

```
In [ ] : import pandas as pd
import sqlalchemy
import matplotlib.pyplot as plt
import plotly.express as px

In [ ] : #Make the connection to the dataset in mysql and read the csv files

engine = sqlalchemy.create_engine('mysql+pymysql://root:root@localhost:3306/abcourse')

athlete_df = pd.read_csv('C:/Users/PC/Desktop/Proyecto/athlete_events.csv')
noc_regions_df = pd.read_csv('C:/Users/PC/Desktop/Proyecto/noc_regions.csv')

First, I wanted to prove or disprove mi hypothesis that USA had a better performance or more medals in Basketball.

In [ ] : USA_Basketball = pd.read_sql(
'''
SELECT noc AS Team, COUNT(Medal) AS Medals
FROM athlete
WHERE Sport = 'Basketball'
GROUP BY noc
ORDER BY Medals DESC
LIMIT 20
''', con = engine
)

plt.figure(figsize=(10, 15))
fig = px.bar(USA_Basketball, y = 'Medals', x = 'Team', text = 'Medals')
fig.update_layout(uniformtext_amsize = 8, uniformtext_mode = 'hide', xaxis_tickangle=45)
fig.show()

<Figure size 1000x1500 with 0 Axes>

I see that USA has more medals than other teams but that's not necessary that they had a good performance through the years, so I decided to see how many medals of each category USA and other teams had

Basketball_Medals = pd.read_sql(
'''
SELECT
    noc AS Team,
    sum(CASE WHEN Medal = 'Bronze' THEN 1 ELSE 0 END) AS Bronze,
    sum(CASE WHEN Medal = 'Silver' THEN 1 ELSE 0 END) AS Silver,
    sum(CASE WHEN Medal = 'Gold' THEN 1 ELSE 0 END) AS Gold
FROM
    athlete
WHERE Sport = 'Basketball'
GROUP BY noc
ORDER BY Gold DESC
LIMIT 10
''', con = engine
)
Basketball_Medals

Out[ ] : Team  Bronze  Silver  Gold
0  USA      36.0    24.0  281.0
1  URS      48.0    50.0   48.0
2  ARG      12.0     0.0   12.0
3  YUG      24.0    48.0   12.0
4  EUN       0.0     0.0   12.0
5  EGY       0.0     0.0    0.0
6  ITA       0.0     0.0    0.0
7  JPN       0.0     0.0    0.0
8  CUB      12.0     0.0    0.0
9  RUS      36.0     0.0    0.0
```

Now it was time to prove or disprove my 2nd hypothesis that Brazil had better performance/more medals in Football

```
In [ ] : Football_Medals = pd.read_sql(
'''
SELECT noc AS Team, COUNT(Medal) AS Medals
FROM athlete
WHERE Sport = 'Football'
GROUP BY noc
ORDER BY Medals DESC
LIMIT 20
''', con = engine
)

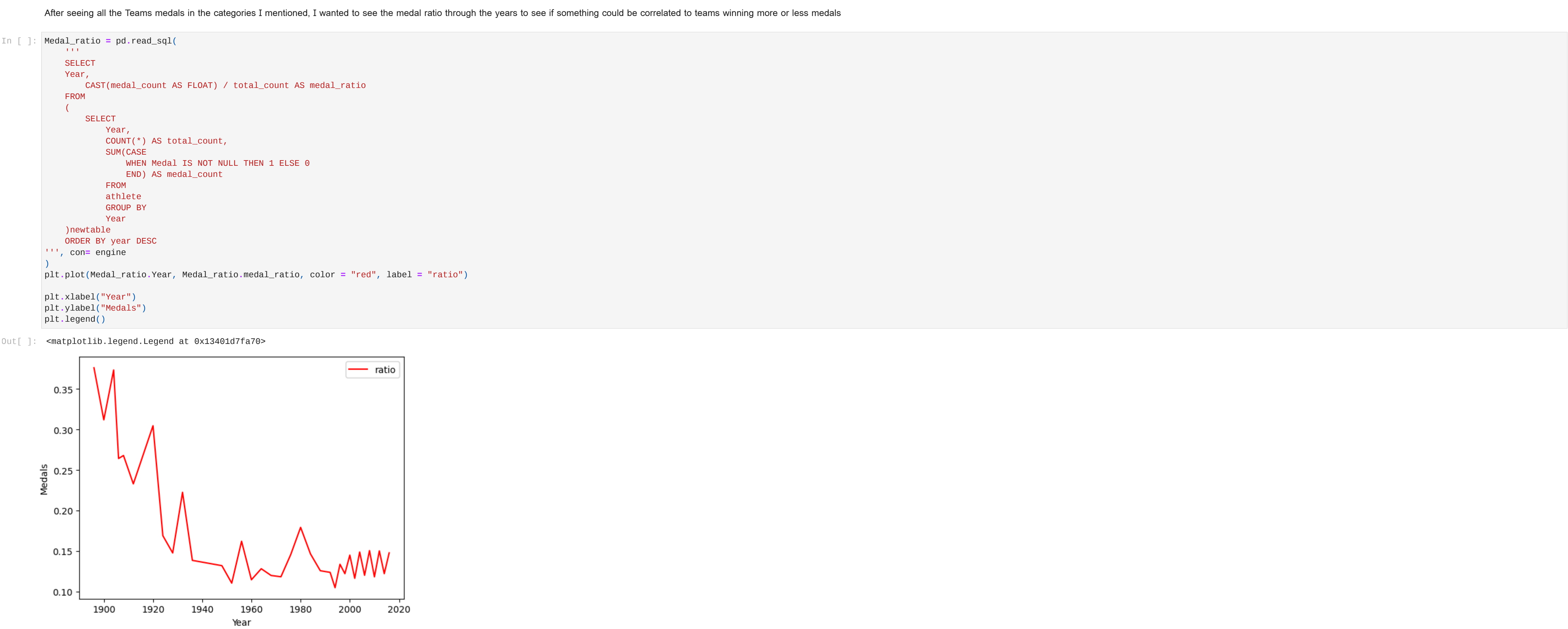
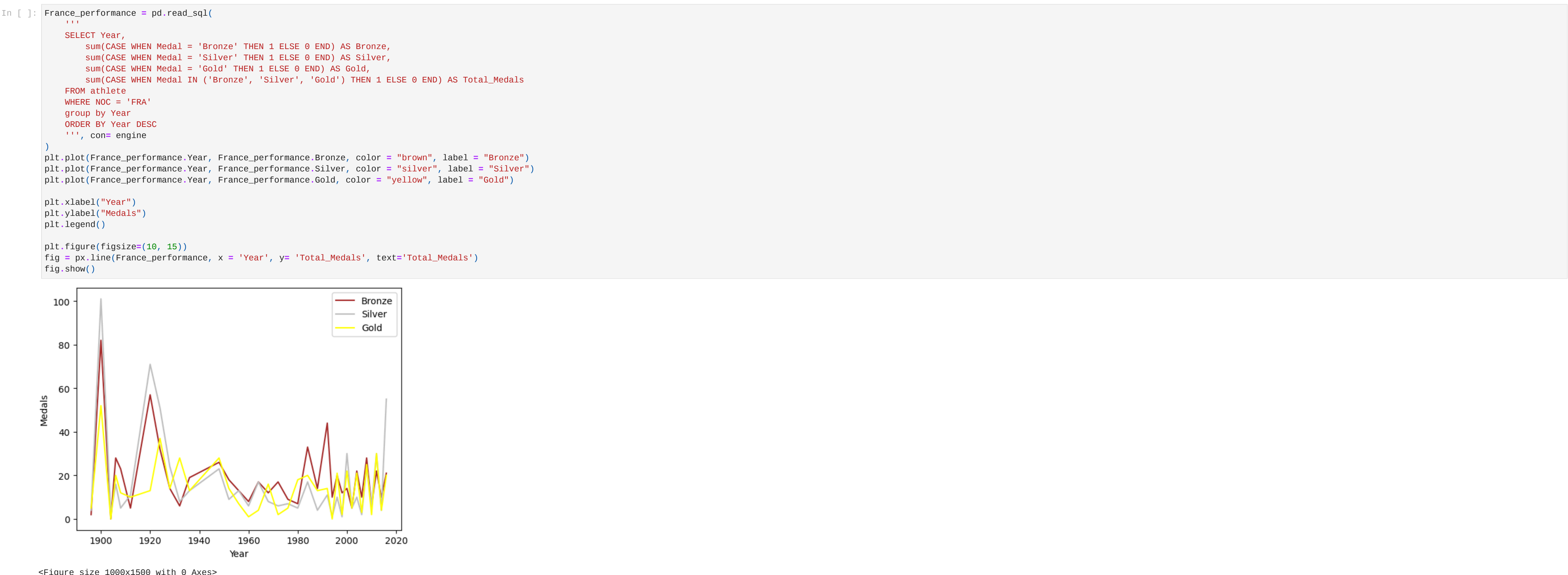
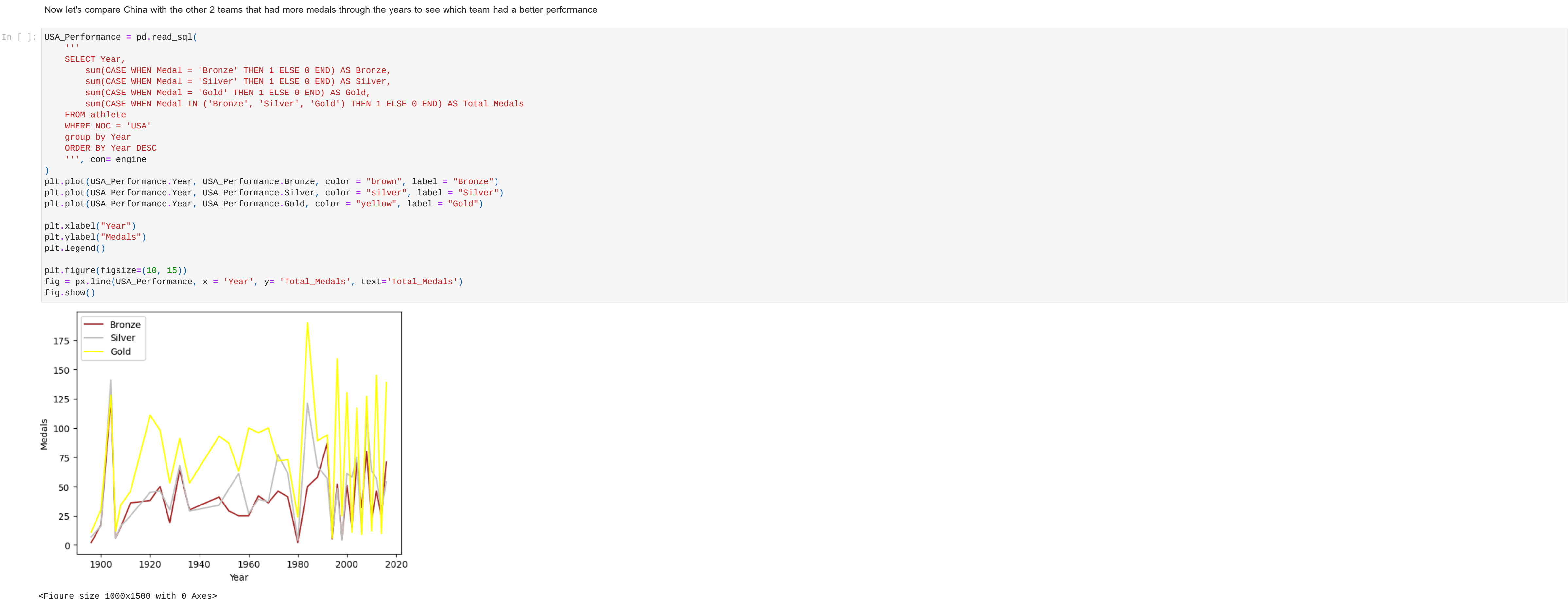
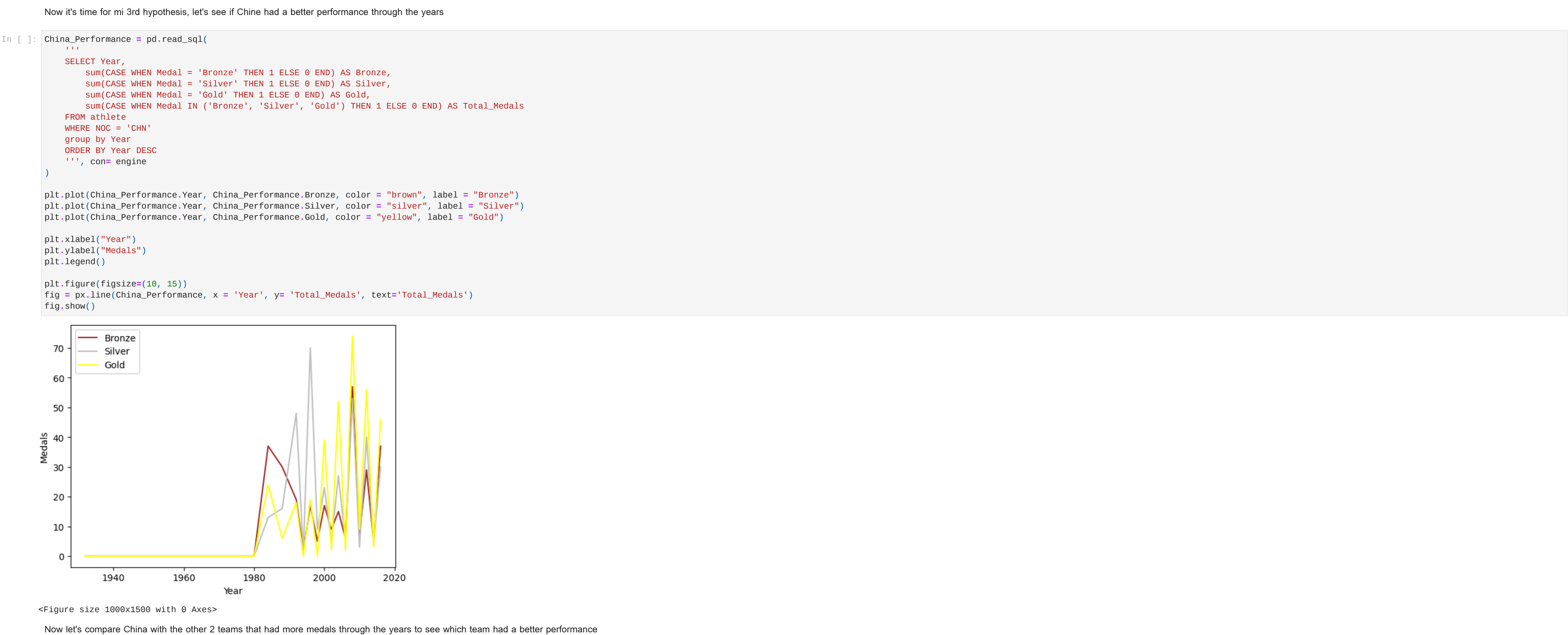
plt.figure(figsize=(10, 15))
fig = px.bar(Football_Medals, y = 'Medals', x = 'Team', text = 'Medals')
fig.update_layout(uniformtext_amsize = 8, uniformtext_mode = 'hide', xaxis_tickangle=45)
fig.show()

<Figure size 1000x1500 with 0 Axes>

With this graph we could see that Brazil has more medals but now we need to see how many of those are gold, silver and bronze

Football_Performance = pd.read_sql(
'''
SELECT
    noc AS Team,
    sum(CASE WHEN Medal = 'Bronze' THEN 1 ELSE 0 END) AS Bronze,
    sum(CASE WHEN Medal = 'Silver' THEN 1 ELSE 0 END) AS Silver,
    sum(CASE WHEN Medal = 'Gold' THEN 1 ELSE 0 END) AS Gold
FROM
    athlete
WHERE Sport = 'Football'
GROUP BY noc
ORDER BY Gold DESC
LIMIT 10
''', con = engine
)
Football_Performance

Out[ ] : Team  Bronze  Silver  Gold
0  USA      12.0    24.0   66.0
1  HUN       6.0    17.0   46.0
2  GBR       1.0     8.0   36.0
3  URS      51.0     0.0   36.0
4  ARG       0.0    34.0   34.0
5  URU       0.0     0.0   31.0
6  CMR       0.0     0.0   18.0
7  GER      69.0    17.0   18.0
8  NOR      30.0     0.0   17.0
9  BRA      34.0    85.0   17.0
```



let's take a look at my last hypothesis, that People with age > 35 have more medals than age < 35

```
In [ ] : AgeAthletes = pd.read_sql(
'''
SELECT
    COUNT(distinct Name) as Athletes
FROM athlete
WHERE Age > 35
)AS Older
COUNT(distinct Name) as Younger
FROM athlete
WHERE Age < 35
''', con= engine
)
AgeAthletes

Out[ ] : Older  Younger
0  8854  122148

In [ ] : YoungerMedals = pd.read_sql(
'''
SELECT
    COUNT(Medal) AS Medals_total,
    SUM(CASE WHEN Medal = 'Gold' THEN 1 ELSE 0 END) AS Gold,
    sum(CASE WHEN Medal = 'Silver' THEN 1 ELSE 0 END) AS Silver,
    sum(CASE WHEN Medal = 'Bronze' THEN 1 ELSE 0 END) AS Bronze
FROM athlete
WHERE Age < 35
''', con= engine
)
YoungerMedals

Out[ ] : Medals_total  Gold  Silver  Bronze
0  36245  12270.0  11877.0  12098.0

In [ ] : OlderMedals = pd.read_sql(
'''
SELECT
    COUNT(Medal) AS Medals_total,
    SUM(CASE WHEN Medal = 'Gold' THEN 1 ELSE 0 END) AS Gold,
    sum(CASE WHEN Medal = 'Silver' THEN 1 ELSE 0 END) AS Silver,
    sum(CASE WHEN Medal = 'Bronze' THEN 1 ELSE 0 END) AS Bronze
FROM athlete
WHERE Age > 35
''', con= engine
)
OlderMedals

Out[ ] : Medals_total  Gold  Silver  Bronze
0  2321  778.0  788.0  755.0
```

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?

- 1.- I found out that besides Brazil had more medals in Football, they're in 9th place in Teams with more gold medals
- 2.- I found out that China actually won medals after 1984
- 3.- I found out that the medal ratio fluctuated at first and then it stabilized

Did you prove or disprove any of your initial hypotheses? If so, which one and what do you plan to do next?

I prove mi first hypothesis that USA has better performance in Basketball than other teams. My other 3 hypothesis were disproved. In the second hypothesis, brazil had more medals in Football but there are 9 teams that have more gold medals, so we could say that they dont have better performance through the years than other teams. Mi 3rd hypothesis was also disprove because despite China in recent years has become a competitive team, before 1980 they had a bad performance therefore less medals than other teams. In my last hypothesis, the younger athletes that won a medal are the 29% of the total of younger athletes and in the other hand, the older athletes that won a medal are 28% of the total of older athletes, so my hypothesis were disproved

What additional questions are you seeking to answer?

