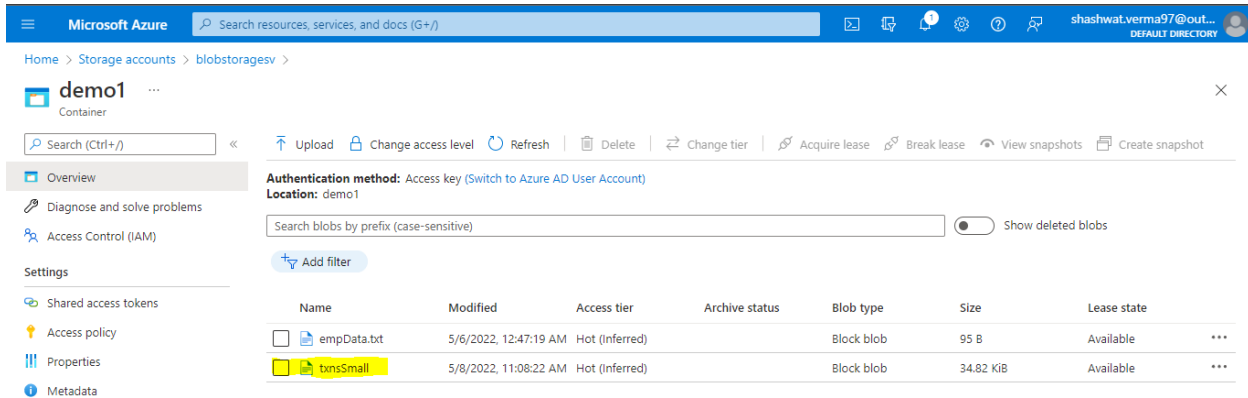## Lab 3: Collect all data whose category is Gymnastics (Databricks)
### Dataset: txnsSmall
### BlobStorage → ADLS → Databricks → ADLS → Synapse
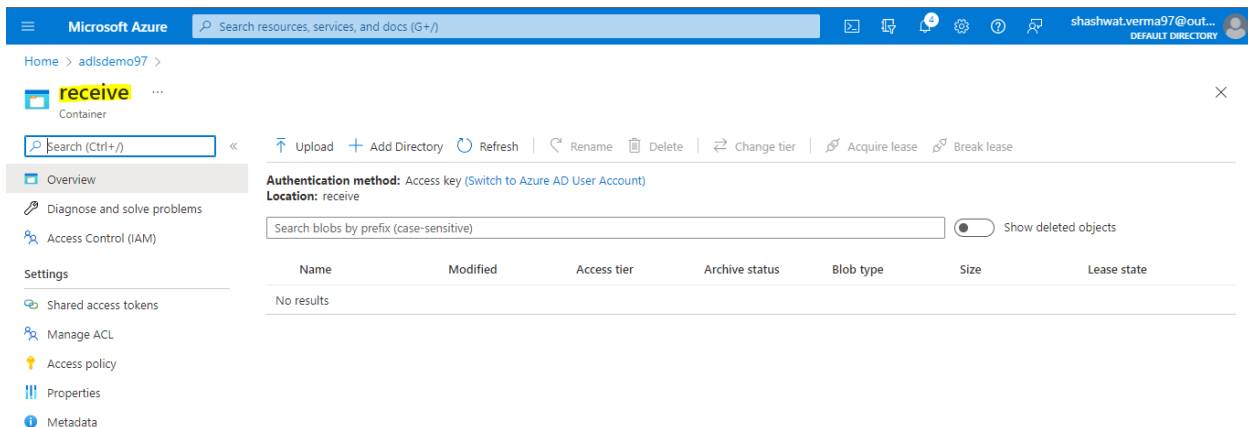
**1.** Upload 'txnsSmall' dataset in blob storage container 'demo1'.



**2.** We will use the 'receive' container in the adls storage 'adlsdemo97' to receive the copy data from blob storage.
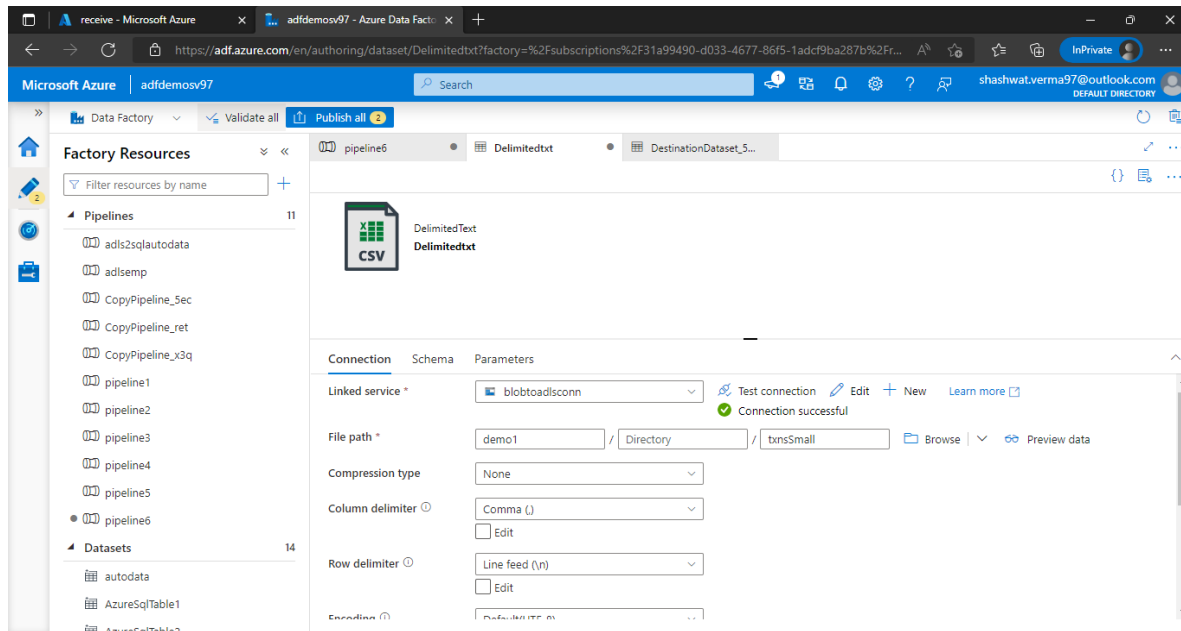


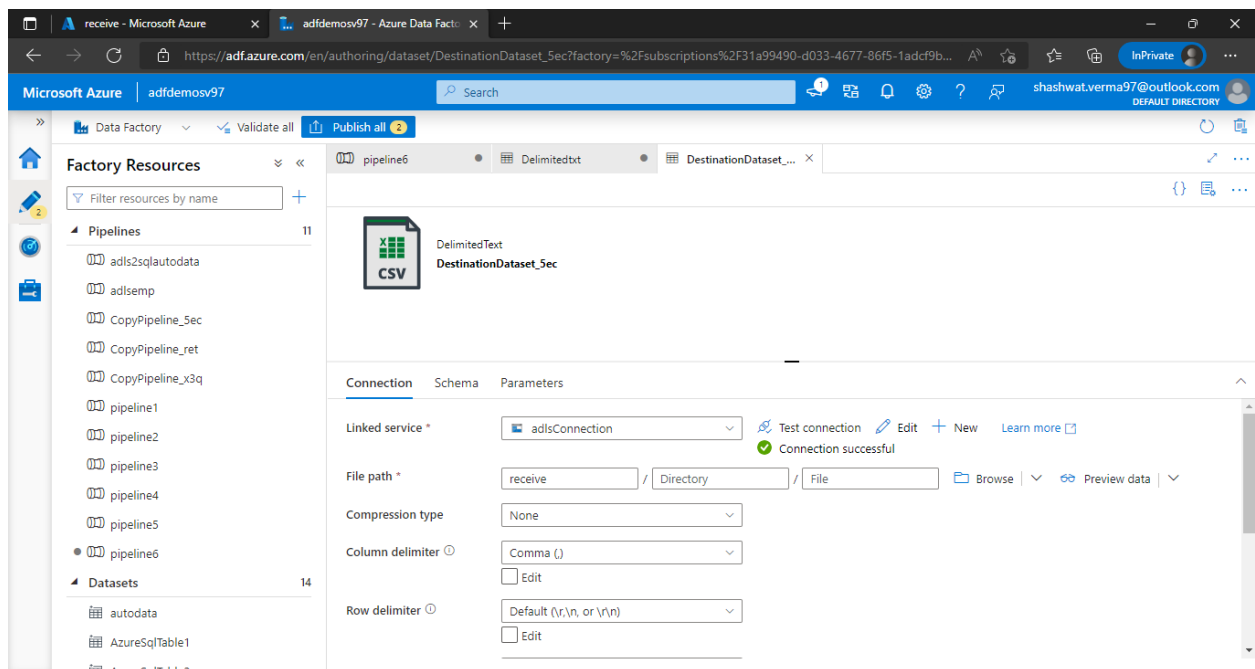**3.** Created pipeline 'blob2adls' to copy the data from blob to adls

Source configuration:
We will use previously created source configuration 'Delimitedtxt' as source for the pipeline

## Sink configuration:

We will use previously created 'DestinationDataset_5ec' as sink configuration which uses 'receive' container as destination location to store copy data

**4.** Created a python notebook in Azure Databricks - 'filterGymnastics' in 'mytestcluster9706'. To store filtered output obtained through this notebook, we created a container 'databricksoutput' in ADLS storage 'adlsdemo97'

Source Code:

```
#Using Access Key for adls storage
spark.conf.set("fs.azure.account.auth.type.adlsdemo97.dfs.core.windows.net","SharedKey")
spark.conf.set("fs.azure.account.key.adlsdemo97.dfs.core.windows.net","IgWwisxE9bIf
Egv5v521p5nE9gb2dX8rl9ZsFDiJAXRv4Y2n+IXderAkHeMJHYLbc1hOFiRPcwZT+AStz
QuYeA==")

# Create a Schema Programatically and Assign the same during DataFrame Creation
from pyspark.sql.types import StructType,StructField
from pyspark.sql.types import IntegerType,StringType,DoubleType,LongType,DateType

schemaFortxnsData = StructType([
    StructField("txnid", IntegerType(),True),
    StructField("txndate", StringType(),True),
    StructField("custid", LongType(),True),
    StructField("amount", DoubleType(),True),
    StructField("category", StringType(),True),
    StructField("subcategory", StringType(),True),
    StructField("city", StringType(),True),
    StructField("state", StringType(),True),
    StructField("txntype", StringType(),True),
])

# Assigning the Schema
txnsDataDF =
spark.read.schema(schemaFortxnsData).option("delimiter",',').csv('abfss://receive@adls
demo97.dfs.core.windows.net/')

txnsDataDF.filter("category =
'Gymnastics'").repartition(1).write.option("header",True).csv("abfss://databricksoutput@
adlsdemo97.dfs.core.windows.net/gymnasticsoutput")
```

filterGymnastics — Data... | adfdemosv97 - Azure | how to connect adf pi... | Authentication using D...

Microsoft Azure    Databricks                                                    Portal    shashwat.verma97@outlook.com

## filterGymnastics (Python)

Free trial ends in 3 days. Upgrade to Premium in Azure Portal

⊘ mytestcluster9706

**Cmd 1**

```python
1  #Using Access Key for adls storage
2  spark.conf.set("fs.azure.account.auth.type.adlsdemo97.dfs.core.windows.net","SharedKey")
3  spark.conf.set("fs.azure.account.key.adlsdemo97.dfs.core.windows.net","IgWwisxE9bIfEgv5v521p5nE9gb2dX8rl9ZsFDiJAXRv4Y2n+IXderAkHeMJHYLbc1hOFiRPcwZT+AStzQu
   YeA==")
```

Command took 0.44 seconds -- by shashwat.verma97@outlook.com at 5/8/2022, 1:52:27 PM on mytestcluster9706

**Cmd 2**

```python
1  # Create a Schema Programatically and Assign the same during DataFrame Creation
2  from pyspark.sql.types import StructType,StructField
3  from pyspark.sql.types import IntegerType,StringType,DoubleType,LongType,DateType
4
5  schemaFortxnsData = StructType([
6      StructField("txnid", IntegerType(),True),
7      StructField("txndate", StringType(),True),
8      StructField("custid", LongType(),True),
9      StructField("amount", DoubleType(),True),
10     StructField("category", StringType(),True),
11     StructField("subcategory", StringType(),True),
12     StructField("city", StringType(),True),
13     StructField("state", StringType(),True),
14     StructField("txntype", StringType(),True),
15  ])
```

---

gymnasticsoutpu... | User Settings - D... | filterGymnastics | adfdemosv97 - A... | how to insert hea... | Spark Write Data

Microsoft Azure    Databricks                                                    Portal    shashwat.verma97@outlook.com

## filterGymnastics (Python)

Free trial ends in 3 days. Upgrade to Premium in Azure Portal

⊘ mytestcluster9706

```python
15  ])
```

Command took 0.05 seconds -- by shashwat.verma97@outlook.com at 5/8/2022, 2:27:08 PM on mytestcluster9706

**Cmd 3**

```python
1  # Assigning the Schema
2  txnsDataDF = spark.read.schema(schemaFortxnsData).option("delimiter",',').csv('abfss://receive@adlsdemo97.dfs.core.windows.net/')
```

▶ ⊞ txnsDataDF: pyspark.sql.dataframe.DataFrame = [txnid: integer, txndate: string ... 7 more fields]

Command took 0.15 seconds -- by shashwat.verma97@outlook.com at 5/8/2022, 2:49:32 PM on mytestcluster9706

**Cmd 4**

```python
1  txnsDataDF.filter("category =
   'Gymnastics'").repartition(1).write.option("header",True).csv("abfss://databricksoutput@adlsdemo97.dfs.core.windows.net/gymnasticsoutput")
```
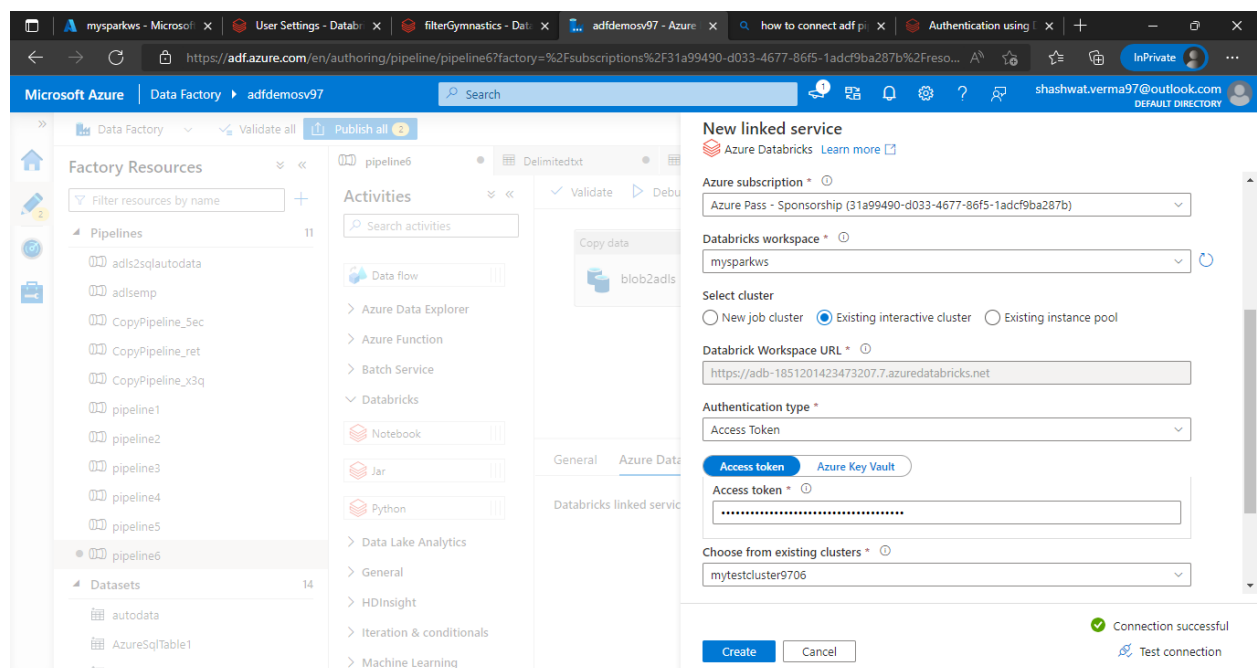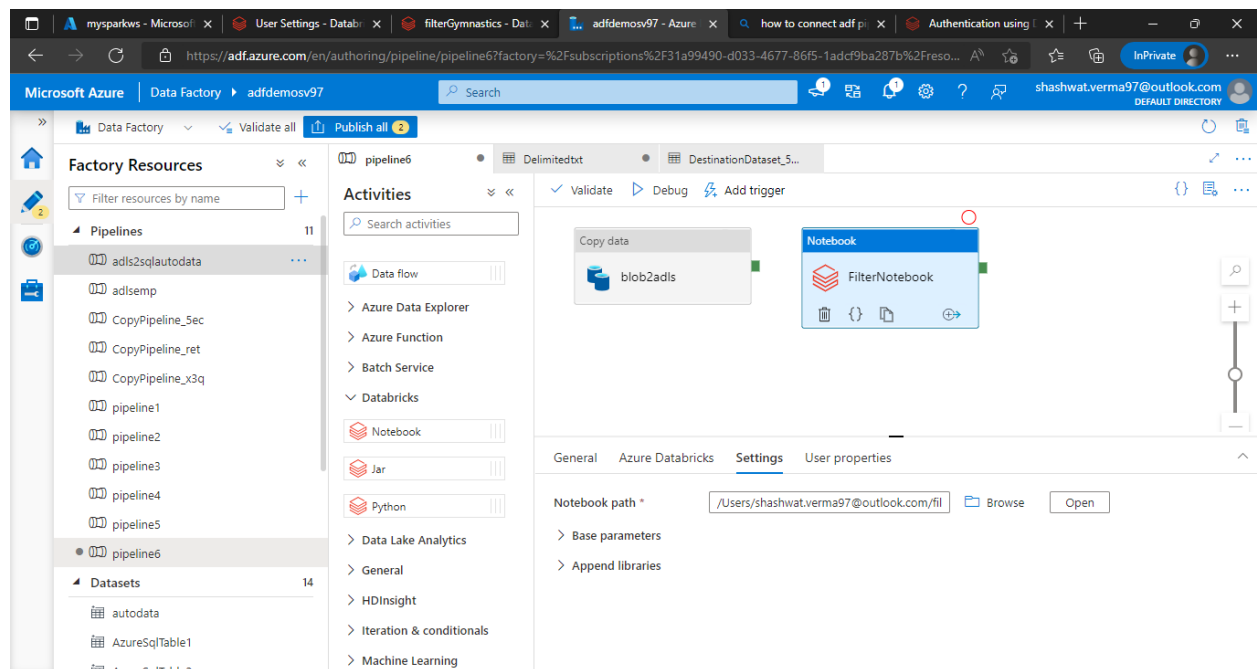
▶ (2) Spark Jobs

Command took 0.80 seconds -- by shashwat.verma97@outlook.com at 5/8/2022, 3:48:11 PM on mytestcluster9706

**Cmd 5**

```python
1  # txnsDataDF.filter("category = 'Gymnastics'").show(5)
```

▶ (1) Spark Jobs

*Untitled - Notepad

**5.** After creating the notebook we will add a Databricks notebook to the pipeline with following configurations:



**6.** Now we have to fetch this filtered data from ADLS container 'databricksoutput/gymnasticsoutput' to Synapse. For that another copy data configuration is created.
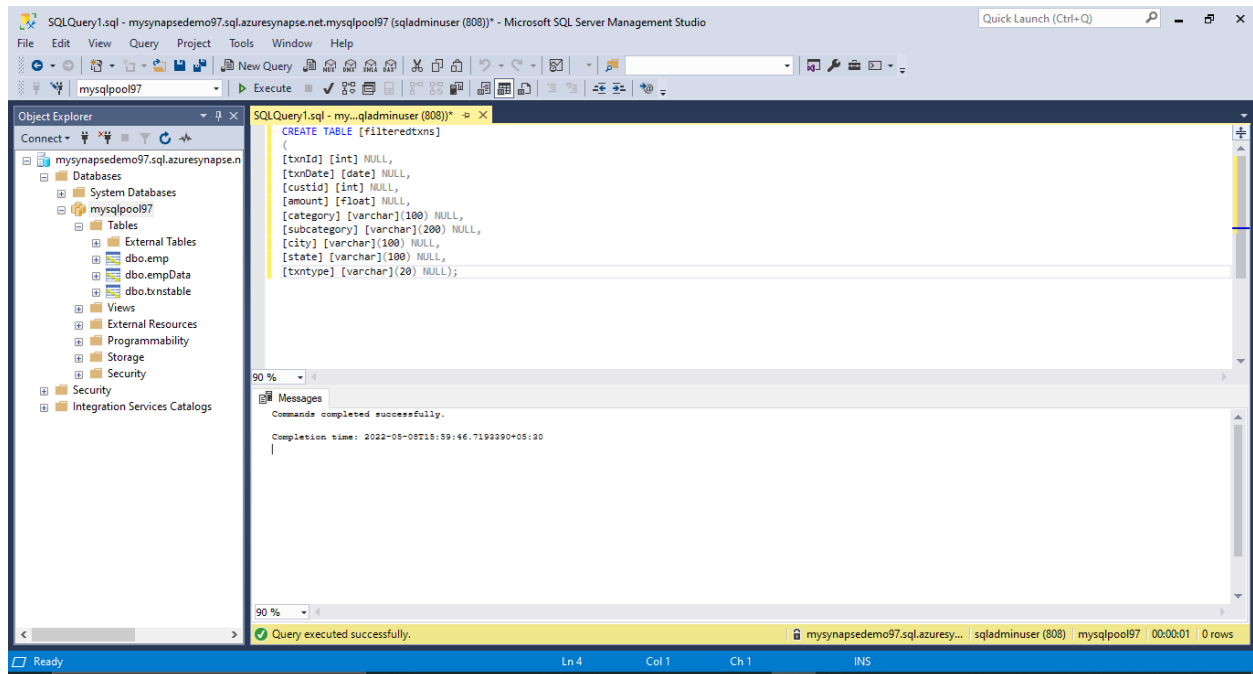
## Source Configurations:

## Sink Configurations:

a. Created a table 'filteredtxns' in Synapse SQLPool



b. Selected Synapse SQLPool table dbo.filteredtxns

Mapping:

<u>Settings:</u> In settings we enabled staging at a container named 'staging' in ADLS storage account 'adlsdemo97'.





**7.** Finally after validating the pipeline with no errors we published the pipeline and then triggered it.

As shown above, it ran successfully.

To confirm we will run query in SSMS as "select * from [filteredtxns];" to fetch all records: