# Project of Datawarehousing

Master degree in computer science - UNIGE - 2021/2022

Pastorino Edoardo
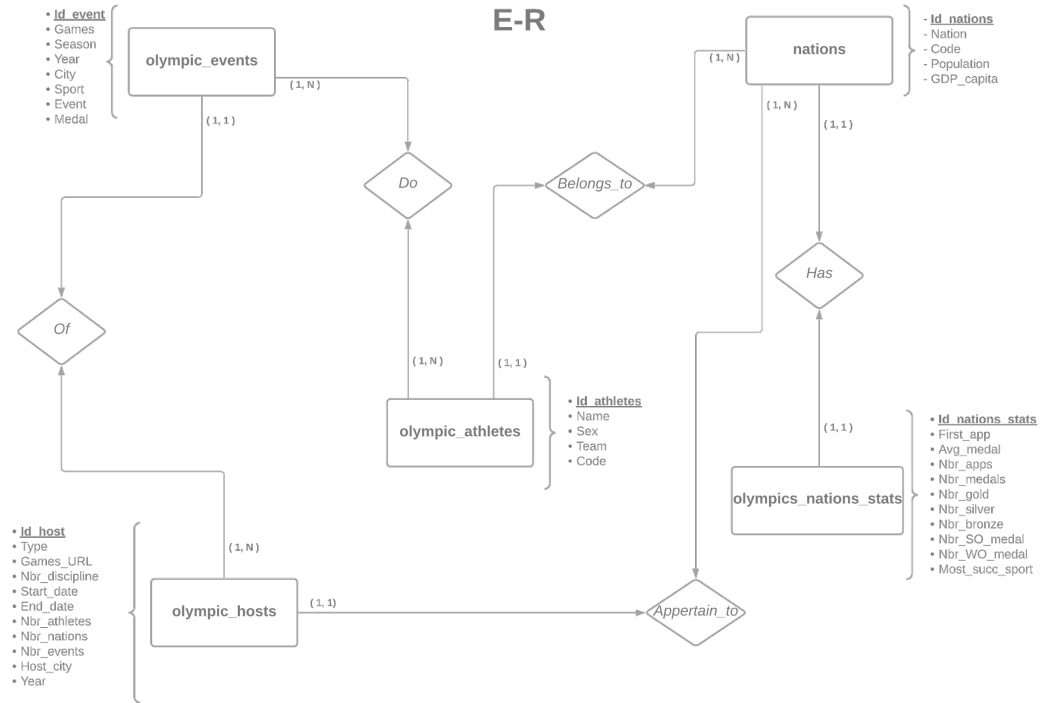Kozy-Korpesh Tolep

# Task 1: Inspection and Profiling

# Hypothesis

1) Do countries with higher GDP / Capita perform better at the Olympics?
2) Is the number of women winners low in the Iraq, Iran, Afghanistan Nations?
3) How is the trend about the winning of medals in the United States over the years?
4) Is it true that the countries with higher populations perform better at the Olympics?
5) Do states that win a lot of medals in the summer olympics also win a lot of medals in the winter olympics?
6) What is the sport in which the Nordic European people (female and male) win more?
7) What are the ten best nations in terms of medal's average?
8) Is it true that the rich olympic host has a higher number of disciplines and athletes?
9) What is the sport with the max number of events in 2016 and in 1896?

# Datasets

1) Dataset1:
https://www.kaggle.com/heesoo37/120-years-of-olympic-history-athletes-and-results/version/2
2) Dataset 2:
https://www.kaggle.com/amirba/olympic-sports-and-medals-18962021
3) Dataset3:
https://www.kaggle.com/piterfm/olympic-games-hosts
4) Dataset4:
https://data.world/johayes13/summer-winter-olympic-games

# Integration and E-R Schema

For integrating the different data sources we decided to remove redundant and unuseful attributes (columns) for the analysis, and sometimes remove an entire table. We have also changed the name of the columns for avoiding mismatch.
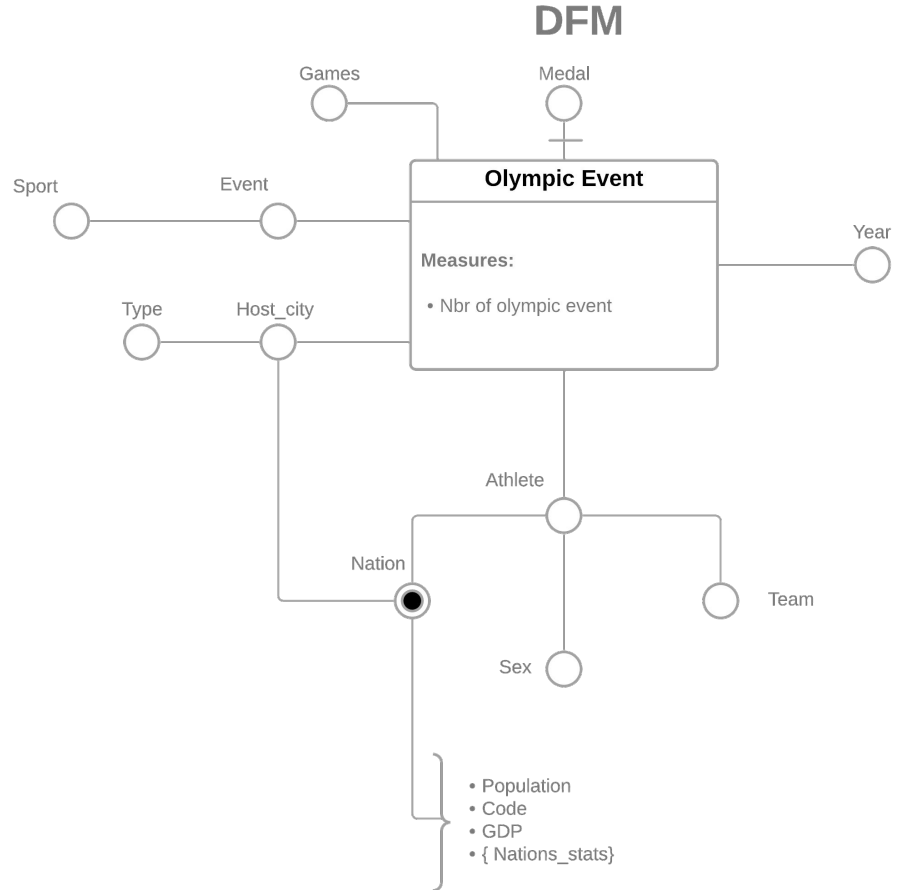
# Task 2: Conceptual Design

# Dimensional Fact Model



According to the DFM conceptual model we have identified the Olympic Events as the fact of interest in the analysis.

We have inserted one measure, the classical count of instance of the fact.

We have identified 6 dimensions: Event, Athlete, Medal, Year, Host City and games.

# Workload

# Data Volume

Each component of this workload refers to the business questions, and it is written following the syntax related to the aggregation patterns used for the relative query.

1) Olympic_participation( Athlete.Nation, Athlete.Nation.GDP, Athlete.Nations.Nbr_medals )
   For each nation shows the GDP and number of medals, ordered by GDP

   . . .

271000 olympic events

201 nations

50 Host city and 2 types

135000 athletes of 2 sex and 1184 teams

765 event, 3 kind of medals, 66 sports and 51 games

40 years, approximately 20 days of duration of each olympics
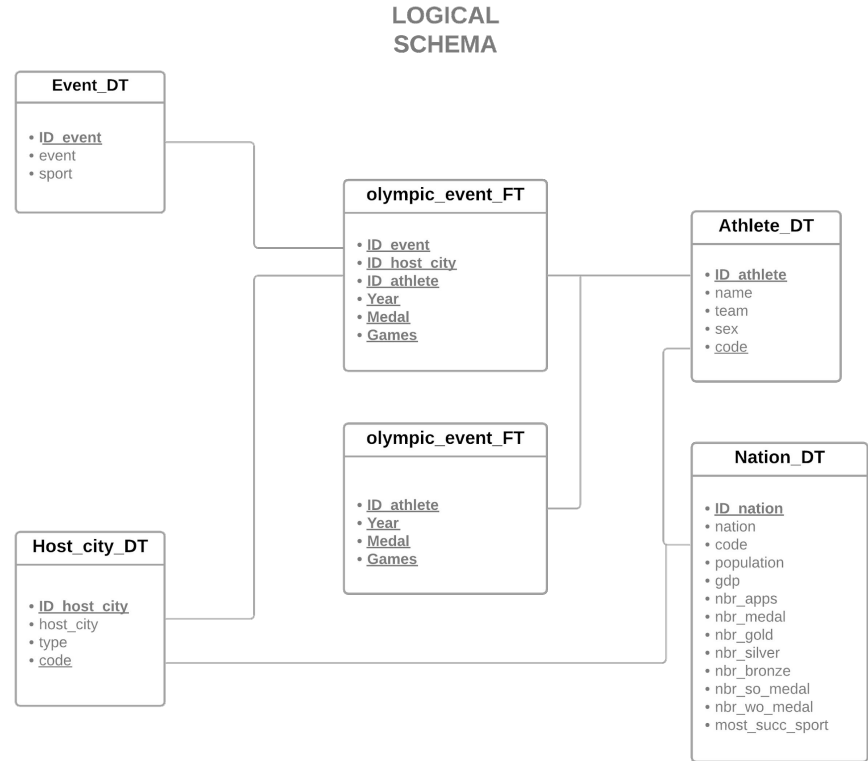
# Task 3: Logical Design

# Logical Schema

For the logical schema at first we have imported the csv files in the tables in postgresql regarding the E-R Schema.
After we have filled the dimension tables and primary fact table and a secondary fact table.
 The degenerate dimensions are represented directly in the fact table.
The logical schema in question is a snowflake schema



**Event_DT**
- ID_event
- event
- sport

**olympic_event_FT**
- ID_event
- ID_host_city
- ID_athlete
- Year
- Medal
- Games

**Athlete_DT**
- ID_athlete
- name
- team
- sex
- code

**olympic_event_FT**
- ID_athlete
- Year
- Medal
- Games

**Host_city_DT**
- ID_host_city
- host_city
- type
- code

**Nation_DT**
- ID_nation
- nation
- code
- population
- gdp
- nbr_apps
- nbr_medal
- nbr_gold
- nbr_silver
- nbr_bronze
- nbr_so_medal
- nbr_wo_medal
- most_succ_sport

# Task 4: Olap Queries

# Queries for creating tables

Ex. of query for creating dimension table

```
CREATE MATERIALIZED VIEW Athlete_DT
(
    ID_athlete,
    Name,
    Team,
    Sex,
    Code

)

AS

SELECT
    "LogicalSchema_Project_DW".olympic_athletes.Id_athlete,
    "LogicalSchema_Project_DW".olympic_athletes.Name,
    "LogicalSchema_Project_DW".olympic_athletes.Team,
    "LogicalSchema_Project_DW".olympic_athletes.Sex,
    "LogicalSchema_Project_DW".olympic_athletes.Code
FROM
    "LogicalSchema_Project_DW".olympic_athletes
```

Ex. of  query for creating E-R's table

```
CREATE TABLE olympic_athletes
(
    ID_athlete serial PRIMARY KEY,
    Name varchar(256),
    Sex varchar(10),
    Team varchar(256),
    Code varchar(10)
);
```

# Queries for the Workload and OLAP Extension

Example of workload query

SELECT nation, gdp, nbr_medal FROM "LogicalSchema_Project_DW".nation_dt  where gdp is not null order by Gdp desc

Example of OLAP extension query

Select team, year, count(medal) as nbr_medal, RANK() over (PARTITION BY team ORDER BY count(medal) desc ) as ranking
FROM
    "LogicalSchema_Project_DW".olympic_events_ft2
    join "LogicalSchema_Project_DW".athlete_dt on

"LogicalSchema_Project_DW".olympic_events_ft2.id_athlete = "LogicalSchema_Project_DW".athlete_dt.id_athlete
    where (medal != 'NA' and Team= 'Italy') or (medal != 'NA' and Team= 'Germany') or (medal != 'NA' and Team= 'France')
Group by (team,year)

Task 5:Hive/SparkSQL

# Example of workload query

```
data = datasetDF.select('Nation', 'GDP per Capita',
'Nbr_medal ').filter(col('GDP per
Capita').isNotNull()).sort(desc('GDP per Capita')).show()
```

```
+--------------+-------------+-------------+
|        Nation|Nbr_SO_medal|Nbr_WO_medal|
+--------------+-------------+-------------+
|     Argentina|           74|            0|
|     Australia|          497|           15|
|       Austria|           87|          232|
|       Belarus|           78|           18|
|       Belgium|          148|            6|
|        Brazil|          129|            0|
|      Bulgaria|          218|            6|
|        Canada|          302|          199|
|         China|          546|           62|
|          Cuba|          226|            0|
|Czech Republic|           56|           31|
|       Denmark|          194|            1|
|      Ethiopia|           54|            0|
|       Finland|          303|          167|
|        France|          716|          124|
|       Germany|          615|          240|
|United Kingdom|          851|           32|
|        Greece|          116|            0|
|       Hungary|          491|            7|
|          Iran|           69|            0|
+--------------+-------------+-------------+
only showing top 20 rows
```
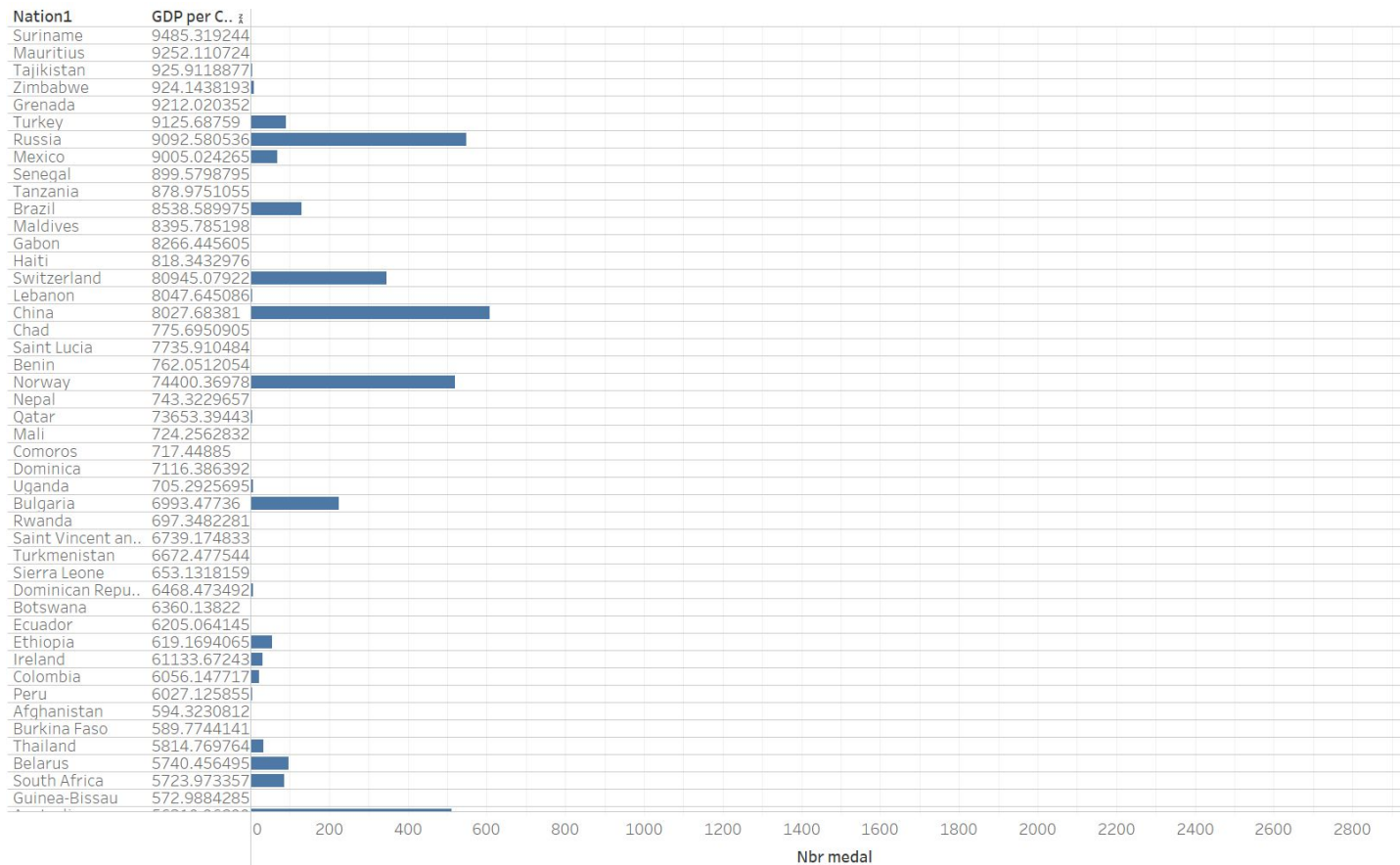
Result of the query
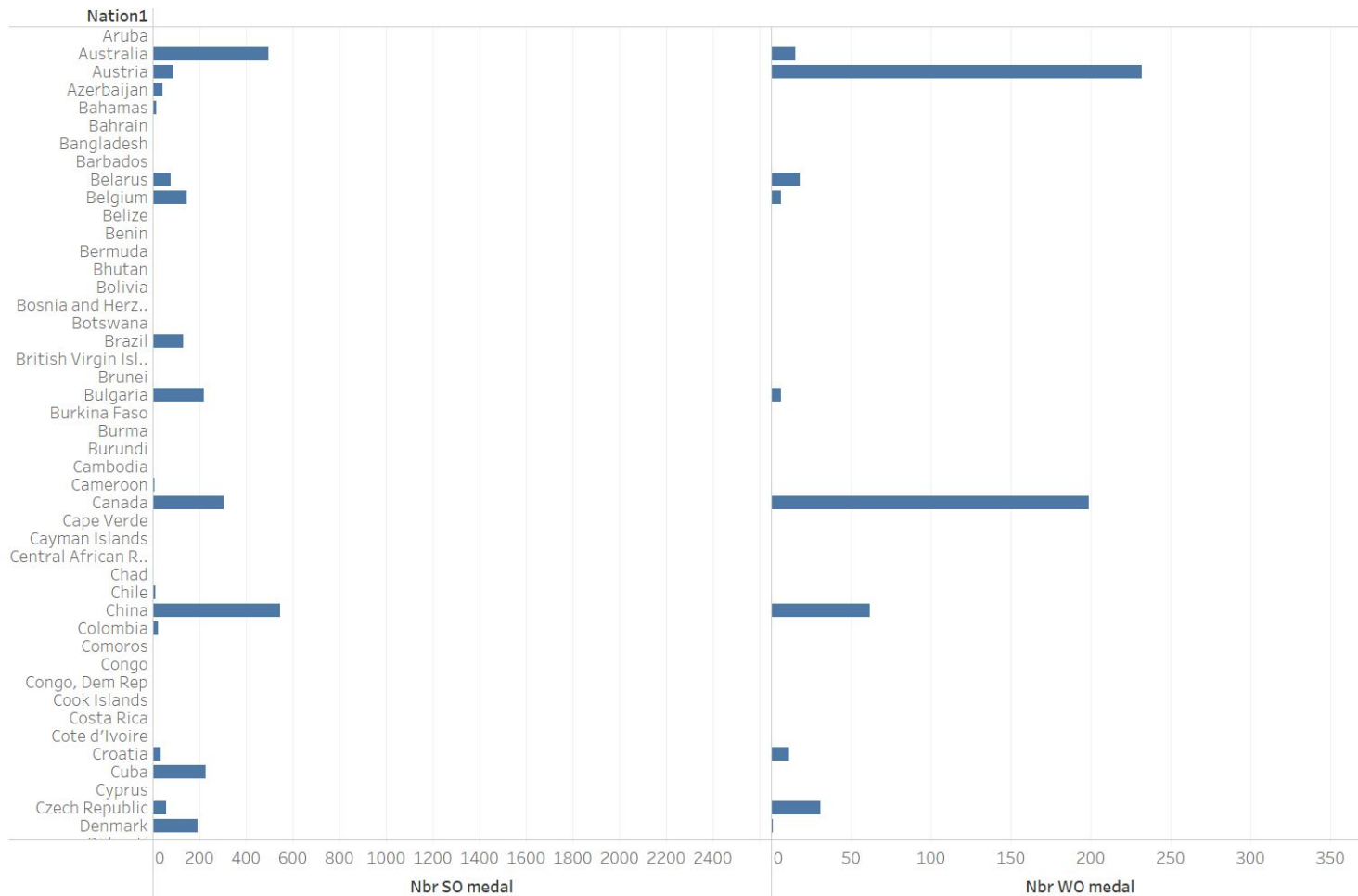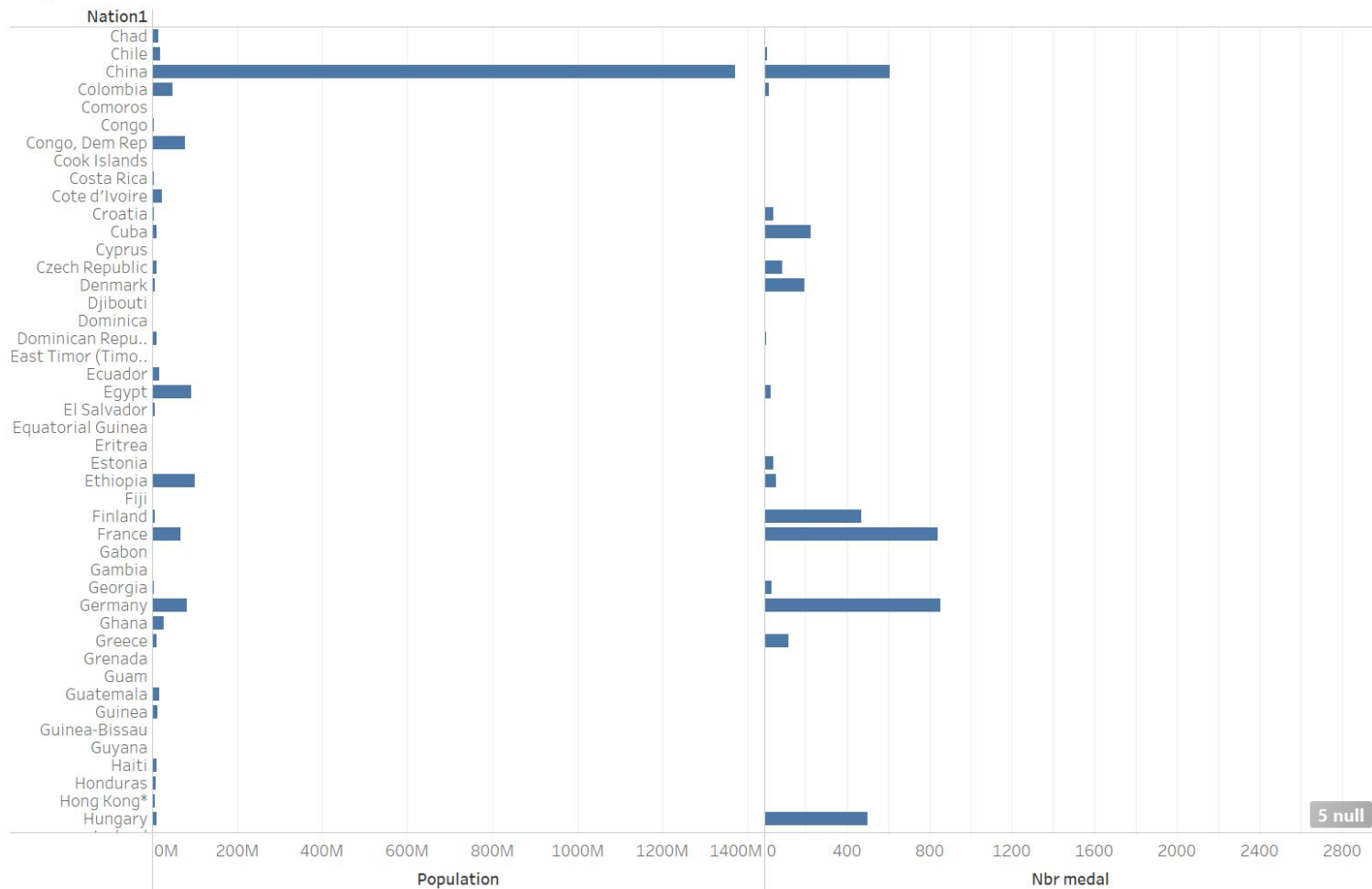
# Task 6: Tableau Report

Report 1

Report 3

Report 4

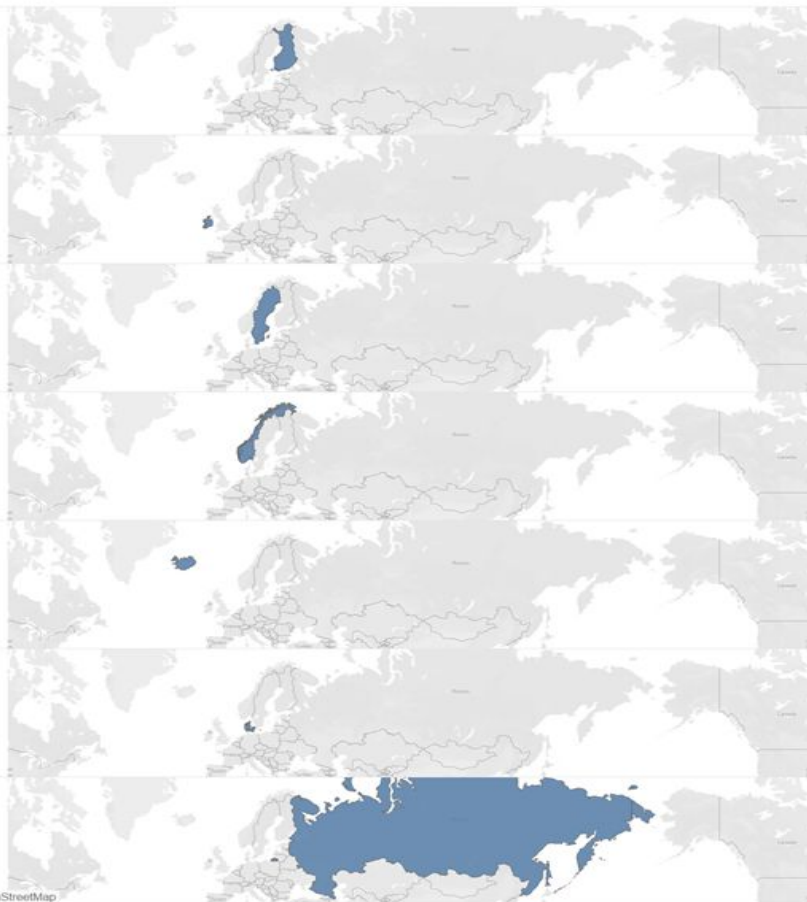# Report 5



Foglio 5

Most s..

Athleti..

Boxing

Cross country skiing

Cross Country Skiing

Handb..

Sailing

Wrestli..

© 2022 Mapbox © OpenStreetMap

Basato su mappa su Longitudine (generata) e Latitudine (generata) suddiviso per Most succ sport. Vengono mostrati i dettagli per Nation1.La vista è filtrata su Nation1, che mantiene 7 di 198 membri.

# Final Considerations

# Critical Analysis

About Tableau Prep application we have noted that it isn't intuitive the way for making the correspondent v-lookup function on the CSV in the table of Tableau Prep and also we didn't find a method for insert a serial number as the primary key of a table.

# Hours Effort

The effort for finishing the entire project in terms of hours is approximately 50 hours.

The first task was about 14 hours of work.

The second task was also about 12 hours of work.

The third task required 8 hours, instead the task number four required approximately 6 hours of work.

The last part, the task number six required about 10 hours.

# Organization of the work

We have tried to split the amount of work based on our engagements (other exam and personal appointment). We often spoke in chat and occasionally also via meetings. We have work on shared files and directories in Google Drive.

The first task, relative to the inspection and profiling of the data, it has been done by Edoardo Pastorino and checked and correct by Kozy-Korpesh Tolep, the same for the task number two, relative to conceptual design and task number four, the one relative to the OLAP Sql queries.

Instead the 5th task has been done by Kozy-Korpesh Tolep alone and only after checked by Edoardo Pastorino.

The task number three and six, the final task, they have been designed and done together.

# THE END

THANK YOU FOR YOUR ATTENTION