Imran Iqbal      Microsoft Azure Profile Project

This is a Microsoft Azure Profile Project created by me and the idea is to upload Covid.xlsx in
the datalake storage gen 2 and create a pipeline in Azure synapse to convert xlsx format to
parquet.

(Parquet is optimized to work with complex data in bulk and features different ways for
efficient data compression and encoding types. This approach is best especially for those
queries that need to read certain columns from a large table)

Then using Azure synapse, Apache spark pool is created that would allocate resource to notebook
which would be used to read the parquet file , create database and tables to perform aggregation
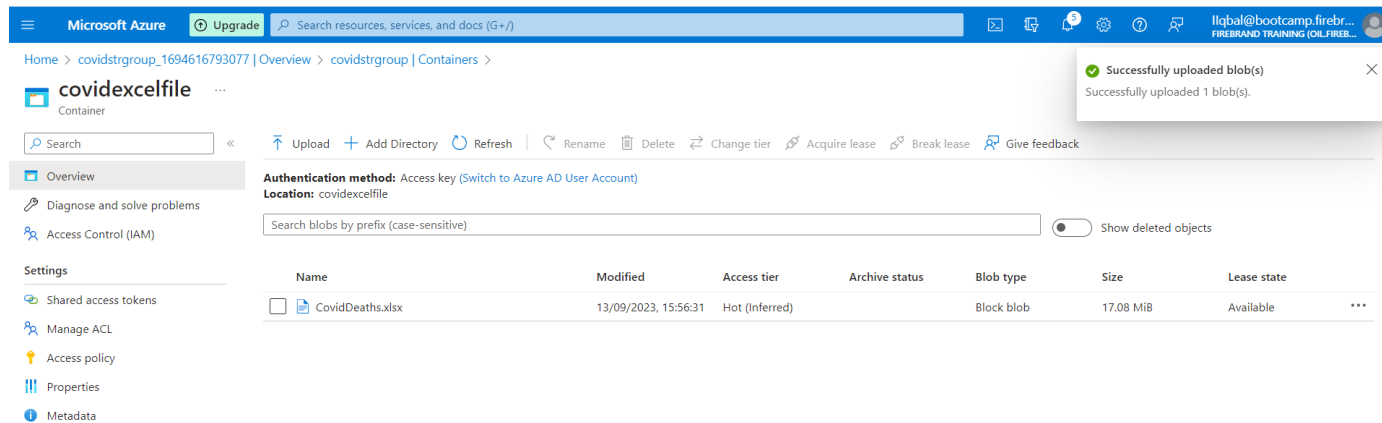functions on the file which is then saved on lake database.

First created a storage account and enabled it to be Azure data lake Gen 2.

Imran Iqbal       Microsoft Azure Profile Project

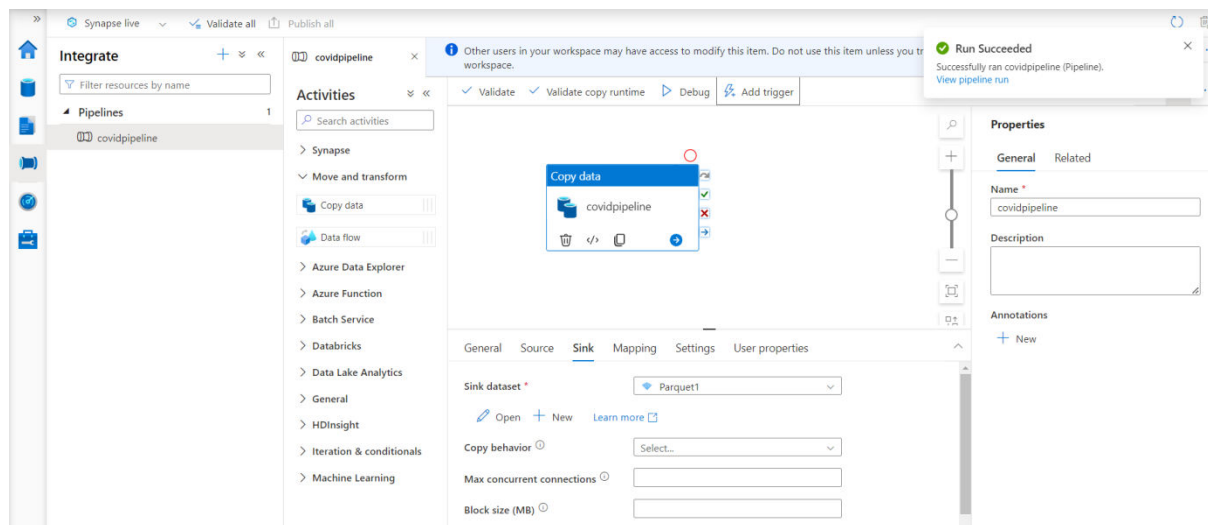Uploaded Covid.xlsx which is  excel format into a folder



Then created a synapse space to run a pipeline which converts excel format of a file to parquet as parquet is columnar based format and best used for query purpose.
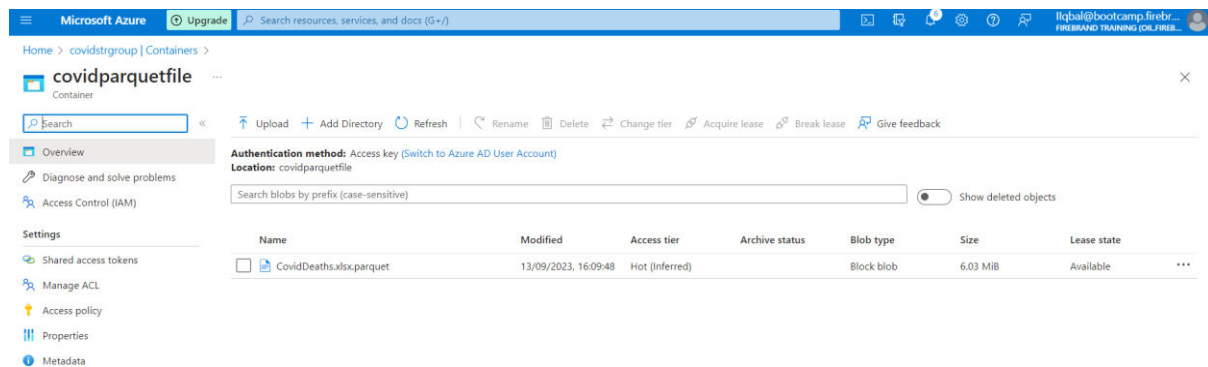
To accomplish this Synapse workspace is created and assigned myself (the Storage Blob Data Contributor role on the Data Lake Storage Gen2 account to interactively query it in the workspace).

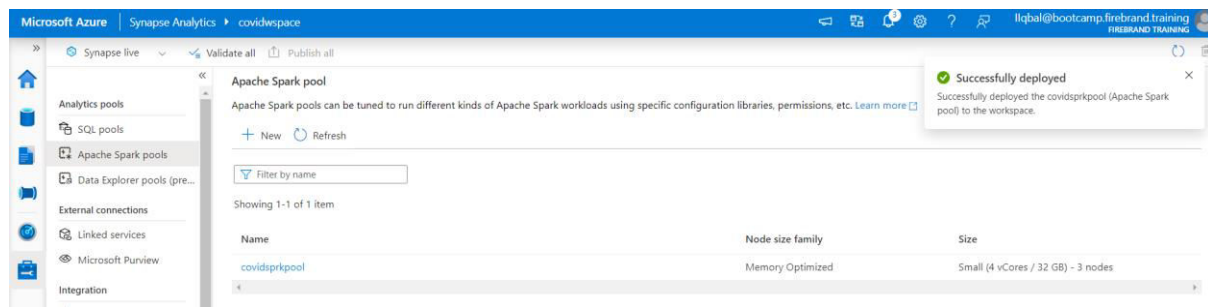Imran Iqbal    Microsoft Azure Profile Project

Furthermore, run a pipeline successfully to change excel file format to parquet



Parquet file is loaded in parquet folder as a result of pipeline trigger

Imran Iqbal     Microsoft Azure Profile Project

Memory optimized Apache spark pool created  in order to use the notebook



Parquet file loaded to dataframe in order to explore and use aggregate function

All the data from CovidDeaths parquet file is loaded and displayed into the table using %pyspark



Schema is displayed using pyspark for CovidDeath file

This data from dataframe is loaded into spark database by creating a database called "coviddb" and a table called coviddetails is created and data is loaded into the table
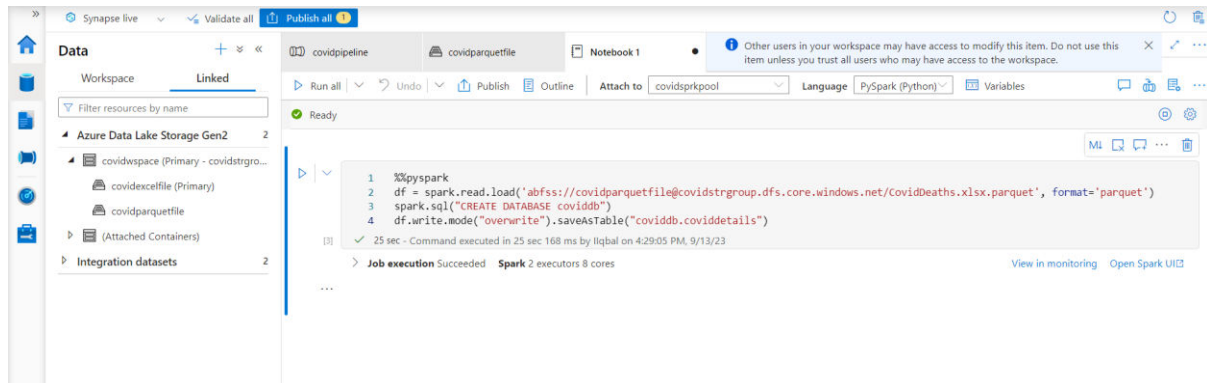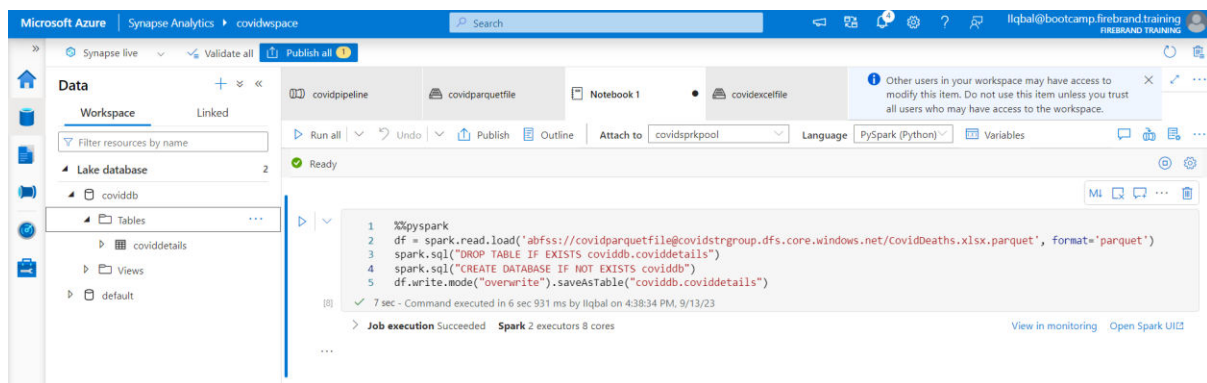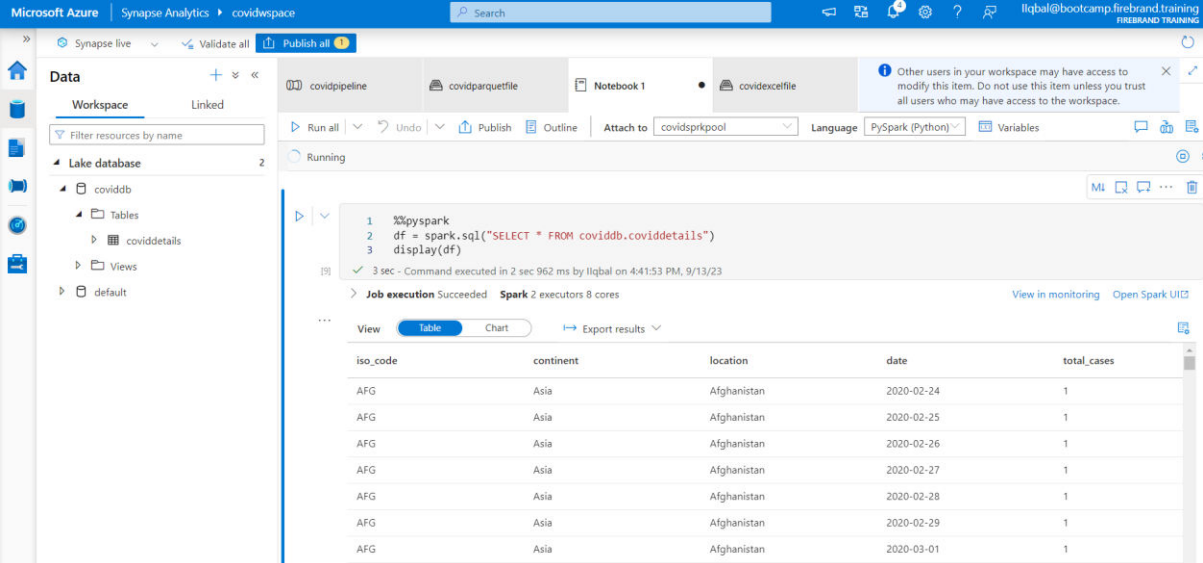


Table coviddetails created successfully in lakedatabase

Imran Iqbal     Microsoft Azure Profile Project

Coviddetails table in the coviddb database is loaded using in notebook

Aggregation is applied to display which continents have most deaths and Africa is the top continent which was displayed using group by and order by function.
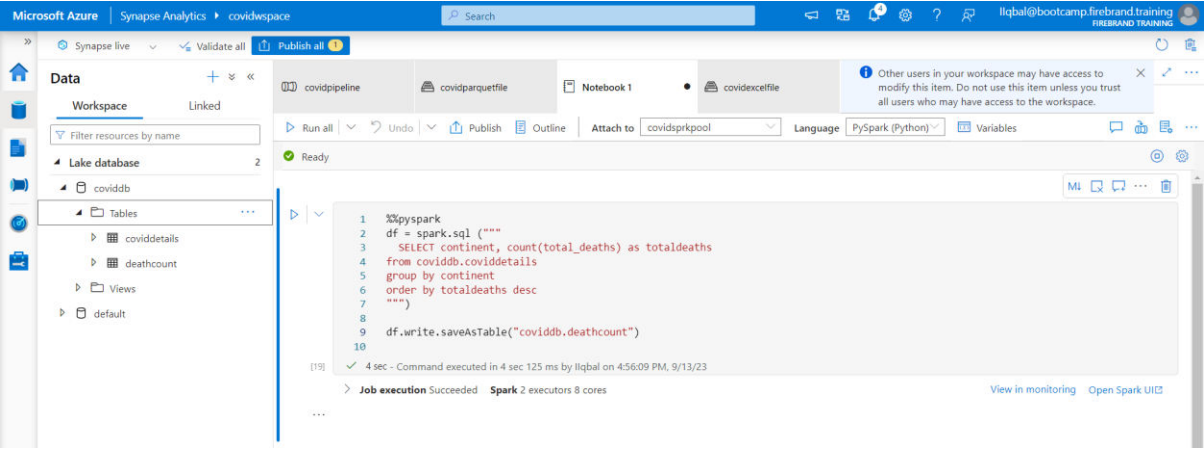


Aggregated table deathcount is created and saved in the database as shown