

# Create an Azure Databricks workspace ...



- Basics**
- Networking
- Encryption
- 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 \*

Azure for Students

▼

Resource group \*

covid-reporting-rg

▼

[Create new](#)

## Instance Details

Workspace name \*

covid-reporting-dbk-ws

✓

Region \*

UK South

▼

Pricing Tier \*

Standard (Apache Spark, Secure with Azure AD)

▼

Managed Resource Group name

Enter name for managed resource group



## covid-reporting-cluster

☐ Multi node ☒ Single node

Access mode

Single user access

Single user



Chinta, Avinash



### Performance

Databricks runtime version

Runtime: 12.2 LTS (Scala 2.12, Spark 3.3.2)



☐ Use Photon Acceleration

Node type

Standard\_D4a\_v4

16 GB Memory, 4 Cores



☒ Terminate after  minutes of inactivity

### Tags

Add tags

Key

Value

Add

> Automatically added tags

Create Cluster

Cancel

UI | [JSON](#)

### Summary

1 Driver

16 GB Memory, 4 Cores

Runtime

12.2.x-scala2.12

Standard\_D4a\_v4

0.75 DBU/h

Compute > UI preview [Provide feedback](#)

# covid-reporting-cluster

[More](#)[Terminate](#)[Edit](#)[Configuration](#)[Notebooks \(0\)](#)[Libraries](#)[Event log](#)[Spark UI](#)[Driver logs](#)[Metrics](#)[Apps](#)[Spark cluster UI - Master](#)☐ Multi node ☒ Single node

Access mode

Single user access

Single user

Chinta, Avinash

## Performance

Databricks Runtime Version

12.2 LTS (includes Apache Spark 3.3.2, Scala 2.12)

☐ Use Photon Acceleration 

Node type

Standard\_DS3\_v2

14 GB Memory, 4 Cores

☒ Terminate after  minutes of inactivity 

## Tags

No custom tags

&gt; Automatically added tags

▶ Advanced options

[UI](#) | [JSON](#)

## Summary

1 Driver

14 GB Memory, 4 Cores

Runtime

12.2.x-scala2.12

Standard\_DS3\_v2

0.75 DBU/h



mount\_storage

Python

File Edit View Run Help Last edit was 7 minutes ago Provide feedback

Run all

covid-reporting-cluster

Schedule

Share

#Linking the datalake storage to workspace by using access keys

Cmd 4

```
1 storage_account_name = "ercovidreportingdatalake"
2 storage_account_key = "9S43wAC0w0hKl9Nsu6USZtFKuAo5kmjyVZA+qh30QVZPWxUVRyqmR5wBsyaovbVDMX7ZrRk8qxLD+Ast3rQ0/A=="
3
4 spark.conf.set(
5     f"fs.azure.account.key.ercovidreportingdatalake.dfs.core.windows.net",
6     f"9S43wAC0w0hKl9Nsu6USZtFKuAo5kmjyVZA+qh30QVZPWxUVRyqmR5wBsyaovbVDMX7ZrRk8qxLD+Ast3rQ0/A==")
7
8 container_name = "raw" # make sure to use the right container
9 df_raw_population = spark.read.csv(f"abfss://raw@ercovidreportingdatalake.dfs.core.windows.net/population", sep=r'\t', header=True)
```

(1) Spark Jobs

df\_raw\_population: pyspark.sql.dataframe.DataFrame = [indic\_de,geo\time: string, 2008 : string ... 11 more fields]

Command took 8.63 seconds -- by avinashchinta@my.unt.edu at 7/2/2023, 1:11:24 PM on covid-reporting-cluster

Cmd 5

1 display(df\_raw\_population.head(5))

(3) Spark Jobs

Table +

	indic_de,geo\time	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	
1	PC_Y0_14,AD	14.6	14.5	14.5	15.5	15.5	15.5	:	:	:	:	:	1.
2	PC_Y0_14,AL	24.1	23.3	22.5	21.6	20.7	20.1	19.6	19.0	18.5	18.2	17.7	1'
3	PC_Y0_14,AM	19.0	18.6	18.3	:	:	:	:	19.4	19.6	20.0	20.2	2i

New

Workspace

Repos

Recents

Data

Compute

Workflows

Marketplace

1/3 Tasks Completed

Enable new UI

Menu options

## Workspace

covid

setup

transform

transform

transform\_population\_data

Home

Run all

covid-reporting-cluster

Schedule

Share

Python

xUVRyqmR5wBsyaoVbVDMX7ZrRk8qx1D+ASt3rQ0/A=="

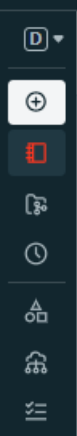
.net",  
VDMX7ZrRk8qx1D+ASt3rQ0/A==")

alake.dfs.core.windows.net/population", sep=r'\t', header=True)

... 11 more fields]

-reporting-cluster

2012	2013	2014	2015	2016	2017	2018	2
15.5	15.5	:	:	:	:	:	1.
20.7	20.1	19.6	19.0	18.5	18.2	17.7	1'



## Replace storage account name with your storage account name before executing.

Cmd 3

```
1 from pyspark.sql.functions import *
```

Cmd 4

Connect datalake storage and read csv file

Cmd 5

```
1 storage_account_name = "ercovidreportingdatalake"
2 storage_account_key = "9S43wAC0w0hKl9Nsu6USZtFKuAo5kmjyVZA+qh30QVZPWxUVRyqmR5wBsyaoVbVDMX7ZrRk8qx1D+AST3rQ0/A=="
3
4 spark.conf.set(
5     f"fs.azure.account.key.ercovidreportingdatalake.dfs.core.windows.net",
6     f"9S43wAC0w0hKl9Nsu6USZtFKuAo5kmjyVZA+qh30QVZPWxUVRyqmR5wBsyaoVbVDMX7ZrRk8qx1D+AST3rQ0/A=="
7 )
8 container_name = "raw" # make sure to use the right container
9 df_raw_population = spark.read.csv(f"abfss://raw@ercovidreportingdatalake.dfs.core.windows.net/population", sep=r'\t', header=True)
```

Cmd 6

## Read the population data & create a temp view

Cmd 7

```
1
2 df_raw_population = df_raw_population.withColumn('age_group', regexp_replace(split(df_raw_population['indic_de,geo\\time'], ',')[0], 'PC_', ''))
3 df_raw_population = df_raw_population.select(col("country_code").alias("country_code"),
4     col("age_group").alias("age_group"),
5     col("2019 ").alias("percentage_2019"))
6 df_raw_population.createOrReplaceTempView("raw_population")
```

Cmd 8

## Pivot the data by age group

Cmd 9

```
1 # Create a data frame with pivoted percentages
2 df_raw_population_pivot = spark.sql("SELECT country_code, age_group, cast(regexp_replace(percentage_2019, '[a-z]', '') AS decimal(4,2)) AS percentage_2019 FROM raw_population WHERE length(country_code) = 2").groupBy("country_code").pivot("age_group").sum
3 df_raw_population_pivot.createOrReplaceTempView("raw_population_pivot")
```





- General
- Factory settings
- Connections
- Linked services
- Integration runtimes
- Microsoft Purview
- Source control
- Git configuration
- ARM template
- Author
- Triggers
- Global parameters
- Data flow libraries
- Security
- Credentials
- Customer managed key
- Outbound rules
- Managed private endpoints
- Workflow orchestration manager
- Apache Airflow

## Linked services

Linked service defines the connection information to a data store or compute. [Learn more](#)

+ New

Filter by name

Annotations : Any

Showing 1 - 3 of 3 items

Name ↑↓	Type ↑↓	Related ↑↓
ls_azblob_ercovidreportingsa	Azure Blob Storage	2
ls_azdl_ercovidreportingdatalake	Azure Data Lake Storage Gen2	9
ls_http_ecdc_data	HTTP	1

## New linked service

Azure Databricks [Learn more](#)

Connect via integration runtime \* ⓘ

AutoResolveIntegrationRuntime

Account selection method \*

☒ From Azure subscription ☐ Enter manually

Azure subscription \* ⓘ

Azure for Students (50381b31-77bb-4cc2-a6f1-156c0ef3f834)

Databricks workspace \* ⓘ

covid-reporting-dbk-ws

Select cluster

☐ New job cluster ☒ Existing interactive cluster ☐ Existing instance pool

Databrick Workspace URL \* ⓘ

https://adb-7843688801986019.19.azuredatabricks.net

Authentication type \*

Access Token

Access token

Azure Key Vault

Access token \* ⓘ

Existing cluster ID \* ⓘ

Add workspace and access token to list options

Annotations

+ New

Create

Back

Test connection

Cancel



1/3



# User Settings

Access tokens

Git integration

Editor settings

Email preferences

Language settings

Personal access tokens can be used for secure authentication to the [Databricks API](#) instead of passwords.

Generate new token

Comment	Creation	Expiration	
demo	2023-07-02 13:54:41 CDT	2023-07-22 13:54:41 CDT	





General

Factory settings

Connections

Linked services

Integration runtimes

Microsoft Purview

Source control

Git configuration

ARM template

Author

Triggers

Global parameters

Data flow libraries

Security

Credentials

Customer managed key

Outbound rules

Managed private endpoints

Workflow orchestration manager

Apache Airflow

## Linked services

Linked service defines the connection information to a data store or compute. [Learn more](#)

+ New

Filter by name

Annotations : Any

Showing 1 - 3 of 3 items

Name ↑↓	Type ↑↓	Related ↑↓
ls_azblob_ercovidreportingsa	Azure Blob Storage	2
ls_azdl_ercovidreportingdatalake	Azure Data Lake Storage Gen2	9
ls_http_ecdc_data	HTTP	1

## New linked service

Azure Databricks [Learn more](#)

### Account selection method \*

☒ From Azure subscription ☐ Enter manually

### Azure subscription \* ⓘ

Azure for Students (50381b31-77bb-4cc2-a6f1-156c0ef3f834)

### Databricks workspace \* ⓘ

covid-reporting-dbk-ws

### Select cluster

☐ New job cluster ☒ Existing interactive cluster ☐ Existing instance pool

### Databrick Workspace URL \* ⓘ

https://adb-7843688801986019.19.azure.databricks.net

### Authentication type \*

Access Token

Access token

Azure Key Vault

### Access token \* ⓘ

.....

### Choose from existing clusters \* ⓘ

covid-reporting-cluster

### Annotations

+ New

> Parameters

> Advanced ⓘ

Create

Back

Test connection

Cancel

Data Factory

Validate all

Publish all

Preview experience

Off

Factory Resources

Filter resources by name

Pipelines

5

pl\_process\_pop\_data

pl\_cases\_deaths\_data\_processed

pl\_ingest\_ecdc\_data

pl\_ingest\_population\_data

pl\_processed\_hospitaladmissions\_data

Change Data Capture (preview)

0

Datasets

12

Data flows

2

Power Query

0

Activities

Search activities

Move & transform

Synapse

Azure Data Explorer

Azure Function

Batch Service

Databricks

Notebook

Jar

Python

Data Lake Analytics

General

HDInsight

Iteration & conditionals

Machine Learning

Power Query

pl\_process\_pop\_data

Validate

Debug

Add trigger

Notebook

exec-pop-transfor  
m

General

Azure Databricks

Settings

User properties

Databricks linked service \*

ls\_dbk\_covid\_cluster

Test connection

Edit

New

Properties

General

Related

Name \*

pl\_process\_pop\_data

Description

Annotations

New

Data Factory

Validate all

Publish all 1

Preview experience

Off

Factory Resources

Filter resources by name

Pipelines5

pl\_process\_pop\_data

pl\_cases\_deaths\_data\_processed

pl\_ingest\_ecdc\_data

pl\_ingest\_population\_data

pl\_processed\_hospitaladmissions\_data

Change Data Capture (preview)0

Datasets12

Data flows2

Power Query0

Activities

Search activities

Move & transform

Synapse

Azure Data Explorer

Azure Function

Batch Service

Databricks

Notebook

Jar

Python

Data Lake Analytics

General

HDInsight

Iteration & conditionals

Machine Learning

Power Query

Validate

Debug

Add trigger

Notebook

exec-pop-transform

General

Azure Databricks

Settings

User properties

Notebook path \*

/covid/transform/transform\_population\_...

Browse

Open

Base parameters

Append libraries

Properties

General

Related

Name \*

pl\_process\_pop\_data

Description

Annotations

+ New

1:58 PM

>>

<<

Home

Runs

Pipeline runs

Trigger runs

Change Data Capture (previ...

Runtime & sessions

Integration runtimes

Data flow debug

Notifications

Alerts & metrics

Rerun

Refresh

Update pipeline

List

Gantt

Notebook

exec-pop-transfor  
m

+

-

100%

Activity runs

Pipeline run ID 6d31e9d6-80bf-4e21-98aa-932caa574725

All status

Export to CSV

Showing 1 - 1 items

Activity name	Status	Activity type	Run start	Duration	Log	Integration runtime	User properties	Activity run ID
exec-pop-transform	Succeeded	Notebook	7/2/2023, 2:06:53 PM	35s		AutoResolveIntegrator		3cd97361-9c10-4c35-af0c-9ec76

EXPLORED

Search for resources

[Collapse all](#)

[Refresh all](#)

Quick Access

Emulator & Attached

Storage Accounts

Azure for Students (AvinashChinta@my.unt.edu)

Storage Accounts

cs71003200143e43ad6

dbstorageocb4jzhwww7hy (ADLS Gen2)

dlspractise (ADLS Gen2)

ercovidreportingdatalake (ADLS Gen2)

Blob Containers

\$logs

lookup

processed

raw

File Shares

Queues

Tables

ercovidreportingsa

Blob Containers

\$logs

configs

population

scripts

File Shares

Queues

Tables

Disks

cloud-shell-storage-southcentralus

covid-reporting-rg

databricks-rg-dbs-handson-bj47lln3b

NetworkWatcherRG

rg-avinash-handson

processed x

\$logs x

population x

configs x

lookup x

scripts x

Upload

Download

Open

Preview

New Folder

Select All

Copy

Paste

Clone

Rename

Move

Manage ACLs

Properties

Delete

Undelete

Folder Statistics

Refresh

← → ∨ ↑

Active blobs (default)

processed > population

Filter by prefix (case-sensitive)

Name	Access Tier	Access Tier Last Modified	Last Modified	Blob Type	Content Type	Size
_SUCCESS	Hot (inferred)		7/2/2023 2:07 PM	Block Blob	application/octet-stream	0 B
_committed_210491613074435726	Hot (inferred)		7/2/2023 2:07 PM	Block Blob	application/octet-stream	111 B
_started_210491613074435726	Hot (inferred)		7/2/2023 2:07 PM	Block Blob	application/octet-stream	0 B
part-00000-tid-210491613074435726-82273a7e-bfa2-4a29-b412-838dae369b66-16-1-c000.csv	Hot (inferred)		7/2/2023 2:07 PM	Block Blob	application/octet-stream	2.70 K

Showing 1 to 4 of 4 cached items

Activities

Clear completed

Clear successful

✓

Transfer of 'C:\Users\Dheeraji Bajjuri\Downloads\testing\testing\000000\_0' to 'processed/ecdc/testing/' complete: 1 item transferred (used SAS, discovery completed)

Started at: 7/1/2023 7:35 PM, Duration: 5 seconds

[Copy AzCopy Command to Clipboard](#)

✓

Successfully created folder 'ecdc/testing'

✓

Transfer of 'C:\Users\Dheeraji Bajjuri\Downloads\covid\_transform\_testing.hql' to 'scripts/hql/' complete: 1 item transferred (used SAS, discovery completed)

Started at: 7/1/2023 7:32 PM, Duration: 7 seconds

[Copy AzCopy Command to Clipboard](#)

## ADF\_er-covid-reporting-adf\_pl\_process\_pop\_data\_exec-pop-transform\_3cd97361-9c10-4c35-af0c-9ec765c36c0f run

[Delete job run](#)

## Output

[Hide code](#)[Export as HTML](#)

## Task run details

Job ID	283344686361047
Task run ID	665
Run as	Chinta, Avinash
Started	2023-07-02 14:06:54 CDT
Ended	2023-07-02 14:07:15 CDT
Duration	20s
Status	Succeeded



## Notebook

[/covid/transform/transform\\_population\\_data](/covid/transform/transform_population_data) 

## Compute

0702-022001-r0adqqmj



## Permissions

No permissions

```
df_processed_population = spark.sql("""SELECT c.country,
c.country_code_2_digit,
c.country_code_3_digit,
c.population,
p.Y0_14 AS age_group_0_14,
p.Y15_24 AS age_group_15_24,
p.Y25_49 AS age_group_25_49,
p.Y50_64 AS age_group_50_64,
p.Y65_79 AS age_group_65_79,
p.Y80_MAX AS age_group_80_max
FROM raw_population_pivot p
JOIN dim_country c ON p.country_code = country_code_2_digit
ORDER BY country""")
```

df\_processed\_population: pyspark.sql.dataframe.DataFrame = [country: string, country\_code\_2\_digit: string ... 8 more fields]

Command took 0.19 seconds

## Write output to the processed mount point

```
df_processed_population.write.format("com.databricks.spark.csv").option("header","true").option("delimiter", "," ).mode
("overwrite").save("abfss://processed@ercovidreportingdatalake.dfs.core.windows.net/population")
```

Command took 3.66 seconds

