

Olympic Data Analysis Using Azure and Power BI

Dhanya R Sangolli

A20541093

Introduction

This report delves into a thorough analysis of Olympic Games data, with a specific focus on assessing athlete performances and the distribution of medals. In this endeavor, we utilize the robust capabilities of Microsoft Azure for data storage, processing, and orchestration, coupled with Power BI for visualization. Our objective is to glean actionable insights from the wealth of historical Olympic records available.

Microsoft Azure: Serving as the cornerstone of our data infrastructure, Azure offers unparalleled scalability and reliability. We leverage Azure Blob Storage as a resilient solution for housing both raw and transformed data, ensuring its accessibility and security. Azure Data Factory plays a pivotal role in streamlining the orchestration of data pipelines, enabling the seamless movement and processing of datasets. Furthermore, Apache Spark on Azure Databricks empowers us to execute large-scale data transformations efficiently, thereby facilitating the extraction of valuable insights from the intricate Olympic datasets.

Medallion Architecture: Our approach adopts the Medallion architecture, comprising three distinct layers: Bronze, Silver, and Gold. Each layer plays a crucial role in ensuring the integrity and usability of the data.

{Bronze Layer}: Serving as the foundational layer, the Bronze layer houses the raw Olympic data, directly ingested into Azure Blob Storage. This layer forms the bedrock upon which subsequent processing relies.

Silver Layer: Positioned as an intermediate stage, the Silver layer is dedicated to cleaning and transforming the data. Here, we meticulously remove duplicate records, fill in missing values, and standardize text fields to ensure consistency. Additionally, we standardize date formats for uniformity across datasets and organize categorical data for enhanced analytical capabilities.

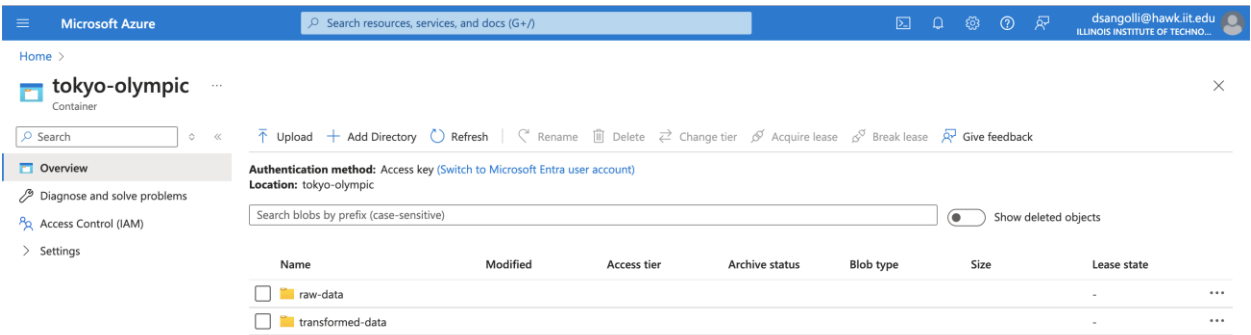
Gold Layer: At the apex of the architecture lies the Gold layer, which comprises structured data in the form of dimension and fact tables. Optimized for in-depth analysis and reporting, dimension tables capture static attributes such as athlete details, coach information, and team profiles, facilitating efficient data slicing and dicing. Conversely, fact tables aggregate measures such as medal counts by country and gender participation across sports, enabling insightful analysis and trend identification.

By harnessing the capabilities of Microsoft Azure and Power BI, this project aims to unlock invaluable insights from Olympic data. The structured Medallion architecture ensures scalability and clarity in data processing and analysis, thereby empowering informed decision-making in the dynamic realm of Olympic sports.

Tools and Technologies

Microsoft Azure:

Positioned as the backbone for data storage, processing, and orchestration, Microsoft Azure provides a comprehensive suite of cloud services. Its scalability, reliability, and robustness make it an ideal choice for handling large volumes of data in diverse environments.

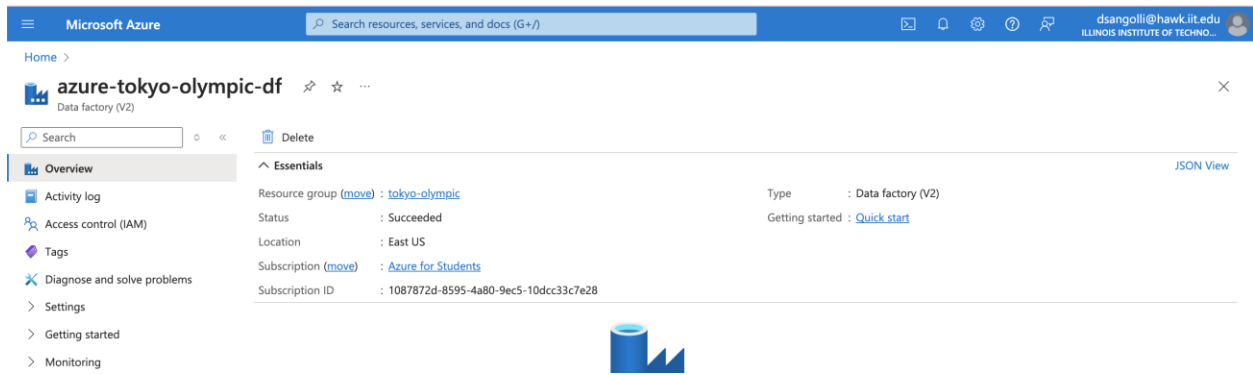


Azure Blob Storage:

Azure Blob Storage serves as a foundational component, offering a scalable solution for storing both raw and transformed data. Its flexibility and durability make it well-suited for accommodating the diverse needs of data storage in this project.

Azure Data Factory:

Azure Blob Storage serves as a foundational component, offering a scalable solution for storing both raw and transformed data. Its flexibility and durability make it well-suited for accommodating the diverse needs of data storage in this project.



Apache Spark on Azure Databricks:

Apache Spark on Azure Databricks empowers data transformation at scale, enabling the efficient processing of large datasets. By leveraging Spark's distributed computing capabilities within the Azure Databricks environment, complex data transformation tasks can be executed with speed and agility.

Jupyter Notebook:

Additionally, Jupyter Notebook has been utilized as a complementary tool for data exploration, analysis, and experimentation. Its interactive and collaborative nature makes it well-suited for prototyping and iterating on data analysis workflows, complementing the capabilities of other tools in the project's toolkit.

Power BI:

Power BI is used to enable the creation of interactive dashboards and reports. With its intuitive interface and powerful visualization features, Power BI facilitates the exploration and communication of insights derived from the analyzed Olympic data.

Data Processing Workflow

Data Ingestion

Historical Olympic records are ingested into the Bronze layer, providing the raw material for subsequent processing.

Data Cleaning and Transformation (Silver Layer)

Cleaning Steps

Duplicate records are removed to ensure data integrity. Missing values are filled, and text fields are standardized for consistency.

Transformation

Date formats are standardized for uniformity across datasets. Categorical data is organized, and new columns are computed to facilitate analytical queries.

Data Structuring (Gold Layer)

Dimension Tables

Static attributes such as athlete details, coach information, and team profiles are extracted and structured into dimension tables. These tables support efficient slicing and dicing of data in reports.

Fact Tables

Aggregate measures such as medal counts by country and gender participation across sports are computed and stored in fact tables. These tables enable insightful analysis and trend identification.

Detailed Descriptions and Steps for Dimension and Fact Tables

Dimension Tables

dim_athletes.csv

Attributes:

- PersonName
- AthleteID
- Country
- Discipline

Description: Contains detailed information about athletes participating in the Olympics.

- **Extract** athlete names, countries, disciplines, genders, birthdates, heights, weights, and medal counts.
- **Deduplicate** records based on unique athlete identifiers.
- **Standardize** text fields and formats for consistency.

- **Compute** aggregate measures such as total medal counts for each athlete.

dim_coaches.csv

Attributes:

- CoachID
- Name
- Country
- Discipline

Description: Stores data related to coaches involved in Olympic sports.

- **Capture** coach names, countries, disciplines, genders, and birthdates.
- **Eliminate** duplicate records based on unique coach identifiers.
- **Ensure** consistency in text fields and formats.
- **Optionally**, calculate coaching tenure or other relevant metrics.

dim_teams.csv

Attributes:

- TeamID
- TeamName
- Discipline
- Country

Description: Provides insights into the teams competing in various sports at the Olympics.

- **Gather** team names, countries, disciplines, and establishment years.
- **Remove** duplicate entries based on unique team identifiers.
- **Standardize** text fields and formats to enhance data consistency.
- **Link** teams to their respective coaches using unique identifiers.

Fact Tables

fact_medal_wins.csv

Attributes:

- FactID
- CountryID

- DisciplineID
- GoldCount
- BronzeCount
- SilverCount
- TotalMedals

Description: Aggregates medal counts by country, reflecting their performance in Olympic events.

- **Group** medal data by country to calculate total gold, silver, and bronze medals.
- **Aggregate** medal counts across disciplines to provide a comprehensive view.
- **Ensure** data accuracy and completeness through rigorous validation.

fact_entries_gender.csv

Attributes:

- EntryID
- DisciplineID
- FemaleCount
- MaleCount
- TotalCount

Description: Summarizes gender participation across different sports at the Olympics.

- **Tabulate** the number of male and female participants in each sport discipline.
- **Aggregate** gender participation data to analyze overall trends.
- **Validate** data integrity to ensure accurate representation.

Visualization and Reporting in Power BI

Dashboards Created

1. Bar Graph: Medal Counts by Country

- **Description:** This bar graph displays the sum of bronze, silver, and gold medals won by each country. It allows users to quickly assess which countries have the highest medal counts.
- **Data Used:** `fact_medal_wins.csv`
- **Fields:**

- X-axis: **Country_Name**
- Y-axis: **BronzeCount, SilverCount, GoldCount**
- **Visualization Type:** Bar Graph

2. Stacked Bar Graph: Coach Count by Country

- **Description:** This stacked bar graph represents the count of coaches by country, segmented by sport discipline. It provides insights into the distribution of coaching resources across different sports within each country.

3. Pie Chart: Participant Count by Team

- **Description:** This pie chart shows the count of participants by team name. It is useful for identifying the teams with the highest number of participants, highlighting prominent teams in the Olympics.

4. Table: Total Medals by Country

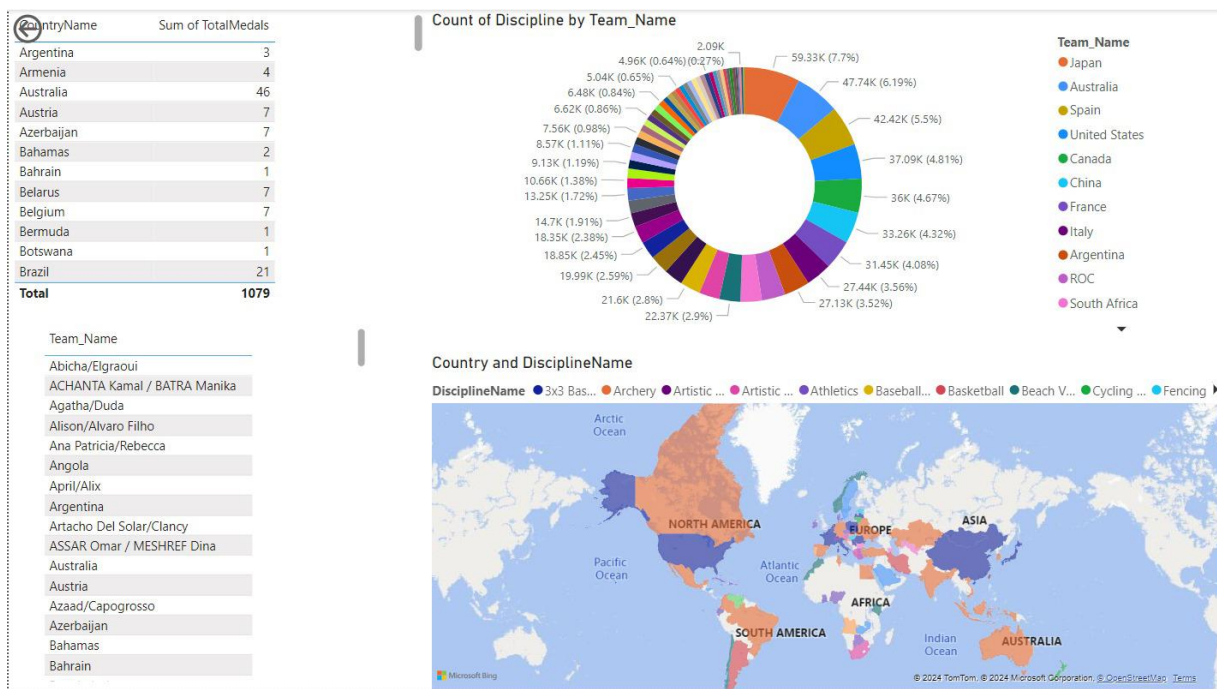
- **Description:** This table displays a summary of total medals by country name, providing a detailed breakdown of medal types (gold, silver, bronze).

5. Donut Chart: Discipline Count by Team

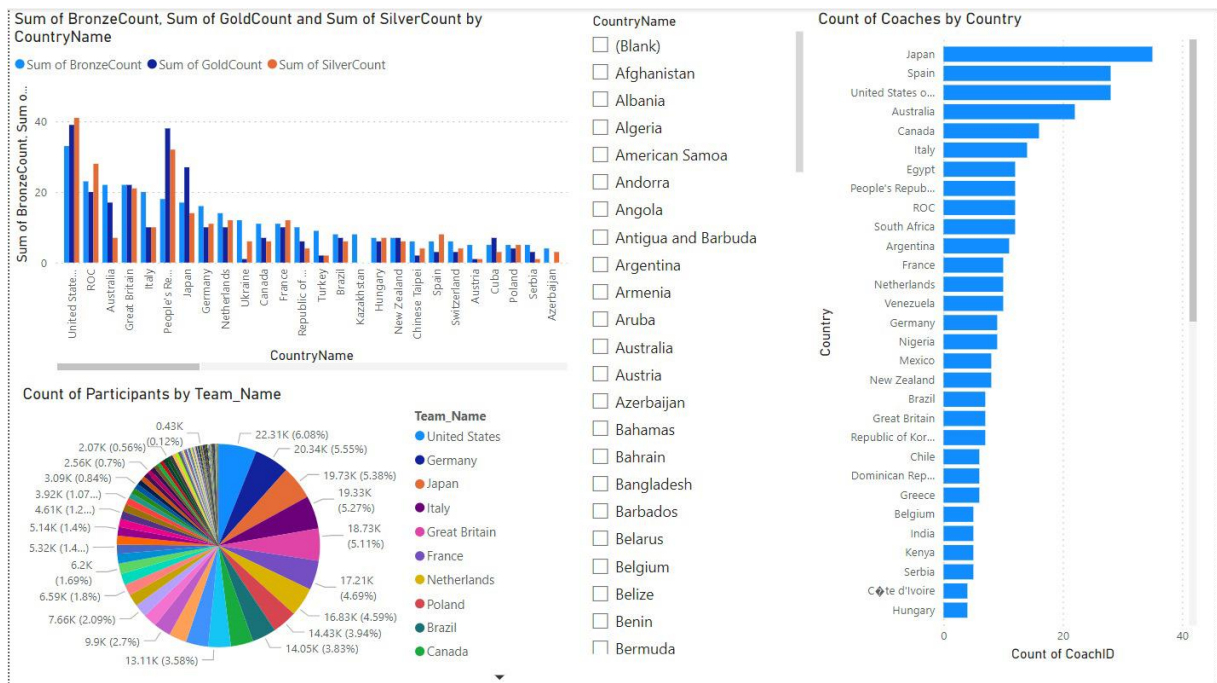
- **Description:** The donut chart visualizes the count of different disciplines represented by each team. It offers a clear view of the diversity of sports that teams are participating in.

6. Filled Map: Discipline and Country Distribution

- **Description:** This filled map provides a geographical visualization of countries and the disciplines they participate in, colored by the number of participants in each discipline. It helps visualize the global distribution of sports participation at the Olympic level.
- Dashboard 1



Dashboard 2



Conclusion

By leveraging Azure and Power BI, this project successfully unlocks valuable insights from Olympic data, demonstrating the effectiveness of cloud and BI technologies in handling large datasets. The structured Medallion architecture ensures scalability and clarity in data processing and analysis, facilitating informed decision-making in the realm of Olympic sports.