

Project Proposal

1. Which client/dataset did you select and why?

Client: SportsStats

SportsStats is a sports analysis firm partnering with local news and elite personal trainers to provide "interesting" insights to help their partners. Insights could be patterns/trends highlighting certain groups/events/countries, etc. for the purpose of developing a news story or discovering key health insights.

The data set being used: Olympics Dataset - 120 years of data

I selected the SportsStats client set as I have always had an interest in sports and being able to use the skills I learned through this course to dive deeper into sports related data is very appealing.

```
In [54]: import pandas as pd
from pandas import sqlqrd
pysqlqr = lambda q: sqlqrd(q, globals())

In [55]: athlete_events = pd.read_csv('C:\Users\Shay\Documents\SQL_Data_Science\Captstone_Project\athlete_events.csv')
#store athlete_events

Stored 'athlete_events' (DataFrame)

In [56]: noc_regions = pd.read_csv('C:\Users\Shay\Documents\SQL_Data_Science\Captstone_Project\noc_regions.csv')

In [57]: athlete_events.head()
```

ID	Name	Sex	Age	Height	Weight	Team	NOC	Games	Year	Season	City	Sport	Event	Medal
0	A Dyang	M	24.0	180.0	80.0	China	CHN	1992 Summer	1992	Summer	Barocone	Basketball	Basketball Men's Basketball	None
1	A Lanius	M	23.0	170.0	60.0	China	CHN	2012 Summer	2012	Summer	London	Judo	Judo Men's Extra-Lightweight	None
2	Gurner Nelson Asby	M	24.0	NaN	NaN	Denmark	DEN	1920 Summer	1920	Summer	Antwerpen	Football	Football Men's Football	None
3	Edgar Lindqvist Aulby	M	34.0	NaN	NaN	Denmark/Sweden	DEN	1900 Summer	1900	Summer	Paris	Tug-Of-War	Tug-Of-War Men's Tug-Of-War	Gold
4	Christine Jacobs Astrha	F	21.0	185.0	82.0	Netherlands	NED	1988 Winter	1988	Winter	Calgary	Speed Skating	Speed Skating Women's 500 metres	None

```
In [58]: noc_regions.head()

NOC      region      notes
0  JPN      Japan
1  JHO      Curacao  Netherlands Antilles
2  ALB      Albania
3  ALG      Algeria
4  AND      Andorra

In [59]: athlete_events.isnull().sum()

ID              0
Name            0
Sex             0
Age            9474
Height        62577
Weight        62875
NOC            940
Year           100
Season         100
City           100
Sport          100
Event          100
Medal          100
dtype: int64

Looking at the null values in each category, we have a significant number of nulls within 'Age', 'Height' and 'Weight'. The nulls under the 'Medal' column is expected as not every athlete is going to receive a medal.
```

```
In [60]: pysqlqr(''SELECT COUNT(*) Total
FROM athlete_events'')
```

Total
671116

First, I am going to split the events into two categories since we are working with data from the Olympics, 'Summer' and 'Winter'.

```
In [61]: summer_games = pysqlqr(''SELECT ID, Name, Sex, Age, Height, Weight, Team, NOC, Games, Year, Season, City, Sport, Event, Medal
WHERE Season = "Summer"')
```

Store 'summer_games' (DataFrame)

```
In [62]: summer_games.head()
```

ID	Name	Sex	Age	Height	Weight	Team	NOC	Games	Year	Season	City	Sport	Event	Medal
0	A Dyang	M	24.0	180.0	80.0	China	CHN	1992 Summer	1992	Summer	Barocone	Basketball	Basketball Men's Basketball	None
1	A Lanius	M	23.0	170.0	60.0	China	CHN	2012 Summer	2012	Summer	London	Judo	Judo Men's Extra-Lightweight	None
2	Gurner Nelson Asby	M	24.0	NaN	NaN	Denmark	DEN	1920 Summer	1920	Summer	Antwerpen	Football	Football Men's Football	None
3	Edgar Lindqvist Aulby	M	34.0	NaN	NaN	Denmark/Sweden	DEN	1900 Summer	1900	Summer	Paris	Tug-Of-War	Tug-Of-War Men's Tug-Of-War	Gold
4	Conelia "Cici" Asken (Stranoud)	F	18.0	168.0	NaN	Netherlands	NED	1932 Summer	1932	Summer	Los Angeles	Athletics	Athletics Women's 100 metres	None

```
In [63]: winter_games = pysqlqr(''SELECT ID, Name, Sex, Age, Height, Weight, Team, NOC, Games, Year, Season, City, Sport, Event, Medal
FROM athlete_events
WHERE Season = "Winter"')
```

Store 'winter_games' (DataFrame)

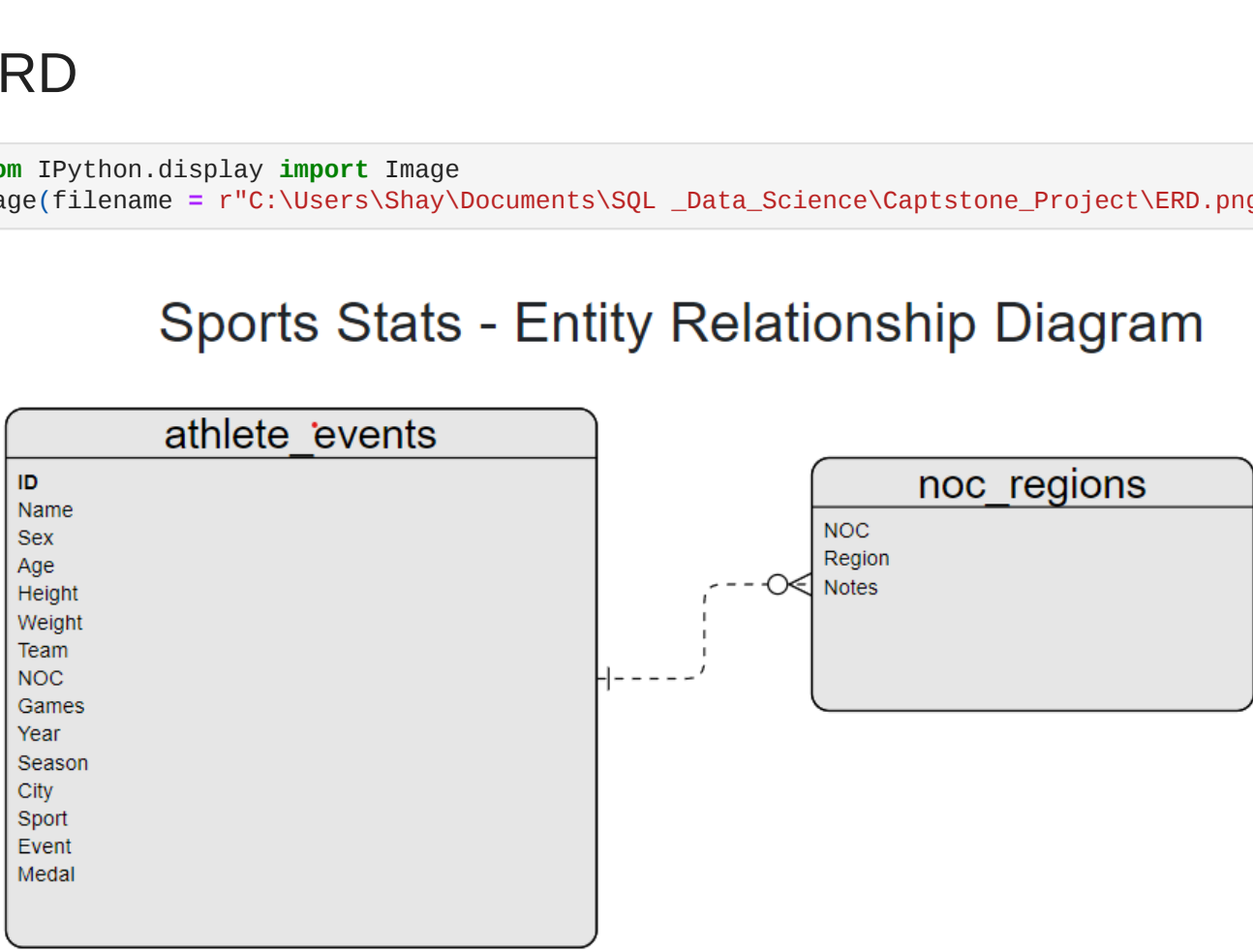
```
In [64]: winter_games.head()
```

ID	Name	Sex	Age	Height	Weight	Team	NOC	Games	Year	Season	City	Sport	Event	Medal
0	Christine Jacobs Astrha	F	21.0	185.0	82.0	Netherlands	NED	1988 Winter	1988	Winter	Calgary	Speed Skating	Speed Skating Women's 500 metres	None
1	Christine Jacobs Astrha	F	21.0	185.0	82.0	Netherlands	NED	1988 Winter	1988	Winter	Calgary	Speed Skating	Speed Skating Women's 1,000 metres	None
2	Christine Jacobs Astrha	F	21.0	185.0	82.0	Netherlands	NED	1988 Winter	1988	Winter	Aberdeen	Speed Skating	Speed Skating Women's 500 metres	None
3	Christine Jacobs Astrha	F	21.0	185.0	82.0	Netherlands	NED	1992 Winter	1992	Winter	Aberdeen	Speed Skating	Speed Skating Women's 1,000 metres	None
4	Christine Jacobs Astrha	F	27.0	185.0	82.0	Netherlands	NED	1994 Winter	1994	Winter	Lillehammer	Speed Skating	Speed Skating Women's 500 metres	None

ERD

```
In [65]: from IPython.display import Image
Image(filename = 'C:\Users\Shay\Documents\SQL_Data_Science\Captstone_Project\ERD.png', width = 800, height = 400)
```

Sports Stats - Entity Relationship Diagram



Questions

Question 1: Is there more men or women competing in the Olympics, how has this changed over time

Question 2: Countries with the most medals and how they compare between Summer and Winter

Question 3: Age distribution, oldest and youngest Olympians

Hypothesis

1: I would think there are more men competing in the Olympics currently but I am very interested in seeing how this has changed overtime and would think in current years it will be a lot more evenly divided.

2: Countries with the most medals are presumed to be the largest countries by population. e.g USA, China etc.

3: I would imagine the age gap would be quite significant and will be interesting to see which athletes are their most divided by.

Approach

1: First step is to review the table to get familiar with the data. I then ran a simple SQL query to get the total counts of males and females that have ever participated. I then declared a variable to store all of the info on males and another variable to store all of the info on females using CAST, CASE and SUM functions to tally up the total number of males and females so I can visualize it.

I created a simple line chart to take a look at the data and see how the male/female participation rates have changed over time, once I plotted the graph it was difficult to read and understand as the numbers significantly dropped every 2 years for the winter Olympics and then spiked back up for the Summer Olympics so it was difficult to grasp the actual difference between males and females.

I then split this out into two categories, participation for Summer Olympics and for Winter Olympics. Once I plotted the graphs broken out it was much easier to read and its easy to tell that there was a huge gap in gender participation in the early years but since the 1980 there has been a steady climb in female participants so eventually this should become very close.

2: I wanted to see the countries with the most medals won over the entire Olympics period (this dataset covers 120 years of data). As expected in my hypothesis I saw a significant lead over other countries.

I then broke it down to see the Summer and Winter figures. Once again as expected for Summer games it was USA leading the pack. The winter games were led by Canada with USA coming in second and then followed by Scandinavian countries such as Norway, Sweden and Finland which is expected for the Winter games.

3: This question was a bit easier but it was something that I was interested in knowing. I simply assigned a variable of ages and ran an SQL query to return all the ages and the count of each of these ages and then grouped by ages. I then plotted this onto a line chart as its always easier to get an understanding of your data when you can visualize it.

Question 1:

Question 1: Is there more men or women competing in the Olympics, how has this changed over time

```
In [66]: pysqlqr(''SELECT Sex, COUNT(Sex)
FROM athlete_events
WHERE Year >= 2000
GROUP BY Sex'')
```

Sex	COUNT(Sex)
0	F 3040
1	M 49218

As expected there are recently more men competing in the Olympics compared to women but lets take a look at how this has changed over time

```
In [67]: male = pysqlqr(''SELECT Year,
CAST(SUM(CASE WHEN Sex = "M" THEN 1 ELSE 0 END) AS INT) Male
FROM athlete_events
GROUP BY Year'')
```

male.tail()

Store 'male' (DataFrame)

```
In [68]: female = pysqlqr(''SELECT Year,
CAST(SUM(CASE WHEN Sex = "F" THEN 1 ELSE 0 END) AS INT) Female
FROM athlete_events
GROUP BY Year'')
```

female.tail()

Store 'female' (DataFrame)

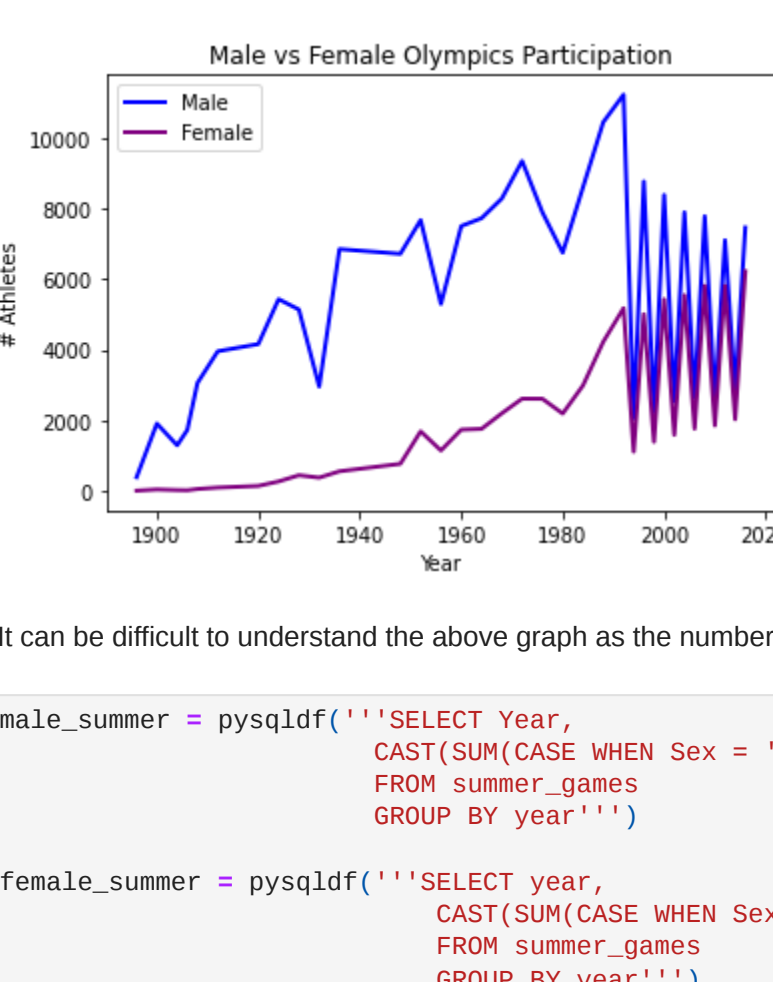
```
In [69]: pip install matplotlib

Requirement already satisfied: matplotlib in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (3.5.1)
Requirement already satisfied: packaging<20.0,>=20.0 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from matplotlib) (21.3)
Requirement already satisfied: fonttools<=4.22.0,>=4.22.0 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from matplotlib) (4.23.1)
Requirement already satisfied: pillow<=6.2.0,>=6.2.0 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from matplotlib) (9.0.3)
Requirement already satisfied: python-dateutil<=2.7,>=2.7 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from matplotlib) (2.8.2)
Requirement already satisfied: pyparsing<=2.2.1,>=2.2.1 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from matplotlib) (2.0.7)
Requirement already satisfied: kiwisolver<=1.3.1,>=1.3.1 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from matplotlib) (1.4.0)
Requirement already satisfied: cycler<=0.10,>=0.10 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from matplotlib) (0.11.0)
Requirement already satisfied: numpy<=1.17,>=1.17 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from matplotlib) (1.22.3)
Requirement already satisfied: six<=1.15,>=1.15 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from python-dateutil<=2.7,>=2.7->matplotlib) (1.16.0)
Note: you may need to restart the kernel to use updated packages.
```

```
In [70]: import matplotlib as mpl
import matplotlib.pyplot as plt
import numpy as np

In [71]: plt.plot(male.Year, male.Male, marker='o', color='blue', linewidth=2, label='Male')
plt.plot(female.Year, female.Female, marker='o', color='purple', linewidth=2, label='Female')
plt.legend(loc=2)
plt.xlabel("Year")
plt.ylabel("Participants")
plt.title("Male vs Female Olympics Participation")

In [72]: Text(0.5, 1.0, "Male vs Female Olympics Participation")
```



It can be difficult to understand the above graph as the numbers decreased significantly every 2 years due to the switch between Summer and Winter games. I will split this out into two different graphs, one for Summer and one for Winter and it will be much easier to read and understand.

```
In [73]: male_summer = pysqlqr(''SELECT Year,
CAST(SUM(CASE WHEN Sex = "M" THEN 1 ELSE 0 END) AS INT) Male
FROM summer_games
GROUP BY Year'')
```

female_summer = pysqlqr(''SELECT Year,
CAST(SUM(CASE WHEN Sex = "F" THEN 1 ELSE 0 END) AS INT) Female
FROM summer_games
GROUP BY Year'')

```
In [74]: male_summer.tail()
```

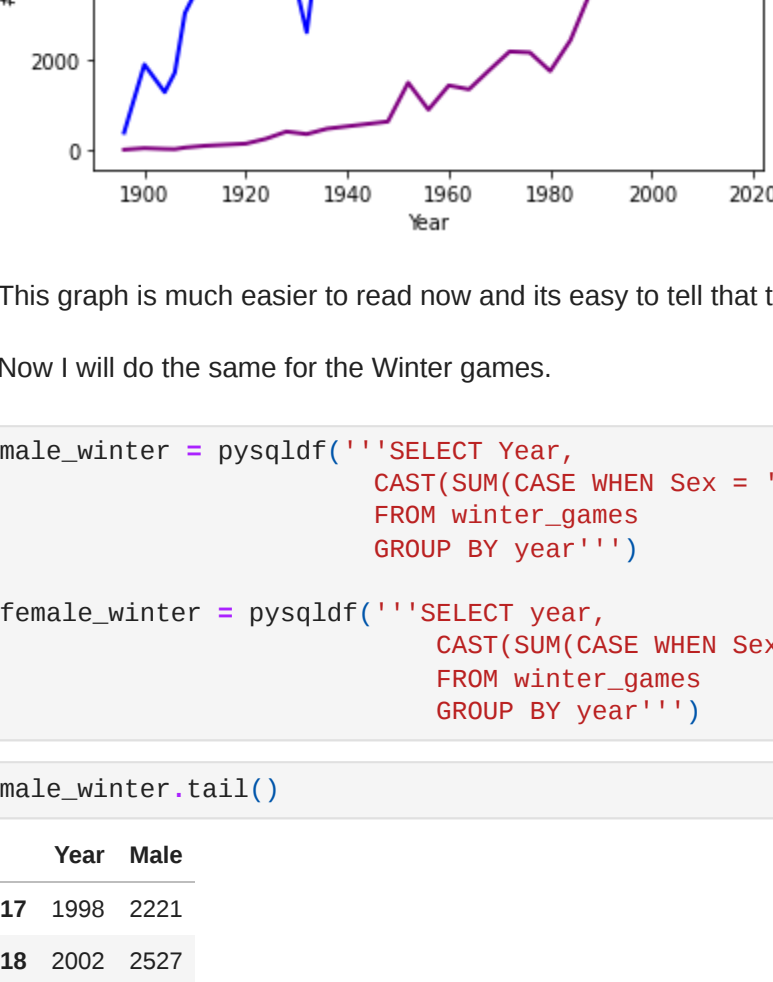
Year	Male
2020	2840
2018	7897
2016	7796
2012	7195
2010	7495

```
In [75]: female_summer.tail()
```

Year	Female
2020	841
2018	5646
2016	5616
2012	5615
2010	6223

```
In [76]: plt.plot(male_summer.Year, male_summer.Male, marker='o', color='blue', linewidth=2, label='Male')
plt.plot(female_summer.Year, female_summer.Female, marker='o', color='purple', linewidth=2, label='Female')
plt.legend(loc=2)
plt.xlabel("Year")
plt.ylabel("Participants")
plt.title("Male vs Female Summer Olympics Participation")

In [77]: Text(0.5, 1.0, "Male vs Female Summer Olympics Participation")
```



This graph is much easier to read now and its easy to tell that there was a huge gap in gender participation in the early years but since the 1980 there has been a steady climb in female participants so eventually this should become very close.

Now I will do the same for the Winter games.

```
In [78]: male_winter = pysqlqr(''SELECT Year,
FROM winter_games
GROUP BY Year'')
```

female_winter = pysqlqr(''SELECT Year,
FROM winter_games
GROUP BY Year'')

```
In [79]: male_winter.tail()
```

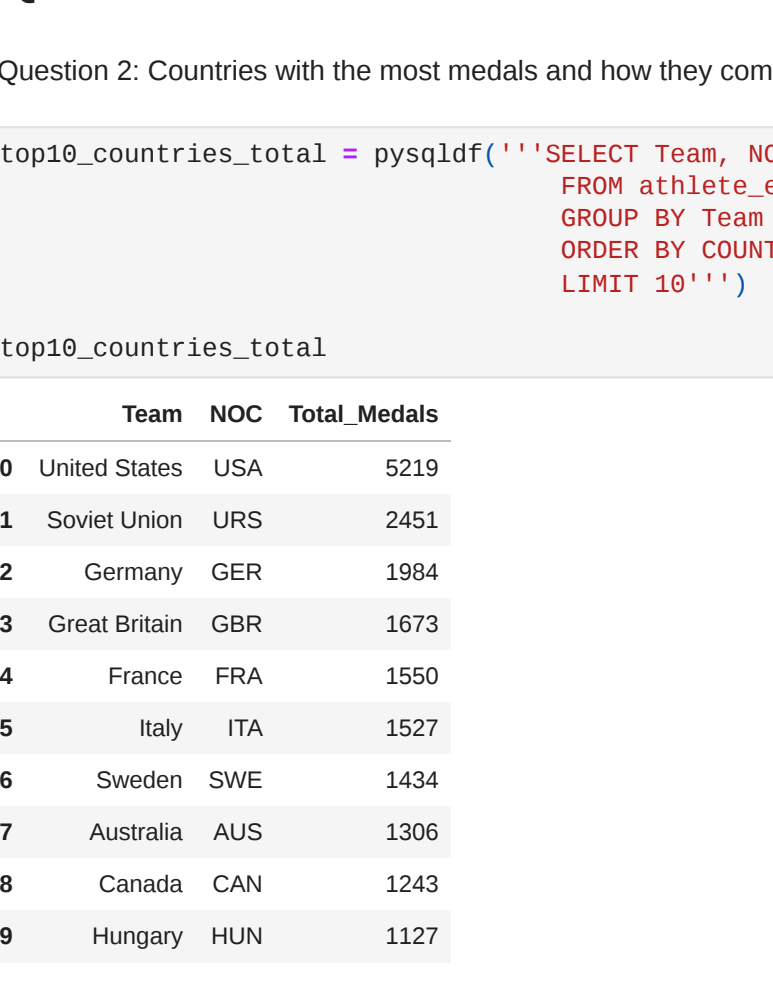
Year	Male
1998	2221
2002	2527
2006	2625
2010	2595
2014	2868

```
In [80]: female_winter.tail()
```

Year	Female
1998	1384
2002	1582
2006	1757
2010	1847
2014	2023

```
In [81]: plt.plot(male_winter.Year, male_winter.Male, marker='o', color='blue', linewidth=2, label='Male')
plt.plot(female_winter.Year, female_winter.Female, marker='o', color='purple', linewidth=2, label='Female')
plt.legend(loc=2)
plt.xlabel("Year")
plt.ylabel("Participants")
plt.title("Male vs Female Winter Olympics Participation")

In [82]: Text(0.5, 1.0, "Male vs Female Winter Olympics Participation")
```



Question 2:

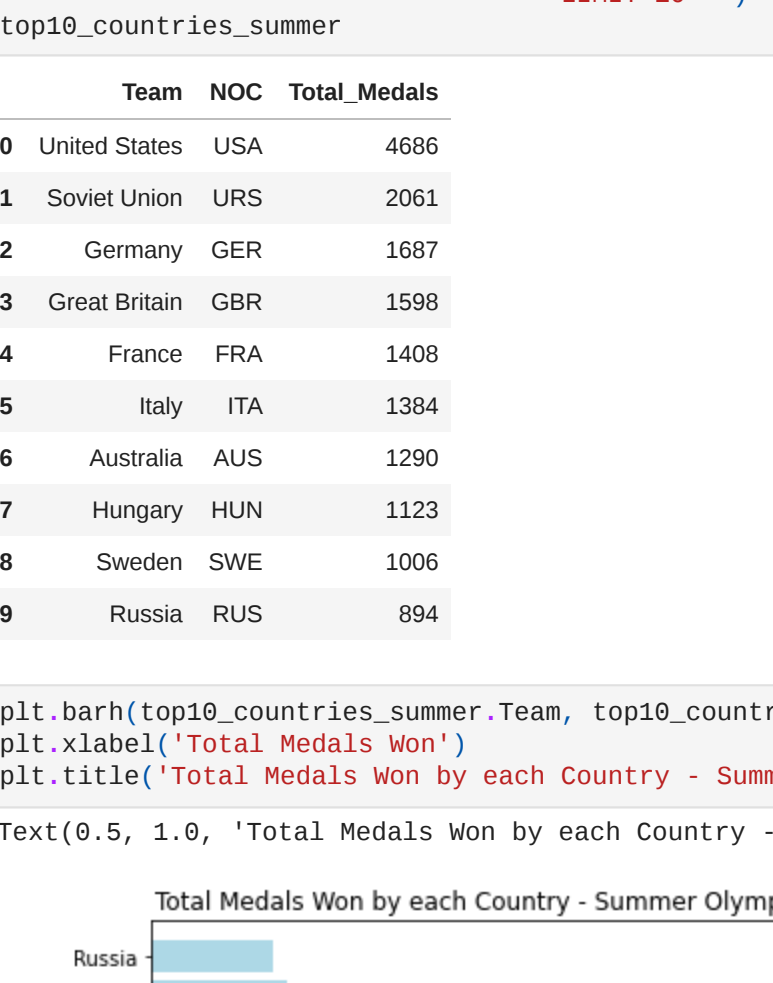
Question 2: Countries with the most medals and how they compare between Summer and Winter

```
In [83]: top18_countries_total = pysqlqr(''SELECT Team, NOC, COUNT(Medal) Total_Medals
FROM summer_games
ORDER BY Team
DESC
LIMIT 10'')
```

Team	NOC	Total_Medals
0	United States	USA 5219
1	Soviet Union	URS 2673
2	Germany	GER 1984
3	Great Britain	GBR 1491
4	France	FRA 1550
5	Italy	ITA 1527
6	Sweden	SWE 1434
7	Australia	AUS 1306
8	Canada	CAN 1243
9	Hungary	HUN 1127

```
In [84]: plt.barh(top18_countries_total.Team, top18_countries_total.Total_Medals, color='lightblue')
plt.xlabel("Total_Medals won")
plt.title("Total Medals won by each Country")

In [85]: Text(0.5, 1.0, "Total Medals won by each Country - Summer Olympics")
```

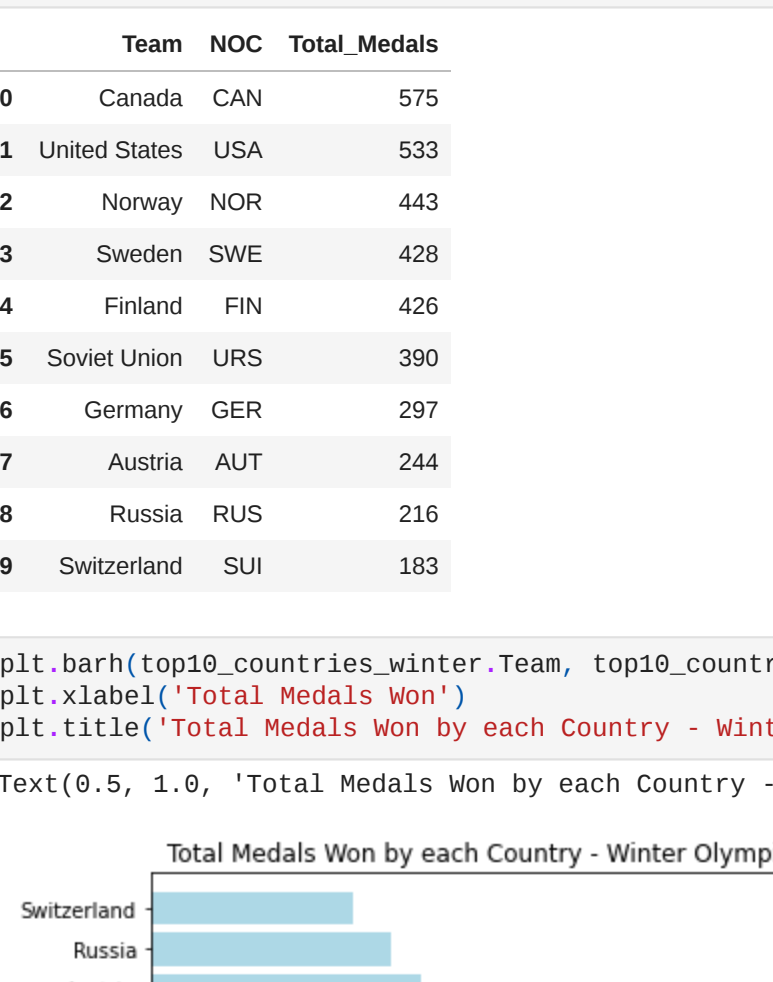


```
In [86]: top18_countries_summer = pysqlqr(''SELECT Team, NOC, COUNT(Medal) Total_Medals
FROM summer_games
ORDER BY Team
DESC
LIMIT 10'')
```

Team	NOC	Total_Medals
0	United States	USA 4686
1	Soviet Union	URS 2081
2	Germany	GER 1687
3	Great Britain	GBR 1598
4	France	FRA 1408
5	Italy	ITA 1286
6	Australia	AUS 1290
7	Hungary	HUN 1123
8	Sweden	SWE 1006
9	Russia	RUS 894

```
In [87]: plt.barh(top18_countries_summer.Team, top18_countries_summer.Total_Medals, color='lightblue')
plt.xlabel("Total Medals won")
plt.title("Total Medals Won by each Country - Summer Olympics")

In [88]: Text(0.5, 1.0, "Total Medals won by each Country - Summer Olympics")
```

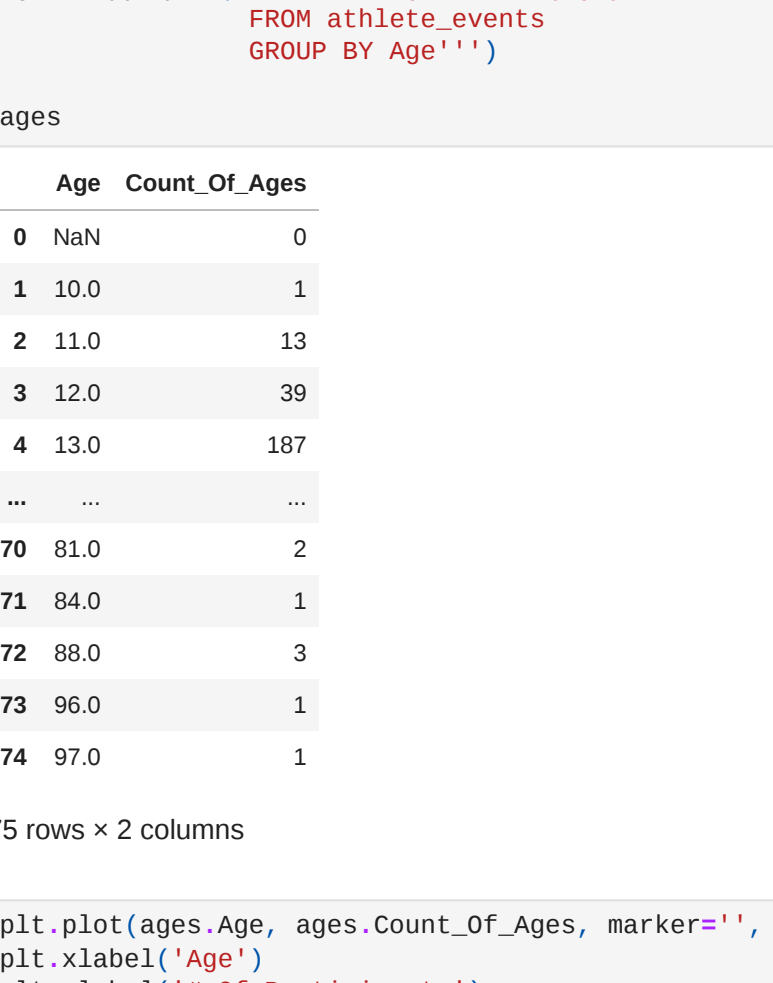


```
In [89]: top18_countries_winter = pysqlqr(''SELECT Team, NOC, COUNT(Medal) Total_Medals
FROM winter_games
ORDER BY Team
DESC
LIMIT 10'')
```

Team	NOC	Total_Medals
0	Canada	CAN 575
1	United States	USA 619
2	Norway	NOR 443
3	Sweden	SWE 428
4	Finland	FIN 426
5	Soviet Union	URS 390
6	Germany	GER 397
7	Austria	AUT 344
8	Russia	RUS 216
9	Switzerland	SUI 183

```
In [90]: plt.barh(top18_countries_winter.Team, top18_countries_winter.Total_Medals, color='lightblue')
plt.xlabel("Total Medals won")
plt.title("Total Medals Won by each Country - Winter Olympics")

In [91]: Text(0.5, 1.0, "Total Medals won by each Country - Winter Olympics")
```



Question 3: Age distribution, oldest and youngest Olympians

```
In [89]: ages = pysqlqr(''SELECT Age, COUNT(Age) Count_Of_Ages
FROM athlete_events
GROUP BY Age'')
```

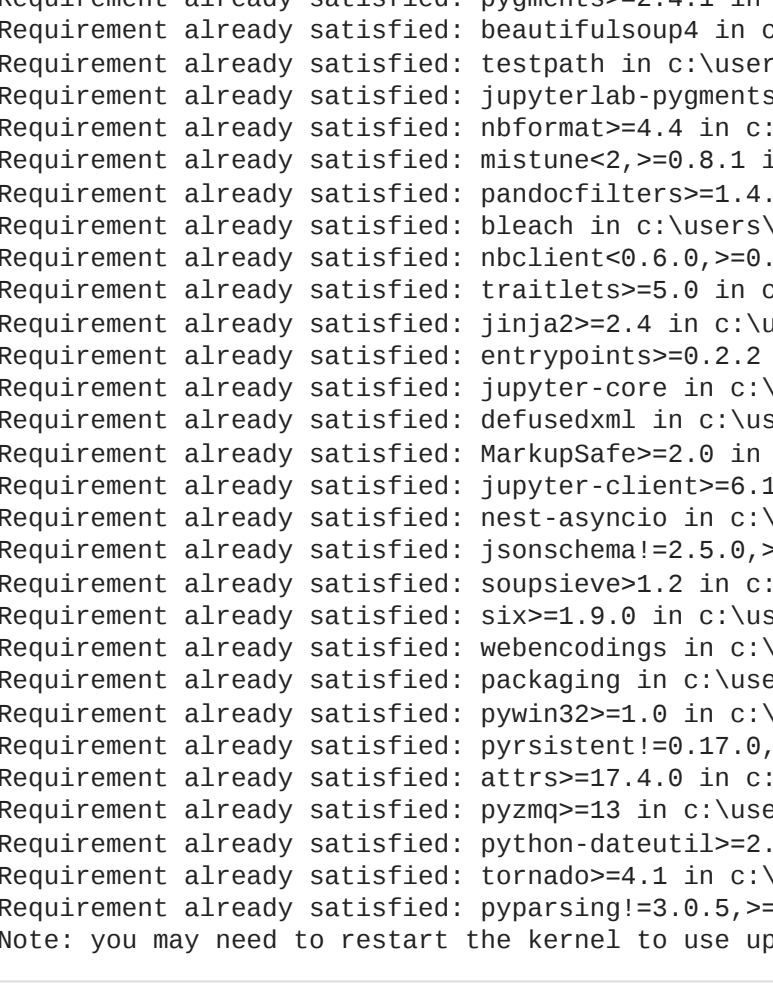
ages

```
In [90]: ages
```

Age	Count_Of_Ages
0	NaN 0
1	10.0 1
2	11.0 13
3	12.0 39
4	13.0 187
...	...
70	84.0 2
71	84.0 1
72	85.0 3
73	86.0 1
74	87.0 1
75	rows = 2 columns

```
In [91]: plt.plot(ages.Age, ages.Count_Of_Ages, marker='o', linewidth = 2)
plt.xlabel("Age")
plt.ylabel("Count of Participants")
plt.title("Count of Participants by Age")

In [92]: Text(0.5, 1.0, "Count of Participants by Age")
```



```
In [93]: avg_age = pysqlqr(''SELECT Average_Age, Average_Age
FROM athlete_events'')
```

avg_age

```
In [94]: avg_age
```

Average_Age	
0	25.55698

```
In [95]: pip install nbconvert

Requirement already satisfied: nbconvert in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (6.4.4)
Requirement already satisfied: pygments>=2.4.1 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from nbconvert) (2.11.2)
Requirement already satisfied: bleach[css]>=3.1.0 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from nbconvert) (4.1.0)
Requirement already satisfied: testpath in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from nbconvert) (0.6.0)
Requirement already satisfied: jupyterlab_pygments in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from nbconvert) (0.2.2)
Requirement already satisfied: pandocfilters<=1.4.1,>=1.4.1 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from nbconvert) (1.5.0)
Requirement already satisfied: nbclient<0.6.0,>=0.5.0 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from nbconvert) (0.6.13)
Requirement already satisfied: nbformat<5.0.0,>=4.4.0 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from nbconvert) (5.1.3)
Requirement already satisfied: traitlets<5.0,>=4.3.3 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from nbconvert) (5.1.1)
Requirement already satisfied: jinja2<=2.4,>=2.4 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from nbconvert) (3.0.3)
Requirement already satisfied: ipynb<=0.5.0,>=0.5.0 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from nbconvert) (0.7.1)
Requirement already satisfied: MarkupSafe<=2.0,>=2.0 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from jinja2<=2.4,=>nbconvert) (2.1.1)
Requirement already satisfied: IPython<=8.0,>=8.0 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from nbclient<0.6.0,>=0.5.0->nbconvert) (7.1.2)
Requirement already satisfied: jmespath<=0.9.4,>=0.9.4 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from nbformat<5.0.0,>=4.4.0->nbconvert) (4.0.0)
Requirement already satisfied: defusedxml in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from bleach->nbconvert) (0.6.0)
Requirement already satisfied: packaging in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from bleach->nbconvert) (23.0)
Requirement already satisfied: pynb<=1.0,>=1.0 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from jupyter-core->nbconvert) (3.0.1)
Requirement already satisfied: pyzmq<=18.17.0,>=18.17.0 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from jupyter-client<6.0,>=6.0->nbconvert) (17.1.2)
Requirement already satisfied: pyzmq<=18.17.0,>=18.17.0 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from jupyter-client<6.0,>=6.0->nbconvert) (17.1.2)
Requirement already satisfied: pyzmq<=18.17.0,>=18.17.0 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from jupyter-client<6.0,>=6.0->nbconvert) (17.1.2)
Requirement already satisfied: tornado<=4.1,>=4.1 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from jupyter-client<6.0,>=6.0->nbconvert) (6.1)
Requirement already satisfied: pygments<=2.8,>=2.8 in c:\users\shay\appdata\local\programs\python\python310\lib\site-packages (from packaging->nbconvert) (3.0.7)
Note: you may need to restart the kernel to use updated packages.
```

```
In [ ]:
In [ ]:
In [ ]:
In [ ]:
```