanced-python/main/WEEK-8%20Assignment/covid19\_cases.csv )

In [ ]: 1. Read in the covid19\_cases.csv file.  
2. Create a date column using the data in the dateRep column and the pd.to\_datetime() function.  
3. Set the date column as the index and sort the index.  
4. Replace occurrences of United\_States\_of\_America and United\_Kingdom with USA and UK, respectively.  
5. Using the countriesAndTerritories column, filter the data down to Argentina, Brazil, China, Colombia, India, Italy, Mexico, Peru, Russia, Spain, Turkey, the UK, and the USA.  
6. Pivot the data so that the index contains the dates, the columns contain the country names, and the values are the case counts in the cases column. Be sure to fill in NaN values with 0.

countriesAndTerritories	Argentina	Brazil	China	Colombia	India	Italy	Mexico	Peru	Russia	Spain	Turkey	UK	USA
date													
2020-01-01	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
2020-01-02	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
2020-01-03	0.0	0.0	17.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
2020-01-04	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
2020-01-05	0.0	0.0	15.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
...	...	...	...	...	...	...	...	...	...	...	...	...	...
2020-09-14	10778.0	14768.0	29.0	7355.0	92071.0	1456.0	4408.0	6787.0	5449.0	27404.0	1527.0	3330.0	33871.0
2020-09-15	9056.0	15155.0	22.0	5573.0	83809.0	1008.0	3335.0	4241.0	5509.0	9437.0	1716.0	2621.0	34841.0
2020-09-16	9908.0	36653.0	24.0	6698.0	90123.0	1229.0	4771.0	4160.0	5529.0	11193.0	1742.0	3103.0	51473.0
2020-09-17	11893.0	36820.0	7.0	7787.0	97894.0	1452.0	4444.0	6380.0	5670.0	11291.0	1771.0	3991.0	24598.0
2020-09-18	11674.0	36303.0	44.0	7568.0	96424.0	1583.0	3182.0	5698.0	5762.0	14389.0	1648.0	3395.0	43567.0

262 rows × 13 columns

In [1]: import pandas as pd  
url="https://raw.githubusercontent.com/svkarthik86/Advanced-python/main/WEEK-8%20Assignment/covid19\_cases.csv"  
covid\_data=pd.read\_csv(url)  
covid\_data.dateRep=pd.to\_datetime(covid\_data[["year", "month", "day"]])  
covid\_data.pivot(index="dateRep",columns="countriesAndTerritories",values="cases").fillna(0)  
covid\_data.pivot.columns=covid\_data.pivot.columns.str.replace("United\_States\_of\_America", "USA")  
covid\_data.pivot.columns=covid\_data.pivot.columns.str.replace("United\_Kingdom", "UK")  
covid\_data.pivot.filter(["Argentina", "Brazil", "China", "Colombia", "India", "Italy", "Mexico", "Peru", "Russia", "Spain", "Turkey", "UK", "USA"])

countriesAndTerritories	Argentina	Brazil	China	Colombia	India	Italy	Mexico	Peru	Russia	Spain	Turkey	UK	USA
dateRep													
2020-01-01	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
2020-01-02	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
2020-01-03	0.0	0.0	17.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
2020-01-04	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
2020-01-05	0.0	0.0	15.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
...	...	...	...	...	...	...	...	...	...	...	...	...	...
2020-09-14	10778.0	14768.0	29.0	7355.0	92071.0	1456.0	4408.0	6787.0	5449.0	27404.0	1527.0	3330.0	33871.0
2020-09-15	9056.0	15155.0	22.0	5573.0	83809.0	1008.0	3335.0	4241.0	5509.0	9437.0	1716.0	2621.0	34841.0
2020-09-16	9908.0	36653.0	24.0	6698.0	90123.0	1229.0	4771.0	4160.0	5529.0	11193.0	1742.0	3103.0	51473.0
2020-09-17	11893.0	36820.0	7.0	7787.0	97894.0	1452.0	4444.0	6380.0	5670.0	11291.0	1771.0	3991.0	24598.0
2020-09-18	11674.0	36303.0	44.0	7568.0	96424.0	1583.0	3182.0	5698.0	5762.0	14389.0	1648.0	3395.0	43567.0

262 rows × 13 columns

Problem-2

In [ ]: In order to determine the case totals per country efficiently, we need the aggregation skills , so the ECDC data in the covid19\_cases.csv file has been aggregated for us and saved in the covid19\_total\_cases.csv file. It contains the total number of case per country. Use this data to find the 20 countries with the largest COVID-19 case totals. Hints:  
(Get covid19\_total\_cases.csv file using this  
link: https://raw.githubusercontent.com/svkarthik86/Advanced-python/main/WEEK-8%20Assignment/covid19\_total\_cases.csv)

In [ ]: \* When reading in the CSV file, pass in index\_col='index'.  
\* Note that it will be helpful to transpose the data before isolating the countries.

index	cases
USA	6724667
India	5308014
Brazil	4495183
Russia	1091186
Peru	756412
Colombia	750471
Mexico	688954
South_Africa	657627
Spain	640040
Argentina	601700
Chile	442827
France	428696
Iran	416198
UK	385936
Bangladesh	345805
Saudi_Arabia	328720
Iraq	311690
Pakistan	305031
Turkey	299810
Italy	294932

In [2]: covid\_cases=pd.read\_csv("https://raw.githubusercontent.com/svkarthik86/Advanced-python/main/WEEK-8%20Assignment/covid19\_total\_cases.csv",index\_col="index").T  
covid\_cases.sort\_values(by="cases", ascending=False)[:20]

index	cases
USA	6724667
India	5308014
Brazil	4495183
Russia	1091186
Peru	756412
Colombia	750471
Mexico	688954
South_Africa	657627
Spain	640040
Argentina	601700
Chile	442827
France	428696
Iran	416198
UK	385936
Bangladesh	345805
Saudi_Arabia	328720
Iraq	311690
Pakistan	305031
Turkey	299810
Italy	294932