

## Pre-Report for the Summer Olympics Dataset:

### Index

Sl.No	Report	Page No
1	Problem Statement	1
2	Data Requirements	1
3	Data Collection	1
4	Data Validation	2
5	Data Cleaning	2
6	Tools	2
7	Dashboard	3
8	Storytelling	3

#### 1. Problem Statement

**"Develop a comprehensive dashboard to analyze and visualize key insights from the Summer Olympics dataset."**

The goal is to create a dashboard that provides insights into athlete demographics, event participation, and medal distributions to understand historical trends and highlight areas of dominance or growth.

#### 2. Data Requirements

- The datasets required include:

**Athlete Events Data:** Detailed records of athletes, their demographic information, events participated in, and medals won.

- Key Features:**

**Demographics:** Age, Gender (Sex), Height, Weight.

**Participation:** Team, NOC, Games, Year, Season, City, Sport, Event.

**Achievements:** Medal.

- NOC Regions Data:** Mapping of NOC codes to regions for geographical analysis.
- Key Features:** NOC, region.

The data must be:

- Sufficient in size to analyze trends (current dataset contains over 271,000 rows).
- Clean and complete, with no missing values in critical columns.
- Merged effectively to integrate region-level insights.

#### 3. Data Collection

The datasets provided are:

- Athlete Events Data:** Covers athlete participation across multiple Summer Olympics, including 271,116 rows and 15 features.

- **NOC Regions Data:** Maps 230 NOC codes to their respective regions.

#### 4. Data Validation

To ensure data quality:

##### Check Completeness:

- Confirm no missing values in critical columns like Name, Sex, Sport, Event, and Medal.

##### Validate Ranges:

- Ensure Age is within a realistic range (e.g., 10–100 years).
- Validate Height (e.g., 100–250 cm) and Weight (e.g., 30–200 kg).

##### Consistency in Categories:

- Verify consistent values for categorical fields like Sex (M, F) and Season (Summer only).

##### NOC Mapping:

- Cross-check that all NOC codes in the Athlete Events dataset map correctly to regions.

#### 5. Data Cleaning

Key steps include:

##### Handling Missing Values:

- Drop or impute missing values for non-critical fields like Height and Weight.

##### Standardizing Numerical Features:

- Normalize variables like Height and Weight for analysis.

##### Merging Datasets:

- Combine Athlete Events and NOC Regions datasets on the NOC column.

##### Encoding:

- Encode categorical variables (e.g., Medal) for analysis and visualization.

#### 6. Tools

##### Python Libraries:

- Pandas and Numpy for data manipulation.
- Matplotlib and Seaborn for static visualizations.

##### Dashboard Tools:

- Tableau or Power BI for interactive dashboards with filters and drilldowns.

## 7. Dashboard

The dashboard will include:

### Key Metrics:

- Total number of athletes, events, and medals.
- Medal counts by country and year.

### Demographics:

- Distribution of athlete age, gender, height, and weight.

### Participation Trends:

- Trends in athlete participation over the years.
- Popular sports/events by the number of participants.

### Medal Insights:

- Medal distribution by region and sport.
- Comparison of top-performing countries.

### Interactive Features:

- Filters for Year, Region, Sport, and Gender.

## 8. Storytelling

### Presenting the Data:

- Highlight the growth of the Olympics in terms of participation and global representation.
- Showcase the dominance of specific countries in certain sports.
- Use data to explore demographic trends among athletes over time.

### Visual Storytelling:

- **Trends:** Medal counts and athlete participation by year and sport.
- **Insights:** Correlation between athlete demographics (e.g., age, height, weight) and medal success.
- **Regional Highlights:** Top-performing regions and their dominant sports.

### Actionable Insights:

- Identify underrepresented regions or sports to target for future development.
- Highlight key factors contributing to athletic success in specific events.