

Descriptive Statistics and Understanding the Data

```
from pandasql import sqldf
pqsqlf = lambda q: sqldf(q, globals())

istore -> athlete_events
istore -> summer_games
istore -> winter_games
istore -> female
istore -> male
```

```
In [189]:
athlete_events.head()
```

```
Out[189]:
```

ID	Name	Sex	Age	Height	Weight	Team	NOC	Games	Year	Season	City	Sport	Event	Medal	
0	1	A Djang	M	24.0	180.0	80.0	China	CHN	1992 Summer	1992	Summer	Barcelona	Basketball	Basketball Men's Basketball	NaN
1	2	A Lamusi	M	23.0	170.0	60.0	China	CHN	2012 Summer	2012	Summer	London	Judo	Judo Men's Extra-Lightweight	NaN
2	3	Gunnar Nielsen Asby	M	24.0	NaN	NaN	Denmark	DEN	1920 Summer	1920	Summer	Antwerpen	Football	Football Men's Football	NaN
3	4	Edgar Lindqvist Asby	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	5	Christine Jacobs Athink	F	21.0	185.0	82.0	Netherlands	NED	1988 Winter	1988	Winter	Calgary	Speed Skating	Speed Skating Women's 500 metres	NaN

```
In [190]:
summer_games.head()
```

```
Out[190]:
```

ID	Name	Sex	Age	Height	Weight	Team	NOC	Games	Year	Season	City	Sport	Event	Medal	
0	1	A Djang	M	24.0	180.0	80.0	China	CHN	1992 Summer	1992	Summer	Barcelona	Basketball	Basketball Men's Basketball	None
1	2	A Lamusi	M	23.0	170.0	60.0	China	CHN	2012 Summer	2012	Summer	London	Judo	Judo Men's Extra-Lightweight	None
2	3	Gunnar Nielsen Asby	M	24.0	NaN	NaN	Denmark	DEN	1920 Summer	1920	Summer	Antwerpen	Football	Football Men's Football	None
3	4	Edgar Lindqvist Asby	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	8	Connelia "Cor" Aalton (Seamond)	F	18.0	188.0	NaN	Netherlands	NED	1922 Summer	1922	Summer	Los Angeles	Artletics	Artletics Women's 100 metres	None

```
In [191]:
winter_games.head()
```

0	5	Christine Jacobs Astrink	F	21.0	185.0	82.0	Netherlands	NED	1988 Winter	1988	Winter	Calgary	Speed Skating	Speed Skating Women's 500 metres	None
1	5	Christine Jacobs Astrink	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	5	Christine Jacobs Astrink	F	25.0	185.0	82.0	Netherlands	NED	1992 Winter	1992	Winter	Albertville	Speed Skating	Speed Skating Women's 500 metres	None
3	5	Christine Jacobs Astrink	F	25.0	185.0	82.0	Netherlands	NED	1992 Winter	1992	Winter	Albertville	Speed Skating	Speed Skating Women's 1,000 metres	None
4	5	Christine Jacobs Astrink	F	27.0	185.0	82.0	Netherlands	NED	1994 Winter	1994	Winter	Lillehammer	Speed Skating	Speed Skating Women's 500 metres	None

```
In [192]: #Calculating the ratios for Men vs Women
# Total Ratio
total_ratio = pysqlor('''SELECT Sex,
COUNT(*) Total_Count,
COUNT(*) * 100.00 / SUM(COUNT(*)) OVER () AS ratio
FROM athlete_events
GROUP BY sex''')
```

```
summer_ratio = pysqlor("SELECT
COUNT(*) Total_Count,
COUNT(*) * 100.00 / SUM(COUNT(*)) OVER () AS ratio
FROM summer_games
GROUP BY sex")

#winter_games_ratio
winter_ratio = pysqlor("SELECT Sex,
COUNT(*) Total_Count,
COUNT(*) * 100.00 / SUM(COUNT(*)) OVER () AS ratio
FROM winter_games
GROUP BY sex")
```

```
total_ratio.head()
```

```
In [193]:
```

```
Out[193]:
```

	Sex	Total_Count	ratio
0	F	74922	27.487127

In [192]:

#Calculating the ratios for Men vs Women
Total Ratio
total_ratio = pqsqlf("""SELECT Sex,
COUNT(*) Total_Count,
COUNT(*) * 100.00 / SUM(COUNT(*) OVER ()) AS ratio
FROM athlete_events
GROUP BY Sex""")

#Summer games ratio
summer_ratio = pqsqlf("""SELECT Sex,
COUNT(*) Total_Count,
COUNT(*) * 100.00 / SUM(COUNT(*) OVER ()) AS ratio
FROM summer_games
GROUP BY Sex""")

#Winter games ratio
winter_ratio = pqsqlf("""SELECT Sex,
COUNT(*) Total_Count,
COUNT(*) * 100.00 / SUM(COUNT(*) OVER ()) AS ratio
FROM winter_games
GROUP BY Sex""")

In [193]:

total_ratio.head()

Out[193]:

Sex	Total_Count	ratio
0	F	74822 27.487127
1	M	196594 72.512873

In [194]:

summer_ratio.head()

Out[194]:

Sex	Total_Count	ratio
0	F	59443 26.707113
1	M	163109 73.292887

In [195]:

winter_ratio.head()

Out[195]:

Sex	Total_Count	ratio
0	F	15079 31.049749
1	M	32485 68.950251

```

COUNT(*) Participants,
SUM(CASE WHEN Medal IS NOT NULL THEN 1 ELSE 0 END) AS Total_Medals,
SUM(CASE WHEN Medal = 'Gold' THEN 1 ELSE 0 END) AS Gold_Medals,
SUM(CASE WHEN Medal = 'Silver' THEN 1 ELSE 0 END) AS Silver_Medals,
SUM(CASE WHEN Medal = 'Bronze' THEN 1 ELSE 0 END) AS Bronze_Medals
FROM athlete_events
GROUP BY Year'''

```

Year	Participants	Total_Medals	Gold_Medals	Silver_Medals	Bronze_Medals
0 1896	300	143	62	43	38
1 1900	1936	604	201	228	175
2 1904	1301	486	173	163	150
3 1906	1733	458	157	156	145
4 1908	3101	831	294	281	256
5 1912	4040	941	326	315	300
6 1920	4292	1308	493	448	367
7 1924	5693	962	332	319	311
8 1928	5574	823	275	267	281
9 1932	3321	739	261	246	232
10 1936	7401	1025	348	347	330
11 1948	7480	987	330	332	325
12 1952	9358	1033	351	336	347
13 1956	6434	1043	353	342	348
14 1960	9235	1058	359	342	357
15 1964	9480	1215	408	406	401

In [198]:

summer_averages.head()

Out[198]:

Sex	Average_Weight	Average_Height	Average_Age	
0	F	60.067844	168.156025	23.660907
1	M	75.604195	178.901874	26.443944

In [199]:

winter_averages.head()

Out[199]:

25	1998	3605	440	145	145	150
26	2000	33821	2004	663	661	680
27	2002	4109	478	162	157	159
28	2004	13443	2001	664	660	677
29	2006	4382	526	176	175	175
30	2008	13602	2048	671	667	710
31	2010	4402	520	174	175	171
32	2012	12620	1941	632	630	679

Next we are going to dive into the medals awarded for the olympians and I will divide it into total medals, summer medals and winter medals. I am interesting in seeing how the medals are divided out, is there a significant area where males are being awarded the most medals? and vice versa is there a specific area where females are dominating medals.

Check the amount of medals won as a % this time

```

In [206]: #Shows the % of medals won compared to the total participants e.g
# 2024 - 14.77% of participants received a medal
# 2020 - 4.65% of participants received a gold medal
# 2020 - 4.78% of participants received a Silver medal
# 2020 - 6.13% of participants received a Bronze medal
total_ratio_medals = pysqlf("""SELECT Year,
CAST(Total_Medals AS FLOAT) * 100 / Participants AS Total_Medals,
CAST(Gold_Medals AS FLOAT) * 100 / Participants AS Gold_Medals,
CAST(Silver_Medals AS FLOAT) * 100 / Participants AS Silver_Medals,
CAST(Bronze_Medals AS FLOAT) * 100 / Participants AS Bronze_Medals
FROM
(
    SELECT Year,
    COUNT(*) Participants,
    SUM(CASE WHEN Medal IS NOT NULL THEN 1 ELSE 0 END) AS Total_Medals,
    SUM(CASE WHEN Medal = 'Gold' THEN 1 ELSE 0 END) AS Gold_Medals,
    SUM(CASE WHEN Medal = 'Silver' THEN 1 ELSE 0 END) AS Silver_Medals,
    SUM(CASE WHEN Medal = 'Bronze' THEN 1 ELSE 0 END) AS Bronze_Medals
    FROM athlete_events
    GROUP BY Year
)
""")

```

Check the amount of medals won as a % this time

```
total_ratio_medals.tail()
```

```
Out[262]:
```

	Year	Total_Medals	Gold_Medals	Silver_Medals	Bronze_Medals
26	2008	15.05609	4.933098	4.903691	5.219821
27	2010	11.812812	3.952749	3.975466	3.884598
28	2012	15.023220	4.891641	4.870161	5.255418
29	2014	12.206093	4.130035	4.027906	4.048252
30	2016	14.779369	4.858270	4.785213	5.135885

Calculating the % of medals won for the Summer Olympic Games

```
In [263]:
```

```
#Calculating the % of medals won for the Summer Olympic Games
summer_ratio_medals = pysqlf("""SELECT Year,
CAST(Total_Medals AS FLOAT) * 100 / Participants AS Total_Medals,
CAST(Gold_Medals AS FLOAT) * 100 / Participants AS Gold_Medals,
CAST(Silver_Medals AS FLOAT) * 100 / Participants AS Silver_Medals,
CAST(Bronze_Medals AS FLOAT) * 100 / Participants AS Bronze_Medals
FROM
(
SELECT Year,
COUNT(*) Participants,
SUM(CASE WHEN Medal IS NOT NULL THEN 1 ELSE 0 END) AS Total_Medals,
SUM(CASE WHEN Medal = 'Gold' THEN 1 ELSE 0 END) AS Gold_Medals,
SUM(CASE WHEN Medal = 'Silver' THEN 1 ELSE 0 END) AS Silver_Medals,
SUM(CASE WHEN Medal = 'Bronze' THEN 1 ELSE 0 END) AS Bronze_Medals
FROM athlete_events
GROUP BY Year)
Medals_Table
""")
```

Calculating the % of medals won for the Summer Olympic Games

```

SUM(CASE WHEN Medal IS NOT NULL THEN 1 ELSE 0 END) AS Total_Medals,
SUM(CASE WHEN Medal = "Gold" THEN 1 ELSE 0 END) AS Gold_Medals,
SUM(CASE WHEN Medal = "Silver" THEN 1 ELSE 0 END) AS Silver_Medals,
SUM(CASE WHEN Medal = "Bronze" THEN 1 ELSE 0 END) AS Bronze_Medals
FROM summer_games
GROUP BY Year) Medals_Table
""")
summer_ratio_medals.tail()

```

In [264]:

summer_ratio_medals.tail()

Out[264]:

	Year	Total_Medals	Gold_Medals	Silver_Medals	Bronze_Medals
24	2000	14.499574	4.797948	4.782577	4.920489
25	2004	14.886070	4.929374	4.909618	5.050789
26	2008	15.05609	4.930908	4.909991	5.219821
27	2012	15.022220	4.891641	4.870161	5.255418
28	2016	14.779369	4.858270	4.785213	5.135885

Calculating the % of medals won for the Winter Olympic Games

In [265]:

#Calculating the % of medals won for the Winter Olympic Games
winter_ratio_medals = pqsqlf("""SELECT Year,
CAST(Total_Medals AS FLOAT) * 100 / Participants AS Total_Medals,
CAST(Gold_Medals AS FLOAT) * 100 / Participants AS Gold_Medals,
CAST(Silver_Medals AS FLOAT) * 100 / Participants AS Silver_Medals,
CAST(Bronze_Medals AS FLOAT) * 100 / Participants AS Bronze_Medals
FROM
(
SELECT Year,
COUNT(*) Participants,
SUM(CASE WHEN Medal IS NOT NULL THEN 1 ELSE 0 END) AS Total_Medals,
SUM(CASE WHEN Medal = 'Gold' THEN 1 ELSE 0 END) AS Gold_Medals,
SUM(CASE WHEN Medal = 'Silver' THEN 1 ELSE 0 END) AS Silver_Medals,
SUM(CASE WHEN Medal = 'Bronze' THEN 1 ELSE 0 END) AS Bronze_Medals
FROM winter_games
GROUP BY Year)
Medals_Table
""")

In [266]:

winter_ratio_medals.tail()

Out[266]:

Year	Total_Medals	Gold_Medals	Silver_Medals	Bronze_Medals	
17	1988	12.295270	4.022191	4.022191	4.160888
18	2002	11.633001	3.942565	3.820881	3.869555
19	2006	12.003651	4.016431	3.993610	3.993610
20	2010	11.812812	3.952749	3.975466	3.884598
21	2014	12.206093	4.130035	4.027906	4.048252

In [267]:

#Visualizing the total amount of medals awarded throughout the years
plt.plot(total_ratio_medals.Year, total_ratio_medals.Total_Medals, marker='x', linewidth = 2, label='Total Medals')
plt.xlabel('Year')
plt.ylabel('Total Medals')
plt.title('Total Medals Awarded at the Olympics')
plt.legend(loc=1)
Text(0.5, 1.6, 'Total Medals Awarded at the Olympics')

Out[267]:

This is very interesting seeing how many medals were awarded in the early years compared to now. It does make sense though considering how many participants competed in the early years for example in 1900-1904 there were less than 2,000 people competing compared to current times we are looking at 13,000 olympians.

Similar to my initial research in my proposal there are aggressive spikes throughout the chart which relates to the switch between Summer Games and Winter Games every 2 years. To get a better view I will break this out into two separate charts. Summer and Winter to give us a better visual representation of the data.

Theres an interesting spike in medals around 1980 which I will need to look into further to understand what this spike relates to.

In [268]:

#Visualizing the total amount of medals awarded throughout the years for the Summer Olympic Games
plt.plot(summer_ratio_medals.Year, summer_ratio_medals.Total_Medals, marker='x', linewidth = 2, label='Total Medals')
plt.xlabel('Year')
plt.ylabel('Total Medals')
plt.title('Total Medals Awarded at the Summer Olympics')
plt.legend(loc=1)
Text(0.5, 1.6, 'Total Medals Awarded at the Summer Olympics')

Out[268]:

In [269]:

#Visualizing the total amount of medals awarded throughout the years for the Winter Olympic Games
plt.plot(winter_ratio_medals.Year, winter_ratio_medals.Total_Medals, marker='x', linewidth = 2, label='Total Medals')
plt.xlabel('Year')
plt.ylabel('Total Medals')
plt.title('Total Medals Awarded at the Winter Olympics')
plt.legend(loc=1)
Text(0.5, 1.6, 'Total Medals Awarded at the Winter Olympics')

Out[269]:

In [270]:

#Comparing the medals awarded from the Summer games to winter games
plt.plot(winter_ratio_medals.Year, winter_ratio_medals.Total_Medals, marker='x', linewidth = 2, label='Winter Medals', color='blue')
plt.plot(summer_ratio_medals.Year, summer_ratio_medals.Total_Medals, marker='x', linewidth = 2, label='Summer Medals', color='red')
plt.xlabel('Year')
plt.ylabel('Total Medals')
plt.title('Total Medals Awarded at both Summer/Winter Olympics')
plt.legend(loc=1)
Text(0.5, 1.6, 'Total Medals Awarded at both Summer/Winter Olympics')

Out[270]:

Lets take a closer look at the medal % awarded throughout the year using subplots

In [231]:

fig, ax = plt.subplots(3, figsize=(16,9))
ax[0].plot(total_ratio_medals.Year, total_ratio_medals.Gold_Medals, marker='x', color='yellow', linewidth=3, label='Gold Ratio')
ax[0].plot(total_ratio_medals.Year, total_ratio_medals.Silver_Medals, marker='x', color='grey', linewidth=3, label='Silver Ratio')
ax[0].plot(total_ratio_medals.Year, total_ratio_medals.Bronze_Medals, marker='x', color='brown', linewidth=3, label='Bronze Ratio')
ax[0].legend(loc=1)
ax[0].set_ylabel('Year')
ax[0].set_title('Total Medals Won')
ax[1].plot(summer_ratio_medals.Year, summer_ratio_medals.Gold_Medals, marker='x', color='yellow', linewidth=3, label='Gold Ratio')
ax[1].plot(summer_ratio_medals.Year, summer_ratio_medals.Silver_Medals, marker='x', color='grey', linewidth=3, label='Silver Ratio')
ax[1].plot(summer_ratio_medals.Year, summer_ratio_medals.Bronze_Medals, marker='x', color='brown', linewidth=3, label='Bronze Ratio')
ax[1].legend(loc=1)
ax[1].set_ylabel('Year')
ax[1].set_title('Total Summer Medals Won')
ax[2].plot(winter_ratio_medals.Year, winter_ratio_medals.Gold_Medals, marker='x', color='yellow', linewidth=3, label='Gold Ratio')
ax[2].plot(winter_ratio_medals.Year, winter_ratio_medals.Silver_Medals, marker='x', color='grey', linewidth=3, label='Silver Ratio')
ax[2].plot(winter_ratio_medals.Year, winter_ratio_medals.Bronze_Medals, marker='x', color='brown', linewidth=3, label='Bronze Ratio')
ax[2].legend(loc=1)
ax[2].set_ylabel('Year')
ax[2].set_title('Total Winter Medals Won')
ax[2].tight_layout()

Out[231]:

The total count of males vs females has started to level out in recent years but in the early years it was a huge discrepancy. The total participants for the past 120 years included 74,522 (27.48%) females and a dominating 196,594 (72.51%) being males.

As you can see from the graph below the difference is starting to improve and hopefully in the coming years it will be around 50/50.

In [234]:

plt.plot(male.Year, male.Male, marker='x', color='blue', linewidth=2, label='Male')
plt.plot(female.Year, female.Female, marker='x', color='purple', linewidth=2, label='Female')
plt.legend(loc=1)
plt.xlabel('Year')
plt.ylabel('M vs F at Olympics')
plt.title('Male vs Female Olympics Participation')
Text(0.5, 1.6, 'Male vs Female Olympics Participation')

Out[234]:

The third key point I noticed was how in the early years there were different amounts of medals distributed for each class (gold, silver, bronze). Eventually this leveled out and it looks like the same amount of medals get distributed for gold, silver and bronze now but its interesting to see the difference compared to how it started.

In []:

Submit 2-3 key points you may have discovered about the data, e.g. new relationships? Aha's! Did you come up with additional ideas for other things to review?

I discovered that there was a significant drop in the amount of medals awarded per the percentage of participants. If we look at the table below we can see that in the early years there was over 30% of participants awarded a medal of some kind whereas if we look at todays figures its dropped by around 50% with the amount of participants being awarded a medal ranging from 10% - 15%.

This line chart displays the total number of medals won in the Summer Olympics from 1900 to 2010. The Y-axis represents the number of medals, ranging from 0 to 15.0. The X-axis represents the year. Three data series are plotted: Gold Ratio (yellow line), Silver Ratio (black line), and Bronze Ratio (red line). The Gold Ratio starts at approximately 15.5 in 1900, drops to around 10.5 in 1904, and then fluctuates between 10 and 12 until 1920. After 1920, the Gold Ratio drops significantly, staying below 10 medals per ratio. The Silver Ratio starts at approximately 12.5 in 1900, drops to around 10.5 in 1904, and then fluctuates between 10 and 12 until 1920. After 1920, the Silver Ratio drops significantly, staying below 10 medals per ratio. The Bronze Ratio starts at approximately 10.5 in 1900, drops to around 8.5 in 1904, and then fluctuates between 8 and 10 until 1920. After 1920, the Bronze Ratio drops significantly, staying below 10 medals per ratio.

Year	Gold Ratio	Silver Ratio	Bronze Ratio
1900	15.5	12.5	10.5
1904	10.5	10.5	8.5
1908	11.5	11.5	9.5
1912	10.5	10.5	8.5
1916	11.5	11.5	9.5
1920	12.5	12.5	10.5
1924	10.5	10.5	8.5
1928	9.5	9.5	7.5
1932	9.5	9.5	7.5
1936	9.5	9.5	7.5
1940	9.5	9.5	7.5
1944	9.5	9.5	7.5
1948	9.5	9.5	7.5
1952	9.5	9.5	7.5
1956	9.5	9.5	7.5
1960	9.5	9.5	7.5
1964	9.5	9.5	7.5
1968	9.5	9.5	7.5
1972	9.5	9.5	7.5
1976	9.5	9.5	7.5
1980	9.5	9.5	7.5
1984	9.5	9.5	7.5
1988	9.5	9.5	7.5
1992	9.5	9.5	7.5
1996	9.5	9.5	7.5
2000	9.5	9.5	7.5
2004	9.5	9.5	7.5
2008	9.5	9.5	7.5
2010	9.5	9.5	7.5