

NOC winter_medals

0	USA	124
1	GER	114
2	NOR	92
3	CAN	92
4	RUS	83
5	AUT	73
6	NED	49
7	SWE	47
8	FRA	46
9	SUI	45

```
In [159]: #counting gold medals by nation time range 2002-2014
comp_winter_gold=pyqldf('select nationality,NOC, Sum(medals_won) as gold_winter_medals from comp_gold left join
SQL Syntax/Gold Medals Winter Games/
Select
nationality,NOC,
Sum(Medals_won) as winter_gold_medals
from (
Select
Year,
NOC_ID,
Sport_ID,
count(distinct Medal_ID) as medals_won
from results
where Medal_ID = 1
group by 1,2,3
order by Year
) as comp_winter_gold
left join nationality
on comp_winter_gold.NOC_ID=nationality.NOC_ID
Where Year in (2002, 2006, 2010, 2014)
group by 1
order by Sum(medals_won) desc limit 10;
```

```
In [160]: #gold medal count winter games
comp_winter_gold
Out[160]:
NOC gold_winter_medals
0 GER 41
1 CAN 37
2 USA 38
3 NOR 34
4 RUS 29
5 SUI 20
6 AUT 20
7 NED 18
8 KOR 17
9 SWE 14
```

But in the winter games, we can see a different trend, the USA was number one in the overall Medal count, but on the gold medal level, they were only number 3.

One of the goals of the study was to figure out who was the most successful nation in the period of 2004-2016 (summer games) and if China was a runner-up. The data shows that the US is, in this period, the most successful nation as well, and the biggest competitor for the US in the Olympic summer games is indeed China. Additionally, I wanted to understand why China was a runner-up. Is it because they were less successful in the same events/disciplines as the US and, therefore, automatically runner-up, or did they succeed in fewer events than the US? Were these events different to the US?

To make a fair comparison, I used a time filter of the last four Olympic summer games to see which events/disciplines were the most successful for the US and China.

China's most successful events/disciplines vs US's most successful events/disciplines

China Medal 2004-2016

```
In [161]: china=pyqldf('Select results.Year, sports.Sport, sports.Sport_ID, results.NOC_ID, results.Medal_ID, count(distinct
In [162]: #China's medal count and most successful events between 2004-2016
usa_m_04_16_dis=pyqldf('Select Sport, (sum(gold_medal)+sum(silver_medal)+sum(bronze_medal)) as total_medals, sum(g
SQL Syntax/China Medal 2004-2016/
Select
Sport,
(sum(gold_medal)+sum(silver_medal)+sum(bronze_medal)) as total_medals,
sum(gold_medal) as number_of_gold_medals,
sum(silver_medal) as number_of_silver_medals,
sum(bronze_medal) as number_of_bronze_medals
from (
Select
results.Year,
sports.Sport,
sports.Sport_ID,
results.NOC_ID,
results.Medal_ID,
count(distinct case when results.Medal_ID=1 then 1 else null end) as gold_medal,
count(distinct case when results.Medal_ID=2 then 1 else null end) as silver_medal,
count(distinct case when results.Medal_ID=3 then 1 else null end) as bronze_medal
from results
left join sports
on results.Sport_ID=sports.Sport_ID
where NOC_ID=42
and results.Season_ID != 2
and results.Year in (2004, 2008, 2012, 2016)
group by 1,2,3,4,5
order by results.Year
) as china_m
where Medal_ID != 4
group by 1
order by (sum(gold_medal)+sum(silver_medal)+sum(bronze_medal)) desc;
```

china_dis_m.head(10)
Sport total_medals number_of_gold_medals number_of_silver_medals number_of_bronze_medals
0 Table Tennis Men's Singles 9 3 4 2
1 Table Tennis Women's Singles 8 4 3 1
2 Diving Women's Springboard 8 4 3 1
3 Trampoline Men's Individual 6 2 1 3
4 Diving Women's Platform 6 3 2 1
5 Diving Men's Springboard 6 3 1 2
6 Badminton Women's Singles 6 3 2 1
7 Trampoline Women's Individual 5 1 1 3
8 Shooting Women's Air Rifle, 10 metres 5 2 1 2
9 Diving Men's Platform 5 2 2 1

An overview in which disciplines China succeeded in the last four summer olympics

```
USA Medal 2004-2016
In [164]: usa_m_04_16=pyqldf('Select results.Year, sports.Sport, results.NOC_ID, results.Medal_ID, count(distinct case w
In [165]: #the US's medal count and most successful events between 2004-2016
usa_m_04_16_dis=pyqldf('Select Sport, (sum(gold_medal)+sum(silver_medal)+sum(bronze_medal)) as total_medals, s
SQL Syntax/USA Medal 2004-2016/
Select
Sport,
(sum(gold_medal)+sum(silver_medal)+sum(bronze_medal)) as total_medals,
sum(gold_medal) as number_of_gold_medals,
sum(silver_medal) as number_of_silver_medals,
sum(bronze_medal) as number_of_bronze_medals
from (
Select
results.Year,
sports.Sport,
sports.Sport_ID,
results.NOC_ID,
results.Medal_ID,
count(distinct case when results.Medal_ID=1 then 1 else null end) as gold_medal,
count(distinct case when results.Medal_ID=2 then 1 else null end) as silver_medal,
count(distinct case when results.Medal_ID=3 then 1 else null end) as bronze_medal
from results
left join sports
on results.Sport_ID=sports.Sport_ID
where NOC_ID=217
and results.Year in (2004, 2008, 2012, 2016)
and results.Season_ID != 2
group by 1,2,3,4
order by results.Year
) as usa_m_04_16
where Medal_ID != 4
group by 1
order by (sum(gold_medal)+sum(silver_medal)+sum(bronze_medal)) desc;
```

usa_m_04_16_dis.head(10)
Sport total_medals number_of_gold_medals number_of_silver_medals number_of_bronze_medals
0 Athletics Women's 100 metres Hurdles 8 3 2 3
1 Swimming Men's 200 metres Individual Medley 7 4 2 1
2 Swimming Men's 100 metres Backstroke 7 4 2 1
3 Athletics Men's 400 metres 7 2 2 3
4 Swimming Men's 400 metres Individual Medley 6 3 2 1
5 Swimming Men's 200 metres Backstroke 6 4 1 1
6 Gymnastics Women's Individual All-Around 6 4 2 0
7 Gymnastics Women's Balance Beam 6 1 3 2
8 Beach Volleyball Women's Beach Volleyball 6 3 1 2
9 Swimming Women's 400 metres Freestyle 5 1 2 2

Here we have the most successful event for the USA in the last four summer games.

The next step was to check if there were no disciplines they had in common.

```
In [167]: #merging with an inner join the US medal table with the chinese medal table to see what they have in common
in_comp=merge(usa_m_04_16_dis['Sport'], china_dis_m['Sport'], how='inner')
```

```
In [168]: #events/disciplines in common
in_comp
Out[168]:
Sport
0 Swimming Men's 200 metres Individual Medley
1 Swimming Men's 100 metres Backstroke
2 Gymnastics Women's Individual All-Around
3 Gymnastics Women's Balance Beam
4 Beach Volleyball Women's Beach Volleyball
...
68 Athletics Women's Shot Put
69 Athletics Women's Marathon
70 Athletics Women's Discus Throw
71 Athletics Women's 10,000 metres
72 Archery Men's Individual
73 rows x 1 columns
```

These are the disciplines the USA and China had in common in the last four summer games. Knowing that challenges the hypothesis already, China had 143 medal-winning events, 73 of which they had in common with the US. Therefore I checked the gold, silver and bronze medals to determine which Nation was more successful.

```
In [169]: #merging the US medal table with table of similar events with China
usa_vs_china=usa_m_04_16_dis.merge(in_comp, how='inner')
```

```
In [170]: #merging the chinese medal table with table of common events with the US
china_vs_usa=china_dis_m.merge(in_comp, how='inner')
```

```
In [171]: usa_vs_china.sort_values('Sport', inplace=True)
```

```
In [172]: china_vs_usa.sort_values('Sport', inplace=True)
```

```
In [173]: usa_vs_china.rename(columns={'total_medals':'total_medals_us'}, inplace=True)
```

I merged the US and China table with the table of the events they had in common and eventually added and edited the new columns so that I could compare the both nations with one another.

```
In [174]: #adding the medal count of each event that China and the US have in common to the US table
usa_vs_china['total_medals_us', 'number_of_gold_medals_us', 'number_of_silver_medals_us', 'number_of_bronze_medals_us', 'number_of_gold_medals', 'number_of_silver_medals', 'number_of_bronze_medals']
```

```
In [175]: #realign(reindex the columns)
usa_vs_china.head()
com_usa_china=usa_vs_china.reindex(columns=['Sport', 'total_medals_us', 'total_medals_china', 'number_of_gold_medals', 'number_of_silver_medals', 'number_of_bronze_medals', 'number_of_gold_medals_china', 'number_of_silver_medals_china', 'number_of_bronze_medals_china'])
```

```
In [176]: com_usa_china
Out[176]:
Sport total_medals_us total_medals_china number_of_gold_medals number_of_gold_medals_china number_of_silver_me
72 Archery Men's Individual 1 1 0 0
43 Archery Men's Team 2 1 0 0
7 Athletics Men's 110 metres Hurdles 5 3 1 0
17 Athletics Men's Triple Jump 4 2 2 2
71 Athletics Women's 10,000 metres 1 1 0 0
...
28 Tennis Women's Doubles 2 2 2 0
18 Volleyball Women's Volleyball 3 2 0 0
44 Weightlifting Women's Super-Heavyweight 1 1 0 0
27 Wrestling Women's Flyweight, Freestyle 2 2 0 1
26 Wrestling Women's Middleweight, Freestyle 2 2 0 1
73 rows x 9 columns
```

```
In [177]: #checking if China has any event where they have more medals total than the US
com_usa_china[com_usa_china.total_medals_us>com_usa_china.total_medals_china]
```

```
Out[177]: Sport total_medals_us total_medals_china number_of_gold_medals number_of_gold_medals_china number_of_silver_medals num
China had no event/discipline where they gathered more medals overall than the US.
```

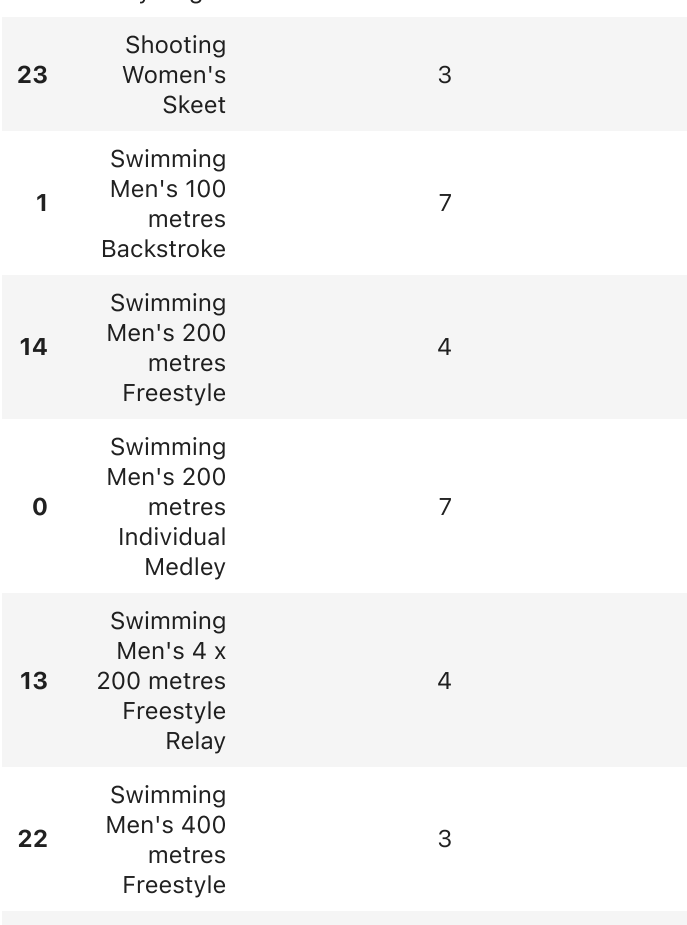
```
In [178]: #checking if the US has any event where they have more medals total than China
com_usa_china[com_usa_china.total_medals_us<com_usa_china.total_medals_china]
```

```
Out[178]: Sport total_medals_us total_medals_china number_of_gold_medals number_of_gold_medals_china number_of_silver_medals num
43 Archery Men's Team 2 1 0 0
7 Athletics Men's 110 metres Hurdles 5 3 1 0
17 Athletics Men's Triple Jump 4 2 2 2
4 Beach Volleyball Women's Beach Volleyball 6 4 3 4
42 Boxing Women's Middleweight 2 1 2 0
41 Diving Men's Platform 2 1 1 1
40 Diving Men's Synchronized Platform 2 1 0 0
39 Diving Men's Synchronized Springboard 2 1 0 0
6 Fencing Women's Sabre, Individual 5 3 2 2
38 Fencing Women's Sabre, Team 2 1 0 0
25 Gymnastics Men's Horizontal Bar 3 2 0 1
37 Gymnastics Men's Individual All-Around 2 1 1 1
36 Gymnastics Men's Team All-Around 2 1 0 0
3 Gymnastics Women's Balance Beam 6 4 1 4
24 Gymnastics Women's Horse Vault 3 2 1 0
2 Gymnastics Women's Individual All-Around 6 4 4 1
16 Gymnastics Women's Team All-Around 4 2 2 1
15 Gymnastics Women's Uneven Bars 4 3 0 2
35 Judo Women's Half-Heavyweight 2 1 2 0
23 Shooting Women's Skeet 3 2 1 2
1 Swimming Men's 100 metres Backstroke 7 5 4 2
14 Swimming Men's 200 metres Freestyle 4 3 1 2
0 Swimming Men's 200 metres Individual Medley 7 5 4 2
13 Swimming Women's 4 x 200 metres Freestyle Relay 4 3 4 2
22 Swimming Men's 400 metres Freestyle 3 2 0 0
5 Swimming Women's 100 metres Backstroke 5 3 3 3
12 Swimming Women's 100 metres Backstroke 4 3 1 2
21 Swimming Women's 100 metres Butterfly 3 2 1 1
20 Swimming Women's 100 metres Freestyle 3 2 1 0
19 Swimming Women's 200 metres Breaststroke 3 2 3 1
11 Swimming Women's 200 metres Individual Medley 4 3 0 1
10 Swimming Women's 4 x 100 metres Freestyle Relay 4 3 2 1
9 Swimming Women's 4 x 200 metres Freestyle Relay 4 3 3 2
8 Swimming Women's 400 metres Individual Medley 4 3 0 1
18 Volleyball Women's Volleyball 3 2 0 0
```

But lets see how each nation performed based on the medal type:

```
USA vs China Gold
In [179]: #checking if China has any event where they have more gold medals than the US
china_g=ien(com_usa_china[com_usa_china.number_of_gold_medals<com_usa_china.number_of_gold_medals_china])
#checking if the US has any event where they have more gold medals than China
us_g=ien(com_usa_china[com_usa_china.number_of_gold_medals>com_usa_china.number_of_gold_medals_china])
gold_diff=pd.DataFrame({'China>USA=Gold':china_g, 'USA>China=Gold':usa_g, index=['Gold']})
plt.xticks(rotation=0)
plt.title('USA vs China Gold', size=16, fontweight='bold')
plt.legend(loc='upper left')
```

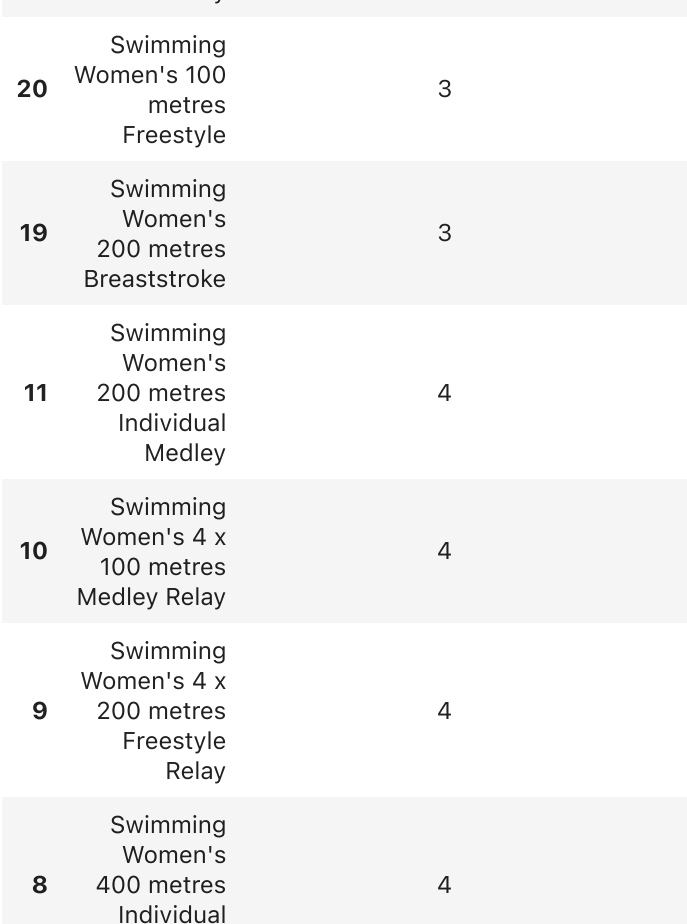
```
Out[179]: <matplotlib.legend.Legend at 0x7f4bec15f190>
```



USA vs China Silver

```
In [180]: #checking if China has any event where they have more silver medals than the US
china_s=ien(com_usa_china[com_usa_china.number_of_silver_medals<com_usa_china.number_of_silver_medals_china])
#checking if the US has any event where they have more silver medals than China
us_s=ien(com_usa_china[com_usa_china.number_of_silver_medals>com_usa_china.number_of_silver_medals_china])
silver_diff=pd.DataFrame({'China>USA=Silver':china_s, 'USA>China=Silver':usa_s, index=['Silver']})
plt.xticks(rotation=0)
plt.title('USA vs China Silver', size=16, fontweight='bold')
```

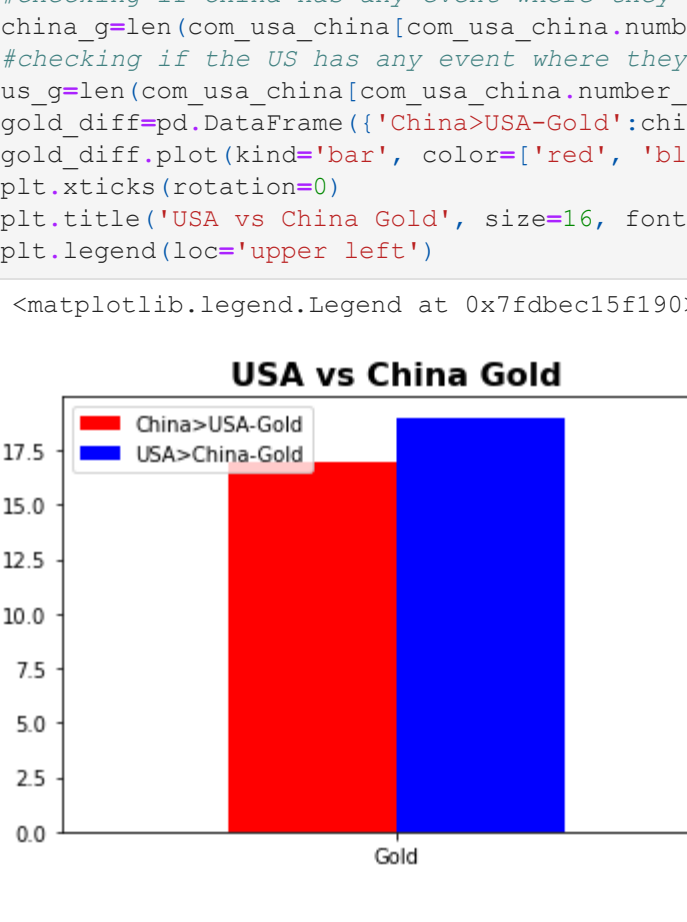
```
Out[180]: Text(0.5, 1.0, 'USA vs China Silver')
```



USA vs China Bronze

```
In [181]: #checking if China has any event where they have more bronze medals than the US
china_b=ien(com_usa_china[com_usa_china.number_of_bronze_medals<com_usa_china.number_of_bronze_medals_china])
#checking if the US has any event where they have more bronze medals than China
us_b=ien(com_usa_china[com_usa_china.number_of_bronze_medals>com_usa_china.number_of_bronze_medals_china])
bronze_diff=pd.DataFrame({'China>USA=Bronze':china_b, 'USA>China=Bronze':usa_b, index=['Bronze']})
plt.xticks(rotation=0)
plt.title('USA vs China Bronze', size=16, fontweight='bold')
```

```
Out[181]: Text(0.5, 1.0, 'USA vs China Bronze')
```



In summary, China had fewer events where they won Olympic medals (143=China; 173=USA), which would support the hypothesis that China was a runner-up because they won fewer medals. Additionally I have discovered that the USA and China had 73 events in common by analysing both profiles. Here, I could observe that the US always came out on top of China in all medal types. What does that mean for the hypothesis?

Hypothesis=> China was a runner-up because it succeed in fewer events than the most successful nation, but these events are unique to the ones of the most successful nation

All in all, China came short of the US and had more common disciplines than unique ones (51% in common with the US). Therefore the hypothesis is false and needs to be rejected.

Most Successful Athlete in Olympic History

```
In [182]: #getting a medal count of the 5 most successful athletes, on an overall medal level
best_ath=pyqldf('select athlete.Name, nationality,NOC, athlete.Sex, count(results.Medal_ID) as number_of_medal
In [183]: best_ath.set_index('Name', inplace=True)
```

```
SQL Syntax/Most Successful Athlete in Olympic History/
select
athlete.Name,
nationality,NOC,
athlete.Sex,
count(results.Medal_ID) as number_of_medals,
count(case when results.Medal_ID=1 then 1 else null end) as number_of_gold_medals,
count(case when results.Medal_ID=2 then 1 else null end) as number_of_silver_medals,
count(case when results.Medal_ID=3 then 1 else null end) as number_of_bronze_medals
from results
left join athlete
on results.Athlete_ID=athlete.Athlete_ID
left join nationality
on results.NOC_ID=nationality.NOC_ID
where results.Medal_ID in (1,2,3)
group by 1,2,3
order by count(results.Medal_ID) desc limit 5;
```

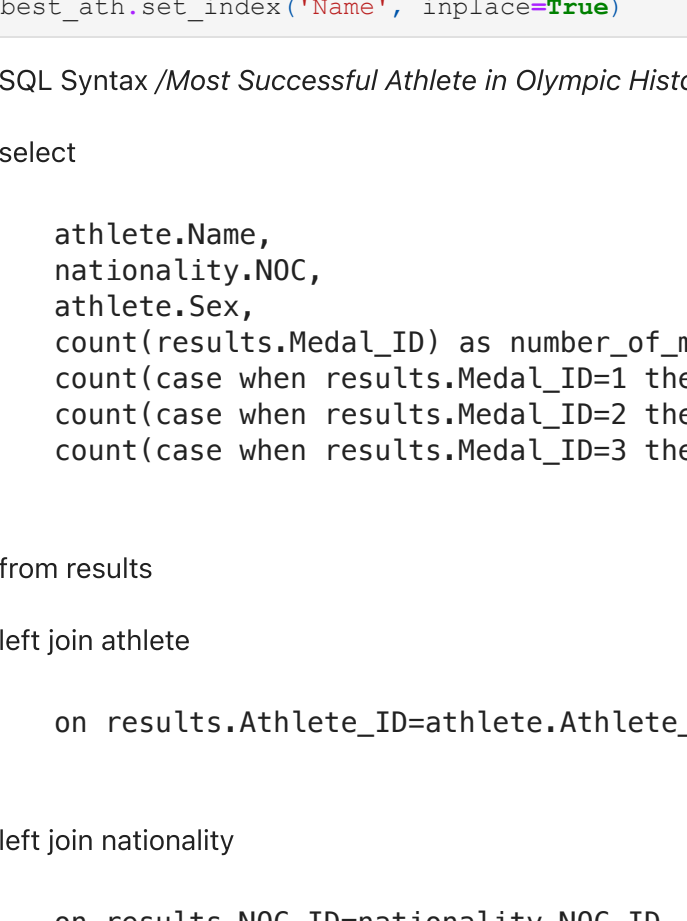
```
In [184]: #table of most successful athletes based on medal count
best_ath
Out[184]:
Name NOC Sex number_of_medals number_of_gold_medals number_of_silver_medals number_of_bronze_medals
Michael Fred Phelps, II USA M 28 23 3 2
Larysa Semenivna Latynina (Diriy-) URS F 18 9 5 4
Nikolay Yefimovich Andrianov URS M 15 7 5 3
Borys Anfinanyevych Shakhlin URS M 13 7 4 2
Edoardo Mangiarotti ITA M 13 6 5 2
```

```
In [185]: best_ath[['number_of_gold_medals', 'number_of_silver_medals', 'number_of_bronze_medals']].plot(kind='bar',
ath_per=pyqldf('select athlete.Name, nationality,NOC, athlete.Sex, count(results.Medal_ID) as number_of_medal
best_ath_gold
Out[186]:
Name NOC Sex number_of_gold_medals
0 Michael Fred Phelps, II USA M 23
1 Raymond Clarence "Ray" Ewry USA M 10
2 Frederick Carlton "Carl" Lewis USA M 9
3 Larysa Semenivna Latynina (Diriy-) URS F 9
4 Mark Andrew Spitz USA M 9
```

```
SQL Syntax/Most successful Gold winning athlete/
select
athlete.Name,
nationality,NOC,
athlete.Sex,
count(results.Medal_ID) as number_of_gold_medals
from results
left join athlete
on results.Athlete_ID=athlete.Athlete_ID
left join nationality
on results.NOC_ID=nationality.NOC_ID
where results.Medal_ID=1
group by 1,2,3
order by count(results.Medal_ID) desc limit 5;
```

```
In [187]: sns.barplot(x=Name, y=number_of_gold_medals, data=best_ath_gold)
plt.xticks(rotation=45)
plt.title('Most Successful Athletes-Gold Medals', size=16, fontweight='bold')
```

```
Out[187]: Text(0.5, 1.0, 'Most Successful Athletes-Gold Medals')
```



Michael Phelps is the most successful athlete in the gold medal segment.

Conversion Rate of most Successful Athletes and Medals won

```
In [188]: #CVR overall Medals of the most successful athletes
ath_per=pyqldf('Select athlete.Name, nationality,NOC, athlete.Sex, count(case when Medal_ID!=4 then 1 else null
SQL Syntax/CVR overall Medals of the most successful athletes/
Select
athlete.Name,
nationality,NOC,
athlete.Sex,
count(case when Medal_ID=4 then 1 else null end) as won_medal,
round(count(case when Medal_ID!=4 then 1 else null end)/count(Medal_ID),2) as CVR_Overall
from results
left join athlete
on results.Athlete_ID=athlete.Athlete_ID
left join nationality
on results.NOC_ID=nationality.NOC_ID
where athlete.Athlete_ID in (86767, 73148, 93760, 16196, 29505)
and results.Sport_ID in (645, 640, 638, 632, 646, 644, 645, 652, 639, 591, 589, 592, 590, 376, 377, 369, 366, 374, 368, 365, 380, 389, 388, 387, 385, 384, 386, 390, 326, 327, 333, 336, 325)
group by 1,2,3;
```

```
In [189]: ath_per.head()
Out[189]:
Name NOC Sex won_medal participation
0 Borys Anfinanyevych Shakhlin URS M 13 24
1 Edoardo Mangiarotti ITA M 13 14
2 Larysa Semenivna Latynina (Diriy-) URS F 18 19
3 Michael Fred Phelps, II USA M 28 30
4 Nikolay Yefimovich Andrianov URS M 15 24
```



```
#adding CVR rate of every athlete
ath_per['CVR_Overall'] = round(ath_per.won_medal/ath_per.participation,4)

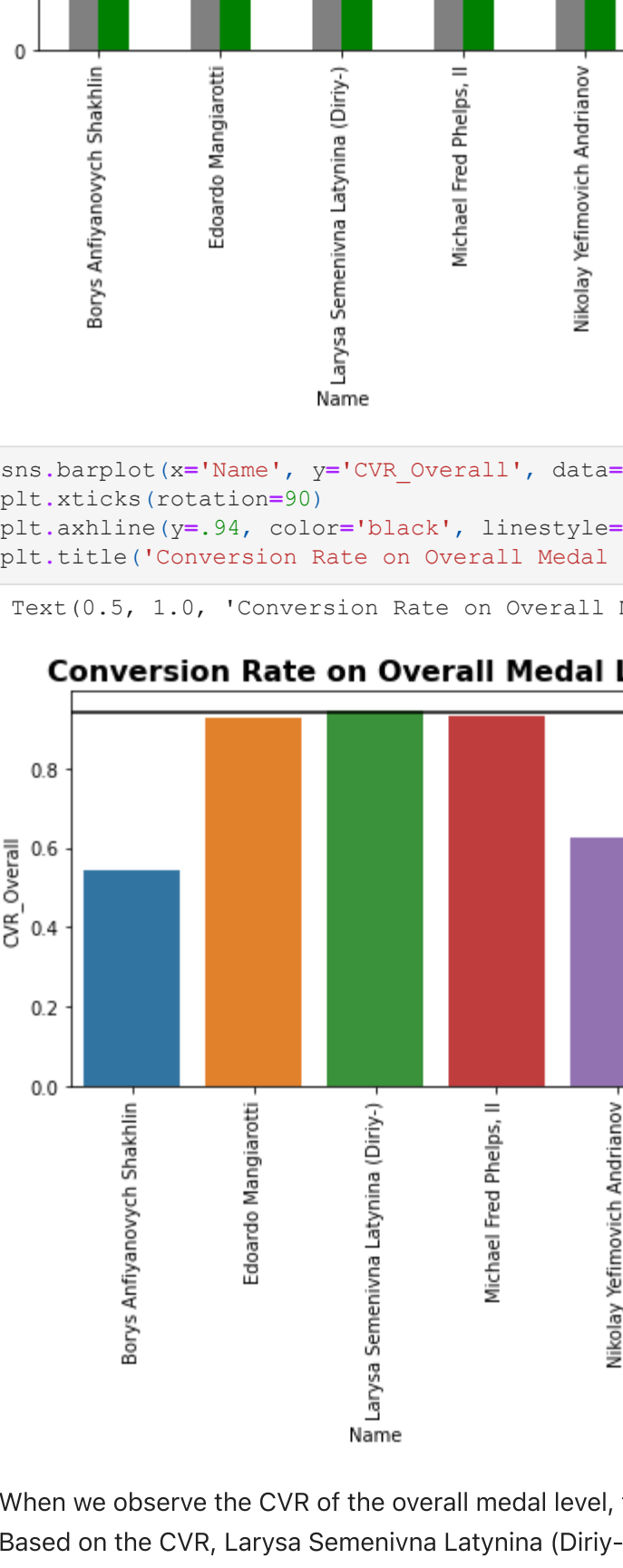
In [191]:
ath_per.set_index('NOC', inplace=True)

In [192]:
#most successful athletes based on medals and CVR
ath_per
```

NOC Sex won_medal participation CVR_Overall						
Borys Anfiyanovych Shakhlin	URS	M	13	24	0.547	
Edoardo Mangiarotti	ITA	M	13	14	0.9286	
Larysa Semenivna Latynina (Diriy-)	URS	F	18	19	0.9474	
Michael Fred Phelps, II	USA	M	28	30	0.9333	
Nikolay Yefimovich Andrianov	URS	M	15	24	0.6250	

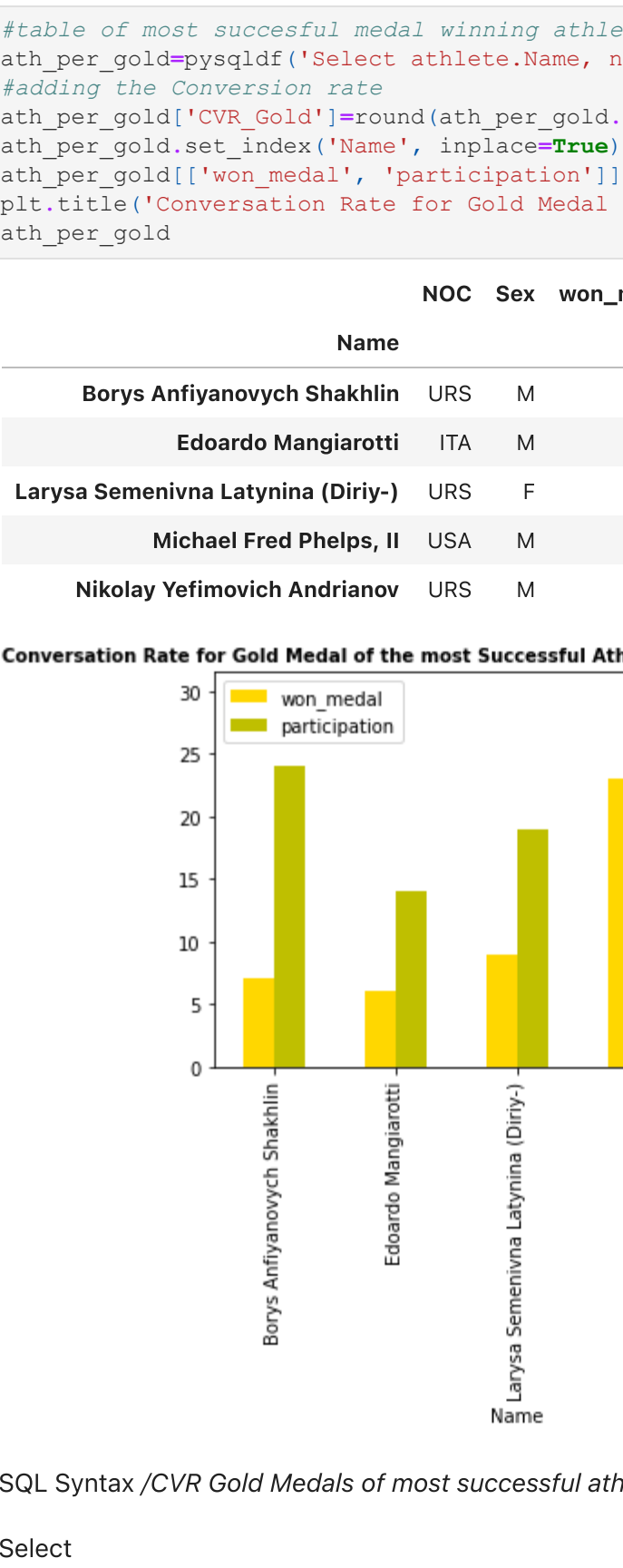
```
In [193]:
ath_per[['participation', 'won_medal']].plot(kind='bar', color=['grey', 'green'])
plt.title('Conversion Rate from Events to Medals', size=16, fontweight='bold')

Out[193]:
Text(0.5, 1.0, 'Conversion Rate from Events to Medals')
```



```
In [194]:
sns.barplot(x='Name', y='CVR_Overall', data=ath_per.reset_index())
plt.xticks(rotation=90)
plt.xlabel(yw=34, color='black', linestyle='-', fontweight='bold')
plt.title('Conversion Rate on Overall Medal Level', size=16, fontweight='bold')

Out[194]:
Text(0.5, 1.0, 'Conversion Rate on Overall Medal Level')
```



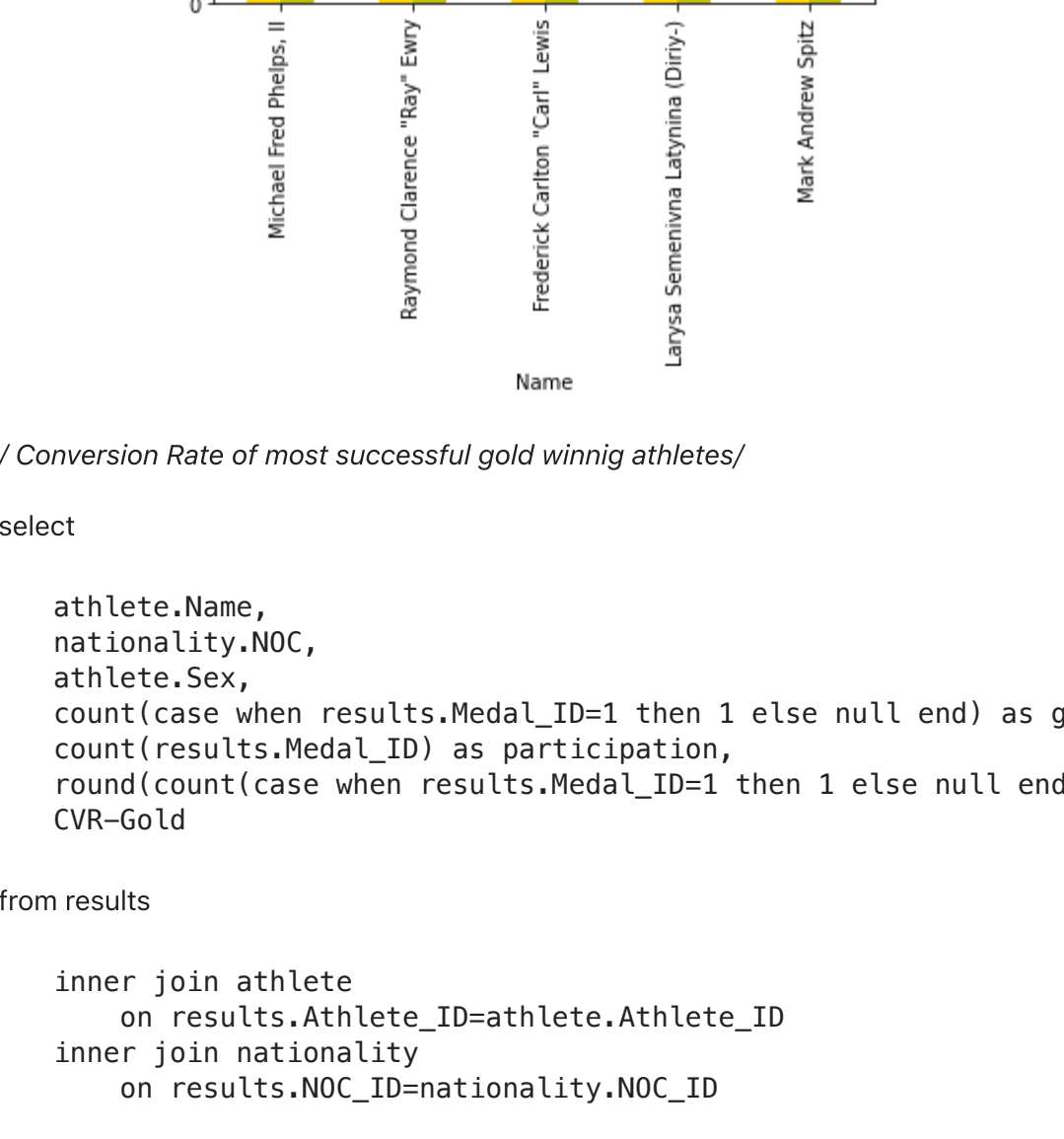
When we observe the CVR of the overall medal level, then we can see that Micheal Phelps is not the most successful athlete. Based on the CVR, Larysa Semenivna Latynina (Diriy-) is the most successful athlete.

Conversion Rate of most Successful Athletes and Gold Medals won

Based on the overall medal count, Michael Phelps is the most successful athlete. But looking at the conversion rate a different picture is observable. In the overall medal conversion rate ranking, Michael Phelps is just a runner-up to Larysa Semenivna Latynina (Diriy-).

```
In [195]:
#table of most successful medal winning athletes and gold medal CVR
ath_per_gold=pyqidf('Select athlete_Name, nationality,NOC, athlete.Sex, count(case when Medal_ID=1 then 1 else 0) as won_medal, count(Medal_ID) as participation')
ath_per_gold.set_index('Name', inplace=True)
ath_per_gold['CVR_Gold']=round(ath_per_gold.won_medal/ath_per_gold.participation, 2)
ath_per_gold[['won_medal', 'participation']].plot(kind='bar', color=['gold', 'y'])
plt.title('Conversion Rate for Gold Medal of the most Successful Athletes based on the Overall Medal Count', fontweight='bold')
ath_per_gold
```

NOC Sex won_medal participation CVR_Gold						
Borys Anfiyanovych Shakhlin	URS	M	7	24	0.29	
Edoardo Mangiarotti	ITA	M	6	14	0.43	
Larysa Semenivna Latynina (Diriy-)	URS	F	9	19	0.47	
Michael Fred Phelps, II	USA	M	23	30	0.77	
Nikolay Yefimovich Andrianov	URS	M	7	24	0.29	



SQL Syntax /CVR Gold Medals of most successful athletes/

```
Select
athlete.Name,
nationality,NOC,
athlete.Sex,
count(case when Medal_ID=1 then 1 else null end) as won_medal,
count(Medal_ID) as participation
round(count(case when Medal_ID=1 then 1 else null end)/count(Medal_ID),2) as CVR-Gold

from results

left join athlete
on results.Athlete_ID=athlete.Athlete_ID
left join nationality
on results.NOC_ID=nationality.NOC_ID

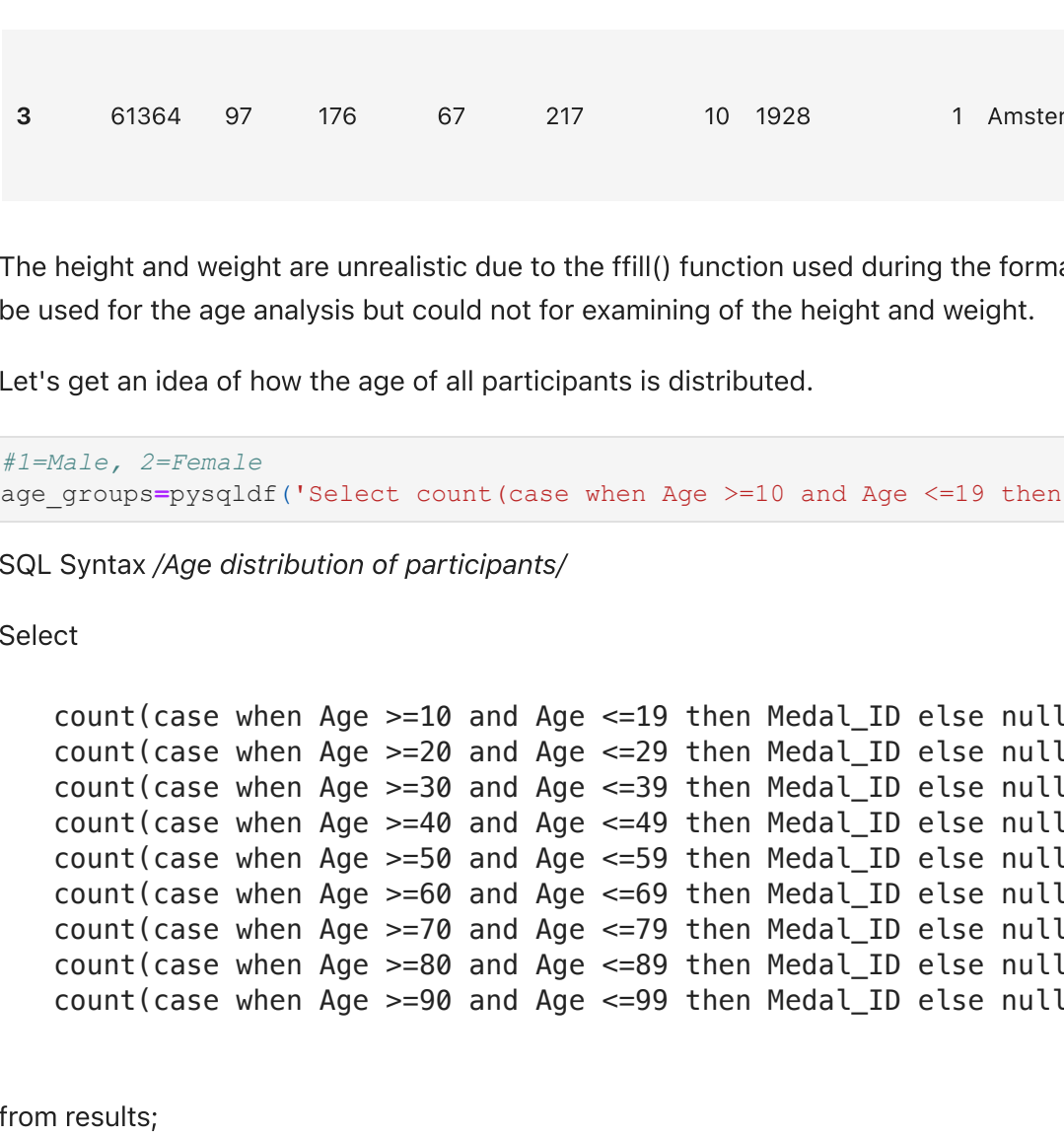
where athlete.Athlete_ID in (86767, 73148, 93760, 16196, 29505)

and results.Sport_ID in (645, 640, 638, 632, 646, 644, 645, 652, 639, 591, 589, 592, 576, 377, 369, 366, 374, 368, 365, 380, 389, 388, 387, 385, 384, 386, 390, 326, 327, 333, 336, 325)

group by 1,2,3;
```

```
In [196]:
#CVR of gold medals, between the most successful gold winning athletes
ath_per_gold=pyqidf('Select athlete_Name, nationality,NOC, athlete.Sex, count(case when results.Medal_ID=1 then 1 else 0) as gold_medals, count(case when results.Medal_ID=2 then 1 else 0) as participation, 2)
ath_per_gold[['CVR_Gold']=round(ath_per_gold.g_gold_medals/ath_per_gold_g.participation, 2)
ath_per_gold[['gold_medals', 'participation']].plot(kind='bar', color=['gold', 'y'])
plt.title('Conversion Rate for Gold Medals of the most Successful Gold Winning Athletes', size=12, fontweight='bold')
ath_per_gold
```

NOC Sex gold_medals participation CVR_Gold						
Michael Fred Phelps, II	USA	M	23	30	0.77	
Frederick Clarence "Carl" Lewis	USA	M	10	10	1.00	
Raymond Clarence "Ray" Ewry	USA	M	9	10	0.90	
Larysa Semenivna Latynina (Diriy-)	URS	F	9	19	0.47	
Mark Andrew Spitz	USA	M	9	12	0.75	



/ Conversion Rate of most successful gold winning athletes/

```
select
athlete.Name,
nationality,NOC,
athlete.Sex,
count(case when results.Medal_ID=1 then 1 else null end) as gold_medals,
count(results.Medal_ID) as participation,
round(count(case when results.Medal_ID=1 then 1 else null end)/count(results.Medal_ID),2) as CVR-Gold

from results

inner join athlete
on results.Athlete_ID=athlete.Athlete_ID
inner join nationality
on results.NOC_ID=nationality.NOC_ID

group by 1,2,3;
```

having count(case when results.Medal_ID=1 then 1 else null end) >0

order by count(case when results.Medal_ID=1 then 1 else null end) desc limit 5

Comparing both gold medal CVR tables, we see two different pictures. We can see in the graph of the athletes with the most overall medals that Michael Phelps has the most gold medals and the best CVR of 77%.

But, when we look at the gold medal CVR of the most successful gold medal athletes, we see that Michael Phelps has the most gold medals but the third best CVR, Raymond Clarence "Ray" Ewry (100%) and Frederick Carlton "Carl" Lewis (90%). At this point, it just depends on the interpretation of these numbers and values. We wanted to see who had the highest count of medals. Therefore Michael Phelps is considered as the most successful athlete.

Hypothesis => The most successful athletes are men

- A men will hold the most gold medals

The Hypothesis is right because Michael Phelps has the highest overall medal and gold medal count

What Role does Age play in Winning a Olympic Medal?

An essential question I wanted to answer in this case study was, if there is a correlation between age and winning an Olympic medal. First, we need to understand what age range we are dealing with.

```
In [197]:
#youngest and oldest participants
pyqidf('select min(Age), max(Age) from results;')

Out[197]:
min(Age) max(Age)
0 10 97
```

The youngest olympian participant was 10 years old, and the oldest 97. First, the dataset seems faulty. Therefore we check what discipline and what year the participants participated.

```
In [198]:
#checking if the weird age range is realistic
pyqidf('select *, sports,Sport from results left join sports on results.Sport_ID=sports.Sport_ID where results.Athlete_ID=1')

Out[198]:
Athlete_ID Age Height Weight NOC_ID Games_ID Year Season_ID City Sport_ID Medal_ID Sex_ID Sport_ID Gym
0 27349 10 176 67 82 1 1896 1 Athens 375 3 1 375 Parallel
```

1	40797	10	176	67	82	1	1896	1	Athina	633	4	1	633	Sw
2	44695	10	176	67	91	1	1896	1	Athina	376	4	1	376	Gym
3	61364	97	176	67	217	10	1928	1	Amsterdam	68	4	1	68	Comp

The height and weight are unrealistic due to the ffill() function used during the formatting phase. Nonetheless, these values could be used for the age analysis but could not for examining of the height and weight.

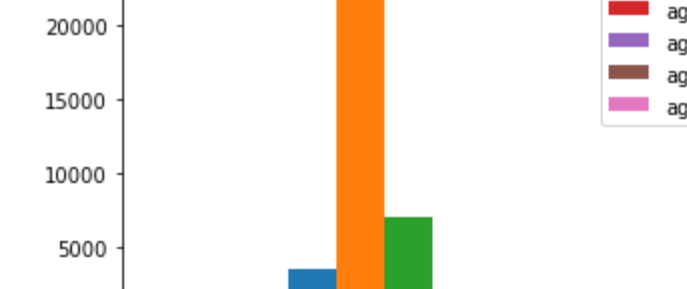
Let's get an idea of how the age of all participants is distributed.

```
In [199]:
#Male, 2=Female
age_groups=pyqidf('Select count(case when Age >=10 and Age <=19 then Medal_ID else null end) as age_10_19, count(case when Age >=20 and Age <=29 then Medal_ID else null end) as age_20_29, count(case when Age >=30 and Age <=39 then Medal_ID else null end) as age_30_39, count(case when Age >=40 and Age <=49 then Medal_ID else null end) as age_40_49, count(case when Age >=50 and Age <=59 then Medal_ID else null end) as age_50_59, count(case when Age >=60 and Age <=69 then Medal_ID else null end) as age_60_69, count(case when Age >=70 and Age <=79 then Medal_ID else null end) as age_70_79, count(case when Age >=80 and Age <=89 then Medal_ID else null end) as age_80_89, count(case when Age >=90 and Age <=99 then Medal_ID else null end) as age_90_99')

from results;
```

age_10_19 age_20_29 age_30_39 age_40_49 age_50_59 age_60_69 age_70_79 age_80_89 age_90_99										
0	33228	183537	44541	7305	1742	602	136	9	2	

```
In [200]:
#Histogram
results.Age.hist()
plt.axline(x=23, color='red')
plt.axline(x=25, color='black')
plt.axline(x=25.628709, color='yellow')
```



age_10_19 age_20_29 age_30_39 age_40_49 age_50_59 age_60_69 age_70_79 age_80_89 age_90_99										
0	33228	183537	44541	7305	1742	602	136	9	2	

As we can see, the Age distribution is right-skewed. This means most of the data is behind the mean, indicating that most participants are below 26 years old.

Next, we wanted to understand the age profile of medal-winning athletes.

```
In [205]:
#age range of medal winning athletes
pyqidf('select min(Age), max(Age) from results where Medal_ID in (1,2,3);')

Out[205]:
min(Age) max(Age)
0 10 73
```

```
In [206]:
#Male, 2=Female
age_groups_medal=pyqidf('Select count(case when Age >=10 and Age <=19 then Medal_ID else null end) as age_10_19, count(case when Age >=20 and Age <=29 then Medal_ID else null end) as age_20_29, count(case when Age >=30 and Age <=39 then Medal_ID else null end) as age_30_39, count(case when Age >=40 and Age <=49 then Medal_ID else null end) as age_40_49, count(case when Age >=50 and Age <=59 then Medal_ID else null end) as age_50_59, count(case when Age >=60 and Age <=69 then Medal_ID else null end) as age_60_69, count(case when Age >=70 and Age <=79 then Medal_ID else null end) as age_70_79, count(case when Age >=80 and Age <=89 then Medal_ID else null end) as age_80_89, count(case when Age >=90 and Age <=99 then Medal_ID else null end) as age_90_99')

SQL Syntax /Age Groups of Medal winning athletes/

Select

count(case when Age >=10 and Age <=19 then Medal_ID else null end) as age_10_19,
count(case when Age >=20 and Age <=29 then Medal_ID else null end) as age_20_29,
count(case when Age >=30 and Age <=39 then Medal_ID else null end) as age_30_39,
count(case when Age >=40 and Age <=49 then Medal_ID else null end) as age_40_49,
count(case when Age >=50 and Age <=59 then Medal_ID else null end) as age_50_59,
count(case when Age >=60 and Age <=69 then Medal_ID else null end) as age_60_69,
count(case when Age >=70 and Age <=79 then Medal_ID else null end) as age_70_79,
count(case when Age >=80 and Age <=89 then Medal_ID else null end) as age_80_89,
count(case when Age >=90 and Age <=99 then Medal_ID else null end) as age_90_99

from results;

where Medal_ID in (1,2,3);
```

age_10_19 age_20_29 age_30_39 age_40_49 age_50_59 age_60_69 age_70_79 age_80_89 age_90_99										
all_athletes	33228	183537	44541	7305	1742	602	136	9	2	

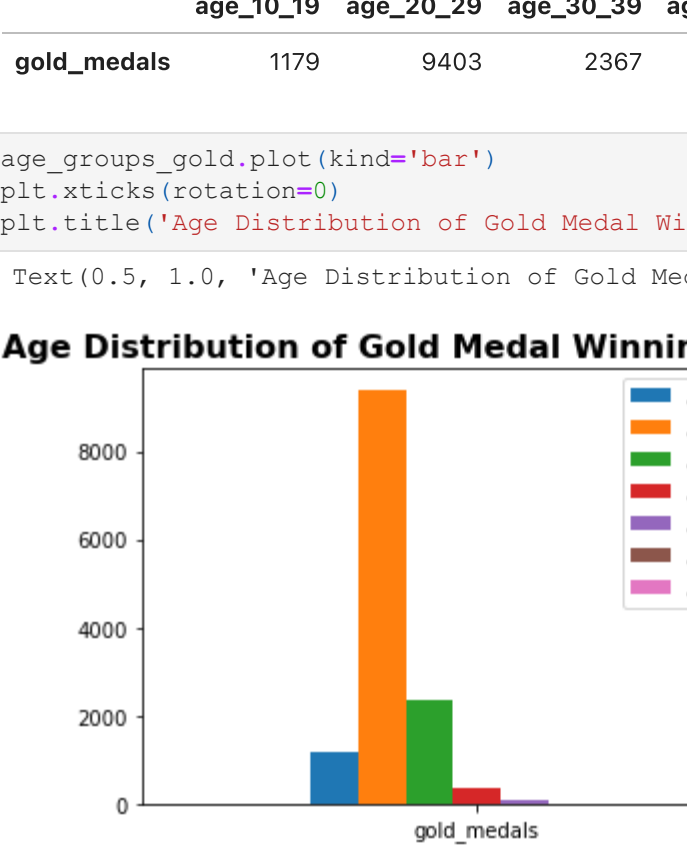
age_groups_medal.index=['medal_winning_athletes']

age_distribution of non-medal winning athletes

age_10_19 age_20_29 age_30_39 age_40_49 age_50_59 age_60_69 age_70_79 age_80_89 age_90_99										
medal_winning_athletes	3481	28002	7003	1039	209	44	5			

```
In [211]:
age_groups.plot(kind='bar')
plt.xticks(rotation=0)
plt.title('Distribution of Age for All Participants', size=16, fontweight='bold')
age_groups_medal.plot(kind='bar')
plt.xticks(rotation=0)
plt.title('Distribution of Age for Medal Winning Participants', size=16, fontweight='bold')

Out[211]:
Text(0.5, 1.0, 'Distribution of Age for Medal Winning Participants')
```



When comparing both graphs, we notice that they look the same when it comes to the distribution of age. The only difference is that there are no medal-winning athletes in their 80s and 90s.

By analyzing the medal count of the age distribution, you could say you are most likely to win a medal in your twenties. But, we would see a similar picture regarding gold, silver and bronze medals?

Gold Medal

```
In [212]:
#age distribution based on gold winning athletes
age_groups_gold=pyqidf('Select count(case when Age >=10 and Age <=19 then Medal_ID else null end) as age_10_19, count(case when Age >=20 and Age <=29 then Medal_ID else null end) as age_20_29, count(case when Age >=30 and Age <=39 then Medal_ID else null end) as age_30_39, count(case when Age >=40 and Age <=49 then Medal_ID else null end) as age_40_49, count(case when Age >=50 and Age <=59 then Medal_ID else null end) as age_50_59, count(case when Age >=60 and Age <=69 then Medal_ID else null end) as age_60_69, count(case when Age >=70 and Age <=79 then Medal_ID else null end) as age_70_79, count(case when Age >=80 and Age <=89 then Medal_ID else null end) as age_80_89, count(case when Age >=90 and Age <=99 then Medal_ID else null end) as age_90_99')

SQL Syntax /Age group distribution of Gold Medals/

Select

count(case when Age >=10 and Age <=19 then Medal_ID else null end) as age_10_19,
count(case when Age >=20 and Age <=29 then Medal_ID else null end) as age_20_29,
count(case when Age >=30 and Age <=39 then Medal_ID else null end) as age_30_39,
count(case when Age >=40 and Age <=49 then Medal_ID else null end) as age_40_49,
count(case when Age >=50 and Age <=59 then Medal_ID else null end) as age_50_59,
count(case when Age >=60 and Age <=69 then Medal_ID else null end) as age_60_69,
count(case when Age >=70 and Age <=79 then Medal_ID else null end) as age_70_79,
count(case when Age >=80 and Age <=89 then Medal_ID else null end) as age_80_89,
count(case when Age >=90 and Age <=99 then Medal_ID else null end) as age_90_99

from results;

where Medal_ID=1;
```

SQL Syntax /Age group distribution of Silver Medals/

```
Select

count(case when Age >=10 and Age <=19 then Medal_ID else null end) as age_10_19,
count(case when Age >=20 and Age <=29 then Medal_ID else null end) as age_20_29,
count(case when Age >=30 and Age <=39 then Medal_ID else null end) as age_30_39,
count(case when Age >=40 and Age <=49 then Medal_ID else null end) as age_40_49,
count(case when Age >=50 and Age <=59 then Medal_ID else null end) as age_50_59,
count(case when Age >=60 and Age <=69 then Medal_ID else null end) as age_60_69,
count(case when Age >=70 and Age <=79 then Medal_ID else null end) as age_70_79,
count(case when Age >=80 and Age <=89 then Medal_ID else null end) as age_80_89,
count(case when Age >=90 and Age <=99 then Medal_ID else null end) as age_90_99

from results;

where Medal_ID=2;
```

SQL Syntax /Age group distribution of Bronze Medals/

```
Select

count(case when Age >=10 and Age <=19 then Medal_ID else null end) as age_10_19,
count(case when Age >=20 and Age <=29 then Medal_ID else null end) as age_20_29,
count(case when Age >=30 and Age <=39 then Medal_ID else null end) as age_30_39,
count(case when Age >=40 and Age <=49 then Medal_ID else null end) as age_40_49,
count(case when Age >=50 and Age <=59 then Medal_ID else null end) as age_50_59,
count(case when Age >=60 and Age <=69 then Medal_ID else null end) as age_60_69,
count(case when Age >=70 and Age <=79 then Medal_ID else null end) as age_70_79,
count(case when Age >=80 and Age <=89 then Medal_ID else null end) as age_80_89,
count(case when Age >=90 and Age <=99 then Medal_ID else null end) as age_90_99

from results;

where Medal_ID=3;
```

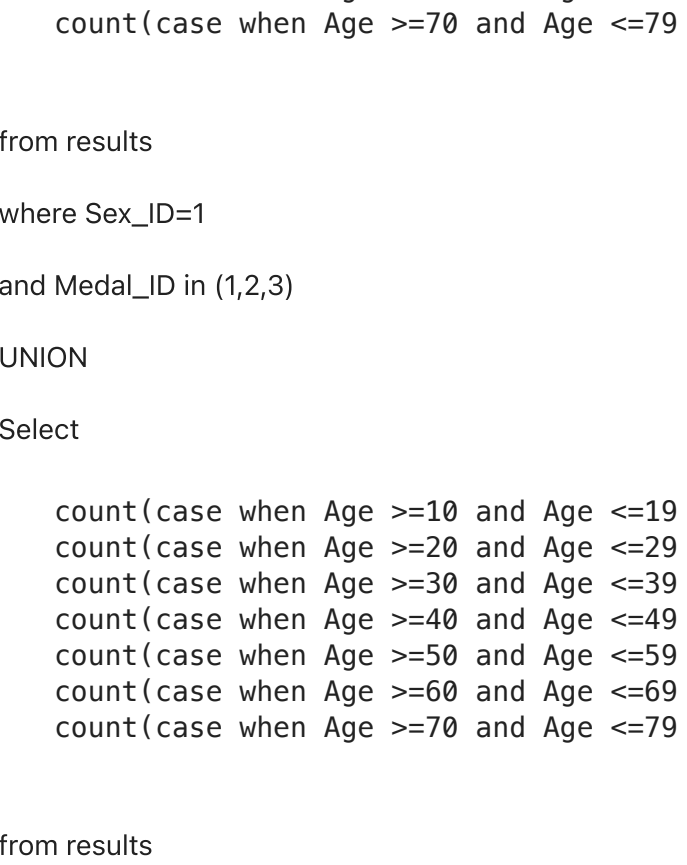
```
In [213]:
age_groups_gold.index=['gold_medals']

In [214]:
#age distribution of gold medal winning athletes
age_groups_gold
```

age_10_19 age_20_29 age_30_39 age_40_49 age_50_59 age_60_69 age_70_79										
gold_medals	1179	9403	2367	343	70	10	0			

```
In [215]:
age_groups_gold.plot(kind='bar')
plt.xticks(rotation=0)
plt.title('Age Distribution of Gold Medal Winning Participants', size=16, fontweight='bold')

Out[215]:
Text(0.5, 1.0, 'Age Distribution of Gold Medal Winning Participants')
```



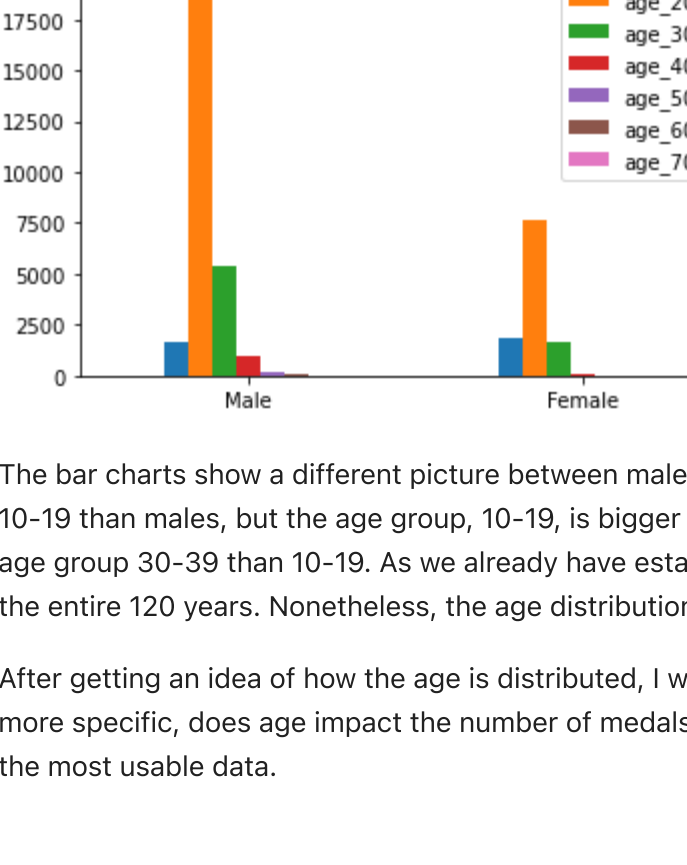
```
In [216]:
age_groups_silver.index=['silver_medals']

In [217]:
#age distribution of silver medal winning athletes
age_groups_silver
```

age_10_19 age_20_29 age_30_39 age_40_49 age_50_59 age_60_69 age_70_79										
silver_medals	1165	9132	2371	350	73	22	3			

```
In [218]:
age_groups_silver.plot(kind='bar')
plt.xticks(rotation=0)
plt.title('Age Distribution of Silver Medal Winning Participants', size=16, fontweight='bold')

Out[218]:
Text(0.5, 1.0, 'Age Distribution of Silver Medal Winning Participants')
```



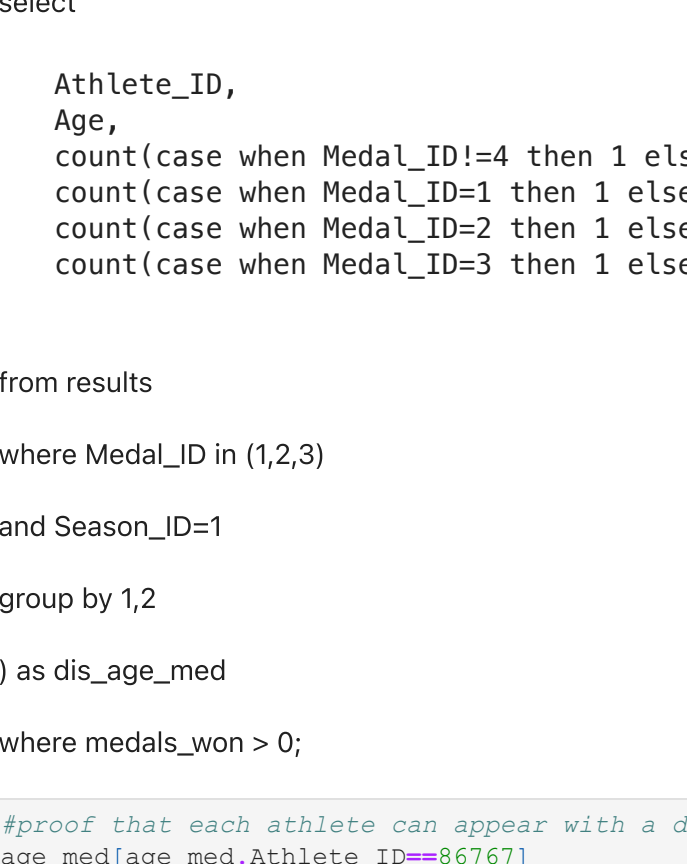
```
In [219]:
age_groups_bronze.index=['bronze_medals']

In [220]:
#age distribution of bronze medal winning athletes
age_groups_bronze
```

age_10_19 age_20_29 age_30_39 age_40_49 age_50_59 age_60_69 age_70_79										
bronze_medals	1137	9467	2265	346	66	12	2			

```
In [221]:
age_groups_bronze.plot(kind='bar')
plt.xticks(rotation=0)
plt.title('Age Distribution of Bronze Medal Winning Participants', size=16, fontweight='bold')

Out[221]:
Text(0.5, 1.0, 'Age Distribution of Bronze Medal Winning Participants')
```



Again we can see a familiar picture for the different medal types. The strongest age group was 20-29. Can we observe a different trend between men and women?

Age Distribution Men vs Women

```
In [222]:
#age distribution by gender
men_women_age=pyqidf('Select count(case when Age >=10 and Age <=19 then Medal_ID else null end) as age_10_19, count(case when Age >=20 and Age <=29 then Medal_ID else null end) as age_20_29, count(case when Age >=30 and Age <=39 then Medal_ID else null end) as age_30_39, count(case when Age >=40 and Age <=49 then Medal_ID else null end) as age_40_49, count(case when Age >=50 and Age <=59 then Medal_ID else null end) as age_50_59, count(case when Age >=60 and Age <=69 then Medal_ID else null end) as age_60_69, count(case when Age >=70 and Age <=79 then Medal_ID else null end) as age_70_79, count(case when Age >=80 and Age <=89 then Medal_ID else null end) as age_80_89, count(case when Age >=90 and Age <=99 then Medal_ID else null end) as age_90_99')

SQL Syntax /Age group distribution between Men and Women/

Select

count(case when Age >=10 and Age <=19 then Medal_ID else null end) as age_10_19,
count(case when Age >=20 and Age <=29 then Medal_ID else null end) as age_20_29,
count(case when Age >=30 and Age <=39 then Medal_ID else null end) as age_30_39,
count(case when Age >=40 and Age <=49 then Medal_ID else null end) as age_40_49,
count(case when Age >=50 and Age <=59 then Medal_ID else null end) as age_50_59,
count(case when Age >=60 and Age <=69 then Medal_ID else null end) as age_60_69,
count(case when Age >=70 and Age <=79 then Medal_ID else null end) as age_70_79,
count(case when Age >=80 and Age <=89 then Medal_ID else null end) as age_80_89,
count(case when Age >=90 and Age <=99 then Medal_ID else null end) as age_90_99

from results;

where Sex_ID=1

and Medal_ID in (1,2,3);

UNION

Select

count(case when Age >=10 and Age <=19 then Medal_ID else null end) as age_10_19,
count(case when Age >=20 and Age <=29 then Medal_ID else null end) as age_20_29,
count(case when Age >=30 and Age <=39 then Medal_ID else null end) as age_30_39,
count(case when Age >=40 and Age <=49 then Medal_ID else null end) as age_40_49,
count(case when Age >=50 and Age <=59 then Medal_ID else null end) as age_50_59,
count(case when Age >=60 and Age <=69 then Medal_ID else null end) as age_60_69,
count(case when Age >=70 and Age <=79 then Medal_ID else null end) as age_70_79,
count(case when Age &
```


In [230]:

```

#Joining concat to to combine the data frames
age_med=pd.concat([age_med|age_med.Age==14].sample(n=50), age_med[age_med.Age==15].sample(n=50),
age_med|age_med.Age==16].sample(n=50), age_med[age_med.Age==17].sample(n=50),
age_med|age_med.Age==18].sample(n=50), age_med[age_med.Age==19].sample(n=50),
age_med|age_med.Age==20].sample(n=50), age_med[age_med.Age==21].sample(n=50),
age_med|age_med.Age==22].sample(n=50), age_med[age_med.Age==23].sample(n=50),
age_med|age_med.Age==24].sample(n=50), age_med[age_med.Age==25].sample(n=50),
age_med|age_med.Age==26].sample(n=50), age_med[age_med.Age==27].sample(n=50),
age_med|age_med.Age==28].sample(n=50), age_med[age_med.Age==29].sample(n=50),
age_med|age_med.Age==30].sample(n=50), age_med[age_med.Age==31].sample(n=50),
age_med|age_med.Age==32].sample(n=50), age_med[age_med.Age==33].sample(n=50),
age_med|age_med.Age==34].sample(n=50), age_med[age_med.Age==35].sample(n=50),
age_med|age_med.Age==36].sample(n=50), age_med[age_med.Age==37].sample(n=50),
age_med|age_med.Age==38].sample(n=50), age_med[age_med.Age==39].sample(n=50),
age_med|age_med.Age==40].sample(n=50), age_med[age_med.Age==41].sample(n=50),
age_med|age_med.Age==42].sample(n=50), age_med[age_med.Age==43].sample(n=50),
age_med|age_med.Age==44].sample(n=50), age_med[age_med.Age==45].sample(n=50),
age_med|age_med.Age==46].sample(n=50)])
age_med
```

Out [236]:

	Athlete_ID	Age	medals_won	gold_won	silver	bronze_won
24624	111140	14	1	1	0	0
3085	14804	14	1	1	0	0
18124	80581	14	2	0	1	1
16245	72433	14	1	1	0	0
23615	106232	14	1	0	1	0
...
24674	111407	46	1	0	0	1
13529	59657	46	1	0	1	0
27345	123157	46	1	0	0	1
24069	108417	46	1	0	1	0
13023	57201	46	1	0	0	1

1650 rows x 6 columns

The sample size was 1650 athletes, that won at least one medal

n=1650

In [232]:

```

#checking if there is a correlation between age and medals won
#monotonic=> we cannot observe a monotonic development
corr_kst=kendalltau(age_med.Age, age_med.medals_won)
corr_kst
```

Out [237]:

KendalltauResult (correlation=-0.0585045306368485, pvalue=0.0035146636686267507)

The correlation was weak and not statistically significant. To ensure we were not missing any trend or possibility that age and the medals won correlate, I have created a scatter plot of athletes that have not won a medal.

I have selected the same age range as I did and randomly chose 50 samples of every age.

In [232]:

```

#getting the table of all non-medal winning athletes and their age
age_no_med=pd.query('Select Athlete_ID, Age, no_medals_won from (select Athlete_ID, Age, count(case when Medal_I
```

In [234]:

```

age_no_med=pd.concat([age_no_med|age_no_med.Age==14].sample(n=50), age_no_med|age_no_med.Age==15].sample(n=50),
age_no_med|age_no_med.Age==16].sample(n=50), age_no_med|age_no_med.Age==17].sample(n=50),
age_no_med|age_no_med.Age==18].sample(n=50), age_no_med|age_no_med.Age==19].sample(n=50),
age_no_med|age_no_med.Age==20].sample(n=50), age_no_med|age_no_med.Age==21].sample(n=50),
age_no_med|age_no_med.Age==22].sample(n=50), age_no_med|age_no_med.Age==23].sample(n=50),
age_no_med|age_no_med.Age==24].sample(n=50), age_no_med|age_no_med.Age==25].sample(n=50),
age_no_med|age_no_med.Age==26].sample(n=50), age_no_med|age_no_med.Age==27].sample(n=50),
age_no_med|age_no_med.Age==28].sample(n=50), age_no_med|age_no_med.Age==29].sample(n=50),
age_no_med|age_no_med.Age==30].sample(n=50), age_no_med|age_no_med.Age==31].sample(n=50),
age_no_med|age_no_med.Age==32].sample(n=50), age_no_med|age_no_med.Age==33].sample(n=50),
age_no_med|age_no_med.Age==34].sample(n=50), age_no_med|age_no_med.Age==35].sample(n=50),
age_no_med|age_no_med.Age==36].sample(n=50), age_no_med|age_no_med.Age==37].sample(n=50),
age_no_med|age_no_med.Age==38].sample(n=50), age_no_med|age_no_med.Age==39].sample(n=50),
age_no_med|age_no_med.Age==40].sample(n=50), age_no_med|age_no_med.Age==41].sample(n=50),
age_no_med|age_no_med.Age==42].sample(n=50), age_no_med|age_no_med.Age==43].sample(n=50),
age_no_med|age_no_med.Age==44].sample(n=50), age_no_med|age_no_med.Age==45].sample(n=50),
age_no_med|age_no_med.Age==46].sample(n=50)])
age_no_med
```

Out [234]:

	Athlete_ID	Age	no_medals_won
31398	30793	14	1
133975	132164	14	1
23002	22661	14	1
23386	23049	14	1
39901	39223	14	1
...
6179	5965	46	2
19915	19484	46	2
125235	123157	46	1
27474	26951	46	2
21287	20905	46	1

1650 rows x 3 columns

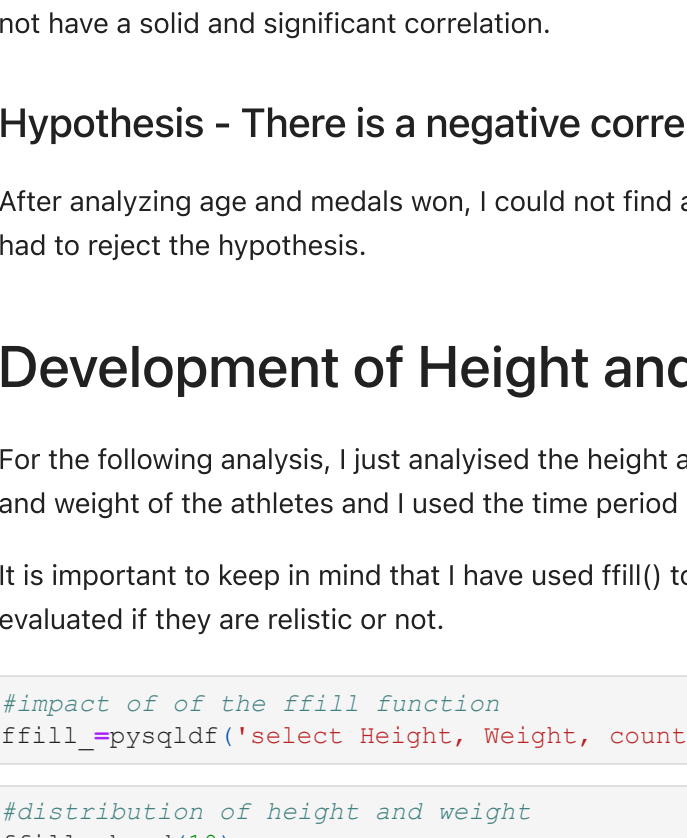
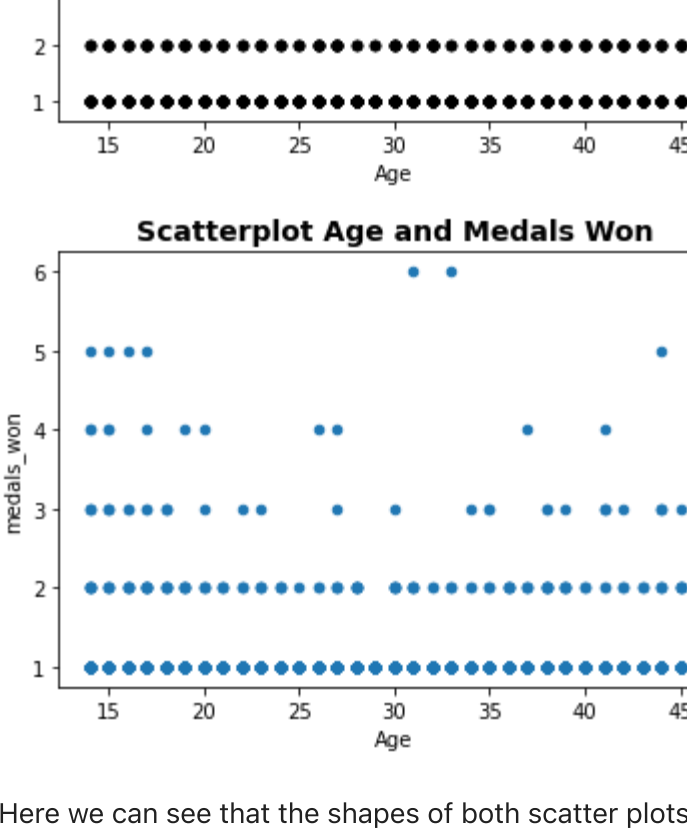
In [235]:

```

age_no_med=plot.scatter(x='Age', y='no_medals_won', color='black')
plt.title('Scatterplot Age and No Medals Won', size=14, fontweight='bold')
age_med=plot.scatter(x='Age', y='medals_won', size=14, fontweight='bold')
plt.title('Scatterplot Age and Medals Won', size=14, fontweight='bold')
```

Out [235]:

Text(0.5, 1.0, 'Scatterplot Age and No Medals Won')



Here we can see that the shapes of both scatter plots looked familiar. Therefore the hypothesis has to be rejected, because we did not have a solid and significant correlation.

Hypothesis - There is a negative correlation between Age and Medals

After analyzing age and medals won, I could not find a statistically significant correlation between these two variables. Therefore I had to reject the hypothesis.

Development of Height and Weight of Olympic Summer Athletes

For the following analysis, I just analysed the height and weight of the summer game athletes. To do so, I used the average height and weight of the athletes and I used the time period of 120 years.

It is important to keep in mind that I have used ffill() to get ride of the na-values, as I already mentioned these values should be evaluated if they are relistic or not.

In [236]:

```

#Impact of of the fill function
ffill_ppyqldf('select Height, Weight, count(Height), count(Weight) from results where Year between 1920 and 19
```

Out [237]:

#distribution of height and weight

	Height	Weight	count(Height)	count(Weight)
0	176	67	38816	38816
1	170	67	346	346
2	173	67	274	274
3	178	67	246	246
4	168	67	244	244
5	167	67	219	219
6	180	67	201	201
7	183	67	195	195
8	175	67	186	186
9	165	67	152	152

As we can see the combination 176 cm and 67 Kg, which were used in the fill function, created a huge amount of data and seemed very unrealistic, therefore they were not included in the analysis.

SQL Syntax /Impact of of the fill function/

Select

Height,
Weight,
count(Height),
count(Weight)

from results

where Year between 1920 and 1960

group by 1,2

order by count(Height) desc;

In [238]:

```

#development of weight and height over 120 years of summer games
dev_ath_ppyqldf('select results.Year, avg(results.Height) as Height, avg(results.Weight) as Weight from results
```

Out [239]:

/Development of Weight and Height over 120 Years/

select

results.Year,
avg(results.Height) as Height,
avg(results.Weight) as Weight

from results

left join sports
on results.Sport_ID=sports.Sport_ID

where Season_ID=1

and Weight=176

and Weight=67

group by 1;

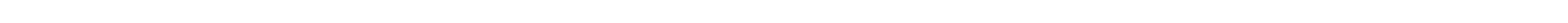
In [239]:

```

dev_ath_plot(x='Year', y='Height', color='black')
dev_ath_plot(x='Year', y='Weight', color='black')
plt.title('Development of Height', size=16, fontweight='bold')
plt.title('Development of Weight', size=16, fontweight='bold')
```

Out [239]:

Text(0.5, 1.0, 'Development of Height')



In [240]:

```

#development of weight and height of Male athletes over 120 Years
dev_ath_ppyqldf('select results.Year, avg(results.Height) as Height, avg(results.Weight) as Weight from resul
```

Out [241]:

/Development of Weight and Height of Male athletes over 120 Years/

select

results.Year,
avg(results.Height) as Height,
avg(results.Weight) as Weight

from results

where results.Sex_ID=1

and Season_ID=1

and Weight=176

and Weight=67

group by 1;

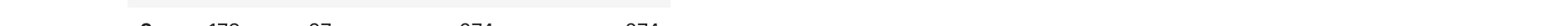
In [241]:

```

dev_ath_m_plot(x='Year', y='Height', color='red')
plt.title('Height Development of Male Athletes', size=16, fontweight='bold')
dev_ath_m_plot(x='Year', y='Weight', color='red')
plt.title('Weight Development of Male Athletes', size=16, fontweight='bold')
```

Out [241]:

Text(0.5, 1.0, 'Height Development of Male Athletes')



In [242]:

```

#development of weight and height of Female athletes over 120 Years
dev_ath_ppyqldf('select results.Year, avg(results.Height) as Height, avg(results.Weight) as Weight from resul
```

Out [243]:

SQL Syntax /Development of Weight and Height of Female athletes over 120 Years/

select

results.Year,
avg(results.Height) as Height,
avg(results.Weight) as Weight

from results

where results.Sex_ID=2

and Season_ID=1

and Weight=176

and Weight=67

group by 1;

In [243]:

```

dev_ath_f_plot(x='Year', y='Height', color='red')
plt.title('Height Development of Female Athletes', size=16, fontweight='bold')
dev_ath_f_plot(x='Year', y='Weight', color='red')
plt.title('Weight Development of Female Athletes', size=16, fontweight='bold')
```

Out [243]:

Text(0.5, 1.0, 'Weight Development of Female Athletes')



The analysis has shown that the height and weight of both genders have increased.

Sources

The Correlation Coefficient (r) - Anon, 2021

URL: <https://sphweb.bumc.bu.edu/otlt/MPH-Modules/PH717-QuantCore/PH717-Module9-Correlation-Regression/PH717-Module9-Correlation-Regression4.html>.