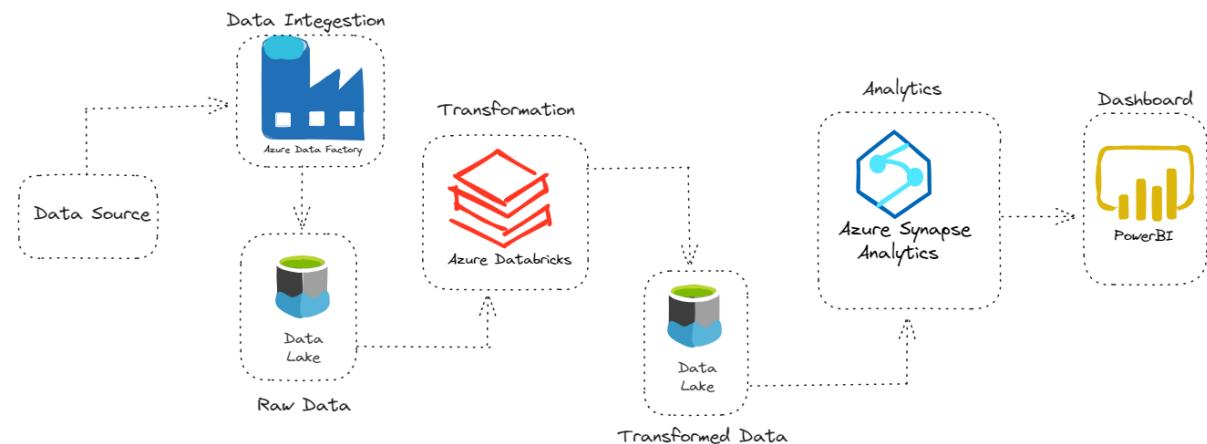


Tokyo 2021 Olympic Data Analytics Project

Pipeline Design Architecture:

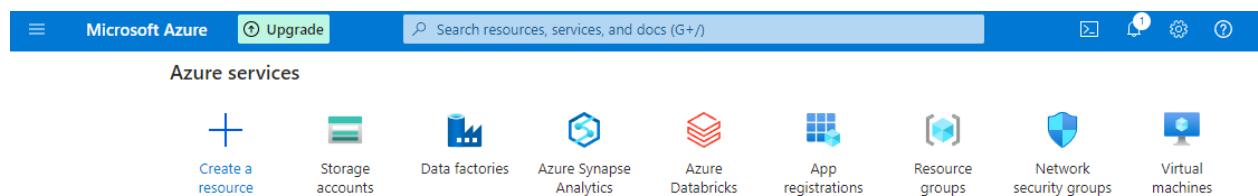


The idea of this project is to analyse Tokyo 2021 Olympics data using data engineering tools and technologies, including Microsoft Azure services such as Azure Data Factory, Data Lake Gen 2, Azure Synapse Analytics and Azure Databricks.

Dataset Used

This dataset contains the details of over 11,000 athletes, with 47 disciplines, along with 743 Teams taking part in the 2021(2020) Tokyo Olympics. This dataset contains the details of the Athletes, Coaches, Teams participating as well as the Entries by gender. It contains their names, countries represented, discipline, gender of competitors, name of the coaches.

Source(Kaggle): <https://www.kaggle.com/datasets/arjunprasadsarkhel/2021-olympics-in-tokyo>



List of Azure Services used for this project

1. **Azure Data Factory:** For data ingestion from GitHub.
2. **Azure Data Lake Storage Gen2:** As the primary data storage solution.
3. **Azure Databricks:** For data transformation tasks.
4. **Azure Synapse Analytics:** To perform in-depth data analytics.

Azure services

Create a resource

tokyo-olympic

Azure Databricks

App registrations

Resource groups

Network security groups

Virtual machines

More

Recent

View

Name

tokyo-olympic-c

tokyo-olympic-c

tokyoolympicda

tokyo-olympic-s

tokyo-olympic

dbstoragergev7

See all

Resource details

Type: Resource group
Location: Southeast Asia
Subscription: Free Trial
Resource group: tokyo-olympic
Status: Succeeded

Type	Last Viewed
Azure Databricks Service	5 minutes ago
Data factory (V2)	28 minutes ago
Storage account	32 minutes ago
Synapse workspace	4 hours ago
Resource group	4 hours ago
Storage account	a day ago

Microsoft Azure

Home >

Resource groups

Default Directory

Create Manage view Refresh Export to CSV Open query Assign tags

Filter for any field... Subscription equals all Location equals all Add filter

No grouping

Name	Subscription	Location
databricks-rg-tokyo-olympic-db-3kcdqqbg62qsu	Free Trial	Southeast Asia
NetworkWatcherRG	Free Trial	East US
synapseworkspace-managedrg-4fbf0656-148e-422d-b647-2a47c9653641	Free Trial	Southeast Asia
tokyo-olympic	Free Trial	Southeast Asia

< Previous Page 1 of 1 Next >

Resource group created.

Microsoft Azure Upgrade Search resources, services, and docs (G+/-) ...

Home > Data factories ...

Default Directory

+ Create Manage view Refresh Export to CSV Open query Assign tags

Filter for any field... Subscription equals all Type equals all Resource group equals all Location equals all Add filter

Showing 1 to 1 of 1 records.

Name	Type	Subscription	Resource group	Location
tokyo-olympic-dfac24	Data factory (V2)	Free Trial	tokyo-olympic	South Central US

< Previous Page 1 of 1 Next >

Microsoft Azure Upgrade Search resources, services, and docs (G+/-) ...

Home > Data factories < Data factories ...

Default Directory

+ Create Manage view ...

Filter for any field... Name ↑↓

tokyo-olympic-dfac24

tokyo-olympic-dfac24 ...

Overview Activity log Access control (IAM) Tags Diagnose and solve problems

Settings Networking Managed identities Properties Locks

Getting started Quick start

Monitoring

Search ...

Delete Subscription ID: 5a2b62fd-8a3e-4a3a-98ad-a90d9b0b98f6

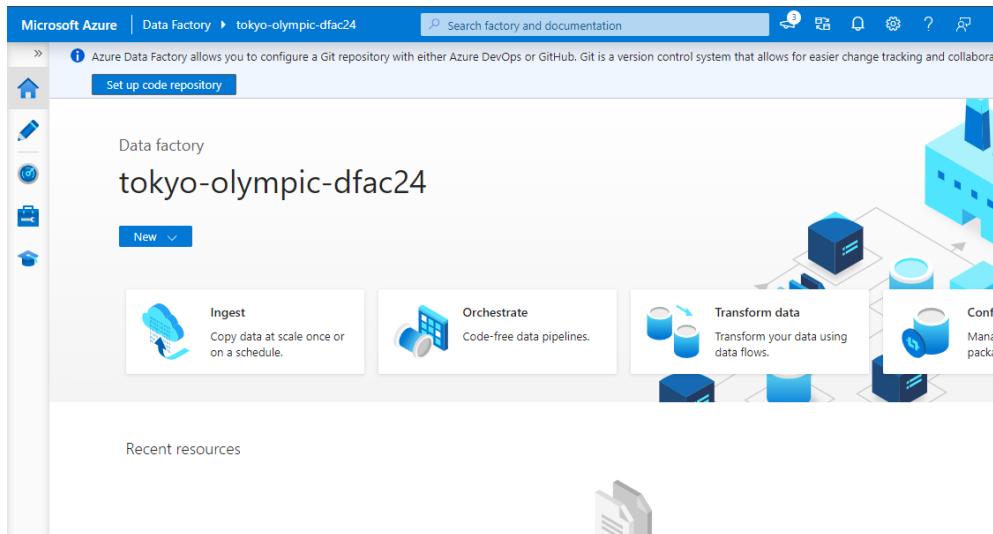
Azure Data Factory Studio

Launch studio

Quick Starts Tutorials Template Gallery Training Modules

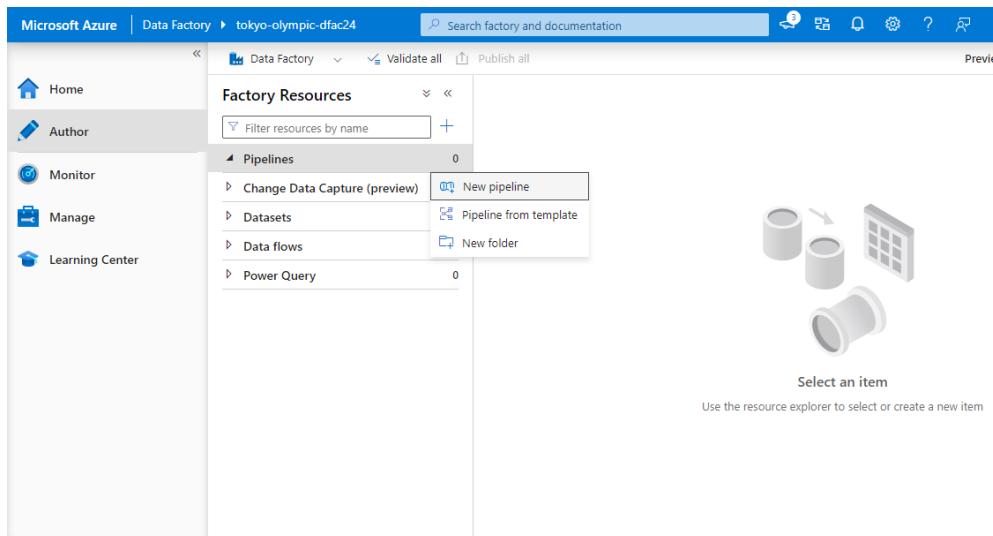
Created the Data Factory Studio for Data injection from GitHub(HTML) Source.

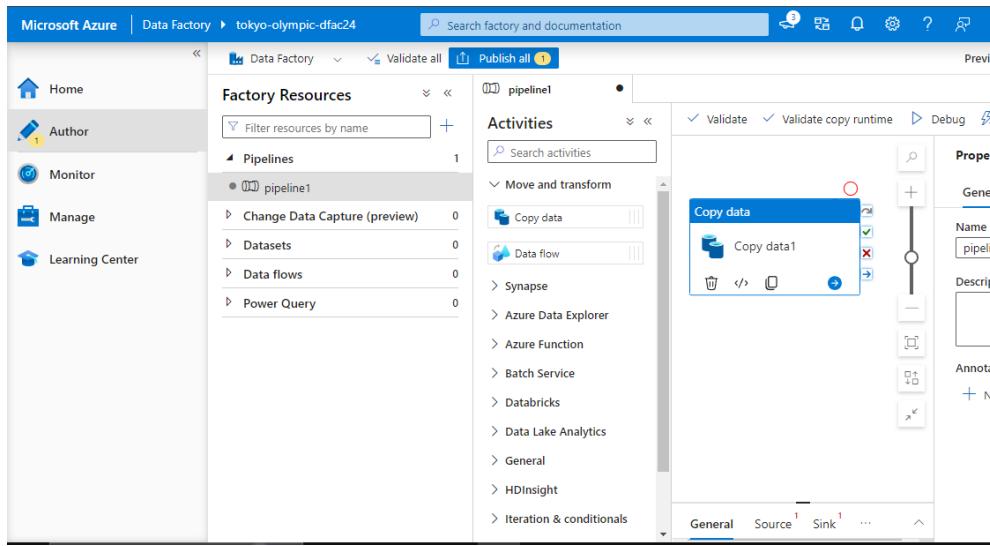




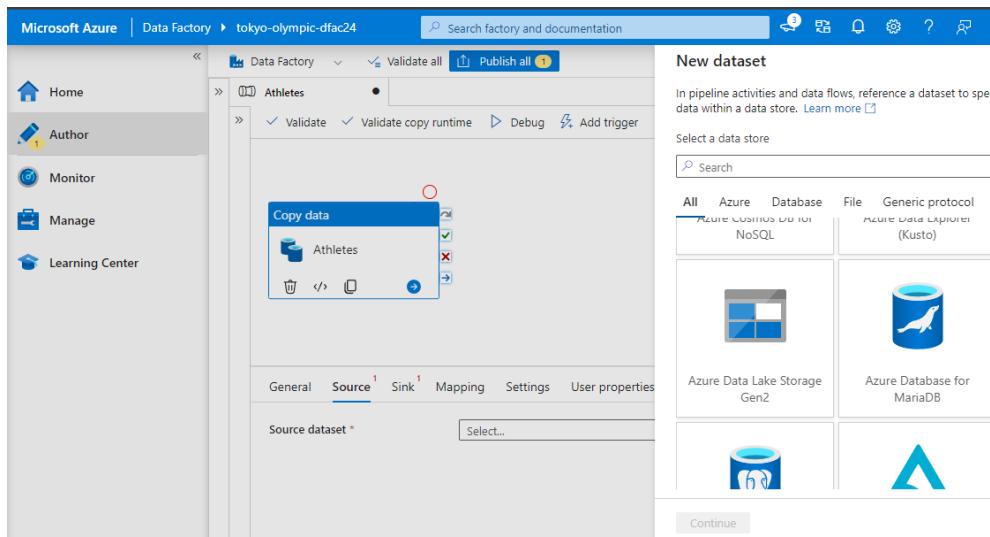
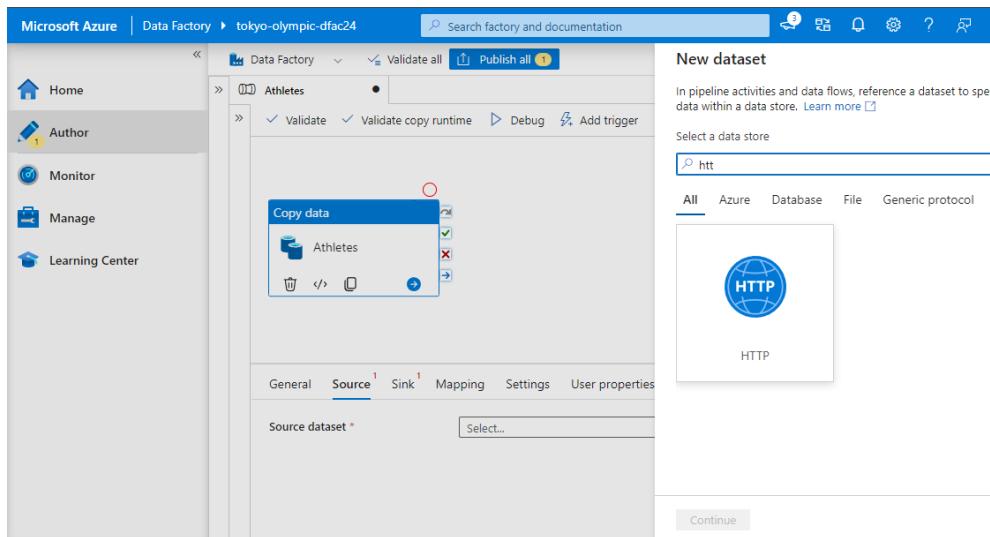
1. Data Ingestion using Azure Data Factory
2. Begin by creating an Azure Data Factory workspace within the previously established resource group.
3. After setting up the workspace, launch the Azure Data Factory Studio.
4. Upload the Tokyo Olympics dataset from kaggle to GitHub.
5. Within the studio, initialize a new data integration pipeline. Now use the task Copy Data to move data efficiently between various supported sources and destinations.
6. Configuring the Data Source with HTTP template as we are using http request to get the data from Github repo.
7. Establishing the Linked Service for source.
8. Configuring the File Format for and setting up the Linked Service Sink.
9. Repeat above steps to load all the datasets.

You can connect all the copy data activity together and run them all at once.





Created the pipeline → Used the Copy Data Activity → For Data Ingestion & pipeline creation.



Microsoft Azure | Data Factory > tokyo-olympic-dfac24

Search factory and documentation

Home Author Monitor Manage Learning Center

Data Factory Validate all Publish all

Athletes

Validate Validate copy runtime Debug Add trigger

Copy data

Athletes

Source Sink Mapping Settings User properties

Source dataset * Select...

New dataset

In pipeline activities and data flows, reference a dataset to specify data within a data store. Learn more

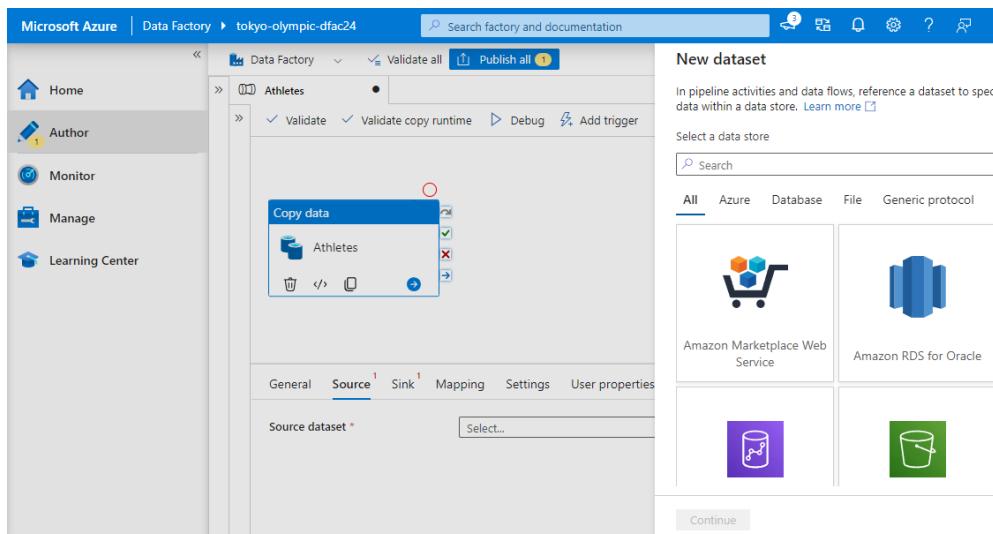
Select a data store

Search

All Azure Database File Generic protocol

 Amazon Marketplace Web Service	 Amazon RDS for Oracle
 Azure Storage	 Azure Blob Storage

Continue



Microsoft Azure | Data Factory > tokyo-olympic-dfac24

Search factory and documentation

Home Author Monitor Manage Learning Center

Data Factory Validate all Publish all

Athletes

Validate Validate copy runtime Debug Add trigger

Copy data

Athletes

Source Sink Mapping Settings User properties

Source dataset * Select...

New dataset

In pipeline activities and data flows, reference a dataset to specify data within a data store. Learn more

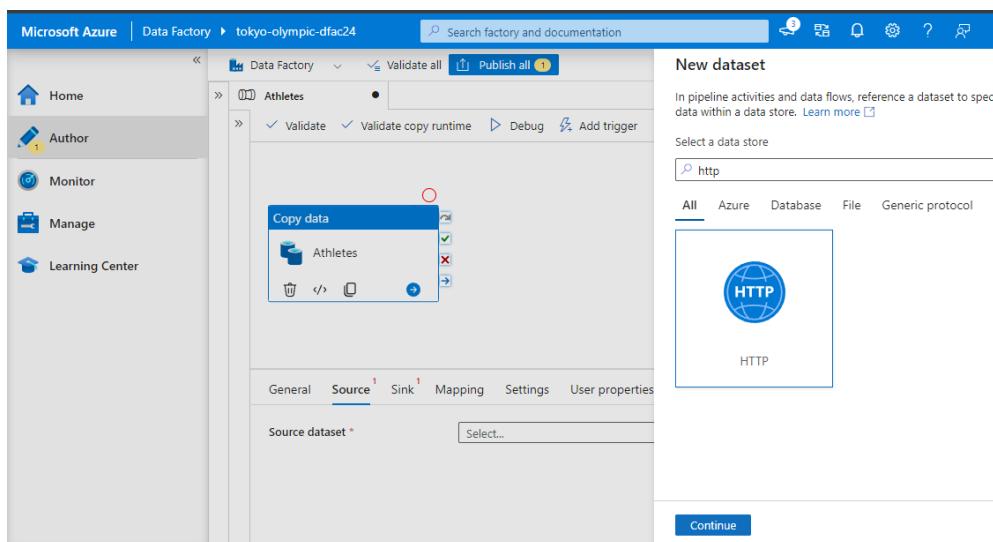
Select a data store

http

All Azure Database File Generic protocol

 HTTP

Continue



Microsoft Azure | Data Factory > tokyo-olympic-dfac24

Search factory and documentation

Home Author Monitor Manage Learning Center

Data Factory Validate all Publish all

Athletes

Validate Validate copy runtime Debug Add trigger

Copy data

Athletes

Source Sink Mapping Settings User properties

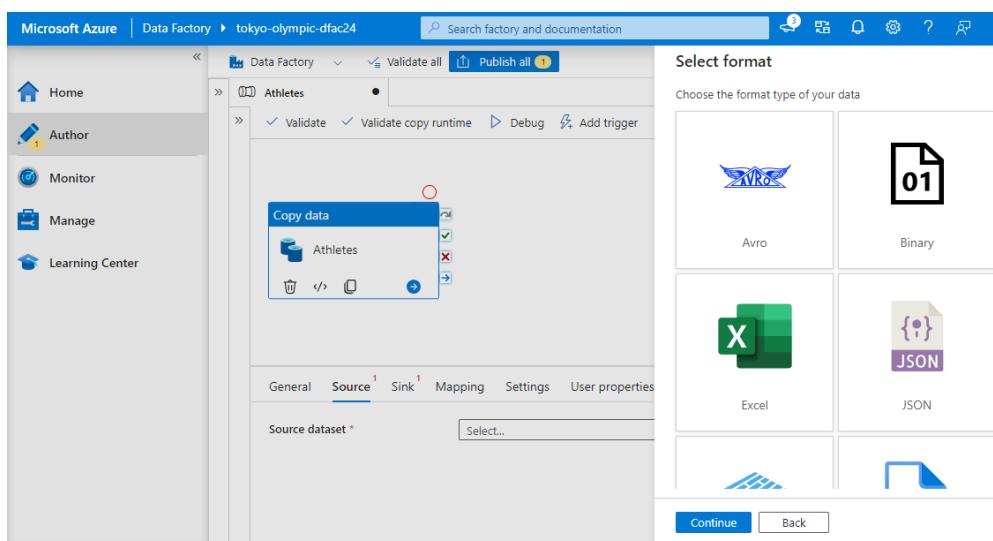
Source dataset * Select...

Select format

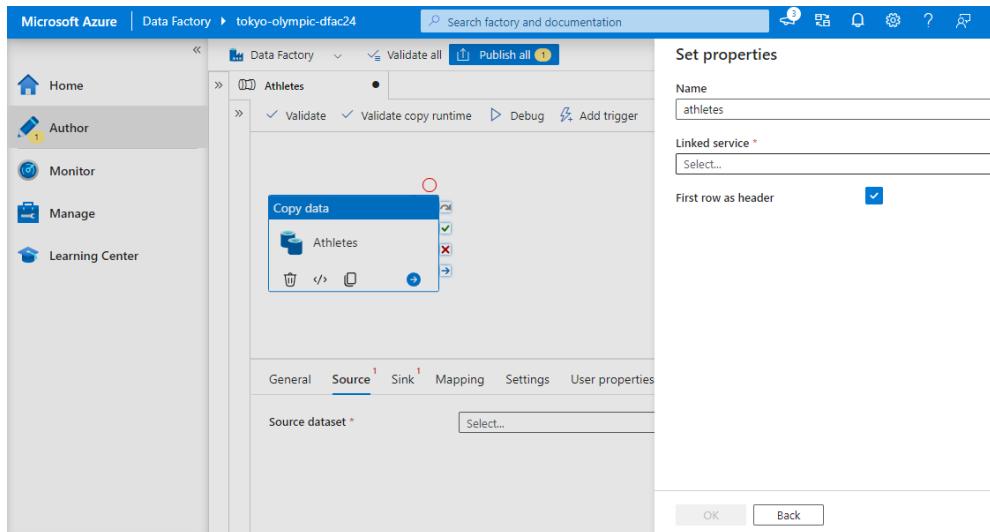
Choose the format type of your data

 Avro	 Binary
 Excel	 JSON
 Delimited Text	 Text

Continue Back



Source file is in CSV. So, I selected the Delimited Text



	PersonName	Country	Discipline
1	AALERUD Katrine	Norway	Cycling Road
2	ABAD Nestor	Spain	Artistic Gymnastics
3	ABADE Alberto	Spain	Basketball
4	ABALDE Tamara	Spain	Basketball
5	ABALDE Luc	France	Handball
6	ABAROA Cesar	Chile	Rowing

Github Data Source Link : <https://raw.githubusercontent.com/darshilparmar/tokyo-olympic-azure-data-engineering-project/main/data/Athletes.csv>

```

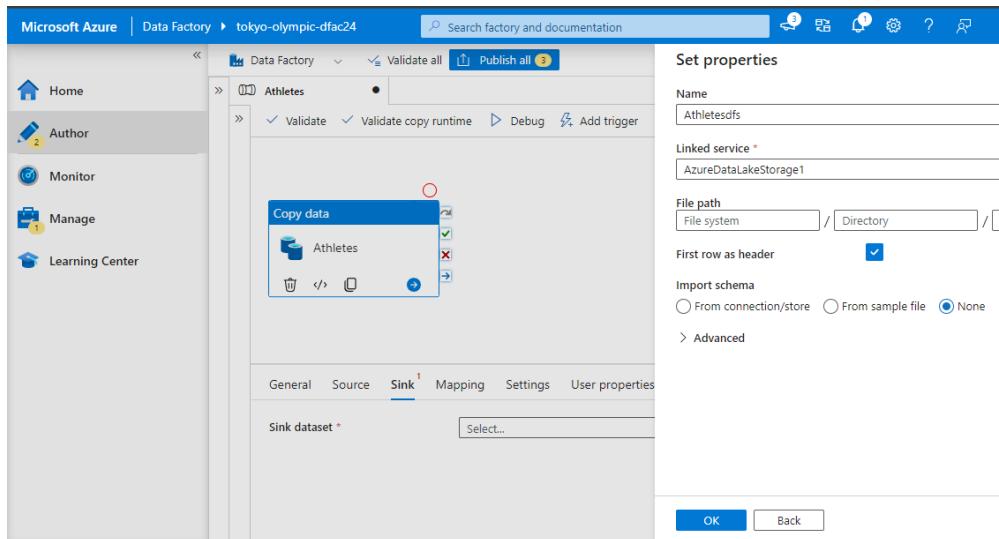
PersonName,Country,Discipline
AALERUD Katrine,Norway,Cycling Road
ABAD Nestor,Spain,Artistic Gymnastics
ABAGNALE Giovanni,Italy,Rowing
ABALDE Alberto,Spain,Basketball
ABALDE Tamara,Spain,Basketball
ABALDE Luc,France,Handball
ABAROA Cesar,Chile,Rowing
ABASS Abobakr,Sudan,Swimming
ABBASALI Hamideh,Islamic Republic of Iran,Karate
ABBASOON Islam,Azerbaijan,Wrestling
ABBINGA Loes,Netherlands,Handball
ABBOT Emily,Australia,Rhythmic Gymnastics
ABBOTT Monica,United States of America,Baseball/Softball
ABDALLA Abubaker Haydar,Qatar,Athletics
ABDALLA Maryam,Egypt,Artistic Swimming
ABDALLAH Shahid,Egypt,Artistic Swimming
ABDALRASOOL Mohamed,Sudan,Judo
ABDEL LATIF Badwa,Egypt,Shooting
ABDEL RAZER Samy,Egypt,Shooting
ABDELRAHMAN Abd El,Egypt,Volleyball
ABDELRAZZZ Farah,Egypt,Table Tennis
ABDELRAZZZ Feryal,Egypt,Karate
ABDELHAMMOUD Mohamed,Egypt,Judo
ABDELNATTALEB Diaeldin Kamal Gouda,Egypt,Wrestling
ABDELRAHMAN Ihab,Egypt,Athletics
ABDELSALAM Mohamed,Egypt,Football
ABDELSALAM Nour,Egypt,Taekwondo
ABDELMWAHED Ahmed,Italy,Athletics
ABDI Bashir,Belgium,Athletics
ABDIRAHMAN Abdi,United States of America,Athletics
ABDUL HADI Farah Ann,Malaysia,Artistic Gymnastics
ABDUL RAHMAN Kirin Tikanah,Singapore,Fencing
ABDUL RAZZAQ Fathimath Nabaha,Maldives,Badminton
ABDULHAMID Saud,Saudi Arabia,Football
ABDULLABBAR Ammar Riad,Germany,Boxing
ABDULLAEV Guomjon,Uzbekistan,Wrestling
ABDULLAEV Muninjon,Uzbekistan,Wrestling
ABDULLAH Rezmat Erwin,Indonesia,Weightlifting
ABDULLAH Ifeaz,Kuwait,Archery

```

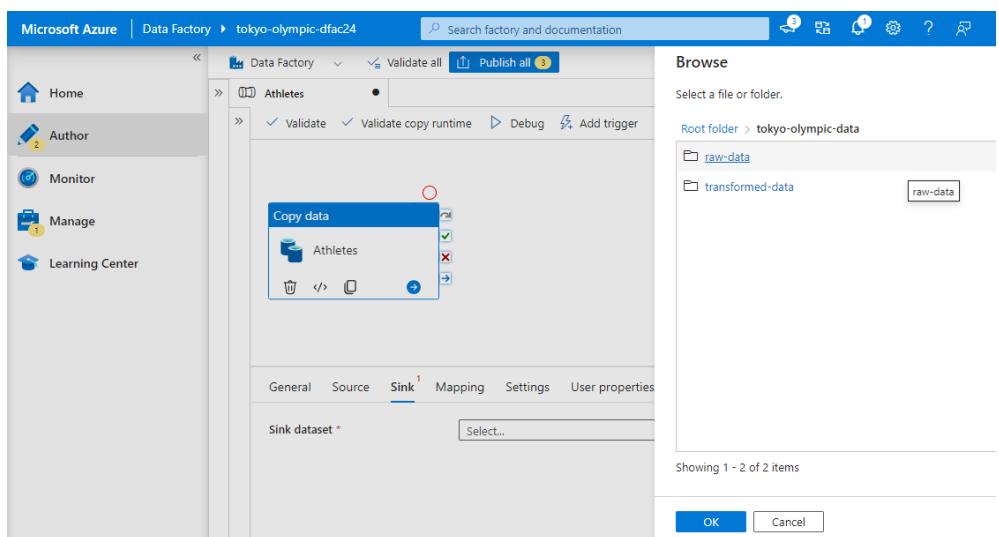
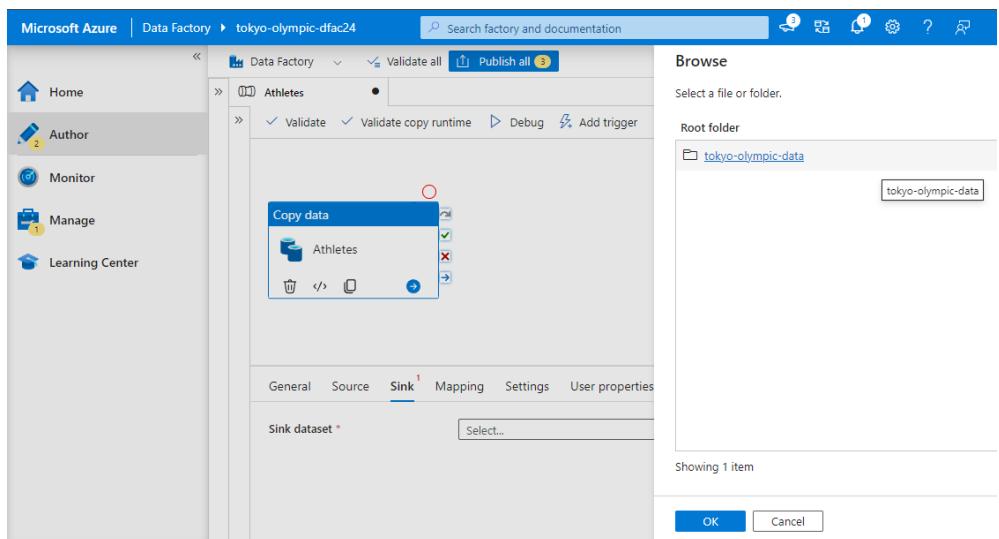
The screenshot shows the Microsoft Azure Data Factory Author interface. On the left sidebar, the 'Author' option is selected. In the main workspace, a pipeline named 'Athletes' is open. A 'Copy data' activity is visible within the pipeline. On the right, a 'New linked service' dialog is displayed for an 'HTTP' type. The 'Name' field is set to 'AthletesHTTP'. The 'Base URL' field contains the value 'https://raw.githubusercontent.com/darshilparmar/tokyo-olympic-data/1.0/Athletes/Athletes.csv'. The 'Authentication type' is set to 'Anonymous'. At the bottom of the dialog are 'Create' and 'Cancel' buttons.

Copy and paste the data source URL and to create the Linked service between the Azure Data factory and GitHub api and ingest the Data from Github HTML API with the help of Linked service.

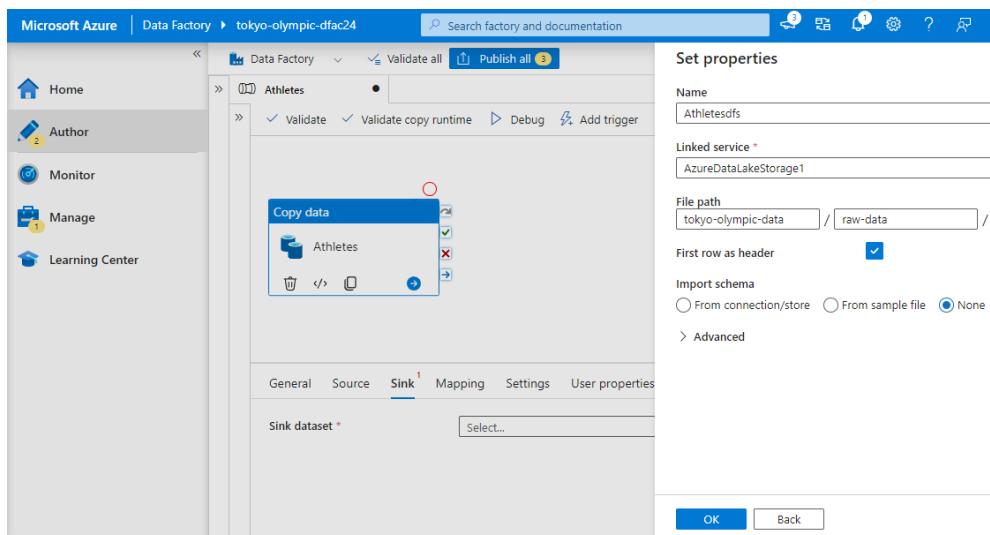
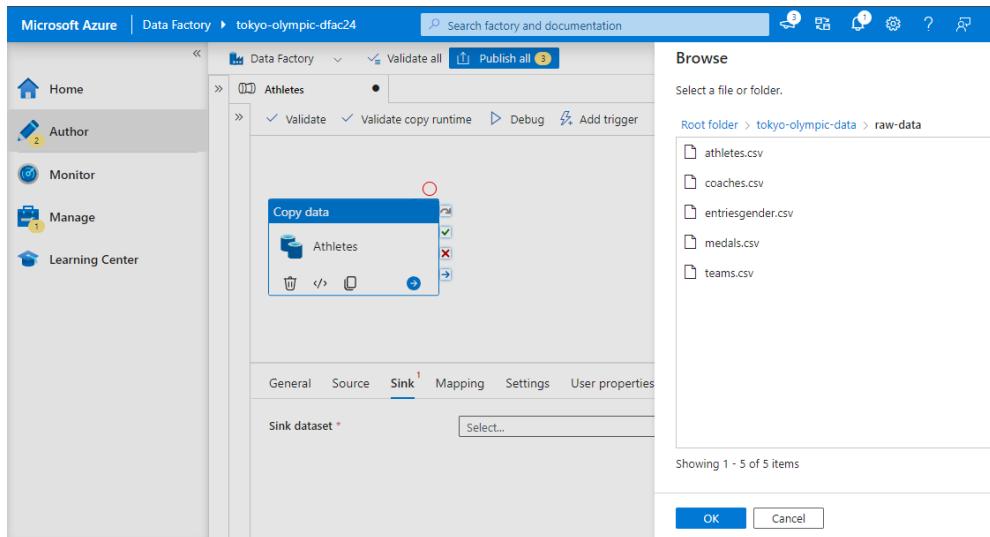
The screenshot shows the Microsoft Azure Data Factory Author interface. The 'Sink' tab is selected for the 'Copy data' activity. A 'New dataset' dialog is open, prompting the user to 'Select a data store'. Below the search bar, tabs for 'All', 'Azure', 'Database', 'File', and 'Generic protocol' are shown. Under the 'All' tab, four options are listed: 'Azure AI Search' (with a magnifying glass icon), 'Azure Blob Storage' (with a cloud icon), 'Azure Data Lake Storage' (with a globe icon), and 'Azure Event Hubs' (with a speech bubble icon). At the bottom of the dialog is a 'Continue' button.



Configure the SINK folder in the Azure Data Lake Gen2 Storage → create containers as Raw DATA, Transformed data.



create containers as Raw DATA, Transformed data.



Microsoft Azure | Data Factory > tokyo-olympic-dfac24

Search factory and documentation

Author (3) Home Monitor Manage Learning Center

Factory Resources Pipelines 1 Olympicdatapipeline

Copy data Athletes

Copy data Coachess

General Source Sink Mapping Settings User properties

Source dataset * Select... New

Olympicdatapipeline Validate Validate copy runtime Debug Add trigger

Publish all

Search

Upload Open in Explorer Delete Move Refresh Open in mobile CLI / PS Feedback

Essentials

Resource group (move) : tokyo-olympic Performance : Standard

Location : southeastasia Replication : Read-access geo-redundant storage (RA-R)

Primary/Secondary Location : Primary: Southeast Asia, Secondary: East Asia Account kind : StorageV2 (general purpose v2)

Subscription (move) : Free Trial Provisioning state : Succeeded

Subscription ID : 5a2b62fd-8a3e-4a3a-98ad-a90d9b0b98f6 Created : 4/1/2024, 1:16:18 PM

Disk state : Primary: Available, Secondary: Available

Tags (edit) : Add tags

Properties Monitoring Capabilities (5) Recommendations (0) Tutorials Tools + SDKs

Data Lake Storage Security

Hierarchical namespace Enabled Require secure transfer for REST API operations Enabled

Default access tier Hot Storage account key access Enabled

Blob anonymous access Disabled Minimum TLS version Version 1.2

Search

Container Change access level Restore containers Refresh Delete Give feedback

Search containers by prefix Show deleted containers

Name	Last modified	Anonymous access level	Lease state
Slogs	1/4/2024, 1:16:50 pm	Private	Available
tokyo-olympic-data	1/4/2024, 1:20:23 pm	Private	Available

Data migration Events Storage browser Data storage Containers File shares Queues Tables Security + networking Networking Access keys Shared access signature Encryption

The screenshot shows the 'Overview' page of a blob container named 'raw-data'. The container is located in the 'tokyo-olympic-data' storage account. The table lists two blobs: 'raw-data' and 'transformed-data'. Both blobs are of type 'Block blob' and have an 'Available' lease state.

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
raw-data				Block blob	-	***
transformed-data				Block blob	-	***

10. After the pipeline completes its execution, navigate to your Azure Data Lake Storage Gen2. Dive into the "raw_data" folder and validate that the files, like "athletes.csv", "medals.csv", etc., are present and populated with the expected data.

The screenshot shows the 'Overview' page of the 'raw-data' container. The container is located in the 'tokyo-olympic-data / raw-data' path. The table lists several CSV files: 'athletes.csv', 'coaches.csv', 'entriesgender.csv', 'medals.csv', and 'teams.csv'. All files are of type 'Block blob' and have an 'Available' lease state. There is also a placeholder entry '[-]'.

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
[-]						***
athletes.csv	2/4/2024, 8:38:27 am	Hot (Inferred)		Block blob	408.68 KiB	Available
coaches.csv	2/4/2024, 8:38:40 am	Hot (Inferred)		Block blob	16.49 KiB	Available
entriesgender.csv	2/4/2024, 8:39:02 am	Hot (Inferred)		Block blob	1.1 KiB	Available
medals.csv	2/4/2024, 8:39:19 am	Hot (Inferred)		Block blob	2.36 KiB	Available
teams.csv	2/4/2024, 8:39:35 am	Hot (Inferred)		Block blob	34.44 KiB	Available

The screenshot shows the 'Overview' page of the 'raw-data' container. The container is located in the 'tokyo-olympic-data' storage account. The table lists the same five CSV files: 'athletes.csv', 'coaches.csv', 'entriesgender.csv', 'medals.csv', and 'teams.csv'. All files are of type 'Block blob' and have an 'Available' lease state.

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
[-]						***
athletes.csv	2/4/2024, 8:38:27 am	Hot (Inferred)		Block blob	408.68 KiB	Available
coaches.csv	2/4/2024, 8:38:40 am	Hot (Inferred)		Block blob	16.49 KiB	Available
entriesgender.csv	2/4/2024, 8:39:02 am	Hot (Inferred)		Block blob	1.1 KiB	Available
medals.csv	2/4/2024, 8:39:19 am	Hot (Inferred)		Block blob	2.36 KiB	Available
teams.csv	2/4/2024, 8:39:35 am	Hot (Inferred)		Block blob	34.44 KiB	Available

Check all the data available in the container.

Search

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

Overview

Diagnose and solve problems

Access Control (IAM)

Settings

Shared access tokens

Manage ACL

Access policy

Properties

Metadata

Authentication method: Access key (Switch to Microsoft Entra user account)
Location: tokyo-olympic-data / transformed-data

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

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
[...]					-	***
athletes					-	***
coaches					-	***
entriesgender					-	***
medals					-	***
teams					-	***

Search

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

Overview

Diagnose and solve problems

Access Control (IAM)

Settings

Shared access tokens

Manage ACL

Access policy

Properties

Metadata

Authentication method: Access key (Switch to Microsoft Entra user account)
Location: tokyo-olympic-data / transformed-data / athletes

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

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
[...]					-	***
_committed_475657...	4/4/2024, 9:02:57 am	Hot (Inferred)		Block blob	112 B	Available
_started_4756574684...	4/4/2024, 9:02:57 am	Hot (Inferred)		Block blob	0 B	Available
_SUCCESS	4/4/2024, 9:02:57 am	Hot (Inferred)		Block blob	0 B	Available
part-00000-tid-4756...	4/4/2024, 9:02:57 am	Hot (Inferred)		Block blob	397.91 KiB	Available

Search

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

Overview

Diagnose and solve problems

Access Control (IAM)

Settings

Shared access tokens

Manage ACL

Access policy

Properties

Metadata

Authentication method: Access key (Switch to Microsoft Entra user account)
Location: tokyo-olympic-data / transformed-data / coaches

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

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
[...]					-	***
_committed_202135...	4/4/2024, 9:02:58 am	Hot (Inferred)		Block blob	112 B	Available
_started_2021355756...	4/4/2024, 9:02:58 am	Hot (Inferred)		Block blob	0 B	Available
_SUCCESS	4/4/2024, 9:02:58 am	Hot (Inferred)		Block blob	0 B	Available
part-00000-tid-2021...	4/4/2024, 9:02:58 am	Hot (Inferred)		Block blob	16.12 KiB	Available

Search «

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

Overview

Diagnose and solve problems

Access Control (IAM)

Settings

Shared access tokens

Manage ACL

Access policy

Properties

Metadata

Authentication method: Access key ([Switch to Microsoft Entra user account](#))
Location: tokyo-olympic-data / transformed-data / entriesgender

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

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
[...]						***
_committed_338725...	4/4/2024, 9:02:59 am	Hot (Inferred)		Block blob	112 B	Available
_started_3387255666...	4/4/2024, 9:02:58 am	Hot (Inferred)		Block blob	0 B	Available
_SUCCESS	4/4/2024, 9:02:59 am	Hot (Inferred)		Block blob	0 B	Available
part-00000-tid-3387...	4/4/2024, 9:02:58 am	Hot (Inferred)		Block blob	1.05 KiB	Available

Search «

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

Overview

Diagnose and solve problems

Access Control (IAM)

Settings

Shared access tokens

Manage ACL

Access policy

Properties

Metadata

Authentication method: Access key ([Switch to Microsoft Entra user account](#))
Location: tokyo-olympic-data / transformed-data / entriesgender

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

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
[...]						***
_committed_338725...	4/4/2024, 9:02:59 am	Hot (Inferred)		Block blob	112 B	Available
_started_3387255666...	4/4/2024, 9:02:58 am	Hot (Inferred)		Block blob	0 B	Available
_SUCCESS	4/4/2024, 9:02:59 am	Hot (Inferred)		Block blob	0 B	Available
part-00000-tid-3387...	4/4/2024, 9:02:58 am	Hot (Inferred)		Block blob	1.05 KiB	Available

Search «

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

Overview

Diagnose and solve problems

Access Control (IAM)

Settings

Shared access tokens

Manage ACL

Access policy

Properties

Metadata

Authentication method: Access key ([Switch to Microsoft Entra user account](#))
Location: tokyo-olympic-data / transformed-data / medals

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

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
[...]						***
_committed_279853...	4/4/2024, 9:02:59 am	Hot (Inferred)		Block blob	112 B	Available
_started_2798538771...	4/4/2024, 9:02:59 am	Hot (Inferred)		Block blob	0 B	Available
_SUCCESS	4/4/2024, 9:02:59 am	Hot (Inferred)		Block blob	0 B	Available
part-00000-tid-2798...	4/4/2024, 9:02:59 am	Hot (Inferred)		Block blob	2.27 KiB	Available

Search «

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

Overview

Diagnose and solve problems

Access Control (IAM)

Settings

Shared access tokens

Manage ACL

Access policy

Properties

Metadata

Authentication method: Access key ([Switch to Microsoft Entra user account](#))
Location: tokyo-olympic-data / transformed-data / teams

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

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
[...]						***
_committed_217174...	4/4/2024, 9:03:00 am	Hot (Inferred)		Block blob	112 B	Available
_started_2171749806...	4/4/2024, 9:03:00 am	Hot (Inferred)		Block blob	0 B	Available
_SUCCESS	4/4/2024, 9:03:00 am	Hot (Inferred)		Block blob	0 B	Available
part-00000-tid-2171...	4/4/2024, 9:03:00 am	Hot (Inferred)		Block blob	33.79 KiB	Available

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

Overview

Diagnose and solve problems Access Control (IAM)

Settings

- Shared access tokens
- Manage ACL
- Access policy
- Properties
- Metadata

Authentication method: Access key ([Switch to Microsoft Entra user account](#))
Location: [tokyo-olympic-data](#) / transformed-data

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

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
[.]					-	***
athletes					-	***
coaches					-	***
entriesgender					-	***
medals					-	***
teams					-	***

Microsoft Azure Upgrade Search resources, services, and docs (G+)

Home > **tokyo-olympic-db** Azure Databricks Service

Search Delete

Overview

Activity log Access control (IAM) Tags Diagnose and solve problems

Settings

- Virtual Network Peering
- Encryption
- Networking
- Security & compliance
- Properties
- Locks

Monitoring

Essentials

Status : Active Managed Resource Group : [databricks-rg-tokyo-olympic](#)
Resource group : [tokyo-olympic](#) URL : <https://adb-76715826667401.azuredatabricks.net>
Location : Southeast Asia Pricing Tier : Premium (+ Role-based access)

Subscription : [Free Trial](#) Subscription ID : [5a2b62fd-8a3e-4a3a-98ad-a90d9b0b90f6](#)

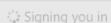
Tags [\(edit\)](#) : [Add tags](#)


[Launch Workspace](#)

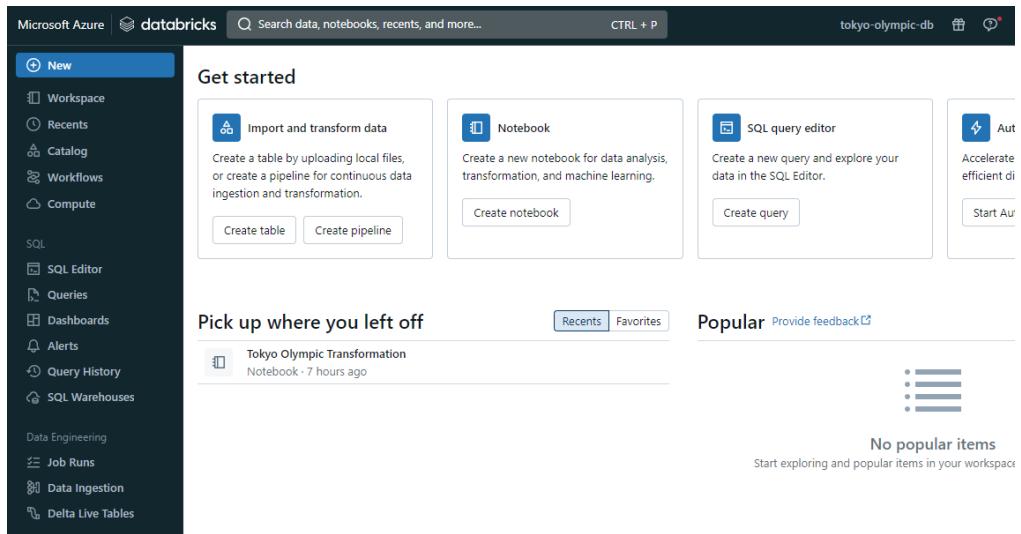
Azure Databricks

Sign In to Databricks

Sign in using Azure Active Directory Single Sign On.

 Signing you in

Contact your site administrator to request access.



Do this transformation by using Data Bricks.

Data Transformation using Azure Databricks

1. Navigate to Azure Databricks within the Azure portal and create a workspace within the previously established resource group and launch it.
2. Configuring Compute in Databricks
3. Create a new notebook within Databricks and rename it appropriately, reflecting its purpose or the dataset it pertains to.
4. Establishing a Connection to Azure Data Lake Storage (ADLS)
5. Using the credentials (Client ID, Tenant ID, Secret), write the appropriate code in the Databricks notebook to mount ADLS.
6. Writing Data Transformations mount ADLS Gen2 to Databricks.
7. Writing Transformed Data to ADLS Gen2.

```

from pyspark.sql.functions import col
from pyspark.sql.types import IntegerType, DoubleType, BooleanType, DataType

```

```

1 Last execution failed
2
3 configs = {"fs.azure.account.auth.type": "OAuth",
4 "fs.azure.account.oauth.provider.type": "org.apache.hadoop.fs.azurebfs.oauth2.ClientCredsTokenProvider",
5 "fs.azure.account.oauth2.client.id": "5f5ff3d8-f05b-4829-bc43-f3b743300ace",
6 "fs.azure.account.oauth2.client.secret": "pu48Q-oqjTg2xsI9V7JHDHmuppTAOzVfSvku0da9",
7 "fs.azure.account.oauth2.client.endpoint": "https://login.microsoftonline.com/1bfcf83e-1bb2-4d03-8384-8466413
8 oauth2/token"
9
10 dbutils.fs.mount(
11 source = "abfss://tokyo-olympic-data@tokyoolympicdatachandru0.dfs.core.windows.net", # contrainer@storageacc
12 mount_point = "/mnt/tokyoolympic",
13 extra_configs = configs
14
15 > java.rmi.RemoteException: java.lang.IllegalArgumentException: requirement failed: Directory already mounted: /
16 mnic: nested exception is:

```

Create the Compute Resource & Mount the data into Data-Bricks.

```

12:59 PM (2s)
%fs
ls "/mnt/tokyoolympic/"

Table + 
path          name        size  modificationTime
1  dbfs:/mnt/tokyoolympic/raw-data/ raw-data/  0   1711957877000
2  dbfs:/mnt/tokyoolympic/transformed-data/ transformed-data/  0   1711957894000

2 rows | 1.83 seconds runtime
Refresh

```

```

01:00 PM (1s)
spark

SparkSession - hive
SparkContext
Spark UI
Version
  v3.3.2
Master
  local[*, 4]
AppName
  Databricks Shell

```

```

01:00 PM (9s)
athletes = spark.read.format("csv").option("header","true").option("InferSchema","true").load("/mnt/tokyoolympic/r
athletes.csv")
coaches = spark.read.format("csv").option("header","true").option("InferSchema","true").load("/mnt/tokyoolympic/r
coaches.csv")
entriesgender = spark.read.format("csv").option("header","true").option("InferSchema","true").load("/mnt/tokyoolympic
entriesgender.csv")

```

Microsoft Azure | databricks | Search data, notebooks, recents, and more... | CTRL + P | tokyo-olympic-db | Balachandar s's Cluster

Tokyo Olympic Transformation | Python | 5

```
athletes = spark.read.format("csv").option("header","true").option("InferSchema","true").load("/mnt/tokyoolympic/raw/athletes.csv")
coaches = spark.read.format("csv").option("header","true").option("InferSchema","true").load("/mnt/tokyoolympic/raw/coaches.csv")
entriesgender = spark.read.format("csv").option("header","true").option("InferSchema","true").load("/mnt/tokyoolympic/raw/entriesgender.csv")
medals = spark.read.format("csv").option("header","true").option("InferSchema","true").load("/mnt/tokyoolympic/raw/medals.csv")
teams = spark.read.format("csv").option("header","true").option("InferSchema","true").load("/mnt/tokyoolympic/raw/teams.csv")
```

(10) Spark Jobs

```
athletes: pyspark.sql.dataframe.DataFrame
  PersonName: string
  Country: string
  Discipline: string
  coaches: pyspark.sql.dataframe.DataFrame = [Name: string, Country: string ... 2 more fields]
  entriesgender: pyspark.sql.dataframe.DataFrame = [Discipline: string, Female: integer ... 2 more fields]
  medals: pyspark.sql.dataframe.DataFrame = [Rank: integer, Team_Country: string ... 5 more fields]
  teams: pyspark.sql.dataframe.DataFrame = [TeamName: string, Discipline: string ... 2 more fields]
```

Microsoft Azure | databricks | Search data, notebooks, recents, and more... | CTRL + P | tokyo-olympic-db | Balachandar s's Cluster

Tokyo Olympic Transformation | Python | 6

```
athletes.show()
```

(1) Spark Jobs

PersonName	Country	Discipline
AALERUD Katrine	Norway	Cycling Road
ABAD Nestor	Spain	Artistic Gymnastics
ABAGNALE Giovanni	Italy	Rowing
ALBALDE Alberto	Spain	Basketball
ALBALDE Tamara	Spain	Basketball
ABALO Luc	France	Handball
ABAROA Cesar	Chile	Rowing
ABASS Abobakr	Sudan	Swimming
ABBASALI Hamideh	Islamic Republic ...	Karate
ABBASOV Islam	Azerbaijan	Wrestling
ABBINGH Lois	Netherlands	Handball
ABBOT Emily	Australia	Rhythmic Gymnastics
ABBOTT Monica	United States of ...	Baseball/Softball
ABDALLA Abubaker ...	Qatar	Athletics
ABDALLA Maryam	Egypt	Artistic Swimming
ABDALLAH Shahd	Egypt	Artistic Swimming

Microsoft Azure | databricks | Search data, notebooks, recents, and more... | CTRL + P | tokyo-olympic-db | Balachandar s's Cluster

Tokyo Olympic Transformation | Python | 7

```
athletes.printSchema()
```

root
|-- PersonName: string (nullable = true)
|-- Country: string (nullable = true)
|-- Discipline: string (nullable = true)

coaches.printSchema()

root
|-- Name: string (nullable = true)
|-- Country: string (nullable = true)
|-- Discipline: string (nullable = true)
|-- Event: string (nullable = true)

9

Microsoft Azure | databricks | Search data, notebooks, recents, and more... | CTRL + P | tokyo-olympic-db | Run all | Balachandar s's Cluster

Tokyo Olympic Transformation Python

```
File Edit View Run Help Last edit was 1 hour ago New cell UI: ON
```

0:00 PM (<1s) coaches.show()

(1) Spark Jobs

Name	Country	Discipline	Event
ABDELMAGID Wael	Egypt	Football	null
ABE Junya	Japan	Volleyball	null
ABE Katsuhiko	Japan	Basketball	null
ADAMA Cherif	Côte d'Ivoire	Football	null
AGEBA Yuya	Japan	Volleyball	null
AIKMAN Siegfried ...	Japan	Hockey	Men
AL SAADI Kais	Germany	Hockey	Men
ALAHMEDA Lonní	Canada	Baseball/Softball	Softball
ALEKNO Vladimir	[Islamic Republic ...]	Volleyball	Men
ALEKSEEV Alexey	ROC	Handball	Women
ALLER CARBALLO Ma...	Spain	Basketball	null
ALSHERI Saad	Saudi Arabia	Football	Men
ALY Kamal	Egypt	Football	null
AMAYA GAITAN Fabian	Puerto Rico	Basketball	null
AMO AGUADO Pablo	Spain	Football	null
ANDONOVSKI Vlatko	[United States of ...]	Football	Women
ANNAN Alyson	Netherlands	Hockey	Women
ARNAU CREUS Xavierl	Japan	Hockey	Women

Microsoft Azure | databricks | Search data, notebooks, recents, and more... | CTRL + P | tokyo-olympic-db | Run all | Balachandar s's Cluster

Tokyo Olympic Transformation Python

```
File Edit View Run Help Last edit was 1 hour ago New cell UI: ON
```

0:00 PM (<1s) entriesgender.show()

(1) Spark Jobs

Discipline	Female	Male	Total
3x3 Basketball	32	32	64
Archery	64	64	128
Artistic Gymnastics	98	98	196
Artistic Swimming	105	0	105
Athletics	969	1072	2041
Badminton	86	87	173
Baseball/Softball	90	144	234
Basketball	144	144	288
Beach Volleyball	48	48	96
Boxing	102	187	289
Canoe Slalom	41	41	82
Canoe Sprint	123	126	249
Cycling BMX Frees...	10	9	19
Cycling BMX Racing	24	24	48
Cycling Mountain ...	38	38	76
Cycling Road	70	131	201
Cycling Track	90	99	189
Diving	72	71	143

Microsoft Azure | databricks | Search data, notebooks, recents, and more... | CTRL + P | tokyo-olympic-db | Run all | Balachandar s's Cluster

Tokyo Olympic Transformation Python

```
File Edit View Run Help Last edit was 2 hours ago New cell UI: ON
```

0:01 PM (<1s) entriesgender.printSchema()

```
root
 |-- Discipline: string (nullable = true)
 |-- Female: integer (nullable = true)
 |-- Male: integer (nullable = true)
 |-- Total: integer (nullable = true)
```

0:01 PM (<1s) entriesgender.withColumn("Female", col("Female").cast(IntegerType()))\n.withColumn("Male", col("Male").cast(IntegerType()))\n.withColumn("Total", col("Total").cast(IntegerType()))

entriesgender: pyspark.sql.dataframe.DataFrame

```
Discipline: string
Female: integer
Male: integer
Total: integer
```

```

root
 |-- Discipline: string (nullable = true)
 |-- Female: integer (nullable = true)
 |-- Male: integer (nullable = true)
 |-- Total: integer (nullable = true)

+-----+-----+-----+
| Discipline|Female|Male|Total|
+-----+-----+-----+
| 3x3 Basketball| 32| 32| 64|
| Archery| 64| 64| 128|
| Artistic Gymnastics| 98| 98| 196|
| Artistic Swimming| 105| 0| 105|
| Athletics| 969| 1072| 2041|
| Badminton| 86| 87| 173|
| Baseball/Softball| 90| 144| 234|
| Basketball| 144| 144| 288|
| Beach Volleyball| 48| 48| 96|
| Boxing| 102| 187| 289|
| Canoe Slalom| 41| 41| 82|
| Canoe Sprint| 123| 126| 249|
| Cycling BMX Frees...| 10| 9| 19|
| Cycling BMX Racing| 24| 24| 48|
| Cycling Mountain ...| 38| 38| 76|
| Cycling Road| 70| 131| 201|
| Cycling Track| 90| 99| 189|

```

Data Transformation part

Rank	Team-Country	Gold	Silver	Bronze	Total
1	Japan	27	14	17	58
2	Great Britain	22	21	22	65
3	ROC	20	28	23	71
4	Australia	17	7	22	46
5	Netherlands	10	12	14	36

Rank	Team-Country	Gold	Silver	Bronze	Total	Rank by Total
1	Japan	27	14	17	58	5
2	Great Britain	22	21	22	65	4
3	ROC	20	28	23	71	3
4	Australia	17	7	22	46	6
5	Netherlands	10	12	14	36	9

Microsoft Azure | databricks Q Search data, notebooks, recents, and more... CTRL + P tokyo-olympic-db Balachandar s's Cluster

Tokyo Olympic Transformation Python

File Edit View Run Help Last edit was 2 hours ago New cell UI: ON

Cell 18: medals.printSchema()

```
root
 |-- Rank: integer (nullable = true)
 |-- Team_Country: string (nullable = true)
 |-- Gold: integer (nullable = true)
 |-- Silver: integer (nullable = true)
 |-- Bronze: integer (nullable = true)
 |-- Total: integer (nullable = true)
 |-- Rank by Total: integer (nullable = true)
```

Cell 19: medals.printSchema()

```
root
 |-- Rank: integer (nullable = true)
 |-- Team_Country: string (nullable = true)
 |-- Gold: integer (nullable = true)
 |-- Silver: integer (nullable = true)
 |-- Bronze: integer (nullable = true)
 |-- Total: integer (nullable = true)
 |-- Rank by Total: integer (nullable = true)
```

Microsoft Azure | databricks Q Search data, notebooks, recents, and more... CTRL + P tokyo-olympic-db Balachandar s's Cluster

Tokyo Olympic Transformation Python

File Edit View Run Help Last edit was 2 hours ago New cell UI: ON

Cell 19: medals.printSchema()

```
root
 |-- Rank: integer (nullable = true)
 |-- Team_Country: string (nullable = true)
 |-- Gold: integer (nullable = true)
 |-- Silver: integer (nullable = true)
 |-- Bronze: integer (nullable = true)
 |-- Total: integer (nullable = true)
 |-- Rank by Total: integer (nullable = true)
```

Cell 20: medals.show()

(1) Spark Jobs

Rank	Team_Country	Gold	Silver	Bronze	Total
1	United States of America	39	41	33	113
2	Great Britain	22	21	22	65
3	ROC	20	28	23	71
4	Australia	17	7	22	46
5	Netherlands	10	12	14	36
6	France	10	12	11	33

Microsoft Azure | databricks Q Search data, notebooks, recents, and more... CTRL + P tokyo-olympic-db Balachandar s's Cluster

Tokyo Olympic Transformation Python

File Edit View Run Help Last edit was 2 hours ago New cell UI: ON

Cell 20: medals.show()

(1) Spark Jobs

Rank	Team_Country	Gold	Silver	Bronze	Total
1	United States of America	39	41	33	113
2	Great Britain	22	21	22	65
3	ROC	20	28	23	71
4	Australia	17	7	22	46
5	Netherlands	10	12	14	36
6	France	10	12	11	33
7	Germany	10	11	16	37
8	Italy	10	10	20	40
9	Canada	7	6	11	24
10	Brazil	7	6	8	21
11	New Zealand	7	6	7	20
12	Cuba	7	3	5	15
13	Hungary	6	7	7	20
14	Republic of Korea	6	4	10	20
15	Poland	4	5	5	14

Microsoft Azure | databricks | Search data, notebooks, recents, and more... | CTRL + P | tokyo-olympic-db | Run all | Balachandar's Cluster

Tokyo Olympic Transformation | Python | 01:02 PM (1s) | 21

```
medals.printSchema()
```

```
root
|-- Rank: integer (nullable = true)
|-- Team_Country: string (nullable = true)
|-- Gold: integer (nullable = true)
|-- Silver: integer (nullable = true)
|-- Bronze: integer (nullable = true)
|-- Total: integer (nullable = true)
|-- Rank by Total: integer (nullable = true)
```

Tokyo Olympic Transformation | Python | 01:02 PM (1s) | 22

```
teams.show()
```

(1) Spark Jobs

	TeamName	Discipline	Country	Event
1	China 3x3 Basketball	People's Republic...	Women	
2	France 3x3 Basketball	France	Women	
3	Italy 3x3 Basketball	Italy	Women	
4	Japan 3x3 Basketball	Japan	Men	
5	Japan 3x3 Basketball	Japan	Women	

Microsoft Azure | databricks | Search data, notebooks, recents, and more... | CTRL + P | tokyo-olympic-db | Run all | Balachandar's Cluster

Tokyo Olympic Transformation | Python | 01:02 PM (1s) | 22

```
teams.show()
```

(1) Spark Jobs

TeamName	Discipline	Country	Event
Belgium 3x3 Basketball	Belgium	Men	
China 3x3 Basketball	People's Republic...	Men	
China 3x3 Basketball	People's Republic...	Women	
France 3x3 Basketball	France	Women	
Italy 3x3 Basketball	Italy	Women	
Japan 3x3 Basketball	Japan	Men	
Japan 3x3 Basketball	Japan	Women	
Latvia 3x3 Basketball	Latvia	Men	
Mongolia 3x3 Basketball	Mongolia	Women	
Netherlands 3x3 Basketball	Netherlands	Men	
Poland 3x3 Basketball	Poland	Men	
ROC 3x3 Basketball	ROC	Men	
ROC 3x3 Basketball	ROC	Women	
Romania 3x3 Basketball	Romania	Women	
Serbia 3x3 Basketball	Serbia	Men	
United States 3x3 Basketball	United States of ...	Women	
Australia Archery	Australia	Men's Team	
Australia Archery	Australia	Mixed Team	

Microsoft Azure | databricks | Search data, notebooks, recents, and more... | CTRL + P | tokyo-olympic-db | Run all | Balachandar's Cluster

Tokyo Olympic Transformation | Python | 01:02 PM (1s) | 23

```
teams.printSchema()
```

```
root
|-- TeamName: string (nullable = true)
|-- Discipline: string (nullable = true)
|-- Country: string (nullable = true)
|-- Event: string (nullable = true)
```

Tokyo Olympic Transformation | Python | 01:02 PM (1s) | 24

```
# Find the top countries with the highest number of gold medals
top_gold_medal_countries = medals.orderBy("Gold", ascending=False).show()
```

(1) Spark Jobs

	TeamName	Discipline	Country	Event	Gold	Silver	Bronze	Total
1	Japan	27	14	17	58	5	5	78
2	Great Britain	22	21	22	65	4	4	91
3	ROC	20	28	23	71	3	3	102
4	Australia	17	7	22	46	6	6	82
5	Netherlands	10	12	14	36	9	9	67
6	France	10	12	11	33	10	10	65
7	Germany	10	11	16	37	8	8	64
8	Italy	10	10	20	40	7	7	67

Microsoft Azure | databricks | Search data, notebooks, recents, and more... | CTRL + P | tokyo-olympic-db | Run all | Balachandar s's Cluster

Tokyo Olympic Transformation Python | 01:02 PM (1s) | 24

```
# Find the top countries with the highest number of gold medals
top_gold_medal_countries = medals.orderBy("Gold", ascending=False).show()
```

(1) Spark Jobs

	Team	Country	Gold	Silver	Bronze	Total
3	Japan	27	14	17	58	5
4	Great Britain	22	21	22	65	4
5	ROC	20	28	23	71	3
6	Australia	17	7	22	46	6
7	Netherlands	10	12	14	36	9
8	France	10	12	11	33	10
9	Germany	10	11	16	37	8
10	Italy	10	10	20	40	7
11	Canada	7	6	11	24	11
12	Brazil	7	6	8	21	12
13	New Zealand	7	6	7	20	13
14	Cuba	7	3	5	15	18
15	Hungary	6	7	7	20	13
16	Republic of Korea	6	4	10	20	13
17	Poland	4	5	5	14	19
18	Czech Republic	4	4	3	11	23
19	Kenya	4	4	2	10	25
20	Norway	4	2	2	8	29

Microsoft Azure | databricks | Search data, notebooks, recents, and more... | CTRL + P | tokyo-olympic-db | Run all | Balachandar s's Cluster | Python

Tokyo Olympic Transformation Python | 01:02 PM (1s) | 25

```
# Find the top countries with the highest number of gold medals
top_gold_medal_countries = medals.orderBy("Gold", ascending=False).select("Team_country", "Gold").show()
```

(1) Spark Jobs

Team	Country	Gold
United States of ...	United States of America	39
People's Republic...	China	38
Japan	Japan	27
Great Britain	United Kingdom	22
ROC	Russia	20
Australia	Australia	17
Netherlands	Netherlands	10
France	France	10
Germany	Germany	10
Italy	Italy	10
Canada	Canada	7
Brazil	Brazil	7
New Zealand	New Zealand	7
Cuba	Cuba	7
Hungary	Hungary	6
Republic of Korea	South Korea	6
Poland	Poland	4

Microsoft Azure | databricks | Search data, notebooks, recents, and more... | CTRL + P | tokyo-olympic-db | Run all | Balachandar s's Cluster

Tokyo Olympic Transformation Python | 01:02 PM (1s) | 26

```
# Calculate the average number of entries by gender for each discipline
average_entries_by_gender = entriesgender.withColumn(
    "Avg_Female", entriesgender['Female'] / entriesgender['Total']
).withColumn(
    "Avg_Male", entriesgender['Male'] / entriesgender['Total']
)
average_entries_by_gender.show()
```

(1) Spark Jobs

Discipline	Female	Male	Total	Avg_Female	Avg_Male
3x3 Basketball	32	32	64	0.5	0.5
Archery	64	64	128	0.5	0.5
Artistic Gymnastics	98	98	196	0.5	0.5
Artistic Swimming	105	0	105	1.0	0.0
Athletics	969	1072	2041	0.4747672709456149	0.5252327290543851
Badminton	86	87	173	0.49710982658959535	0.5028901734104047
Baseball/Softball	90	144	234	0.38461538461538464	0.6153846153846154
Basketball	144	144	288	0.5	0.5
Beach Volleyball	48	48	96	0.5	0.5
Boxing	102	187	289	0.35294117647058826	0.6470588235294118

Do some transformation → Calculate the Average Number of Entries for each discipline.

Microsoft Azure | databricks | Search data, notebooks, recents, and more... | CTRL + P | tokyo-olympic-db | Run all | Balachandar s's Cluster

Tokyo Olympic Transformation Python

```
File Edit View Run Help Last edit was 2 hours ago New cell UI: ON
```

01:02 PM (x1) 27

```
entriesgender.show()
```

(1) Spark Jobs

Discipline	Female	Male	Total
3x3 Basketball	32	32	64
Archery	64	64	128
Artistic Gymnastics	98	98	196
Artistic Swimming	105	0	105
Athletics	969	1072	2041
Badminton	86	87	173
Baseball/Softball	90	144	234
Basketball	144	144	288
Beach Volleyball	48	48	96
Boxing	102	187	289
Canoe Slalom	41	41	82
Canoe Sprint	123	126	249
Cycling BMX Freestyle	10	9	19
Cycling BMX Racing	24	24	48
Cycling Mountain Bike	38	38	76
Curling Road	781	1311	2092

Microsoft Azure | databricks | Search data, notebooks, recents, and more... | CTRL + P | tokyo-olympic-db | Run all | Balachandar s's Cluster

Tokyo Olympic Transformation Python

```
File Edit View Run Help Last edit was 2 hours ago New cell UI: ON
```

Skipped 28

```
athletes.repartition(1).write.mode("overwrite").option("header","true").csv("/mnt/tokyoolympic/transformed-data/athletes")
coaches.repartition(1).write.mode("overwrite").option("header","true").csv("/mnt/tokyoolympic/transformed-data/coaches")
entriesgender.repartition(1).write.mode("overwrite").option("header","true").csv("/mnt/tokyoolympic/transformed-data/entriesgender")
medals.repartition(1).write.mode("overwrite").option("header","true").csv("/mnt/tokyoolympic/transformed-data/medals")
teams.repartition(1).write.mode("overwrite").option("header","true").csv("/mnt/tokyoolympic/transformed-data/teams")
```

Command skipped

01:03 PM (x1) 29

```
medals.show()
```

(1) Spark Jobs

Rank	Team_Country	Gold	Silver	Bronze	Total	Rank by Total
1	United States of America	39	41	33	113	1
2	People's Republic of China	38	32	18	88	2
3	Japan	27	14	17	58	5
4	Great Britain	22	21	22	65	4
5	ROC	20	28	23	71	3
6	Australia	17	7	22	46	6
7	Netherlands	10	12	14	36	9
8	France	10	12	11	33	10
9	Germany	10	11	16	37	8
10	Italy	10	10	20	40	7
11	Canada	7	6	11	24	11
12	Brazil	7	6	8	21	12
13	New Zealand	7	6	7	20	13
14	Cuba	7	3	5	15	18
15	Hungary	6	7	7	20	13
16	Republic of Korea	6	4	10	20	13
17	Poland	4	5	5	14	19

Microsoft Azure | databricks | Search data, notebooks, recents, and more... | CTRL + P | tokyo-olympic-db | Run all | Balachandar s's Cluster

Tokyo Olympic Transformation Python

```
File Edit View Run Help Last edit was 2 hours ago New cell UI: ON
```

01:03 PM (x1) 29

```
medals.show()
```

(1) Spark Jobs

Rank	Team_Country	Gold	Silver	Bronze	Total	Rank by Total
1	United States of America	39	41	33	113	1
2	People's Republic of China	38	32	18	88	2
3	Japan	27	14	17	58	5
4	Great Britain	22	21	22	65	4
5	ROC	20	28	23	71	3
6	Australia	17	7	22	46	6
7	Netherlands	10	12	14	36	9
8	France	10	12	11	33	10
9	Germany	10	11	16	37	8
10	Italy	10	10	20	40	7
11	Canada	7	6	11	24	11
12	Brazil	7	6	8	21	12
13	New Zealand	7	6	7	20	13
14	Cuba	7	3	5	15	18
15	Hungary	6	7	7	20	13
16	Republic of Korea	6	4	10	20	13
17	Poland	4	5	5	14	19

Microsoft Azure Upgrade Search resources, services, and docs (G+)

Home > Microsoft.Azure.SynapseAnalytics-20240404094923 | Overview

Deployment

Search Delete Cancel Redeploy Download Refresh

Overview

Your deployment is complete

Deployment name : Microsoft.Azure.SynapseAnalytics-20... Start time : 4/4/2024, 10:00:15 am
Subscription : Free Trial Correlation ID : 7aa29a7b-7adc-4de2-b431-0120837...
Resource group : tokyo-olympic

Deployment details Next steps

Go to resource group

Give feedback Tell us about your experience with deployment

Cost management Get notified to st and prevent une bill. Set up cost alerts!

Microsoft Defense Secure your app! Go to Microsoft I

Free Microsoft t Start learning to

Work with an ex

Microsoft Azure Upgrade Search resources, services, and docs (G+)

Home > Microsoft.Azure.SynapseAnalytics-20240404094923 | Overview > tokyo-olympic >

tokyo-olympic-sat Synapse workspace

Search New dedicated SQL pool New Apache Spark pool New Data Explorer pool (preview) Refresh Reset SQL admin password

Overview

Activity log Access control (IAM) Tags Diagnose and solve problems

Resource group (move) : tokyo-olympic Networking : Show firewall settings
Status : Succeeded Primary ADLS Gen2 acco... : https://tokyoolympicdata Primary ADLS Gen2 file s... : tokyo-olympic-data
Location : Southeast Asia SQL admin username : sqladminuser
Subscription (move) : Free Trial SQL Microsoft Entra admin : live.com#balachandar2014
Subscription ID : 5a2b62fd-8a3e-4a3a-98ad-a90d9b0b98f6 Dedicated SQL endpoint : tokyo-olympic-sat.sql.azure
Managed virtual network : No Serverless SQL endpoint : tokyo-olympic-sat-ondem Development endpoint : https://tokyo-olympic-sat.d
Managed identity object ... : b167b0c7-15d2-4321-9eb8-3b5da271f540
Workspace web URL : <https://web.azuresynapse.net/workspace=%2bsub...>
Tags (edit) : Add tags

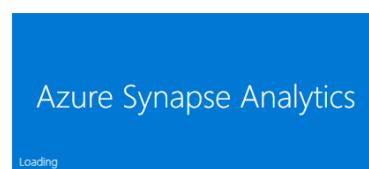
Essentials

Getting started

Open Synapse Studio Start building your fully-integrated analytics solution and unlock new insights.

Read documentation Learn how to be productive quickly. Explore concepts, tutorials, and samples.

Data Loading & Analytics by using Azure Synapse Analytics Studio.

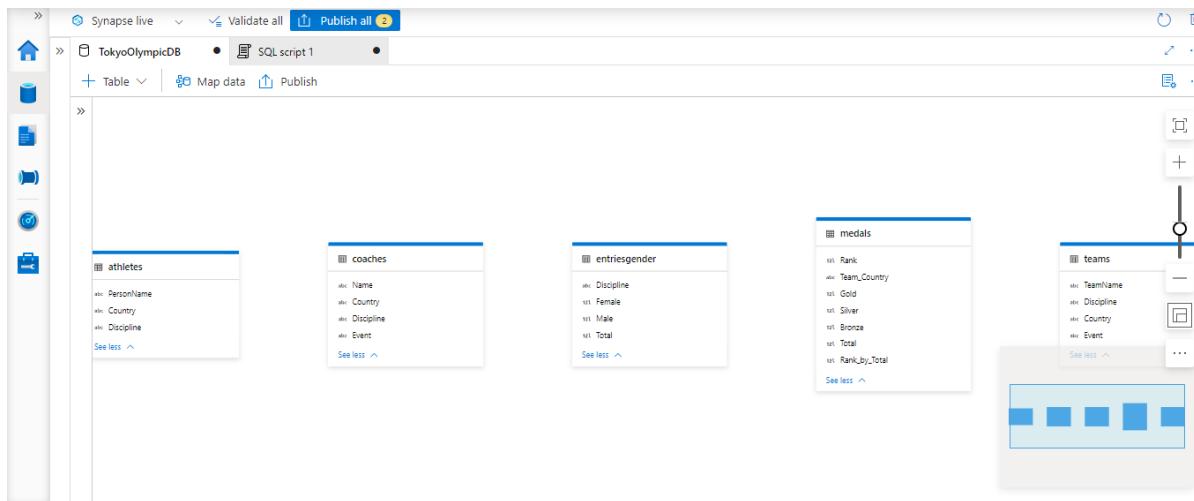


Setting Up and Using Azure Synapse Analytics

1. Creating a Synapse Analytics Workspace.
2. Within Workspace navigate to the "Data" section, choose "Lake Database" and create a Database "TokyoOlympicDB".
3. Creating Table from Data Lake from the Transformed Data folder within your ADLS Gen2 storage.

The screenshot shows the Microsoft Azure Synapse Analytics workspace interface for the 'tokyo-olympic-sat' workspace. The left sidebar includes links for Home, Data, Develop, Integrate, Monitor, and Manage. The main area features a large circular diagram illustrating data flow and analysis. Three cards are displayed: 'Ingest' (Perform a one-time or scheduled data load.), 'Explore and analyze' (Learn how to get insights from your data.), and 'Visualize' (Build interactive reports with Power BI capabilities.). Below these cards, there's a 'Discover more' section with links to 'Knowledge center' and 'Browse partners'.

The screenshot shows the Microsoft Azure Synapse Analytics Data section for the 'TokyoOlympicDB' lake database. The left sidebar remains the same. The main area displays the 'Tables' section, which lists five tables: 'athletes', 'coaches', 'entriesgender', 'medals', and 'teams'. Each table has a detailed view icon next to it, showing columns such as 'PersonName (string)', 'Country (string)', and 'Discipline (string)' for the 'athletes' table.



The screenshot shows the Microsoft Power BI Data Flow interface with a workspace named "Lake database". The workspace contains a table named "TokyoOlympicDB" which includes the following tables: athletes, coaches, entriesgender, medals, and teams. The "Entriesgender" table is expanded to show its columns: Discipline, Female, Male, and Total. A query script is displayed in the center:

```

2 SELECT Country, COUNT(*) AS TotalAthletes
3 FROM athletes
4 GROUP BY Country
5 ORDER BY TotalAthletes DESC;
6
7 --Calculate the total medals won by each country:
8 SELECT
9 Team_Country,
10 SUM(GOLD) Total_Gold,
11 SUM(Silver) Total_Silver,
12 SUM(Bronze) Total_Bronze
13 FROM medals
14 GROUP BY Team_Country
15 ORDER BY Total_Gold DESC;
16
17 --Calculate the average number of entries by gender for each discipline:
18 SELECT Discipline,
19 Avg(Female) Avg_Female,
20 Avg(Male) Avg_Male
21 FROM entriesgender
22 GROUP BY Discipline;

```

The results pane at the bottom shows a message: "00:00:07 Query executed successfully."

The screenshot shows the Microsoft Power BI Data Flow interface with a workspace named "Lake database". The workspace contains a table named "TokyoOlympicDB" which includes the following tables: athletes, coaches, entriesgender, medals, and teams. The "Coaches" table is expanded to show its columns: Name, Country, and Discipline. A query script is displayed in the center:

```

2 SELECT Country, COUNT(*) AS TotalAthletes
3 FROM athletes
4 GROUP BY Country
5 ORDER BY TotalAthletes DESC;
6
7 --Calculate the total medals won by each country:
8 SELECT
9 Team_Country,
10 SUM(GOLD) Total_Gold,
11 SUM(Silver) Total_Silver,
12 SUM(Bronze) Total_Bronze
13 FROM medals
14 GROUP BY Team_Country
15 ORDER BY Total_Gold DESC;
16
17 --Calculate the average number of entries by gender for each discipline:
18 SELECT Discipline,
19 Avg(Female) Avg_Female,
20 Avg(Male) Avg_Male
21 FROM entriesgender
22 GROUP BY Discipline;

```

The results pane at the bottom shows a message: "00:00:07 Query executed successfully."

```
2 SELECT Country, COUNT(*) AS TotalAthletes
3 FROM athletes
4 GROUP BY Country
5 ORDER BY TotalAthletes DESC;
6
7 --Calculate the total medals won by each country:
8 SELECT
9 Team_Country,
10 SUM(GOLD) Total_Gold,
11 SUM(Silver) Total_Silver,
12 SUM(Bronze) Total_Bronze
13 FROM medals
14 GROUP BY Team_Country
15 ORDER BY Total_Gold DESC;
16
17 --Calculate the average number of entries by gender for each discipline:
18 SELECT Discipline,
19 Avg(Female) Avg_Female,
20 Avg(Male) Avg_Male
21 FROM entriesgender
22 GROUP BY Discipline;
```

Results Messages

00:00:07 Query executed successfully.

Calculate the total medals won by each country,

Calculate the average number of entries by gender for discipline by Using sql queries.

```
6 --Calculate the total medals won by each country:
7 SELECT
8 Team_Country,
9 SUM(GOLD) Total_Gold,
10 SUM(Silver) Total_Silver,
11 SUM(Bronze) Total_Bronze
12 FROM medals
13 GROUP BY Team_Country
14 ORDER BY Total_Gold DESC;
```

Results Messages

View Table Chart Export results ▾

Search

Team_Country	Total_Gold	Total_Silver	Total_Bronze
United States of America	39	41	33
People's Republic of China	38	32	18

Properties

General Related (0)

Name * SQL script 1

Description

Type sql script

Size 94 bytes

Results settings per query ▾

First 5000 rows (default)

All rows

00:00:03 Query executed successfully.

Synapse live Validate all Publish all

TokyoOlympicDB SQL script 1

Run Undo Publish Query plan Connect to Built-in Use database TokyoOlympicDB

```
1 --Count the number of Athletes from each Country:
2 SELECT Country, COUNT(*) AS TotalAthletes
3 FROM athletes
4 GROUP BY Country
5 ORDER BY TotalAthletes DESC;
6
```

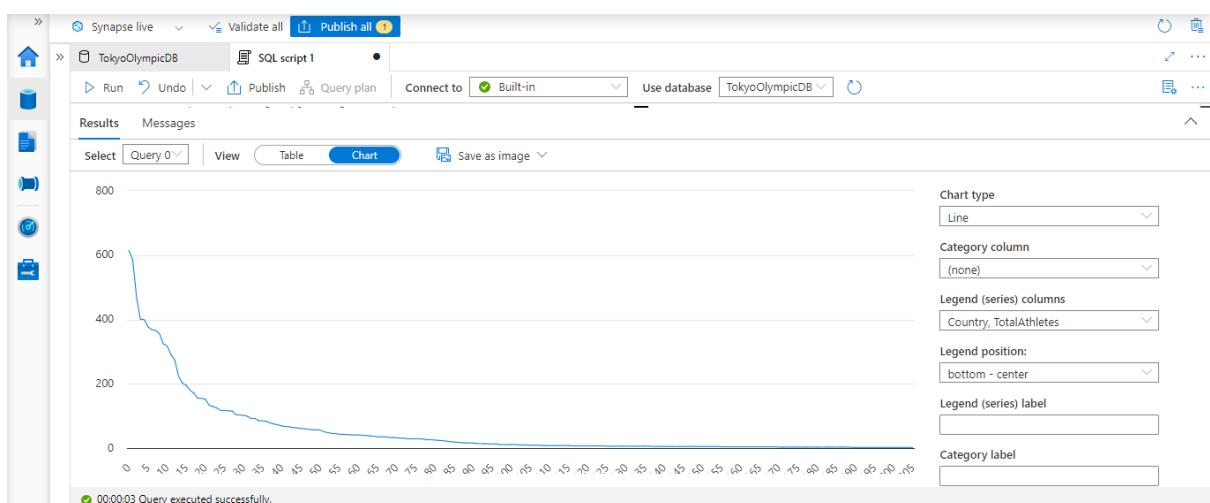
Results Messages

Select Query 0 View Table Chart Export results

Search

Country	TotalAthletes
United States of America	615
Japan	586
Australia	470
People's Republic of China	401
Germany	388

00:00:03 Query executed successfully.



Synapse live Validate all Publish all

TokyoOlympicDB SQL script 1

Run Undo Publish Query plan Connect to Built-in Use database TokyoOlympicDB

```
8 SELECT
9 Team_Country,
10 SUM(GOLD) Total_Gold,
11 SUM(Silver) Total_Silver,
12 SUM(Bronze) Total_Bronze
13 FROM medals
14 GROUP BY Team_Country
15 ORDER BY Total_Gold DESC;
16
```

Results Messages

View Table Chart Export results

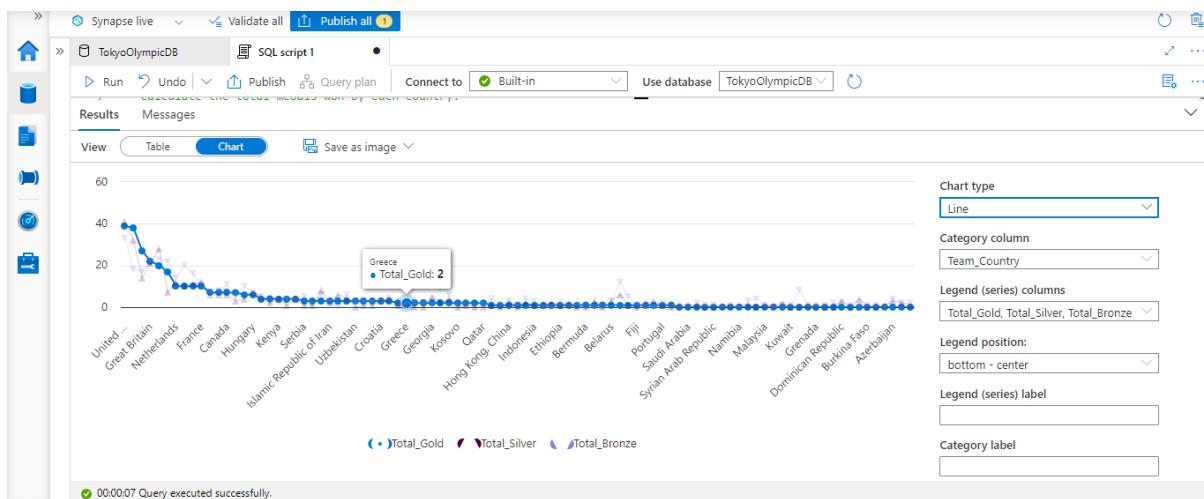
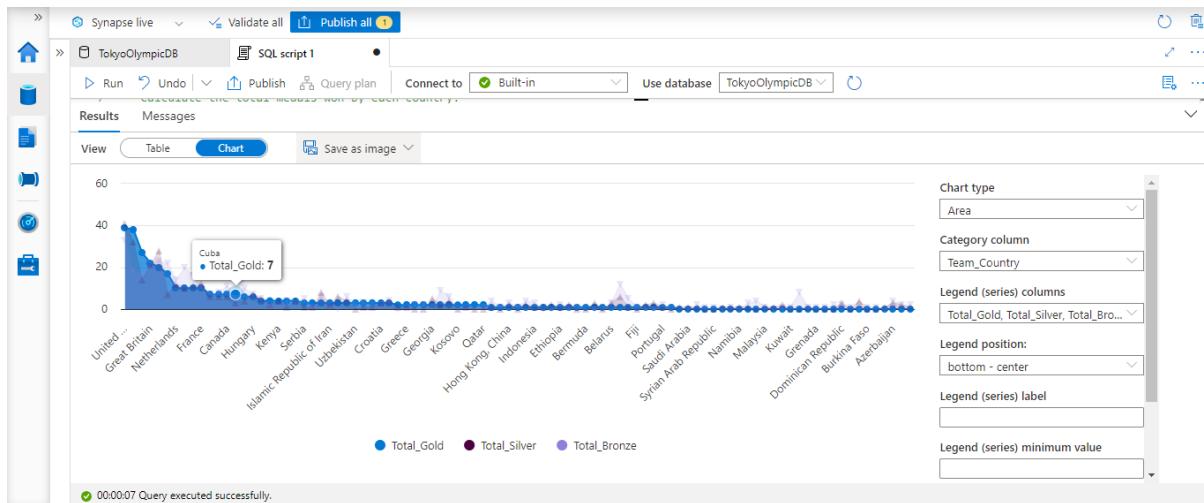
Search

Team_Country	Total_Gold	Total_Silver	Total_Bronze
United States of America	39	41	33
People's Republic of China	38	32	18

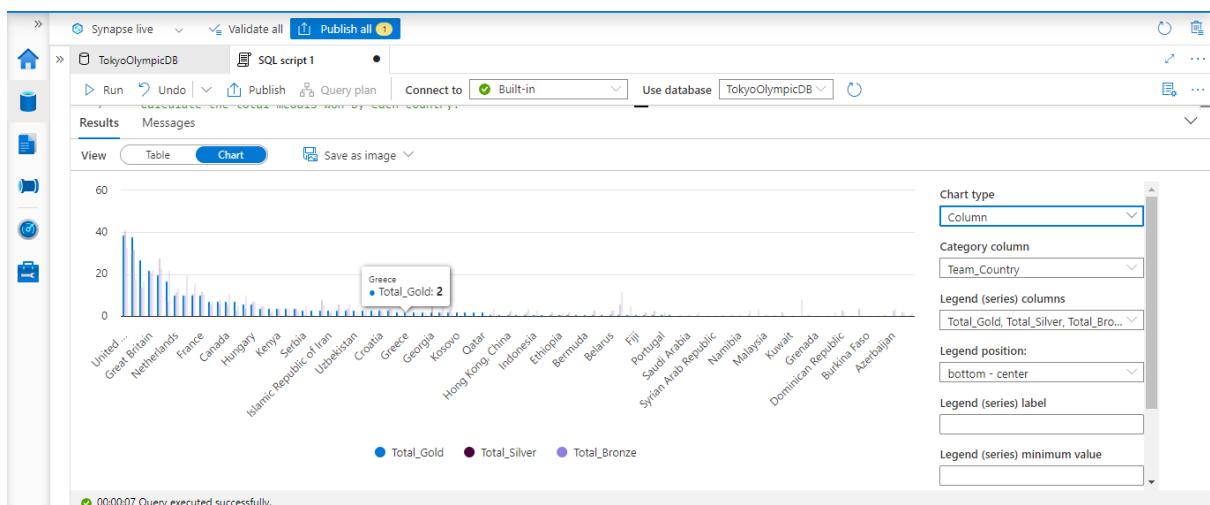
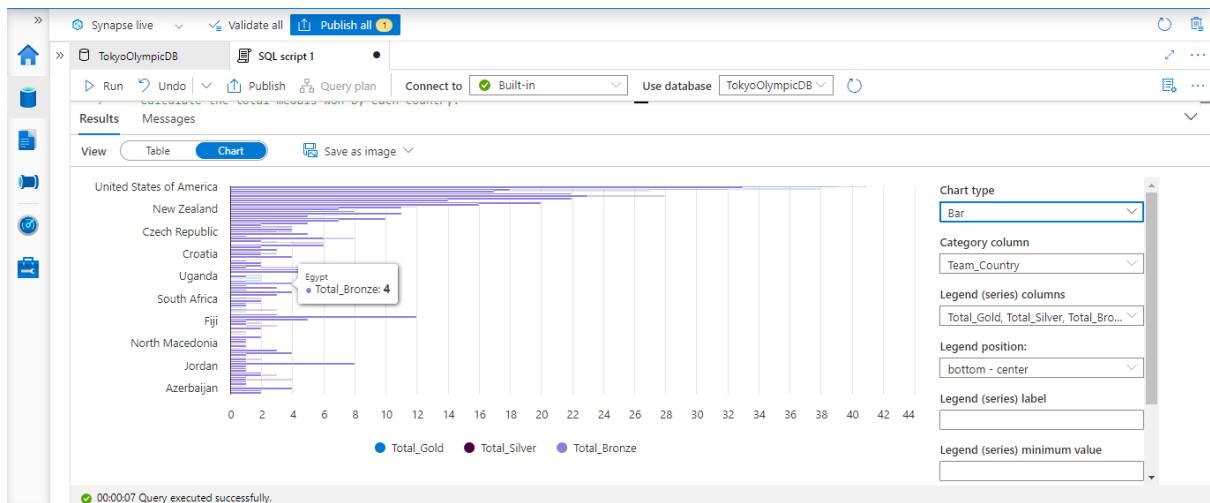
00:00:07 Query executed successfully.

Performing Data Analysis on the Data

Create SQL script to Perform Exploratory data analysis using SQL. You can also use PowerBI to generate your analysis reports.



Visualize the data according to the SQL query in Azure Synapse Analytics.



Synapse live | Validate all | Publish all

TokyoOlympicDB | SQL script 1 | Run | Undo | Publish | Query plan | Connect to: Built-in | Use database: TokyoOlympicDB | Refresh

```

15 ORDER BY Total_Gold DESC;
16
17 --Calculate the average number of entries by gender for each discipline:
18 SELECT Discipline,
19     Avg(Female) Avg_Female,
20     Avg(Male) Avg_Male
21 FROM entriesgender
22 GROUP BY Discipline;

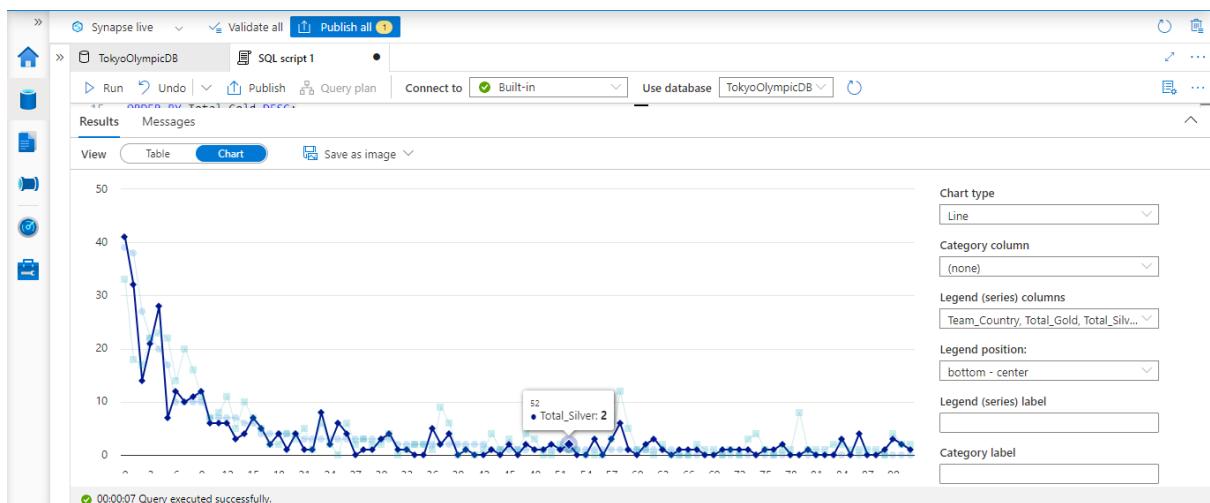
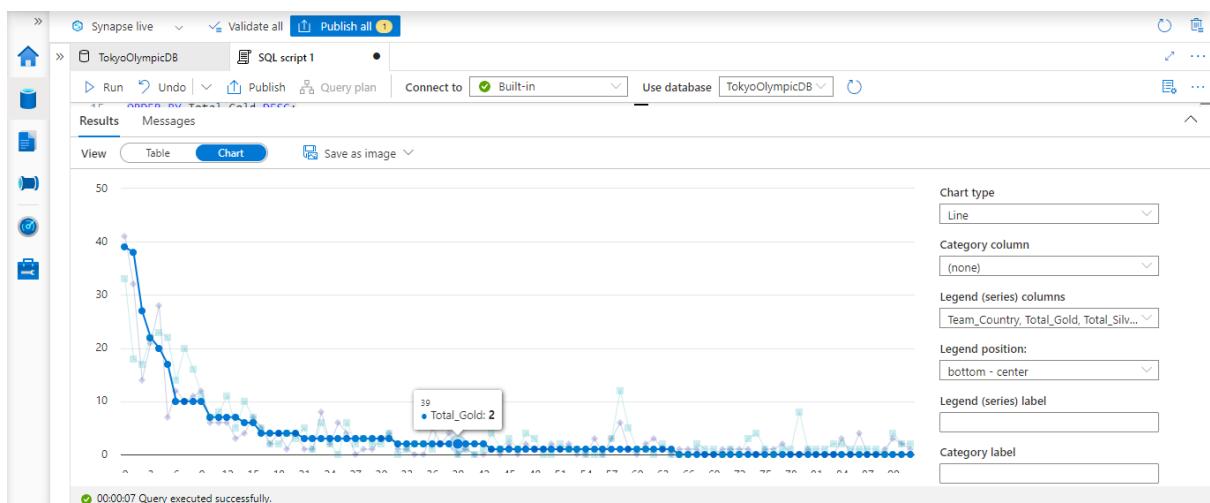
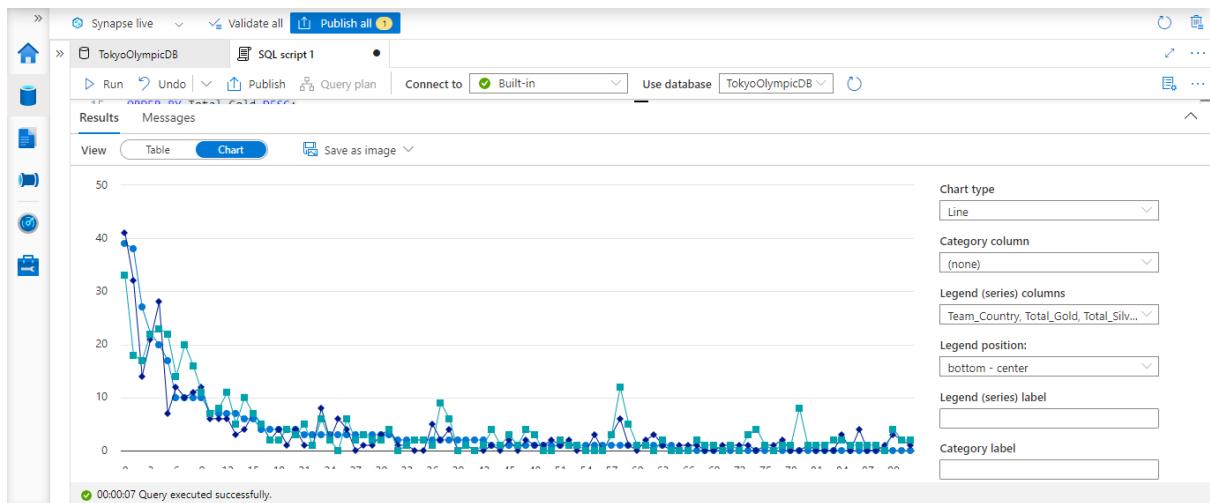
```

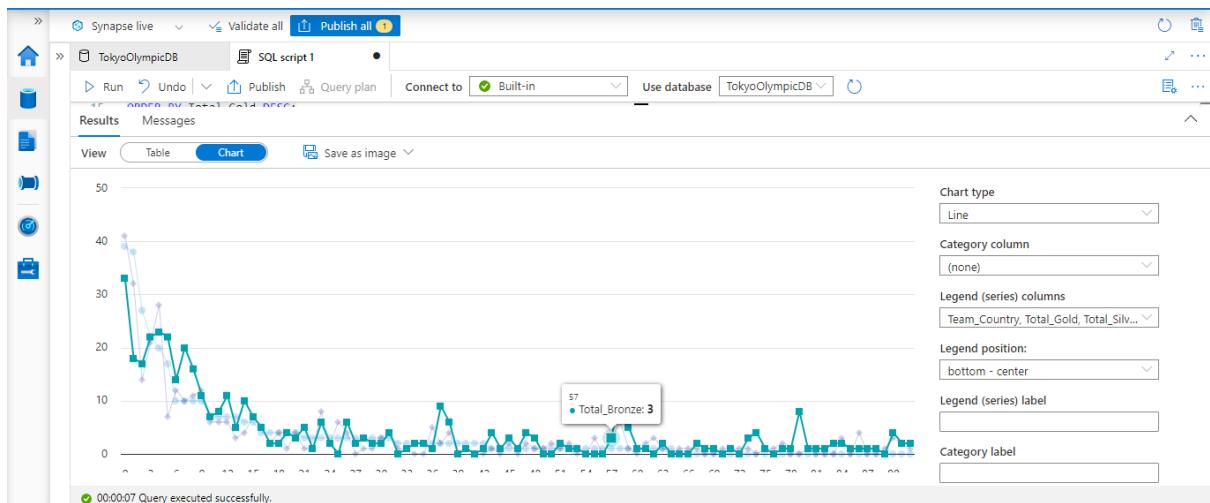
Results | Messages

View: Table | Chart | Export results

Team_Country	Total_Gold	Total_Silver	Total_Bronze
United States of America	39	41	33
People's Republic of China	38	32	18
Japan	27	14	17
...

00:00:07 Query executed successfully.





Untitled - Power BI Desktop

File Home Insert Modeling View Optimize Help

Cut Copy Paste Format painter Get data from multiple sources.

Clipboard

Import data from Excel Import data from SQL Server Paste data into a blank table Use sample data

Add data to your report

Once loaded, your data will appear in the Data pane.

Get data from another source →

Visualizations

Build visual

Filters

Values

Add data fields here

Drill through

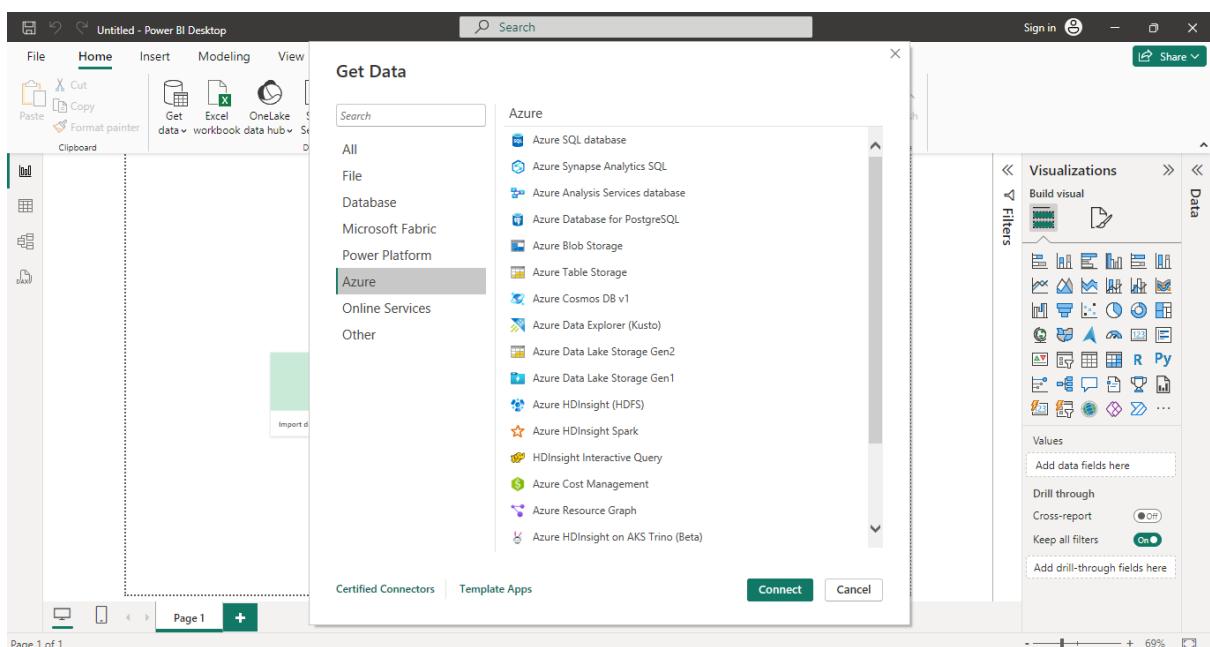
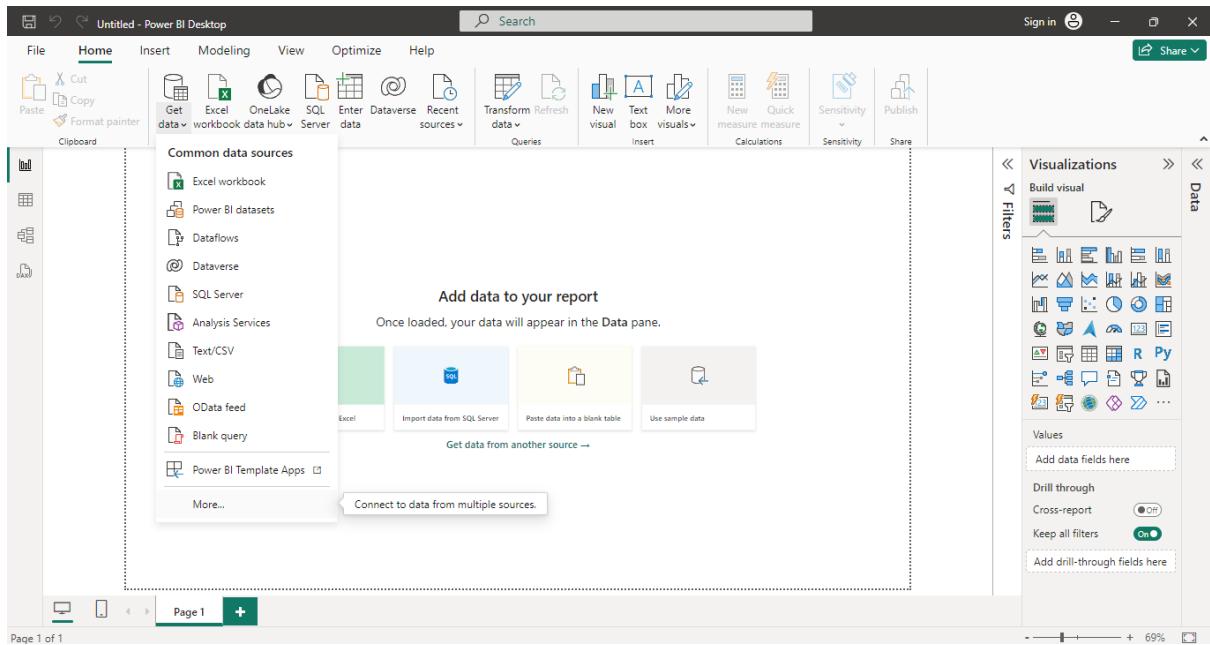
Cross-report

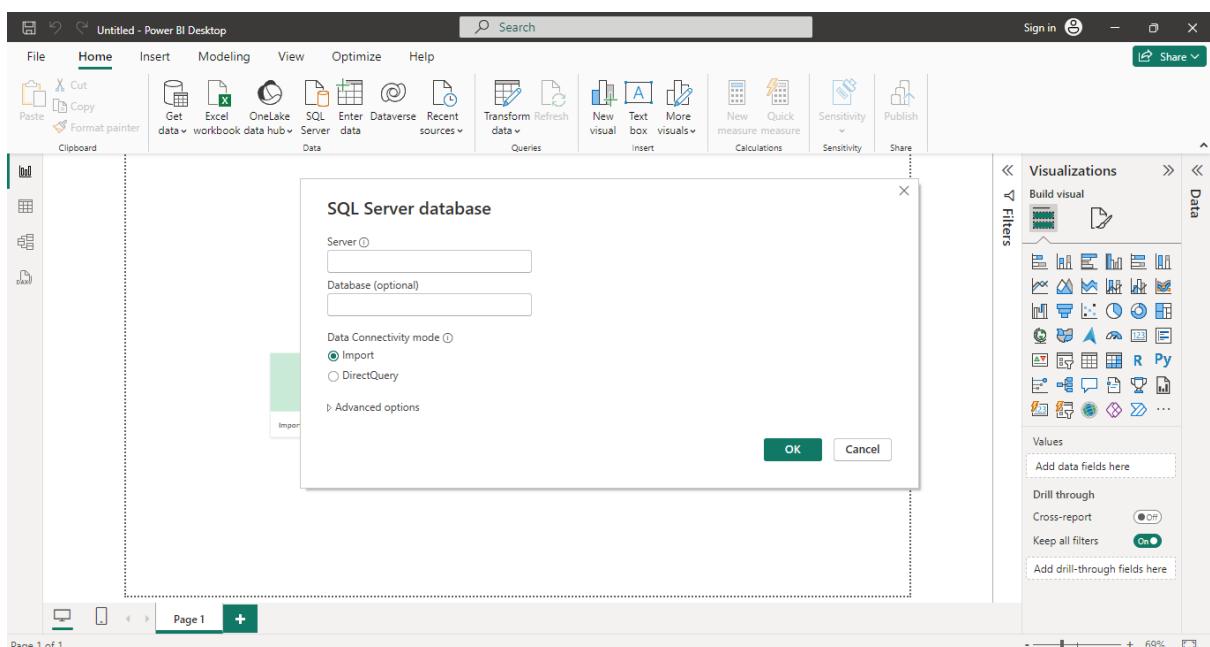
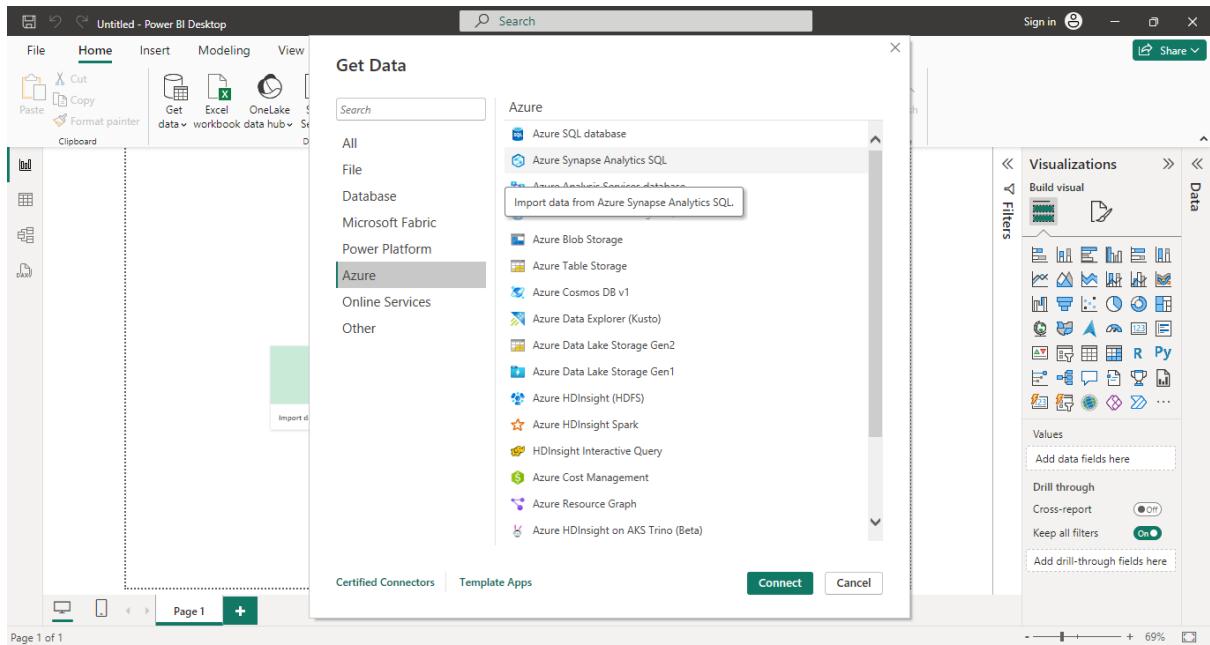
Keep all filters

Add drill-through fields here

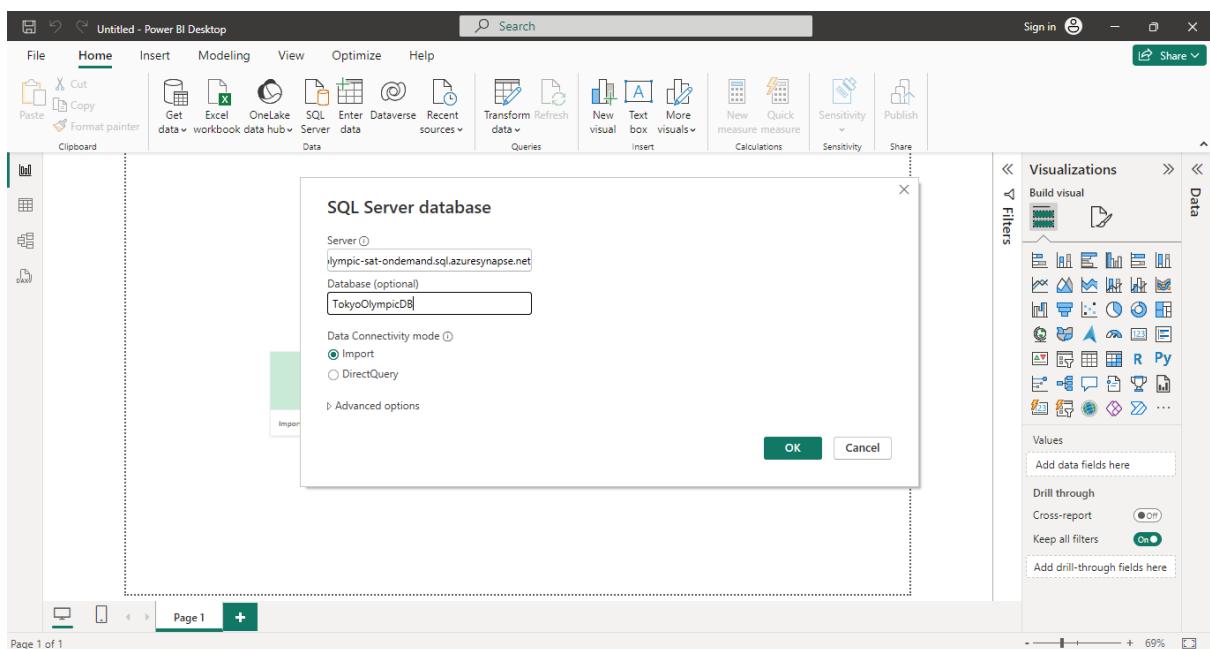
Page 1 +

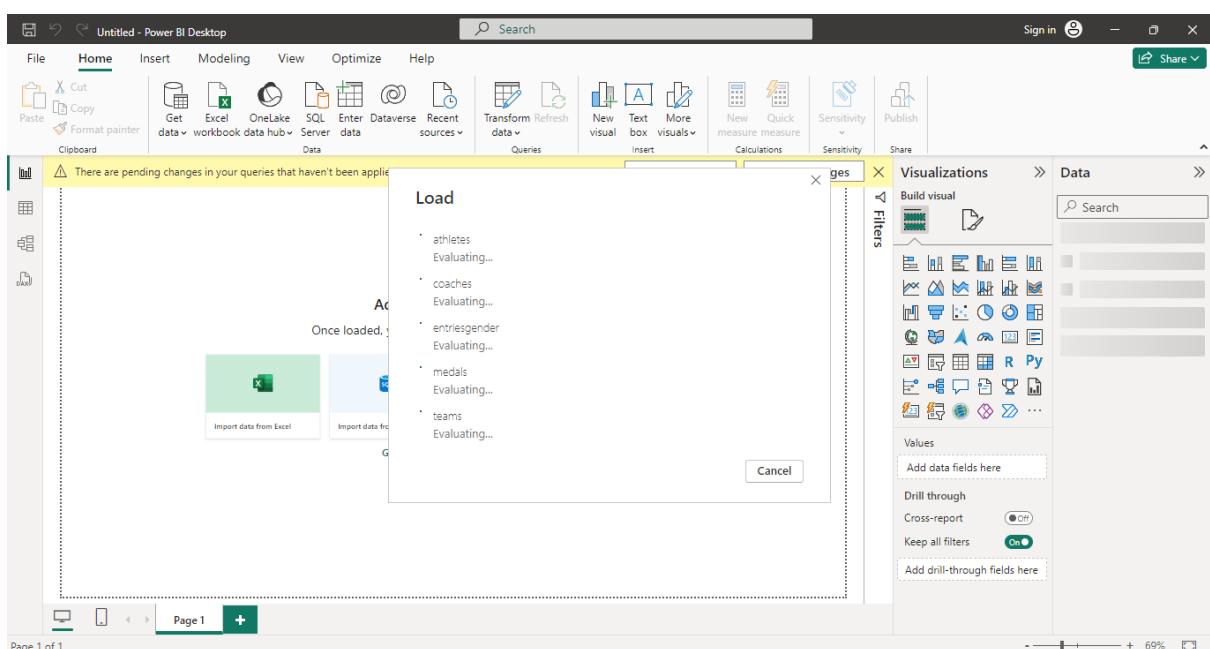
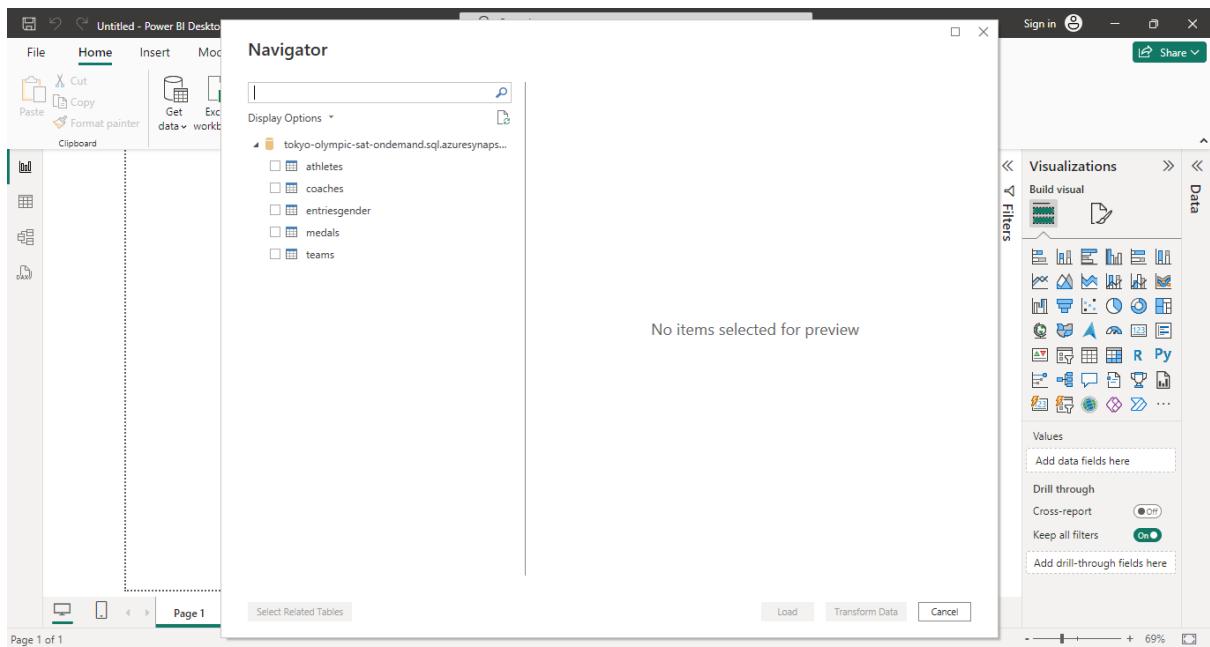
Page 1 of 1

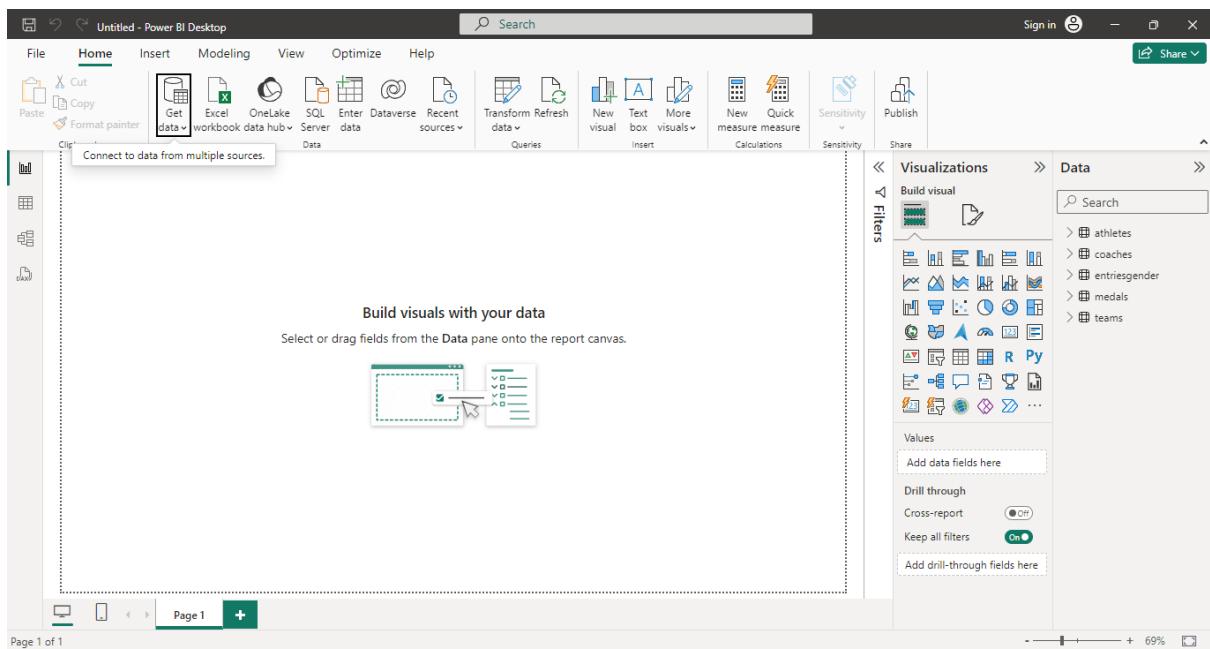
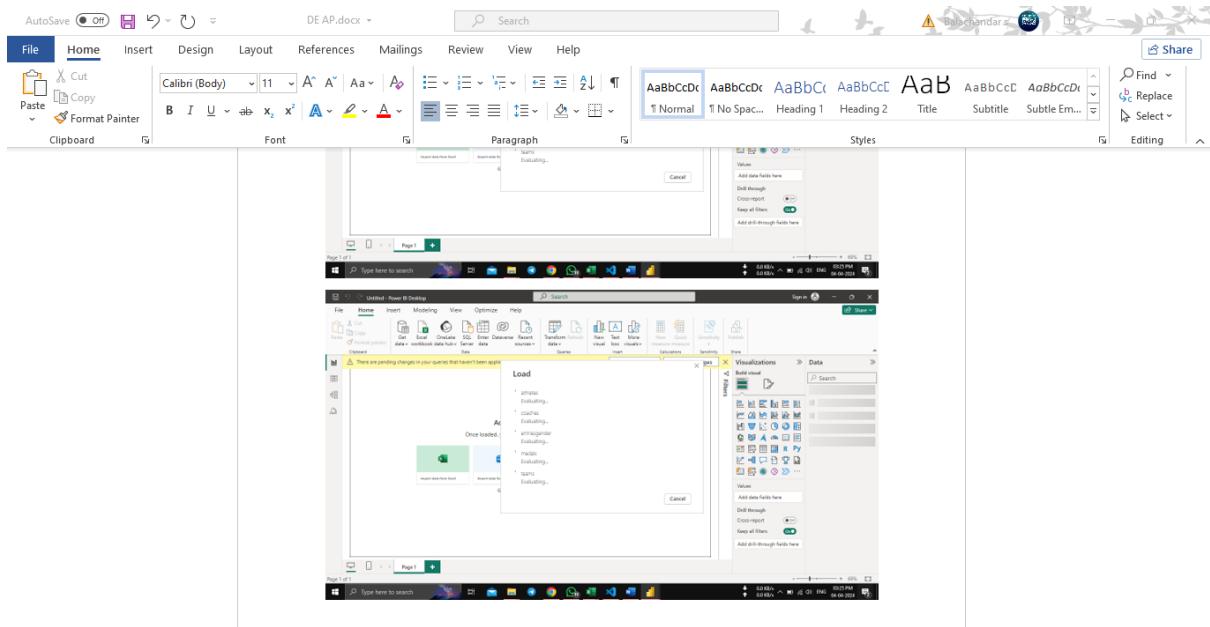




```
1 --Count the number of Athletes from each Country:
2 SELECT Country, COUNT(*) AS TotalAthletes
3 FROM athletes
4 GROUP BY Country
5 ORDER BY TotalAthletes DESC;
6
7 --Calculate the total medals won by each country:
8 SELECT
9 Team_Country,
10 SUM(GOLD) Total_Gold,
11 SUM(Silver) Total_Silver,
12 SUM(Bronze) Total_Bronze
13 FROM medals
14 GROUP BY Team_Country
15 ORDER BY Total_Gold DESC;
16
17 --Calculate the average number of entries by gender for each discipline:
18 SELECT Discipline,
19 Avg(Female) Avg_Female,
20 Avg(Male) Avg_Male
21 FROM entriesgender
```







Untitled - Power BI Desktop

File Home Help Table tools

Name athletes

Structure

Mark as date table Calendars Manage relationships Relationships New measure Quick New measure column New table Calculations

PersonName Country Discipline

PersonName	Country	Discipline
ABDIRAHMAN Abdi	United States of America	Athletics
ALLEN Devon	United States of America	Athletics
ALLMAN Valerie	United States of America	Athletics
ANDERSEN Brooke	United States of America	Athletics
ANDERSON Shae	United States of America	Athletics
AQUILLA Adelaide	United States of America	Athletics
ASHLEY Whitney	United States of America	Athletics
BAILEY JR Aldrich	United States of America	Athletics
BAKER Ronnie	United States of America	Athletics
BASTIEN Steven	United States of America	Athletics
BATTLE Anava	United States of America	Athletics
BEDNAREK Kenneth	United States of America	Athletics
BENARD Chris	United States of America	Athletics
BENJAMIN Rai	United States of America	Athletics
BERRY Gwen	United States of America	Athletics
BOUGARD Erica	United States of America	Athletics
BROMELL Trayvon	United States of America	Athletics
BURKS Quanesha	United States of America	Athletics
BUTTS-TOWNSEND Tyntia	United States of America	Athletics
CARD Kelsey	United States of America	Athletics
CARTER Chris	United States of America	Athletics
CENTROWITZ Matthew	United States of America	Athletics
CHERRY Michael	United States of America	Athletics

Table: athletes (11,085 rows)

Untitled - Power BI Desktop

File Home Help Table tools

Name medals

Structure

Mark as date table Calendars Manage relationships Relationships New measure Quick New measure column New table Calculations

Rank Team_Country Gold Silver Bronze Total Rank_by_Total

Rank	Team_Country	Gold	Silver	Bronze	Total	Rank_by_Total
1	United States of America	39	41	33	113	1
2	People's Republic of China	38	32	18	88	2
3	Japan	27	14	17	58	5
4	Great Britain	22	21	22	65	4
5	ROC	20	28	23	71	3
6	Australia	17	7	22	46	6
7	Netherlands	10	12	14	36	9
8	France	10	12	11	33	10
9	Germany	10	11	16	37	8
10	Italy	10	10	20	40	7
11	Canada	7	6	11	24	11
12	Brazil	7	6	8	21	12
13	New Zealand	7	6	7	20	13
14	Cuba	7	3	5	15	18
15	Hungary	6	7	7	20	13
16	Republic of Korea	6	4	10	20	13
17	Poland	4	5	5	14	19
18	Czech Republic	4	4	3	11	23
19	Kenya	4	4	2	10	25
20	Norway	4	2	2	8	29
21	Jamaica	4	1	4	9	26
22	Spain	3	8	6	17	17
23	Sweden	3	6	0	9	26

Table: medals (93 rows)

Name medals
Storage mode Import
Data refreshed 4/4/2024, 3:26:23 PM

Untitled - Power BI Desktop

File Home Help Table tools

Name teams

Mark as date table Calendars Manage relationships Relationships New Quick measure New measure column New table Calculations

TeamName	Discipline	Country	Event
Belgium	3x3 Basketball	Belgium	Men
China	3x3 Basketball	People's Republic of China	Men
Japan	3x3 Basketball	Japan	Men
Latvia	3x3 Basketball	Latvia	Men
Netherlands	3x3 Basketball	Netherlands	Men
Poland	3x3 Basketball	Poland	Men
ROC	3x3 Basketball	ROC	Men
Serbia	3x3 Basketball	Serbia	Men
Argentina	Basketball	Argentina	Men
Australia	Basketball	Australia	Men
Czech Republic	Basketball	Czech Republic	Men
France	Basketball	France	Men
Germany	Basketball	Germany	Men
Islamic Rep. of Iran	Basketball	Islamic Republic of Iran	Men
Italy	Basketball	Italy	Men
Japan	Basketball	Japan	Men
Nigeria	Basketball	Nigeria	Men
Slovenia	Basketball	Slovenia	Men
Spain	Basketball	Spain	Men
United States	Basketball	United States of America	Men
Azaad/Capogrosso	Beach Volleyball	Argentina	Men
McHugh/Schumann	Beach Volleyball	Australia	Men
Alison/Alvaro Filho	Beach Volleyball	Brazil	Men

Table: teams (743 rows)

Name teams

Storage mode Import Data refreshed 4/4/2024, 3:26:27 PM

Untitled - Power BI Desktop

File Home Help Table tools

Name teams

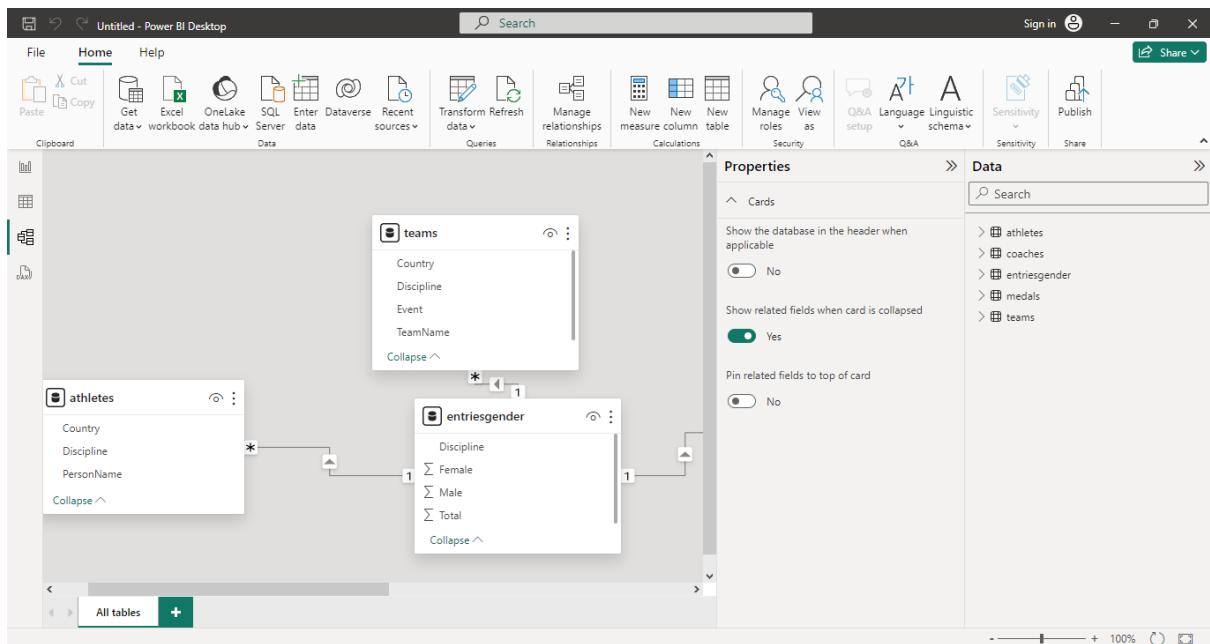
Mark as date table Calendars Manage relationships Relationships New Quick measure New measure column New table Calculations

TeamName	Discipline	Country	Event
Belgium	3x3 Basketball	Belgium	Men
China	3x3 Basketball	People's Republic of China	Men
Japan	3x3 Basketball	Japan	Men
Latvia	3x3 Basketball	Latvia	Men
Netherlands	3x3 Basketball	Netherlands	Men
Poland	3x3 Basketball	Poland	Men
ROC	3x3 Basketball	ROC	Men
Serbia	3x3 Basketball	Serbia	Men
Argentina	Basketball	Argentina	Men
Australia	Basketball	Australia	Men
Czech Republic	Basketball	Czech Republic	Men
France	Basketball	France	Men
Germany	Basketball	Germany	Men
Islamic Rep. of Iran	Basketball	Islamic Republic of Iran	Men
Italy	Basketball	Italy	Men
Japan	Basketball	Japan	Men
Nigeria	Basketball	Nigeria	Men
Slovenia	Basketball	Slovenia	Men
Spain	Basketball	Spain	Men
United States	Basketball	United States of America	Men
Azaad/Capogrosso	Beach Volleyball	Argentina	Men
McHugh/Schumann	Beach Volleyball	Australia	Men
Alison/Alvaro Filho	Beach Volleyball	Brazil	Men

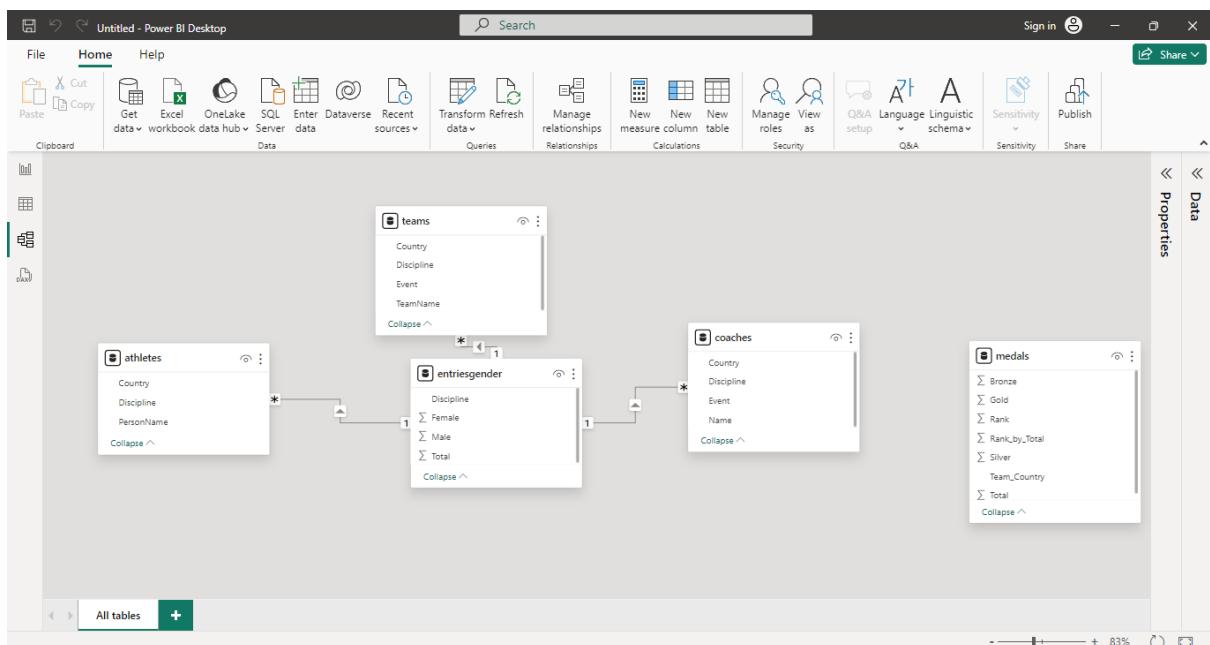
Table: teams (743 rows)

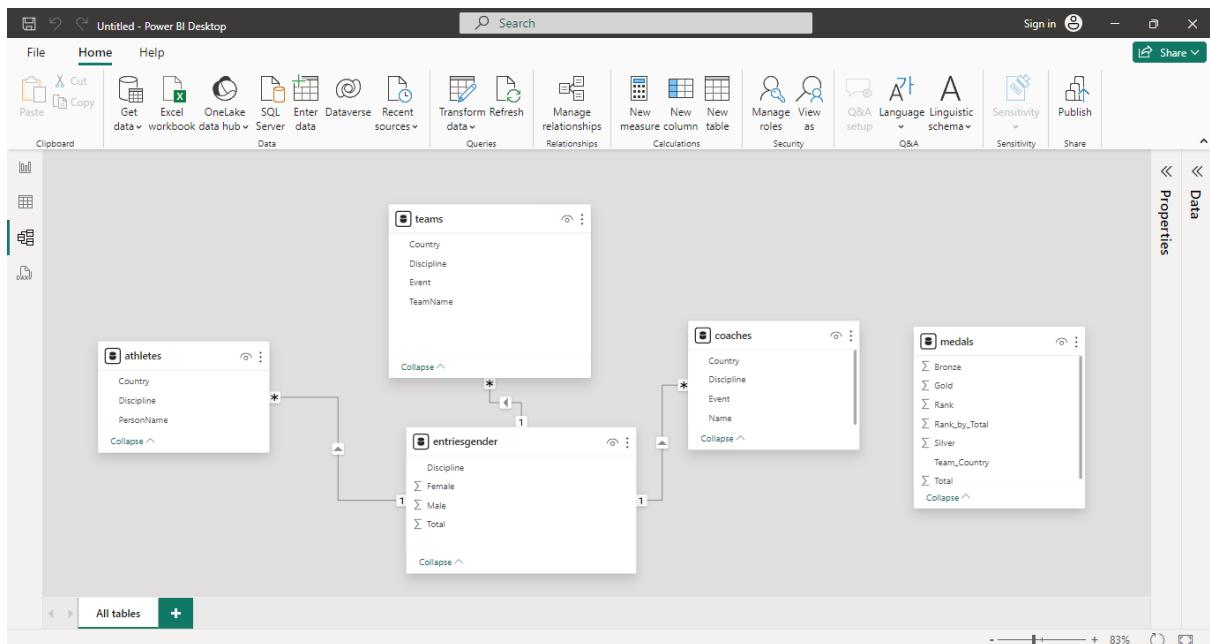
Name teams

Storage mode Import Data refreshed 4/4/2024, 3:26:27 PM



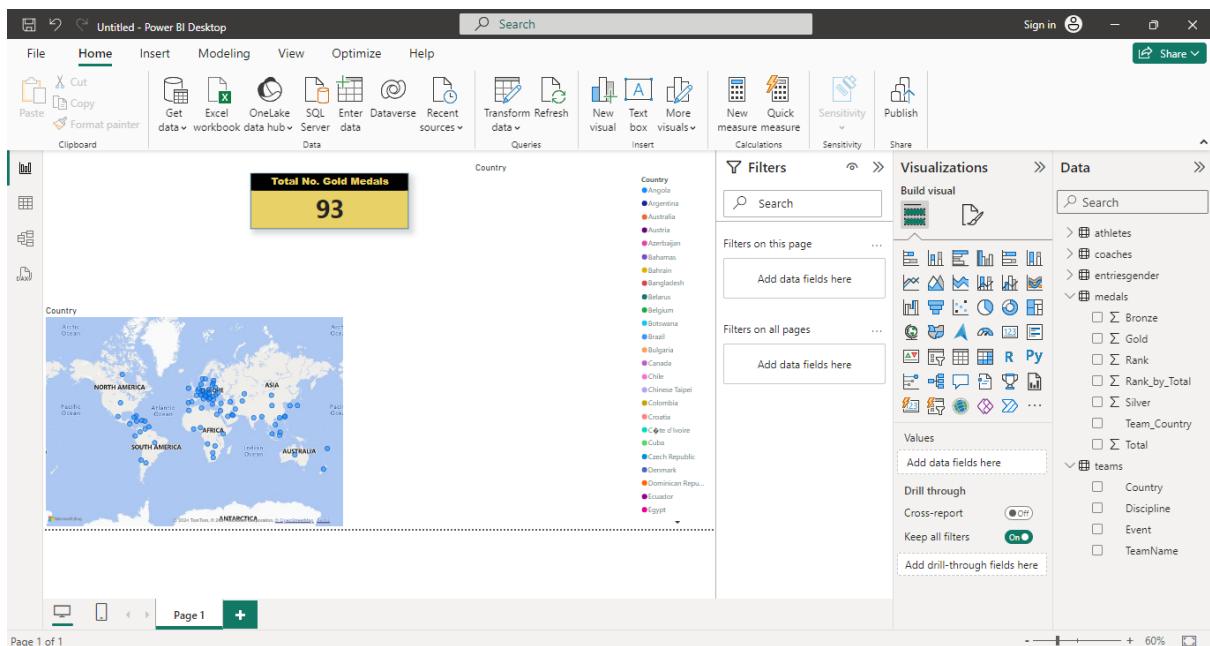
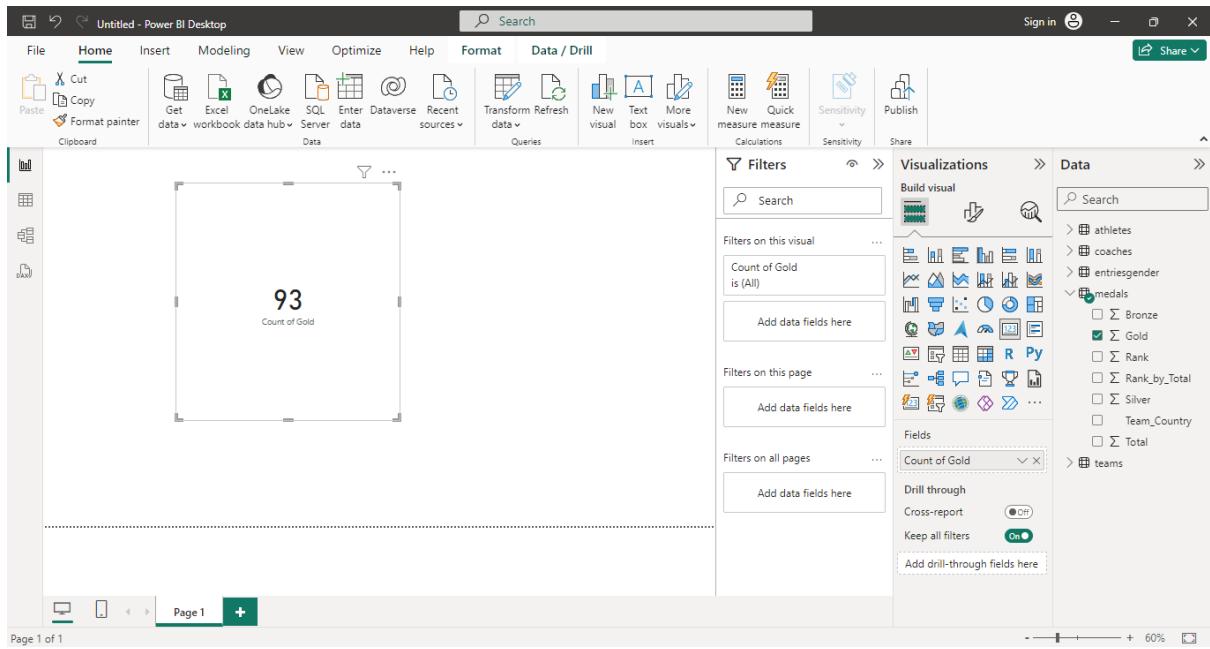
Data Modelling...

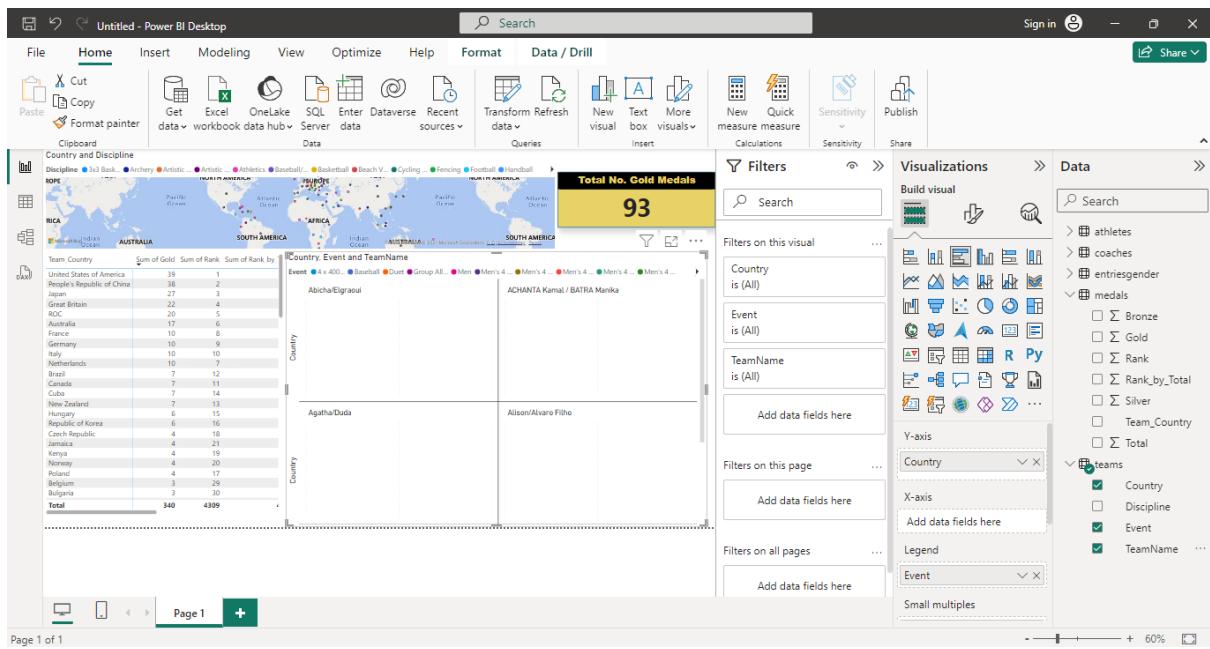




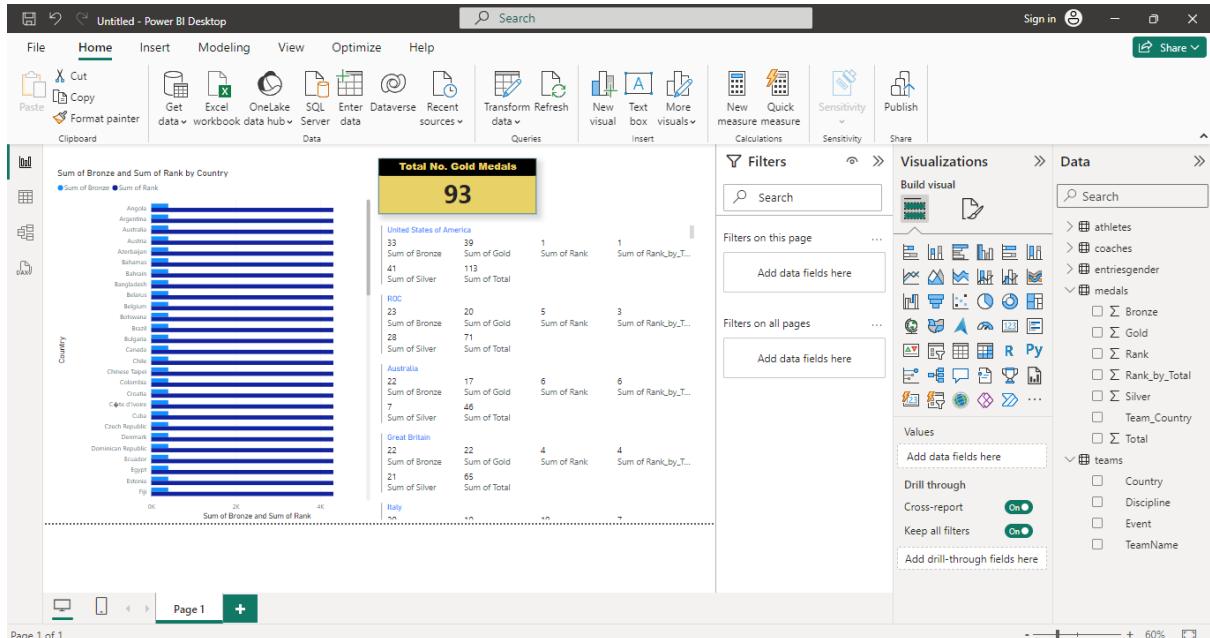
Untitled - Power BI Desktop

The screenshot shows the Power BI Report view. The main area is a blank report canvas with the placeholder text "Build visuals with your data". To the right is the "Data" pane, which lists the five tables used in the data model: 'athletes', 'coaches', 'entriesgender', 'medals', and 'teams'. The "Visualizations" pane is also visible, showing various chart and table icons.





Dashboard Creation by using Microsoft Power BI.



End of Project....

Delete the Project resource and avoid the Over Billing.

Delete Resources

The selected resources along with their related resources and contents will be permanently deleted. If you are unsure of the selected resource dependencies, navigate to the individual resource page to perform the delete operation. More details of the resource dependencies are available in the manage experience.

Resources to be deleted (1)

Name	Resource type
tokyo-olympic-db	Azure Databricks Service

Enter "delete" to confirm deletion *

Delete **Cancel**

Delete a resource group

The following resource group and all its dependent resources will be permanently deleted.

Resource group to be deleted

tokyo-olympic [\[View\]](#)

Dependent resources to be deleted (1)

All dependent resources, including hidden types, are shown

Name	Resource type
tokyo-olympic-db	Azure Databricks Service

Enter resource group name to confirm deletion *

Delete **Cancel**

Thank you!