

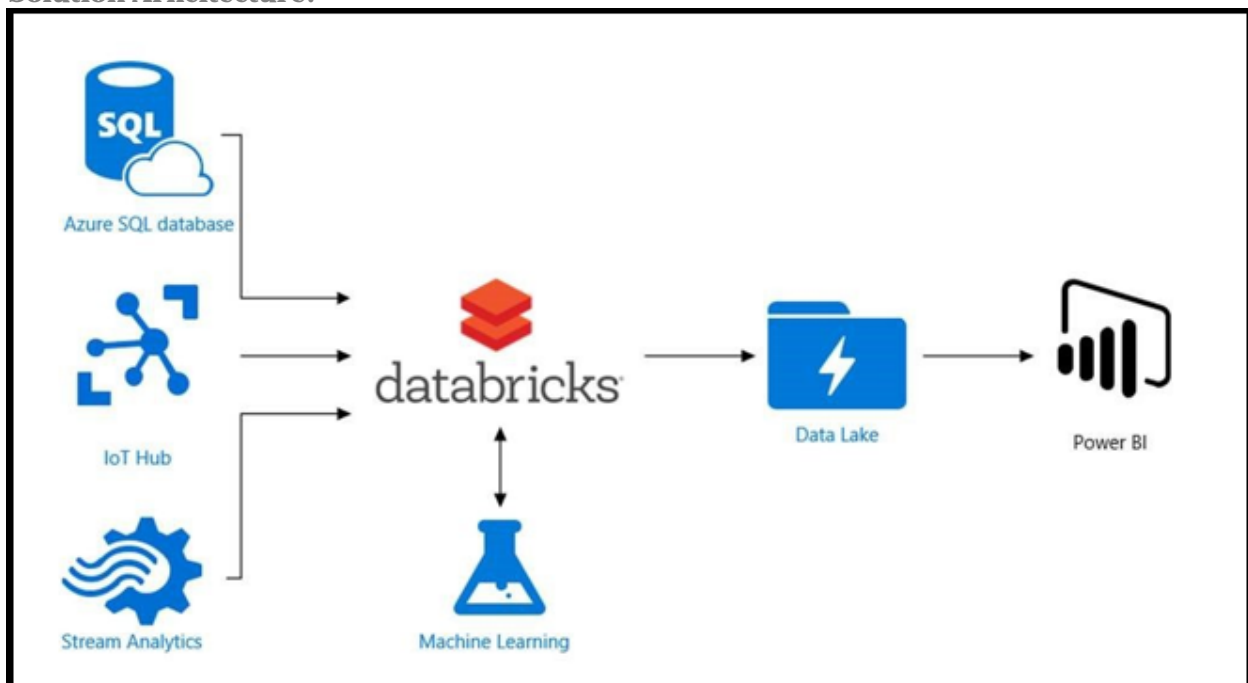
COVID-19 Data Analytics and Reporting with Azure Databricks and Power BI using Python

Problem Statement

In the midst of all of the bad news in COVID era, there are good news reports of the important work done by data engineers and data scientists to get the data of all COVID cases around the world together and provide useful insights which will be useful for pan professionals.

The objective of this article is to focus on a use case that demonstrates the integration between daily changing Source, Azure Databricks and Power BI to deliver insights and data visualizations using a publicly available COVID-19 dataset. While Azure Databricks provides the distributed computing power to process and transform complex datasets, Power BI is a fitting recipient of the transformed dataset that surfaces these insights to business users.

Solution Architecture:



Source

The latest available public data on the geographic distribution of COVID-19 cases worldwide from the [European Center for Disease Prevention and Control \(ECDC\)](#). Each row/entry contains the number of new cases reported per day and per country or region.

Before we start with our exercise, we will need to have the following prerequisites:

1. You need to have an active **Azure Subscription**.

2. **Azure Databricks** – You need to set up both Databricks service and cluster in Azure, you can go over the steps in this article. As shown in this article, we have created a Databricks service.

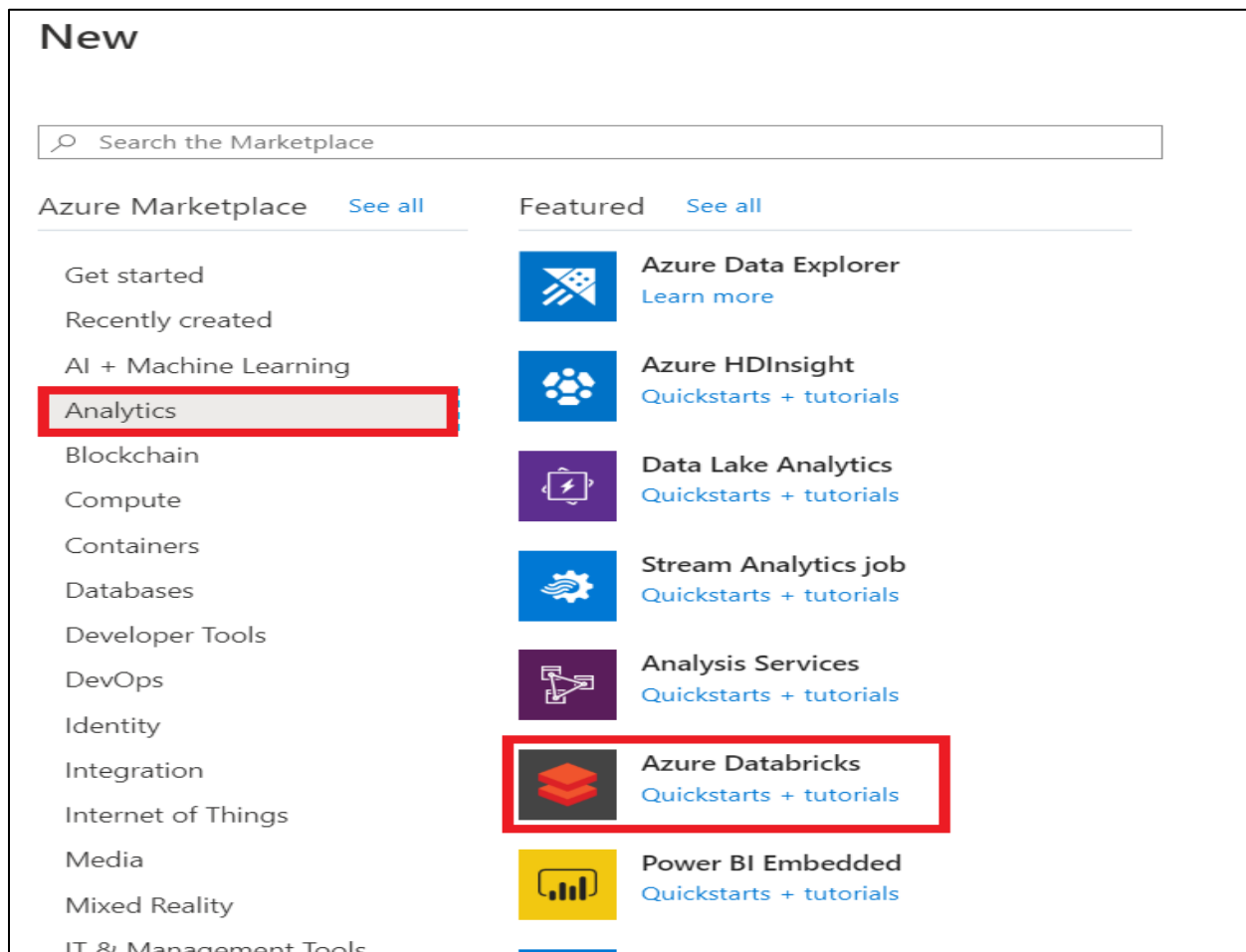
3. Power BI Subscription

Let's go ahead and demonstrate the data load into Databricks using Python notebooks from the source mentioned above.

STEP 1: Create Azure Databricks

The next step is to create a [Databricks Workspace](#). You can think of the workspace like an application that you are installing within Azure, where you will access all of your Databricks assets. Follow these steps to create a workspace:

- On the Azure home screen, click **'Create a Resource'**.
- In the **'Search the Marketplace'** search bar, type **'Databricks'** and you should see **'Azure Databricks'** pop up as an option. Click that option.
- Click **'Create'** to begin creating your workspace.



Use the same resource group you created or selected earlier. Then, enter a workspace name. Remember to always stick to naming standards when creating Azure resources, but for now, enter whatever you would like.

You can keep the location as whatever comes default or switch it to a region closer to you. For the **pricing tier**, select **'Trial'**. Finally, select **'Review and Create'**. We can skip-networking and tags for now which are for more advanced set-ups.

* Basics Networking Tags Review + Create

Project Details

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all your resources.

Subscription *

Resource group * [Create new](#)

Instance Details

Workspace name *

Location *

Pricing Tier *

Review + Create Next : Networking >

- This should bring you to a validation page where you can click **'create'** to deploy your workspace. This will bring you to a deployment page and the creation of the workspace should only take a couple of minutes.
- Once the deployment is complete, click **'Go to resource'** and then click **'Launch Workspace'** to get into the Databricks workspace.

demo1 Azure Databricks Service

Search (Ctrl+/) Delete

Overview Status : Active

Managed Resource Group : databricks-rg-demo1-lukwowaSkcj

Activity log Resource group URL :

Access control (IAM) Location : West US Pricing Tier : Trial (Premium - 14-Days Free DBUs)

Tags Subscription :

Settings Subscription ID :

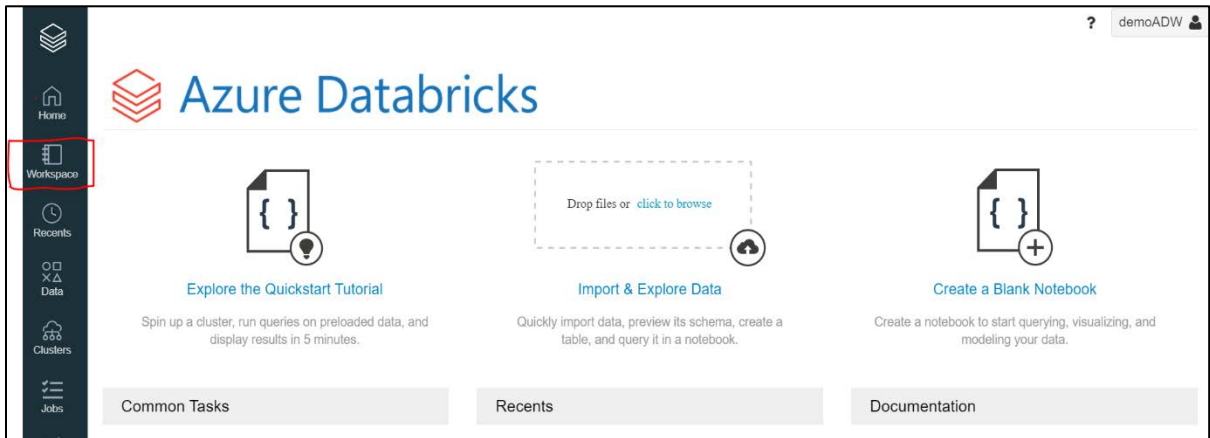
Virtual Network Peerings Tags (change) : [Click here to add tags](#)

Launch Workspace

Upgrade to Premium

Documentation Getting Started Import Data from File Import Data from Azure Storage

- Create a Workspace (Notebook) with Python or your choice of language for performing analysis.



STEP 2: Start Scripting in Notebook

- Import relevant libraries and fetch data from the source according to its data type.

```

1 import seaborn as sns
2 import plotly.express as px
3 %matplotlib inline
4 import pandas as pd
5 import numpy as np
6 %matplotlib inline
7 import matplotlib.pyplot as plt
8
9 df = pd.read_csv("https://pandemicdatalake.blob.core.windows.net/public/raw/covid-19/ecdc_cases/latest/ECDCcases.csv")
10 df.head(1000)

```

Out[1]:

	dateRep	day	month	year	cases	deaths	countriesAndTerritories	geoid	countryterritoryCode	popData2019	continentExp	Cumulative_number_for_14_days_of_COVID-19_cases_per_100000
0	08/11/2020	8	11	2020	126	6	Afghanistan	AF	AFG	38041757.0	Asia	3.656508
1	07/11/2020	7	11	2020	58	2	Afghanistan	AF	AFG	38041757.0	Asia	3.538217
2	06/11/2020	6	11	2020	40	0	Afghanistan	AF	AFG	38041757.0	Asia	3.546103
3	05/11/2020	5	11	2020	121	6	Afghanistan	AF	AFG	38041757.0	Asia	3.745884
4	04/11/2020	4	11	2020	86	4	Afghanistan	AF	AFG	38041757.0	Asia	3.782685

- Digging the data to get useful and relevant analytical data
 - Top Countries with cases and deaths.

```

top countries
4 51
1 import plotly.graph_objects as go
2 import plotly.express as px
3 import matplotlib.pyplot as plt
4
5 df.loc[:, ['countriesAndTerritories', 'cases', 'deaths']].groupby(['countriesAndTerritories']).max().sort_values(by='cases', ascending=False).reset_index()[0:15].style.background_gradient(cmap='rainbow')

```

Out[42]:

countriesAndTerritories	cases	deaths
0 United_States_of_America	184813	4928
1 India	97894	2003
2 France	88852	2004
3 Brazil	89074	1595
4 Spain	55019	1623
5 Russia	45188	745
6 Italy	40902	971
7 Chile	36179	1057
8 United_Kingdom	33470	1224
9 Poland	27875	603
10 Germany	23542	315
11 Belgium	22189	322
12 Switzerland	21842	197
13 Kazakhstan	19285	324
14 Argentina	18326	3351

Command took 0.65 seconds -- by ayush.chauhan@ffl.tech at 11/28/2020, 3:08:11 PM on democluster

- Calculating total cases, deaths, and death rate around the globe.

```

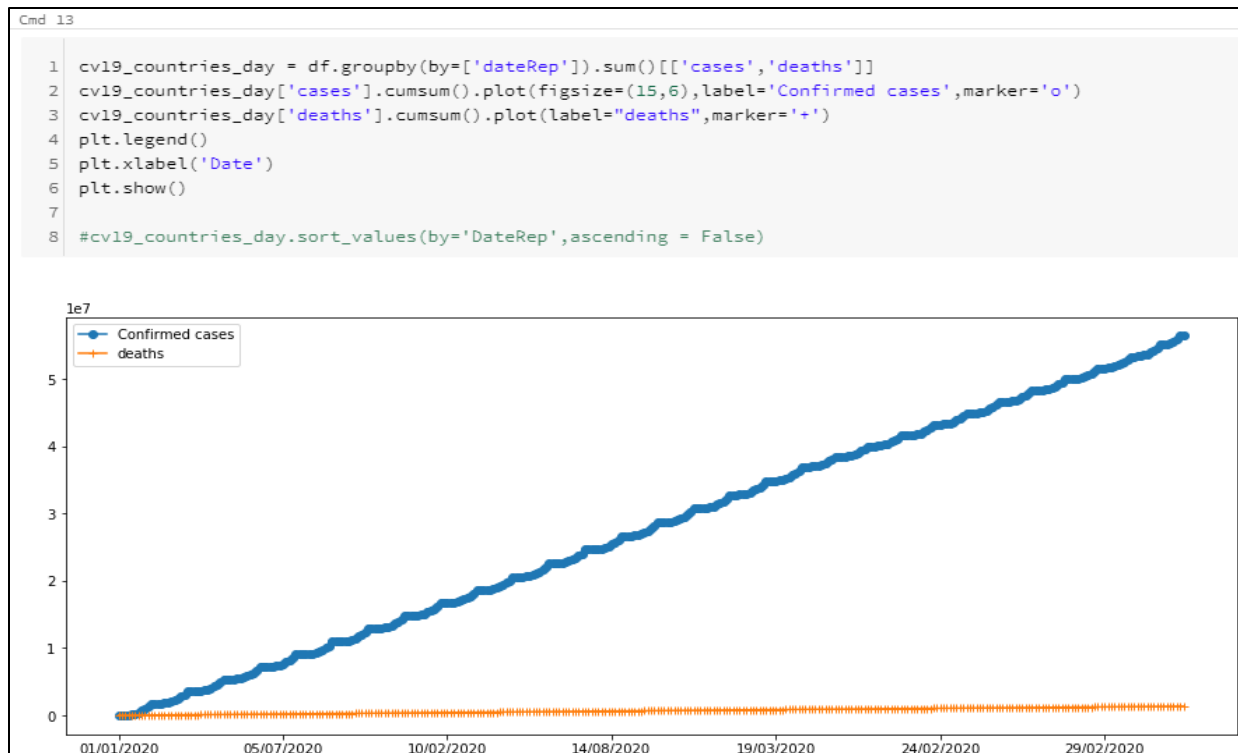
Cmd 8
1 cv19_countries_day = df.groupby(by=['dateRep', 'countriesAndTerritories']).sum()[['cases', 'deaths']]
2
3 Total_confirmed = cv19_countries_day.groupby('dateRep').sum()[['cases', 'deaths']].sum()['cases']
4 Total_deaths = cv19_countries_day.groupby('dateRep').sum()[['cases', 'deaths']].sum()['deaths']
5
6 dicc = {'TotalConfirmed' : Total_confirmed, 'TotalDeaths' : Total_deaths, 'DeathRate' : round((Total_deaths/Total_confirmed)*100,2)}
7 total = pd.DataFrame(dicc, index=['Counter'])[['TotalConfirmed', 'TotalDeaths', 'DeathRate']]
8
9
10 total.style.set_properties(**{
11     'background-color': 'white',
12     'font-size': '20pt',
13     'color' : 'red'
14 })

```

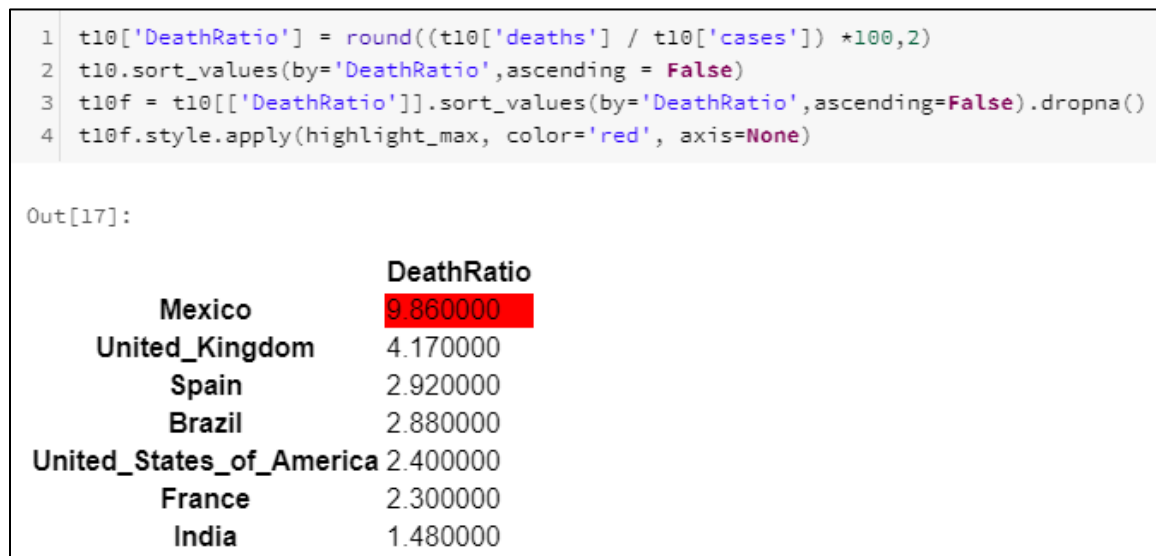
Out[6]:

	TotalConfirmed	TotalDeaths	DeathRate
Counter	56368586	1350713	2.400000

- Visualizing total cases and deaths with a timeline.



- Calculating Death Ratio for top countries.



- Calculating global change rate per day.

1. Global Rate Change for cases per day.



2. Global Rate Change for deaths per day.



- Calculating Change rate for top countries only.

```
Cmd 29
```

```
1 #Calculate change by country for the 15 first
2
3 Impacted_countries = df[['countriesAndTerritories','cases']].sort_values(by=['dateRep','cases'],ascending=False).head(15)['countriesAndTerritories']
4 Impacted_countries
5
6 top_impact = pd.DataFrame()
7
8 for country in Impacted_countries:
9     top_impact[country] = df[df['countriesAndTerritories']==country]['cases']
10
11
12 top_impact = top_impact.reset_index().sort_values(by='dateRep',ascending=True) #true
13
14
15 #Normalize
16
17 #top_impact_norm = top_impact/top_impact.iloc[0] * 100
```

Command took 0.10 seconds -- by ayush.chauhan@ifi.tech at 11/20/2020, 3:08:20 PM on democluster

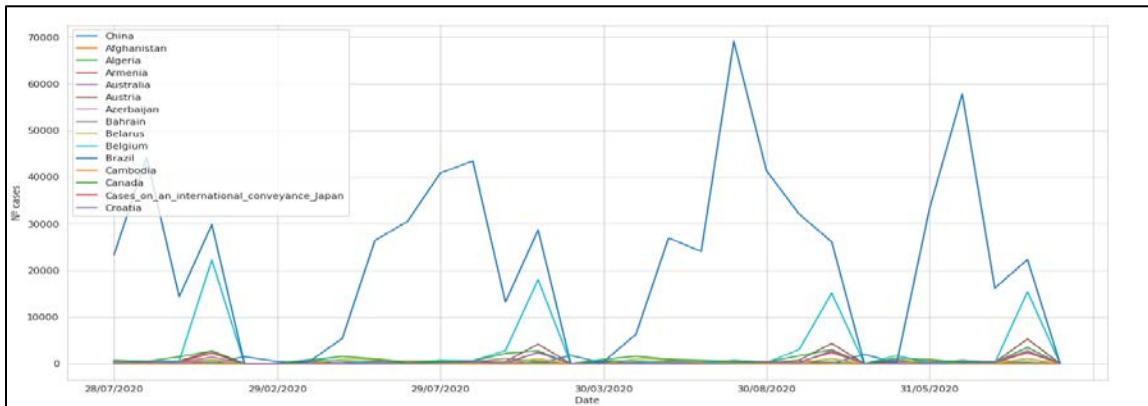
- Top countries Rate Change for cases.

CHANGE COVID RATE

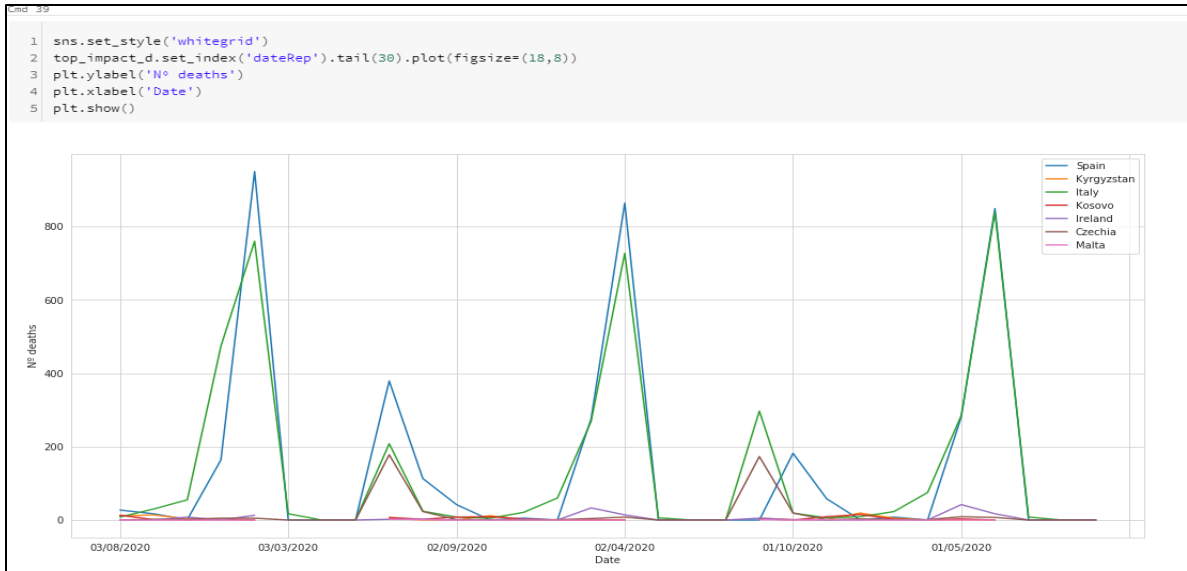
Cmd 26

```
1 covid19_change_global = cv19_countries_day.cumsum()
2 covid19_change_global[['Cases Day','Deaths Day']] = cv19_countries_day[['cases','deaths']]
3 covid19_change_global = covid19_change_global.pct_change(1)
4 covid19_change_global = covid19_change_global.sort_values(by='dateRep',ascending=False)
5 covid19_change_global = covid19_change_global.replace([np.inf, -np.inf], np.nan)
6 covid19_change_global = covid19_change_global.fillna(0)
7 covid19_change_global = round(covid19_change_global*100,2)
8 covid19_change_global = covid19_change_global.reset_index()
9
10 covid19_change_global_d = cv19_countries_day.cumsum()
11 covid19_change_global_d[['Cases Day','Deaths Day']] = cv19_countries_day[['cases','deaths']]
12 covid19_change_global_d = covid19_change_global_d.pct_change(1)
13 covid19_change_global_d = covid19_change_global_d.sort_values(by='dateRep',ascending=False)
14 covid19_change_global_d = covid19_change_global_d.replace([np.inf, -np.inf], np.nan)
15 covid19_change_global_d = covid19_change_global_d.fillna(0)
16 covid19_change_global_d = round(covid19_change_global_d*100,2)
17 covid19_change_global_d = covid19_change_global_d.reset_index()
18
```

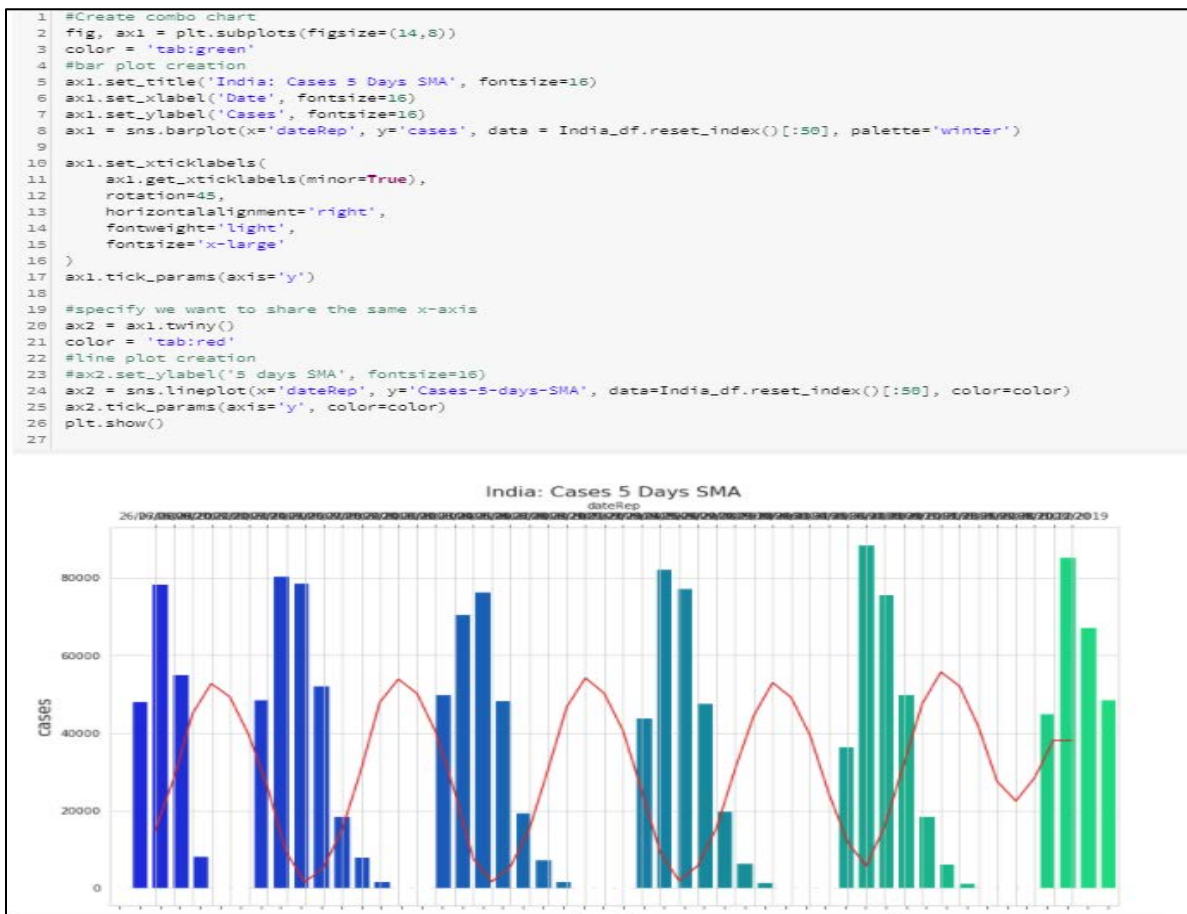
Command took 0.04 seconds -- by ayush.chauhan@ifi.tech at 11/20/2020, 3:08:20 PM on democluster



- Top countries Rate Change for deaths.



- Fetching Country particular data (e.g. INDIA).



- Fetching Continent particular data.

```

AFRICA

1 | africa_df = df[df['continentExp']=='Africa']

Command took 0.10 seconds -- by ayush.chauhan@ifi.tech at 11/20/2020, 3:08:21 PM on democuster

2 | display(africa_df)

(1) Spark Jobs

day month year cases deaths countriesAndTerritories ggold countryterritoryCode popData2019 continentExp Cumulative_number_for_14_days_of_COVID-19_cases_per_100000
1 19 11 2020 1038 20 Algeria DZ DZA 43953054 Africa 25.79578955
2 18 11 2020 1002 18 Algeria DZ DZA 43953054 Africa 24.64865791
3 17 11 2020 910 14 Algeria DZ DZA 43953054 Africa 23.2819967
4 16 11 2020 860 15 Algeria DZ DZA 43953054 Africa 21.84578562
5 15 11 2020 844 14 Algeria DZ DZA 43953054 Africa 20.61874635
6 14 11 2020 867 14 Algeria DZ DZA 43953054 Africa 19.33428462
7 13 11 2020 851 18 Algeria DZ DZA 43953054 Africa 18.06143648
8 12 11 2020 811 16 Algeria DZ DZA 43953054 Africa 16.79655248

```

- Save the relevant data frames as tables.

```

1 | top_countries =spark.createDataFrame(top_countries)

top_countries: pyspark.sql.dataframe.DataFrame = [countriesAndTerritories: string, cases: long ... 1 more fields]

Command took 0.09 seconds -- by ayush.chauhan@ifi.tech at 11/20/2020, 3:08:21 PM on democuster

Cmd 55

1 | top_countries.write.mode("overwrite").saveAsTable("top_countries")

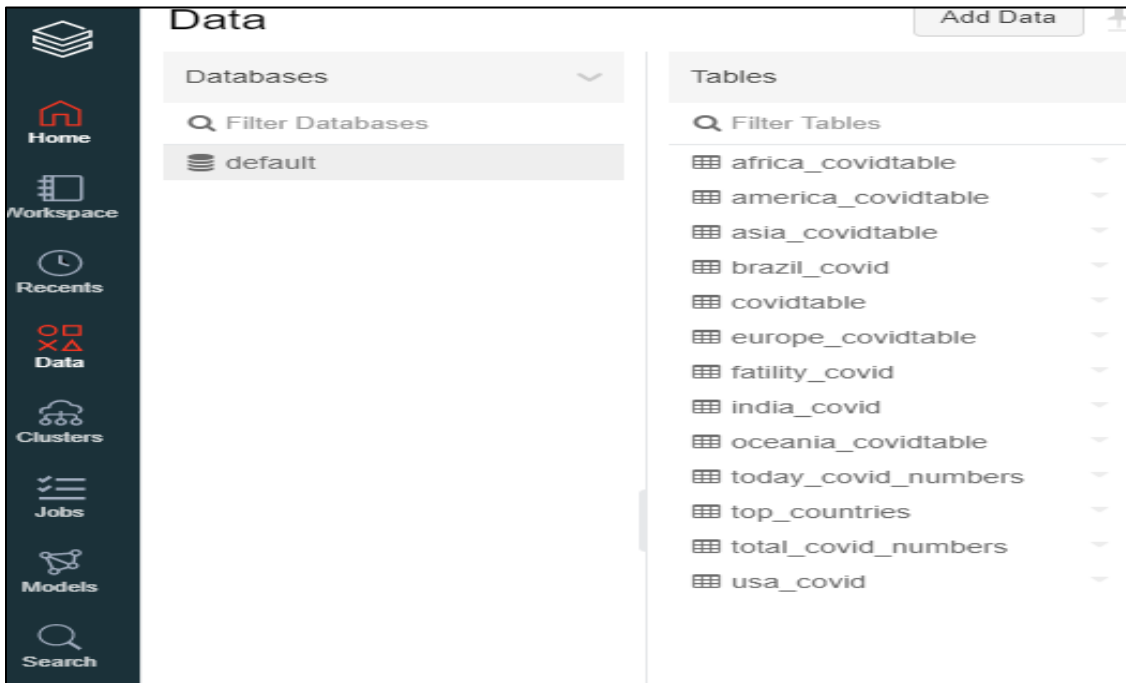
(1) Spark Jobs

Command took 1.11 seconds -- by ayush.chauhan@ifi.tech at 11/20/2020, 3:08:21 PM on democuster

Cmd 56

```

- Check the newly created tables by clicking 'DATA' in the left panel of data bricks workspace.



- Locate and store the credentials from data bricks clusters (Server Hostname & HTTP Path).

Clusters / democloud

democloud 📌 Edit Clone Restart Terminate Delete

Configuration **Notebooks (1)** Libraries Event Log Spark UI Driver Logs Metrics Apps Spar

Standard_DS3_v2 14.0 GB Memory, 4 Cores, 0.75 DBU

Driver Type
Standard_DS3_v2 14.0 GB Memory, 4 Cores, 0.75 DBU

▼ Advanced Options

Azure Data Lake Storage Credential Passthrough ? Available on Azure Databricks Premium [Learn more](#)

Enable credential passthrough for user-level data access

Spark **Tags** Logging Init Scripts JDBC/ODBC Permissions

Server Hostname
[Redacted]

Port
443

Protocol
HTTPS

HTTP Path
[Redacted]

JDBC URL ?

```
jdbc:spark://adb-729461757855538.18.azuredatabricks.net:443/default;transportMode=http;ssl=1;httpPath=sql/protocolv1/o/729461757855538/1120-093248-debut2;AuthMech=3;UID=token;PWD=<personal-access-token>
```

Step 3: Data brick Job Scheduling

The very import scope of this article is to capture the daily changing data, then we must run every other command for every day when a new data will be in our source to avoid this, we will automate the whole process and schedule a job to run our notebook with an interval of a day, so it will check for a new file and apply all the commands accordingly.

- Click on **(calendar like symbol)**, on the top-right corner of your notebook&click on '+New' button.



- Schedule a job according to your need and then click **'ok'**.

Create Schedule

Schedule

Every starting at :

Show Cron Syntax

A new cluster will be created each time this schedule runs. You can modify settings from the [Jobs page](#) once the job is created.

Step 4: Connecting Power Bi to your Databricks.

- Click on 'Get Data' and search for 'Azure data bricks' to add data bricks credentials (Server Hostname & HTTP Path)

Azure Databricks

Server Hostname

HTTP Path

Advanced Options (optional)

Data Connectivity mode Import DirectQuery

- Click on 'DATA' to get a view and validation for the tables in Power BI.

day	month	year	cases	deaths	countriesAndTerritories	geold	countryterritoryCode	popData2019	continentExp	Cumulative_number_for_14_days_of_COVID-
19	2	2020	0	0	Brazil	BR	BRA	211049515	America	
19	1	2020	0	0	Brazil	BR	BRA	211049515	America	
18	2	2020	0	0	Brazil	BR	BRA	211049515	America	
18	1	2020	0	0	Brazil	BR	BRA	211049515	America	
17	2	2020	0	0	Brazil	BR	BRA	211049515	America	
17	1	2020	0	0	Brazil	BR	BRA	211049515	America	
16	2	2020	0	0	Brazil	BR	BRA	211049515	America	
16	1	2020	0	0	Brazil	BR	BRA	211049515	America	
27	2	2020	0	0	Brazil	BR	BRA	211049515	America	
27	1	2020	0	0	Brazil	BR	BRA	211049515	America	
26	1	2020	0	0	Brazil	BR	BRA	211049515	America	
25	2	2020	0	0	Brazil	BR	BRA	211049515	America	
25	1	2020	0	0	Brazil	BR	BRA	211049515	America	

- Start making insightful reports by applying different filters and visuals, a sample report to get an overview.



Here you can select Individual countries or group of them to assign them with a particular date with the help of slicers to view analytics accordingly.

This is the end of this article, we successfully automate the process of Loading, Transforming and Writing the data process from a web source and then migrating the processed data to Power BI, it will also read new data and apply all the transformations on it with an interval that was defined by the developer.

Thank you and feel free to try out some different transformations in data bricks and create some awesome visuals with Power BI.

Challenges Faced:

Bringing ML to production is hard.

Most Business Intelligence (BI) workloads are restricted to a fraction of the data that an organization collects. Most of these are done against data warehouses which aggregate and transform data for specific use cases. This can be great for a lot of common reporting and known analysis, but it can also restrict an organization's ability to form new and valuable insights.

Business Benefits :

It can be used to help Medical Hospitals and Doctors to keep a track of patient history Also with National and International Medical Authorities to analyze the region's covid statistics

Author: Ayush Chauhan, Associate Data Engineer