
Large Scale Data Processing

Project Overview

In today's data-driven world, the ability to efficiently process and analyze large volumes of data is crucial for businesses to gain insights and make informed decisions. This project aims to leverage the power of Azure Databricks and PySpark to perform large-scale data processing tasks, including Extract, Transform, and Load (ETL) operations, on massive datasets. By utilizing Databricks clusters, we ensure scalability and parallel processing capabilities.

About Project

This project leverages Azure Databricks and PySpark for large-scale data processing on a sample CSV file containing employee data. The dataset consists of more than 100,000 records with 11 fields including employee ID, full name, job title, department, gender, ethnicity, age, hire date, annual salary, bonus percentage, country, and exit date.

Architectural Diagram



Key-Components/Requirements of the projects

1. Azure Databricks:

- Azure Databricks provides a cloud-based platform for big data analytics and machine learning. It offers a collaborative environment for data engineers, data scientists, and analysts to work together seamlessly.
- Databricks provides managed Spark clusters, eliminating the need for infrastructure management and allowing teams to focus on data processing tasks.

2. PySpark:

- PySpark is the Python API for Apache Spark, a powerful open-source framework for distributed data processing. PySpark simplifies development tasks by providing a Python interface to Spark's capabilities.
- With PySpark, developers can write concise and expressive code to perform complex data transformations, aggregations, and analytics on large datasets.

3. ETL Operations:

- **Extract:** Data ingestion from various sources such as databases, data lakes, streaming platforms, or external APIs.
- **Transform:** Data transformation tasks including cleansing, filtering, aggregating, joining, and enriching datasets to prepare them for analysis.
- **Load:** Storing processed data into target systems such as data warehouses, data lakes, or serving layers for downstream consumption.

4. Scalability with Databricks Clusters:

- Databricks clusters dynamically allocate computational resources based on workload requirements, ensuring optimal performance and resource utilization.
- Autoscaling capabilities automatically adjust cluster size to accommodate changes in workload demand, allowing for seamless scalability without manual intervention.
- Databricks Runtime optimizes performance with built-in optimizations, caching, and tuning for various workloads, resulting in faster processing times.

Azure Resources Used for this Project

1. Azure Data Lake Storage Gen2:

- This is where the Transformed data is Loaded. Azure Data Lake Storage Gen2 provides a scalable and secure platform for storing large volumes of data. It enables us to manage, access, and analyse data effectively

2. Azure Blob Storage:

- This is where the raw data is stored. Azure Blob Storage integral to Microsoft Azure's storage service, is a cloud-based solution tailored for managing vast amounts of unstructured data, encompassing both text and binary data. Termed "Blob" for "Binary Large Object," it signifies a compilation of binary data treated as a singular entity within a database.

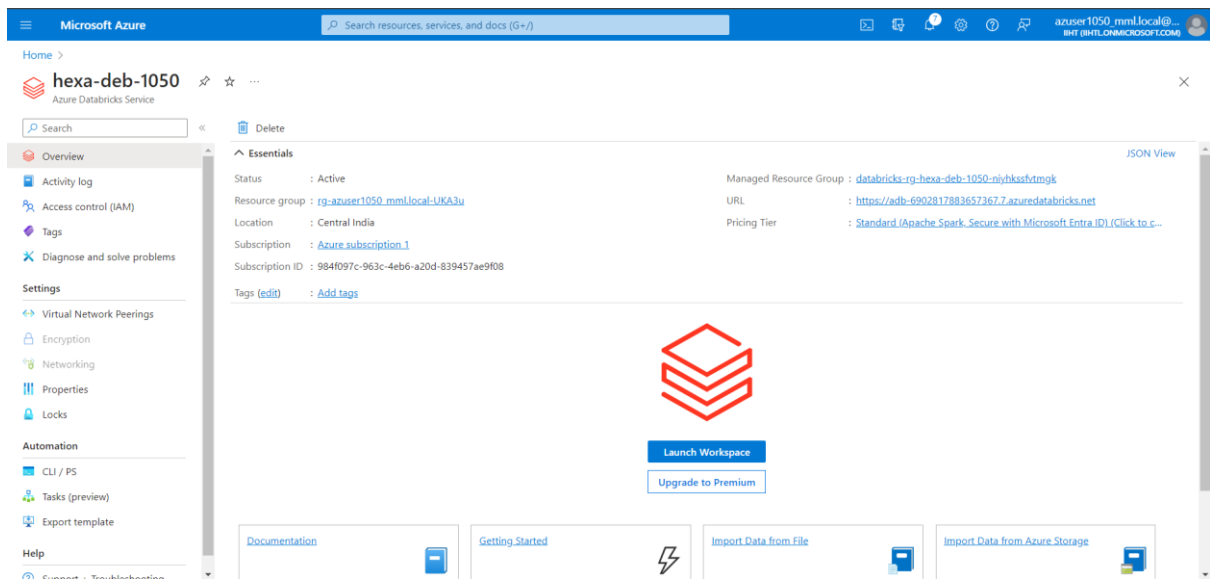
3. Databricks Cluster

- An Azure Databricks cluster process the data depending on the user instructions in the Azure Notebook. It serves as a computational resource facilitating the processing of extensive data and execution of analytics workloads through the Apache Spark platform within the Microsoft Azure cloud.

How It works

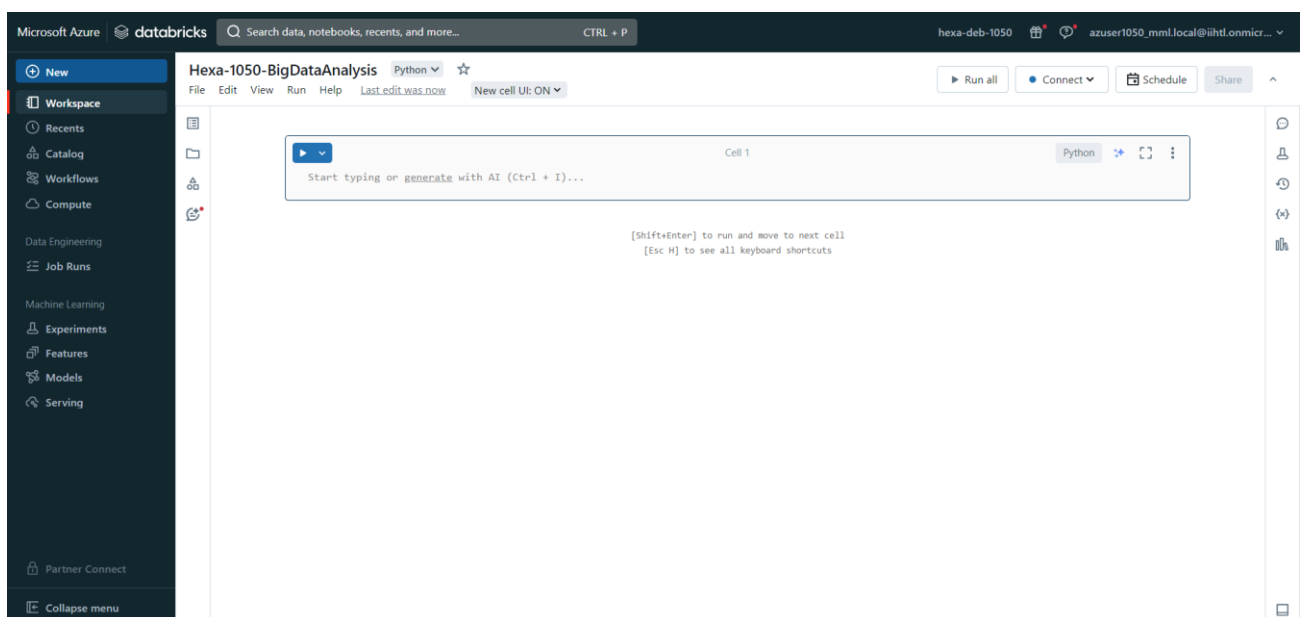
1. Setting Up Azure Databricks Environment:

- Sign in to the Azure portal and create an Azure Databricks workspace. Configure workspace settings, including pricing tier, region, and workspace name



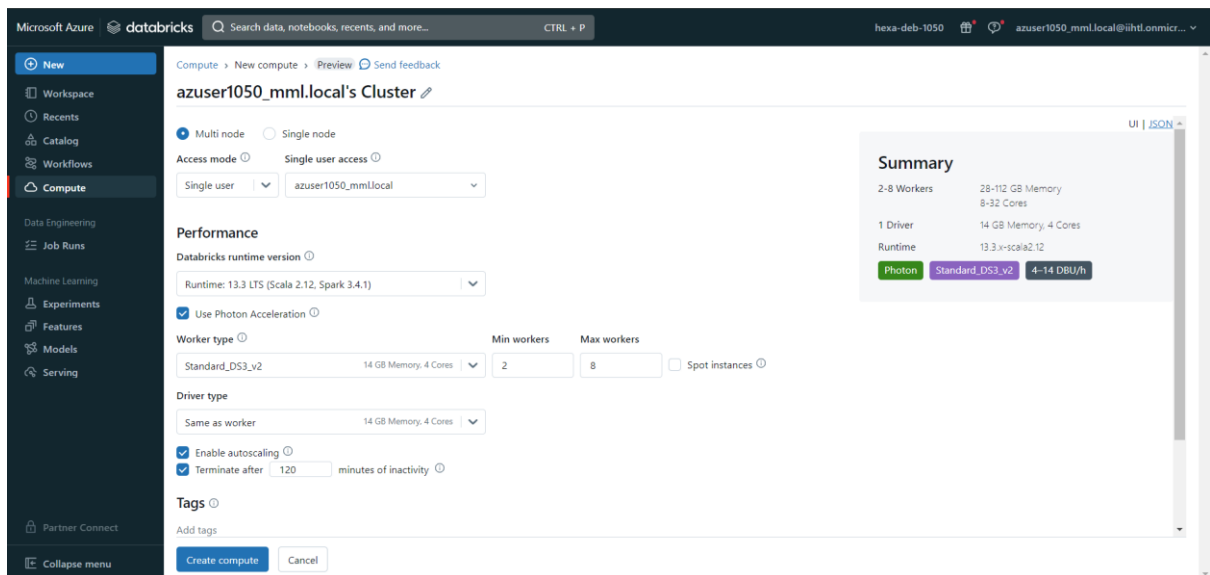
2. Developing PySpark Notebooks:

- Create a new PySpark notebook within the Databricks workspace. Begin writing PySpark code to perform ETL operations, data transformations, and other data processing tasks.



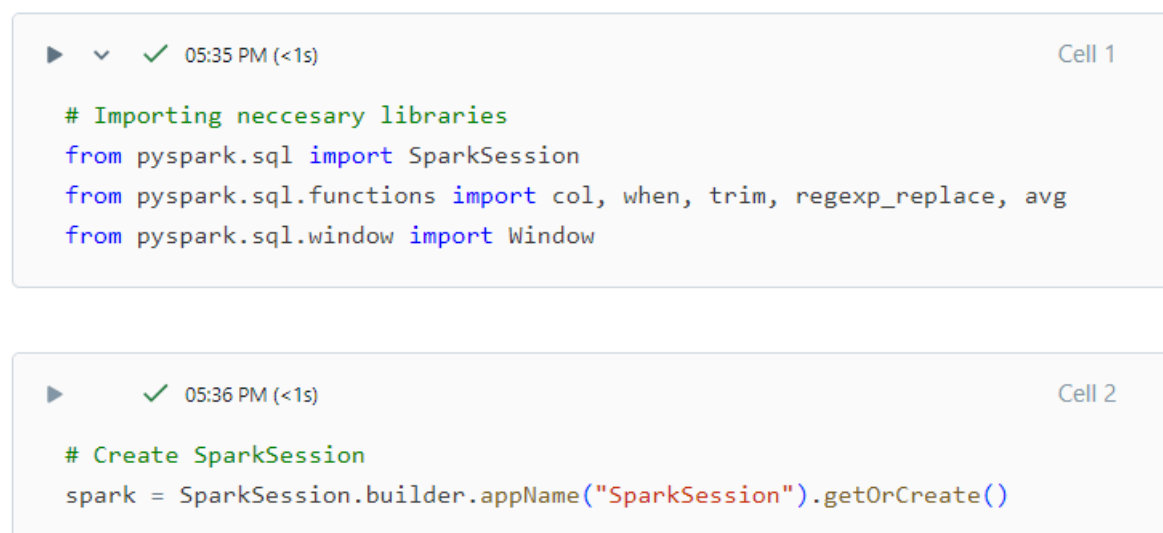
3. Create Cluster and Connecting to notebook

- The cluster is created with 4 working nodes and autoscaling is enabled which automatically adjust cluster size to accommodate changes in workload demand, allowing for seamless scalability without manual intervention.



4. Importing Necessary libraries and Creating Spark Session:

- Use `SparkSession.builder` to configure and create a `SparkSession`. specify the application name using `.appName()` and configure any additional Spark options using `.config()`. Finally, call `.getOrCreate()` to either create a new `SparkSession`



5. Extracting Data from Source storage

- Connecting data source (Azure Blob Storage) by mounting it to the Databricks File System (DBFS) to simplify data access
- It helps to retrieve raw data for processing and analysis within the PySpark environment.

```
05:37 PM (12s) Cell 3

# 1) Extracting the data from blob storage
# Mounting the blob storage with Azure databricks

dbutils.fs.mount(
    source = "wasbs://hexa1050sourcecontainer@hexa1050sourcestorage.blob.core.windows.net",
    mount_point="/mnt/blobStorage1",
    extra_configs={"fs.azure.account.key.hexa1050sourcestorage.blob.core.windows.net":
        "wV1Xp0FR6Njh29V1z/Pzk1KkY/p7JLmpkFwuIBRb55xIgYXPhyx9TWrhMZ0kgb+nnfcEoeQJSgU0+AStdannEQ=="})

True
```

```
05:37 PM (<1s) Cell 4

# Listing the File information to get file path

dbutils.fs.ls('/mnt/blobStorage1')

[FileInfo(path='dbfs:/mnt/blobStorage1/Employee_data.csv', name='Employee_data.csv'
```

```
05:38 PM (10s) Cell 5

# Reading the data of blob storage and converting it to spark RDD

path = "dbfs:/mnt/blobStorage1/Employee_data.csv"
RDD = spark.read.csv(path, header=True,inferSchema=True)

(2) Spark Jobs

RDD: pyspark.sql.dataframe.DataFrame = [EEID: integer, Full Name: string ... 10 more fields]
```

```
05:38 PM (1s) Cell 6

RDD.display()

(1) Spark Jobs
```

	EEID	Full Name	Job Title	Department	Gender	Ethnicity	Age
1	98278	Samaira Raj	Human Resources Manager	Human Resources	Male	White	46
2	19840	Lavanya Hayer	Data Analyst	Research and Development	Female	null	55
3	22487	Shayak Raval	Financial Analyst	Finance	Female	White	29
4	17160	Inaaya Bala	Software Engineer	Research and Development	Male	null	65
5	10385	Aarav Garde	Data Analyst	Research and Development	Male	White	56
6	28297	Romil Keer	Customer Service Representative	Customer Service	null	null	19
7	21123	Pranav Chadha	Human Resources Manager	Human Resources	Male	Hisoanic	46

10,000 rows | Truncated data | 1.10 seconds runtime

Refreshed 7 minutes ago

6. Transforming the raw data

- Utilize PySpark DataFrame transformations and functions to cleanse, transform, and prepare the data for analysis.
- Implement business logic and data processing steps to transform raw dataset up to mark for data analysis purpose.

Transformations done:-

▪ Removing the Duplicate records

```
▶ 05:38 PM (2s) Cell 7

# 2) Transforming Data
# Removing the Duplicate data

print(RDD.count())
RDD=RDD.distinct()
print(RDD.count())

▶ (5) Spark Jobs
  ▶ RDD: pyspark.sql.dataframe.DataFrame = [EEID: integer, Full Name: string ... 10 more fields]
100009
100000
```

▪ Handling anonymous data

```
▶ 05:38 PM (2s) Cell 8

# Removing the anonymous data

print(RDD.count())
RDD = RDD.na.drop("any",subset=["EEID"])
print(RDD.count())

▶ (6) Spark Jobs
  ▶ RDD: pyspark.sql.dataframe.DataFrame = [EEID: integer, Full Name: string ... 10 more fields]
100000
99988
```


■ Removing Extra spaces form the data

05:38 PM (1s) Cell 10

```
# Removing Leading and Trailing spaces from the data

RDD = RDD.withColumn("Full Name", trim("Full Name"))
RDD = RDD.withColumn("Job Title", trim("Job Title"))
RDD = RDD.withColumn("Department", trim("Department"))
RDD = RDD.withColumn("Gender", trim("Gender"))
RDD = RDD.withColumn("Ethnicity", trim("Ethnicity"))
RDD = RDD.withColumn("Country", trim("Country"))

RDD.display()
```

(2) Spark Jobs

RDD: pyspark.sql.dataframe.DataFrame = [EEID: integer, Full Name: string ... 10 more fields]

Table + New result table: OFF

	EEID	Full Name	Job Title	Department	Gender	Ethnicity	Age
1	46400	Riya Grover	Data Analyst	Marketing	Female	Asian	25
2	61085	Ira Wable	Project Manager	Engineering	Female	Asian	18
3	92676	Parinaaz Karpe	Software Engineer	IT	Female	Asian	44
4	38396	Taran Butala	Software Engineer	Research and Development	Female	White	68
5	45876	Jayant Devan	Sales Representative	Sales	Female	Black	24
6	69030	Abram Mani	Software Engineer	IT	Male	Hispanic	44
7	22677	Samiha Vasa	Sales Representative	Sales	Female	null	70

10,000 rows | Truncated data | 0.77 seconds runtime Refreshed 9 minutes ago

■ Filling null values with proper messages and data

05:38 PM (<1s) Cell 11

```
# Filling the null values with proper message

RDD = RDD.na.fill(value="Not Known",subset=["Full Name"])
RDD = RDD.na.fill(value="Not Known",subset=["Job Title"])
RDD = RDD.na.fill(value="Not Known",subset=["Department"])
RDD = RDD.na.fill(value="Prefer Not to say",subset=["Gender"])
RDD = RDD.na.fill(value="Not Known",subset=["Ethnicity"])
RDD = RDD.na.fill(value="Not Known",subset=["Country"])
RDD = RDD.na.fill(value=0,subset=["Bonus %"])
RDD = RDD.withColumn('Hire Date',when(col('Hire Date').isNull(),('No data provided')).otherwise(col('Hire Date')))
RDD = RDD.withColumn('Exit Date',when(col('Exit Date').isNull(),('Currently Working')).otherwise(col('Exit Date')))
```

RDD: pyspark.sql.dataframe.DataFrame = [EEID: integer, Full Name: string ... 10 more fields]

05:38 PM (1s)

Cell 12

RDD.display()

(2) Spark Jobs

Table

New result table: OFF

	EEID	Full Name	Job Title	Department	Gender	Ethnicity	Age
1	46400	Riya Grover	Data Analyst	Marketing	Female	Asian	25
2	61085	Ira Wable	Project Manager	Engineering	Female	Asian	18
3	92676	Parinaaz Karpe	Software Engineer	IT	Female	Asian	44
4	38396	Taran Butala	Software Engineer	Research and Development	Female	White	68
5	45876	Jayant Devan	Sales Representative	Sales	Female	Black	24
6	69030	Abram Mani	Software Engineer	IT	Male	Hispanic	44
7	22677	Samiha Vasa	Sales Representative	Sales	Female	Not Known	70

10,000 rows | Truncated data | 0.94 seconds runtime

Refreshed 11 minutes ago

05:39 PM (2s)

Cell 13

Filling the numerical column's null values with proper average values

window_spec = Window.partitionBy()
RDD = RDD.withColumn('Age',when(col('Age').isNull(),avg(col('Age')).over(window_spec)).otherwise(col('Age')))

average_salaries = RDD.groupBy("Country", "Department").avg("Annual Salary")
RDD = RDD.join(average_salaries, ["Country", "Department"], "left").withColumnRenamed("avg(Annual Salary)", "Average Salary")
RDD = RDD.withColumn('Annual Salary', when(col('Annual Salary').isNull(), col('Average Salary')).otherwise(col('Annual Salary'))).drop("Average Salary")

RDD.display()

(5) Spark Jobs

average_salaries: pyspark.sql.dataframe.DataFrame = [Country: string, Department: string ... 1 more field]
RDD: pyspark.sql.dataframe.DataFrame = [Country: string, Department: string ... 10 more fields]

Table

New result table: OFF

	Country	Department	EEID	Full Name	Job Title	Gender	Ethnicity
1	Canada	Marketing	46400	Riya Grover	Data Analyst	Female	Asian
2	UK	Engineering	61085	Ira Wable	Project Manager	Female	Asian
3	USA	IT	92676	Parinaaz Karpe	Software Engineer	Female	Asian
4	Canada	Research and Development	38396	Taran Butala	Software Engineer	Female	White
5	Australia	Sales	45876	Jayant Devan	Sales Representative	Female	Black
6	Canada	IT	69030	Abram Mani	Software Engineer	Male	Hispanic
7	Canada	Sales	22677	Samiha Vasa	Sales Representative	Female	Not Known

10,000 rows | Truncated data | 2.33 seconds runtime

Refreshed 11 minutes ago

- Renaming USA to US to make dataset consistent

▶ 05:39 PM (1s) Cell 14

```
# Renaming USA as US

RDD=RDD.withColumn('Country',when(RDD.Country=='USA',regexp_replace(RDD.Country,'USA','US')).otherwise(RDD.Country))
print(RDD.select("Country").distinct().collect())
```

▶ (3) Spark Jobs

▶ RDD: pyspark.sql.dataframe.DataFrame = [Country: string, Department: string ... 10 more fields]

[Row(Country='US'), Row(Country='UK'), Row(Country='Canada'), Row(Country='Australia')]

- Statistical data

▶ 05:39 PM (1s) Cell 15

```
# Statistical Data for analysis

RDD.describe("Age", "Annual Salary", "Bonus %").display()
```

▶ (5) Spark Jobs

Table ▼ +

	summary ▲	Age ▲	Annual Salary ▲	Bonus % ▲
1	count	99988	99988	99988
2	mean	43.93738248589831	79966.6102574301	5.003393807256835
3	stddev	15.330340571295418	23092.15240664656	2.884433311011505
4	min	18.0	40000.48	0.0
5	max	70.0	119999.56	10.0

⬇ 5 rows | 1.43 seconds runtime

7. Loading Data into Sink Storage

- To store the transformed data we need to create sink storage (Azure Data Lake) and Container
- Connecting data source (ADLS) by mounting it to the Databricks notebook to load the data

The screenshot shows the 'Create a storage account' wizard in the Microsoft Azure portal, specifically the 'Advanced' tab. The wizard is for creating a storage account named 'hexa1050sinkstorage_1708948423023'. The 'Advanced' tab is selected, showing options for 'Enable storage account key access' (checked), 'Default to Microsoft Entra authorization in the Azure portal' (unchecked), 'Minimum TLS version' (set to 'Version 1.2'), and 'Permitted scope for copy operations (preview)' (set to 'From any storage account'). Below these, the 'Hierarchical Namespace' section is visible, with 'Enable hierarchical namespace' checked. At the bottom, there are navigation buttons: 'Review', '< Previous', and 'Next: Networking >', along with a 'Give feedback' link.

The screenshot shows the 'hexa1050sinkcontainer' container overview in the Microsoft Azure portal. The container is located under 'hexa1050sinkstorage_1708948423023'. The 'Overview' tab is selected, showing the 'Authentication method' as 'Access key (Switch to Microsoft Entra user account)' and the 'Location' as 'hexa1050sinkcontainer / transformed_data'. A search bar is present with the text 'Search blobs by prefix (case-sensitive)'. Below the search bar, there is a table with columns: 'Name', 'Modified', 'Access tier', 'Archive status', 'Blob type', 'Size', and 'Lease state'. The table currently shows one entry: a folder named '[-]'. To the left of the table, there is a sidebar with navigation links: 'Overview', 'Diagnose and solve problems', 'Access Control (IAM)', 'Settings', 'Shared access tokens', 'Manage ACL', 'Access policy', 'Properties', and 'Metadata'.

▶ ✓ 05:40 PM (11s)

Cell 16

3) Loading Data in Azure Data Lake

Mounting the sink storage(Azure Data Lake) with Azure databricks

```
dbutils.fs.mount(  
  source = "wasbs://hexa1050sinkcontainer@hexa1050sinkstorage.blob.core.windows.net",  
  mount_point="/mnt/blobStorage2",  
  extra_configs={"fs.azure.account.key.hexa1050sinkstorage.blob.core.windows.net":  
    "Ry48JkTFLSmtAVkGNMnEmNYAzCh23Ejb4LJ41TnN23AABoPuK3ygE3pB+7BBao2+TLMbTCDIs+cn+ASTi2DVWg=="}  
)
```

True

▶ ✓ 05:40 PM (<1s)

Cell 17

```
dbutils.fs.ls("/mnt/blobStorage2")
```

[FileInfo(path='dbfs:/mnt/blobStorage2/transformed_data/', name='transformed_data/', size=0, modificationTime=0)]

▶ ✓ 05:40 PM (2s)

Cell 18

Converting RDD to pandas dataframe

```
pandas_df=RDD.toPandas()
```

▶ (5) Spark Jobs

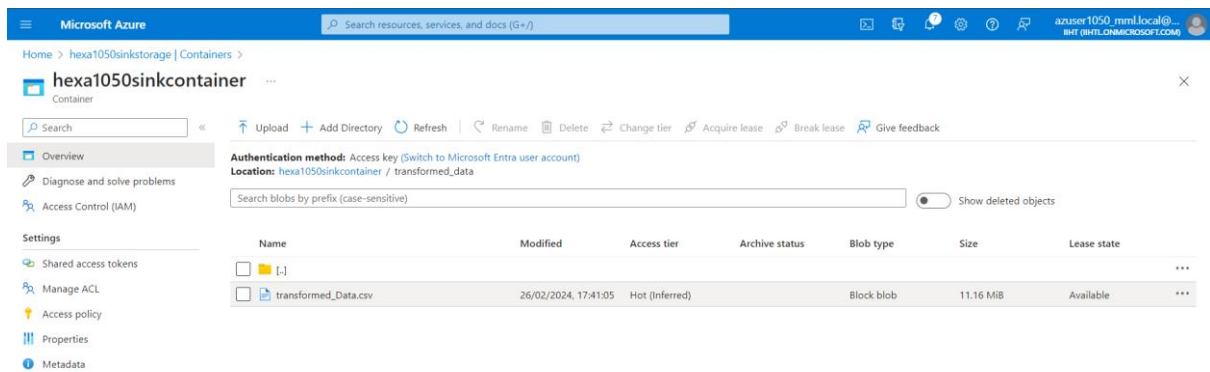
▶ ✓ 05:41 PM (1s)

Cell 19

Loading the transformed data in Azure Data Lake

```
pandas_df.to_csv('/dbfs/mnt/blobStorage2/transformed_data/transformed_Data.csv',index=False)
```

- Data Successfully Loaded



Microsoft Azure

Search resources, services, and docs (G+/I)

Home > hexa1050sinkstorage | Containers >

hexa1050sinkcontainer

Container

Search

Upload Add Directory Refresh Rename Delete Change tier Acquire lease Break lease Give feedback

Authentication method: Access key (Switch to Microsoft Entra user account)

Location: hexa1050sinkcontainer / transformed_data

Search blobs by prefix (case-sensitive) Show deleted objects

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
[...]						...
transformed_Data.csv	26/02/2024, 17:41:05	Hot (Inferred)		Block blob	11.16 MiB	Available

Overview

Diagnose and solve problems

Access Control (IAM)

Settings

Shared access tokens

Manage ACL

Access policy

Properties

Metadata

8. Unmounting the Source and Sink storage

```
05:41 PM (21s)

# Unmounting the source and sink storage

dbutils.fs.unmount("/mnt/blobStorage1")
dbutils.fs.unmount("/mnt/blobStorage2")

/mnt/blobStorage1 has been unmounted.
/mnt/blobStorage2 has been unmounted.

True
```

Raw Data

	A	B	C	D	E	F	G	H	I	J	K	L
1	EEID	Full Name	Job Title	Department	Gender	Ethnicity	Age	Hire Date	Annual Sal	Bonus %	Country	Exit Date
2	98278	Samaira Rao	Human Resources	Human Resources	Male	White	46	#####	52649.82	5.62	Australia	#####
3	19840	Lavanya H	Data Analyst	Research & Development	Female		55	#####	119024.1	6.97	USA	#####
4	22487	Shayak F	Financial Analyst	Finance	Female	White	29	#####	57567.81	7.14	Australia	#####
5	17160	Inaaya B	Software Engineer	Research & Development	Male		65	#####	80599.13	1.13	US	#####
6	10385	Aarav G	Data Analyst	Research & Development	Male	White	56	#####	45077.04	4.89	UK	#####
7	28297	Romil K	Customer Service	Customer Service			19	#####	82680.4	6.21	UK	
8	21123	Pranay Ch	Human Resources	Human Resources	Male	Hispanic	46	#####	102028.9	3.34	USA	#####
9	29382	Trisha S	Lawyer	Legal	Male	Black	59	#####	87641.25	0.09	UK	#####
10	42214	Ira Chaudh	Accountant	Finance	Female	Black	60	#####	79690.1	1.67	Canada	#####
11		Fateh Man	Software Engineer	Engineering	Male		27	#####	67070.3	8.52	UK	#####
12	95098	Ira Dass	Customer Service	Customer Service	Male		70	#####	76975.24	2.8	Canada	
13	15826	Vihaan R	Human Resources	Human Resources	Female		53	#####	62931.12	6.44	Australia	
14	59107	Tushar S	Lawyer	Legal	Male	White	46	#####	56553.87	8.34	UK	
15	80261	Vardaniya	Data Analyst	IT		Asian	69	#####	116655.9	1.69	USA	
16	19118	Riya A	Lawyer	Legal	Female	White	41	#####	49611.44	5.59	Canada	
17		Saanvi K	Human Resources	Human Resources			39	#####	51052.24	7.75	Canada	#####
18	44752	Madhav K	Marketing	Marketing	Male		53	#####	69578.78	3.61	Canada	
19	25242	Tushar Ch	Operations	Operations	Male	Black	67	#####	63484.07	5.03	Canada	#####
20	50360	Miraan R	Financial Analyst	Finance	Male	White	37	#####	55569.24	4.84	USA	#####
21	34275	Aarush A	Customer Service	Customer Service	Male	Hispanic	18	#####	90214.6	1.19	UK	
22	72826	Drishya Ch	Customer Service	Customer Service	Male	Black	62	#####	99760.6	8.73	UK	#####
23	33187	Indranil R	Data Analyst	Research & Development	Male	Black	70	#####	79037.64	0.07	UK	#####
24	30645	Mahika K	Sales Representative	Sales		White	48	#####	69457.92	2.74	USA	#####
25	71220	Nishith L	Data Analyst	IT	Female		32	#####	114713.7	1.97	USA	
26	32553	Hunar D	Data Analyst	Research & Development	Male	White	43	#####	118909.7	8.61	Australia	#####
27	11846	Elakshi D	Operations	Operations	Male		50	#####	56832.93	4.81	Australia	

Transformed Data

	A	B	C	D	E	F	G	H	I	J	K	L	M
1	Country	Department	EEID	Full Name	Job Title	Gender	Ethnicity	Age	Hire Date	Annual Sal	Bonus %	Exit Date	
2	Canada	Marketing	46400	Riya Grove	Data Analyst	Female	Asian	25	#####	75499.29	3.83	#####	
3	UK	Engineering	61085	Ira Wable	Project Manager	Female	Asian	18	#####	110198.8	2.09	#####	
4	US	IT	92676	Parinaaz K	Software Engineer	Female	Asian	44	#####	41026.82	5.31	#####	
5	Canada	Research & Development	38396	Taran But	Software Engineer	Female	White	68	#####	105485.7	8.69	#####	
6	Australia	Sales	45876	Jayant Dev	Sales Representative	Female	Black	24	#####	57866.08	4.18	Currently Working	
7	Canada	IT	69030	Abram Ma	Software Engineer	Male	Hispanic	44	#####	48409.17	8.63	Currently Working	
8	Canada	Sales	22677	Samihha V	Sales Representative	Female	Not Known	70	#####	104495.8	5.83	#####	
9	US	Customer Service	94951	Kabir Dey	Customer Service	Prefer Not to Say	Not Known	40	#####	49184.54	8.37	#####	
10	UK	Sales	30669	Tarini Brah	Sales Representative	Male	White	64	#####	64526.56	9.09	#####	
11	UK	Customer Service	39341	Kismat Kee	Customer Service	Female	Asian	40	#####	110918.4	4.94	Currently Working	
12	UK	Sales	44291	Aradhya C	Sales Representative	Prefer Not to Say	Not Known	29	#####	71443.58	1.7	#####	
13	Canada	Engineering	86184	Stuvan Dui	Software Engineer	Female	Black	69	#####	41516.36	5.26	Currently Working	
14	Canada	Research & Development	20783	Rasha Auri	Data Analyst	Female	White	40	#####	118531.4	2.32	#####	
15	US	Operations	78333	Ritvik Lank	Customer Service	Prefer Not to Say	Not Known	43	#####	90268.73	8.09	Currently Working	
16	Canada	Legal	38418	Ryan Sura	Lawyer	Prefer Not to Say	Not Known	18	#####	46707.41	4.8	#####	
17	US	Operations	45764	Adira Sahn	Customer Service	Male	Not Known	26	#####	78001.94	2.36	#####	
18	Australia	IT	86860	Navya She	Data Analyst	Male	Hispanic	49	#####	71266.08	4.21	#####	
19	US	Legal	45864	Farhan Bal	Lawyer	Male	White	22	#####	116168	5.35	#####	
20	UK	Finance	31359	Anvi Ahluw	Accountant	Prefer Not to Say	Asian	25	#####	101659	9.25	#####	
21	UK	Operations	27549	Emir Tailor	Operations	Prefer Not to Say	Black	64	#####	68657.94	5.17	#####	
22	UK	Engineering	75249	Arnav Mar	Project Manager	Male	Not Known	67	#####	111875.2	8.71	Currently Working	
23	Australia	Marketing	47025	Samar Chh	Marketing	Male	Not Known	60	#####	78734.94	9.19	Currently Working	
24	Australia	Sales	12513	Ranbir Kib	Sales Representative	Prefer Not to Say	Hispanic	30	#####	53666.9	2.53	Currently Working	
25	Australia	Human Resources	80691	Indrans La	Human Resources	Male	White	66	#####	103161.8	0.66	#####	
26	Canada	Operations	88868	Biju Sani	Operations	Male	Hispanic	70	#####	59208.66	7.47	#####	
27	UK	Finance	46756	Dhruv Cha	Financial Analyst	Female	Not Known	32	#####	82979.87	1.98	#####	

Conclusion

In conclusion, this project successfully demonstrated the utilization of Azure Databricks and PySpark for large-scale data processing on a sample CSV file containing employee data. Through the implementation of Extract, Transform, Load (ETL) operations, the dataset was cleaned, transformed, and prepared for further analysis or reporting. Overall, this project serves as a practical demonstration of how Azure Databricks and PySpark can be effectively utilized for large-scale data processing.