# Olympic Historical Dataset (1896 - 2022)

**Min Pak** • **Flora Aka** • **Maria Ardila**

# Questions we found interesting

- Country medal count versus economy (GDP) & Population

- How age has changed and how age is a factor in different events

- Women and men's participation throughout the years

- How height and weight have changed through the years

# Data you used to answer these questions

- Kaggle: Olympic Historical Dataset (1896 - 2022)

  - Olympic_Athlete_Bio.csv

  - Olympic_Athlete_Event_Results.csv

  - Olympic_Games_Medal_Tally.csv

  - Olympic_Results.csv

  - Olympics_Country.csv

  - Olympics_Games.csv

  https://www.kaggle.com/datasets/muhammadehsan000/olympic-historical-dataset-1896-2020?select=Olympic_Athlete_Event_Results.csv

fppt.com

# Additional Data Sources

- World Bank Group
  - data.worldbank.org/indicator/SP.POP.TOTL
    - csv file with the population from different countries from 1960 to current

  - data.worldbank.org/indicator/NY.GDP.MKTP.CD
    - csv file with the GDP of different countries from 1960 to current

# Data exploration and cleanup process

DataFrame with the countries that won medals needed alterations so the data in DataFrames would match

#1960 Olympics change from Soviet Union
medal_df.at[341, 'country'] = 'Russian Federation'

#1980 Olympics change from Soviet Union
medal_df.at[559, 'country'] = 'Russian Federation'

#2020 Olympics change from ROC
medal_df.at[1254, 'country'] = 'Russian Federation'

fppt.com

# Additional Data Cleanup

```
#change East Germany to Germany for 1980
medal_df.at[560, 'country'] = 'Germany'

GDP DataFrame
#fix United Kingdom to Great Britain
gdp_df.at[81, 'Country Name'] = 'Great Britain'

#add GDP to Russia Federation (USSR) for 1980 and 1960
gdp_df.at[202, '1980'] = 1210000000000
gdp_df.at[202, '1960'] = 142400000000
```

# Analysis process

After merging GDP and Medals DataFrames and reorganzing the columns, I plotted the Data.  I first used GDP per capita but did not like the results.  I switched to Total GDP.

I made sample plots of 2020, 2016, 2012, 2008, and 2004 Summer Olympics and I found them to be very similar.

After seeing this, I decided looking at a wider range of years would provide more interesting data.   So I reorganized the Data for 1960, 1980, 2000, and 2020 Summer Olympics.
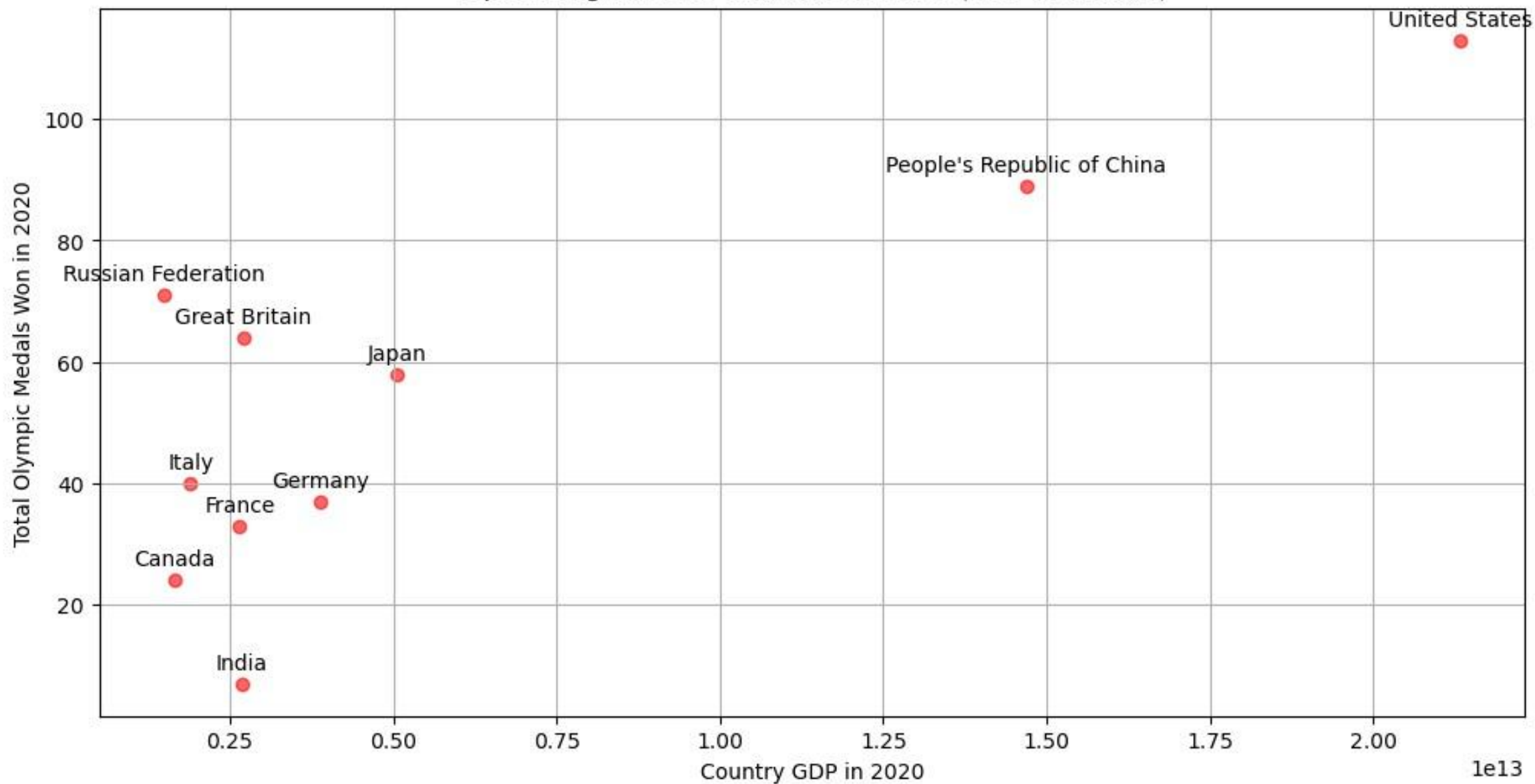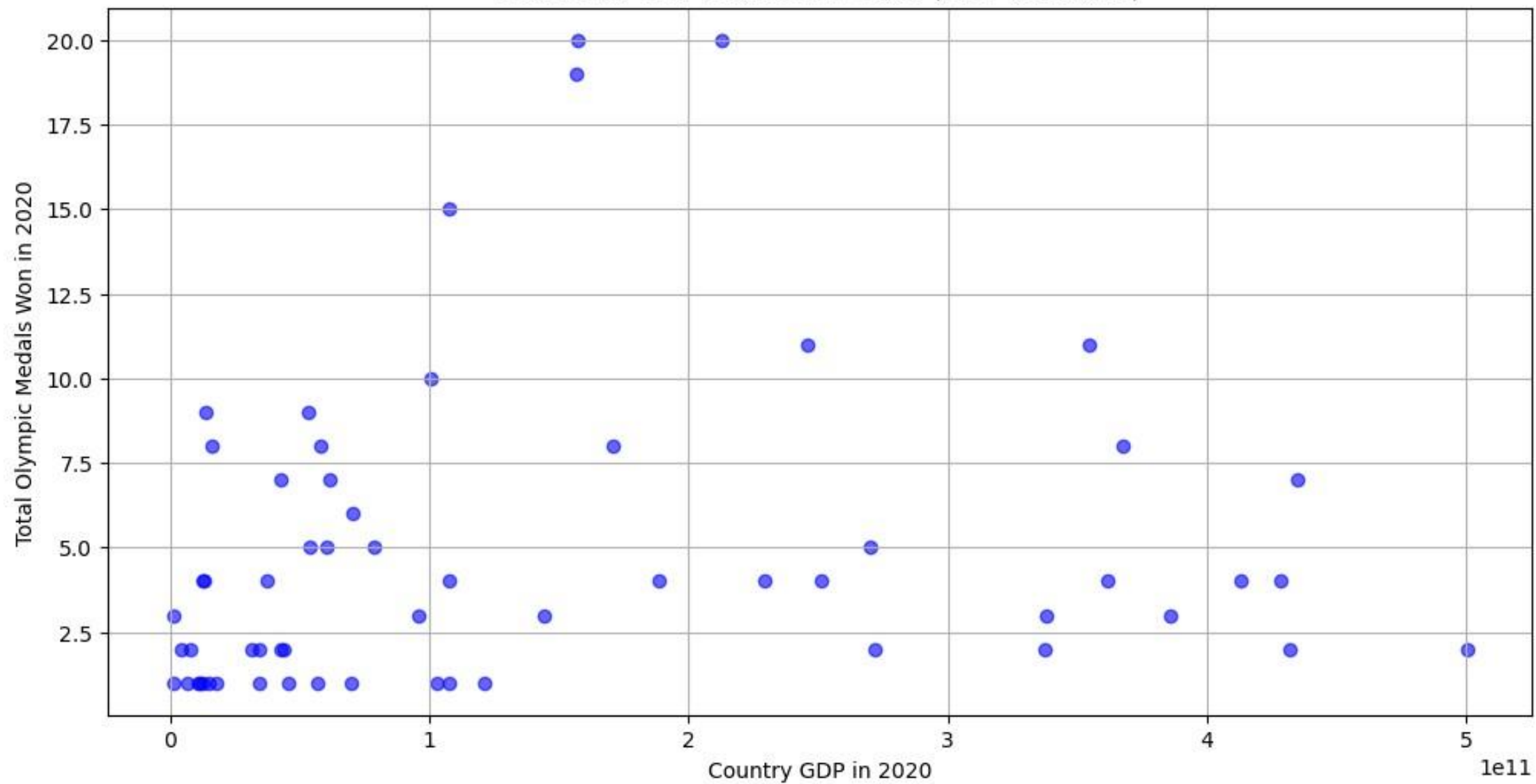
# 2020 Summer Olympics



2020 Summer Olympics - GDP vs Total Medals

2020 Summer Olympics - GDP vs Total Medals

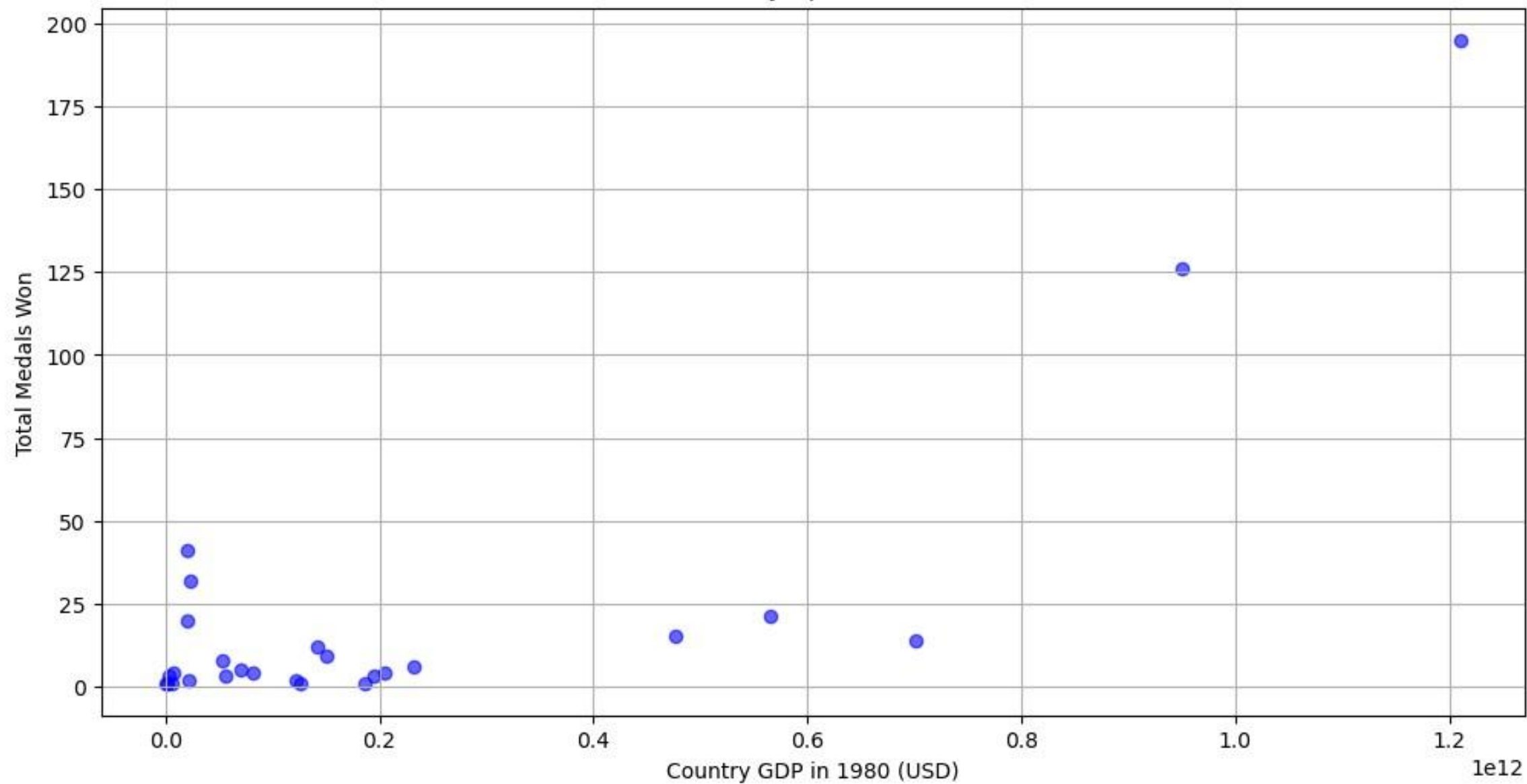Top Ten Highest GDP countries in 2020 (GDP vs Medals)

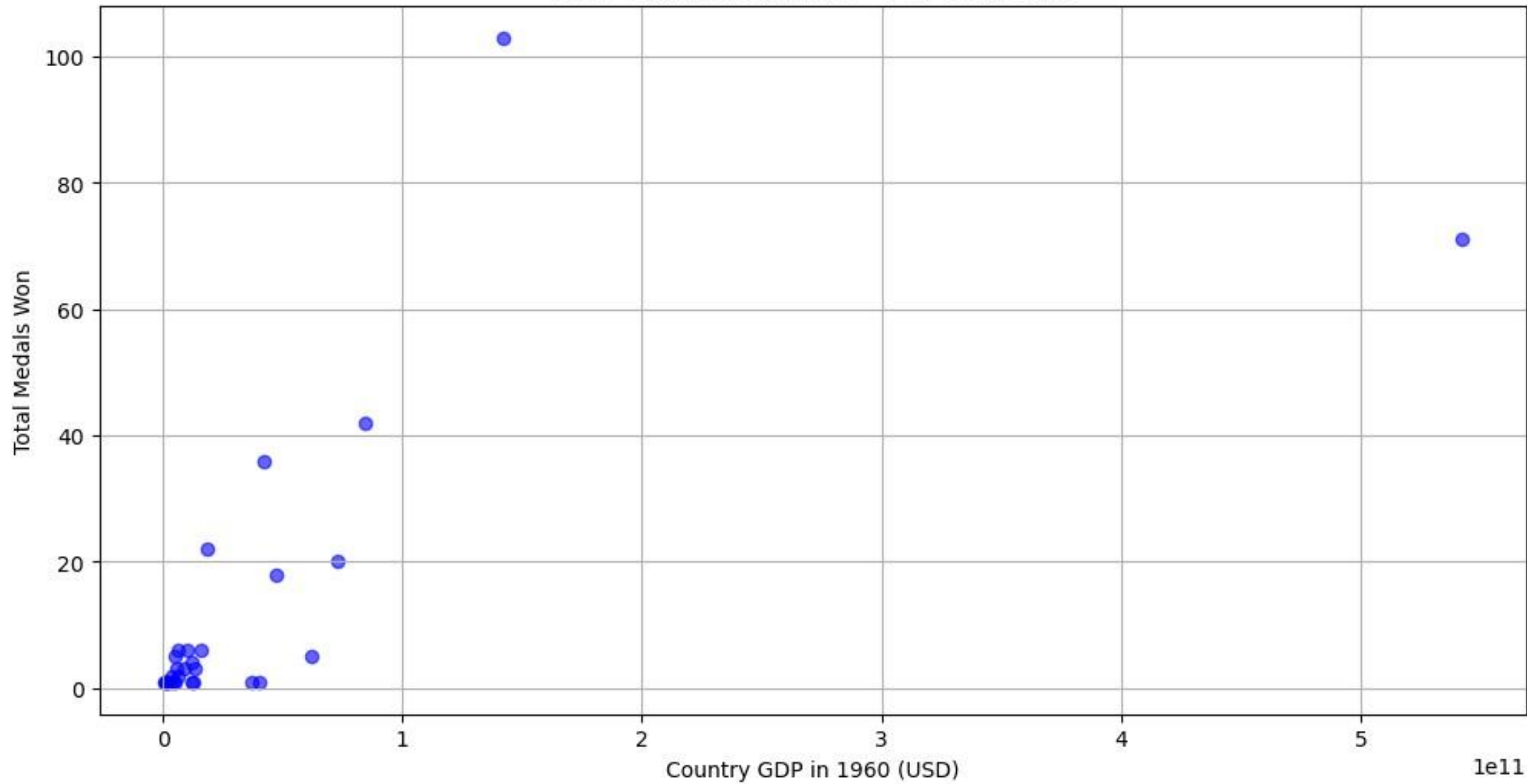Lowest 60 GDP countries in 2020 (GDP vs Medals)
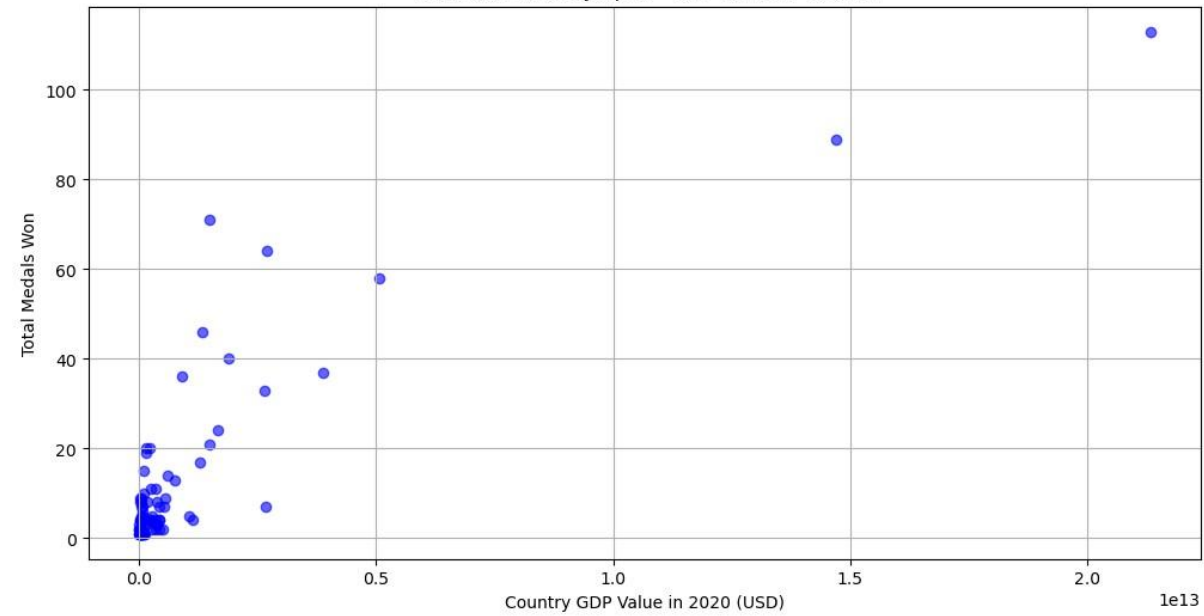
2000 Summer Olympics - GDP vs Medals

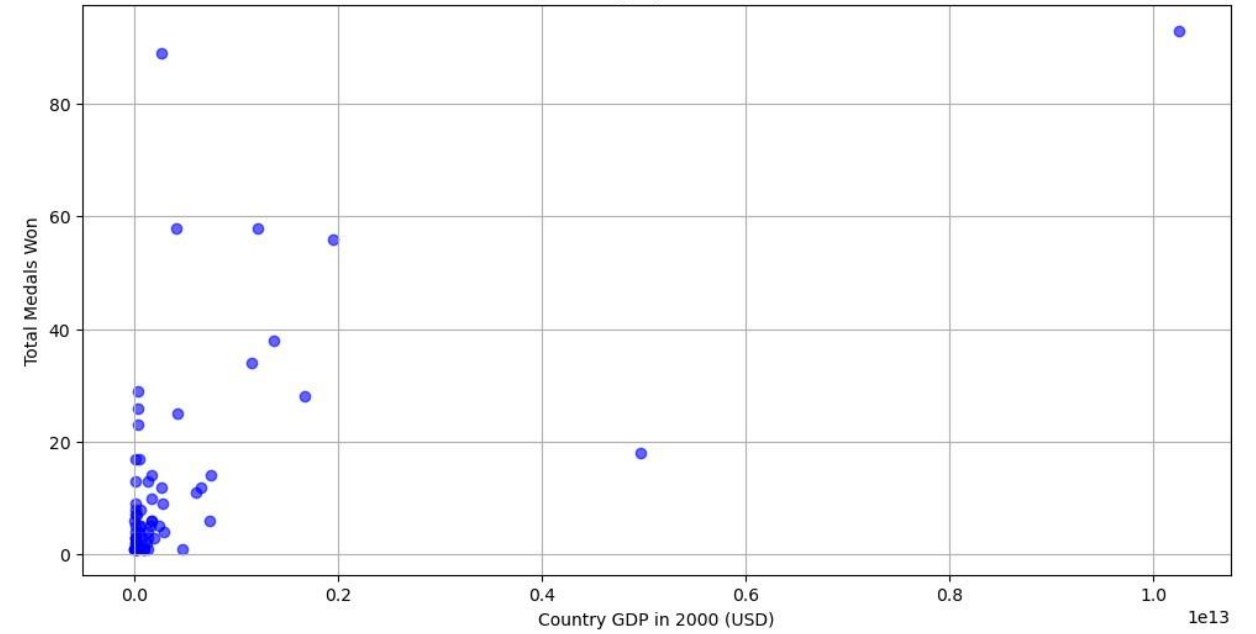1980 Summer Olympics - GDP vs Medals
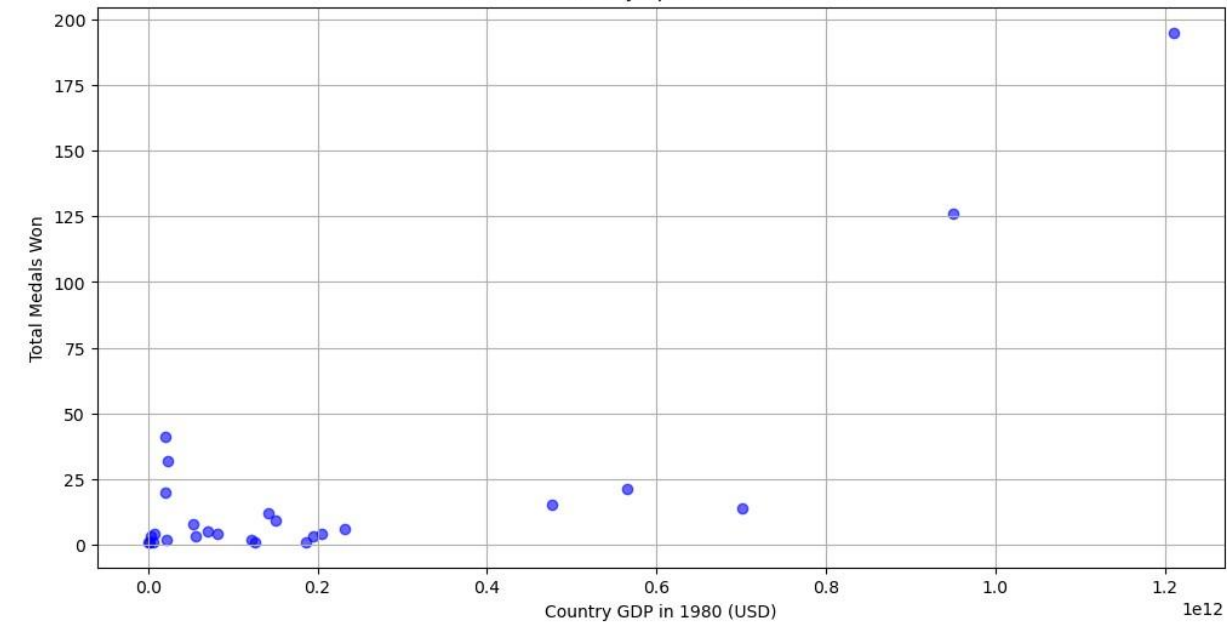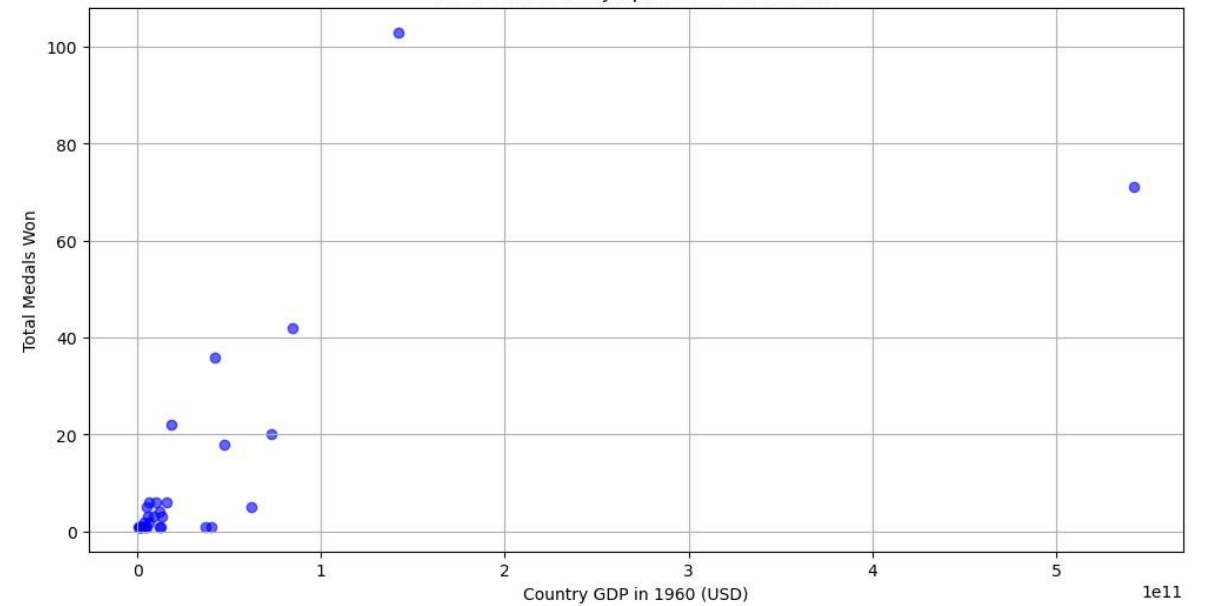
1960 Summer Olympics - GDP vs Medals

**2020 Summer Olympics - GDP vs Total Medals**
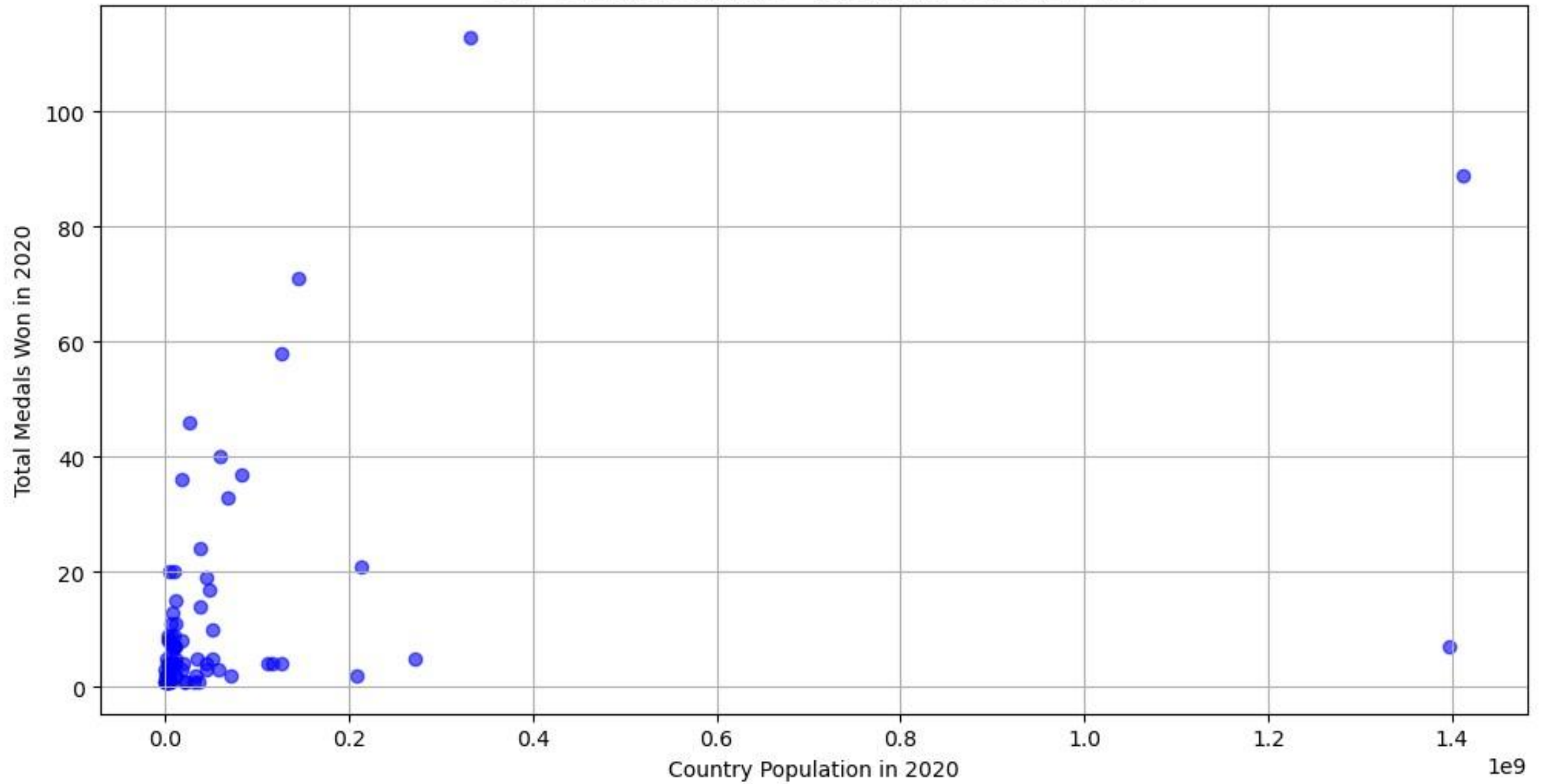**2000 Summer Olympics - GDP vs Medals**
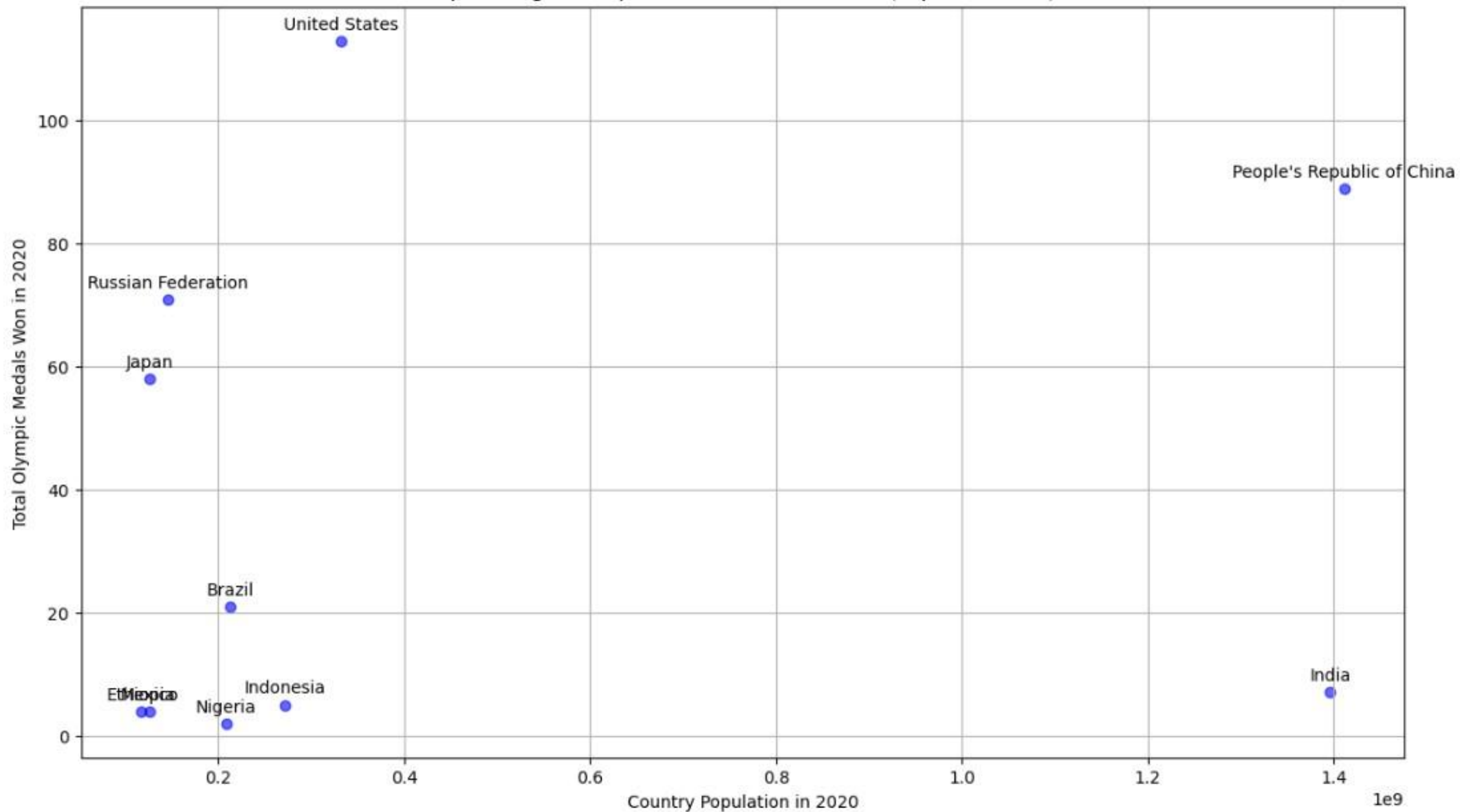**1980 Summer Olympics - GDP vs Medals**
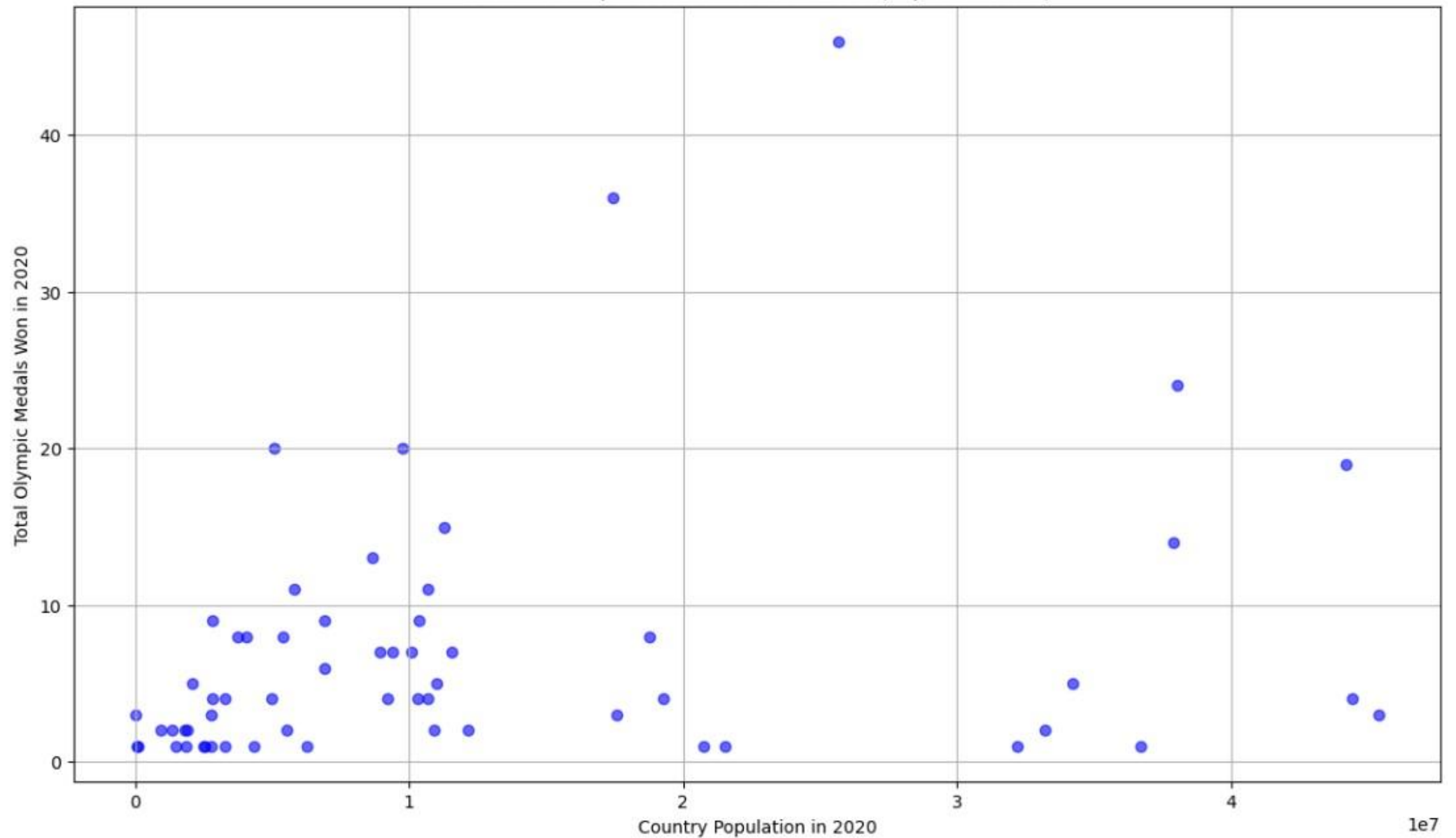**1960 Summer Olympics - GDP vs Medals**

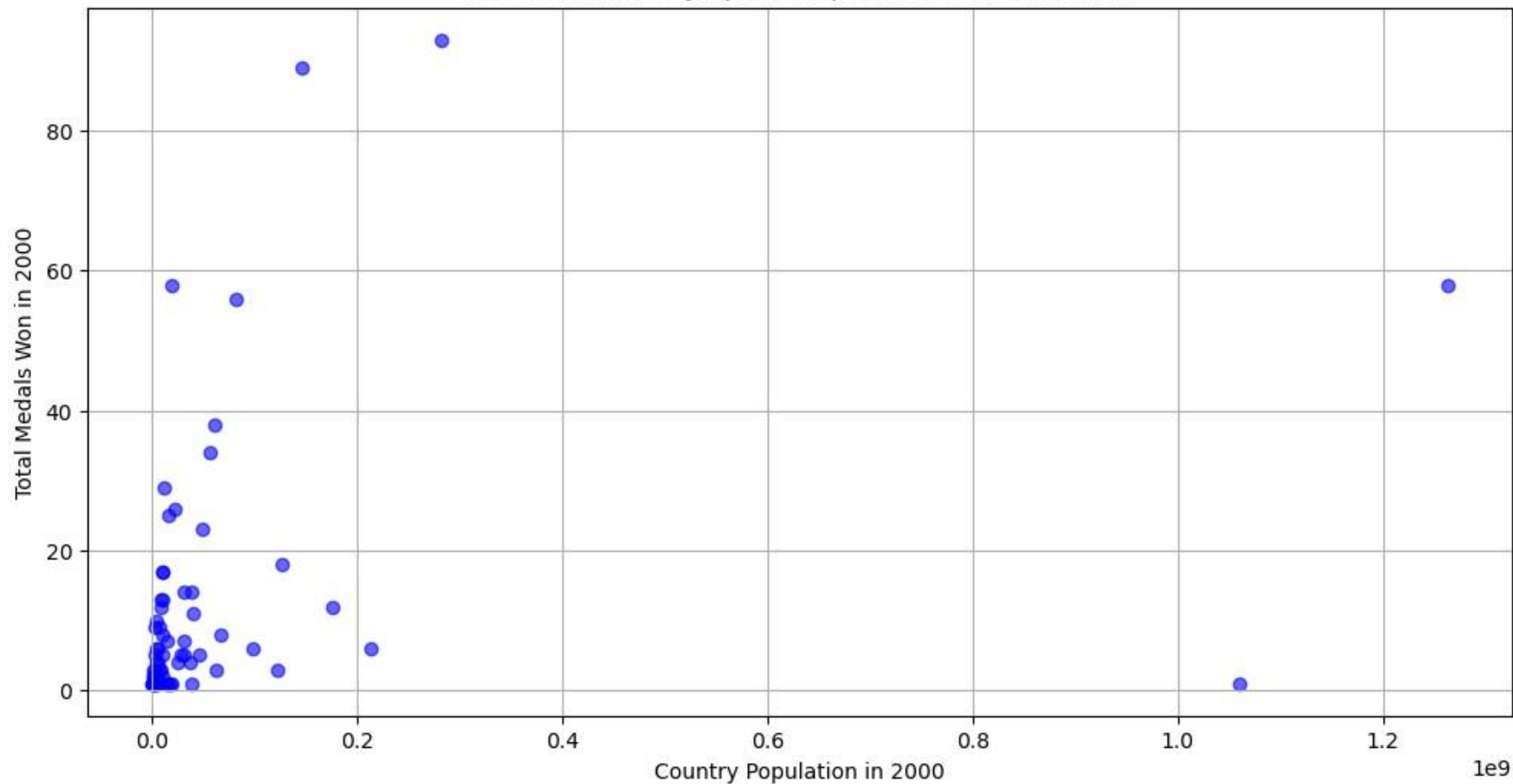2020 Summer Olympics - Population vs Medals Won

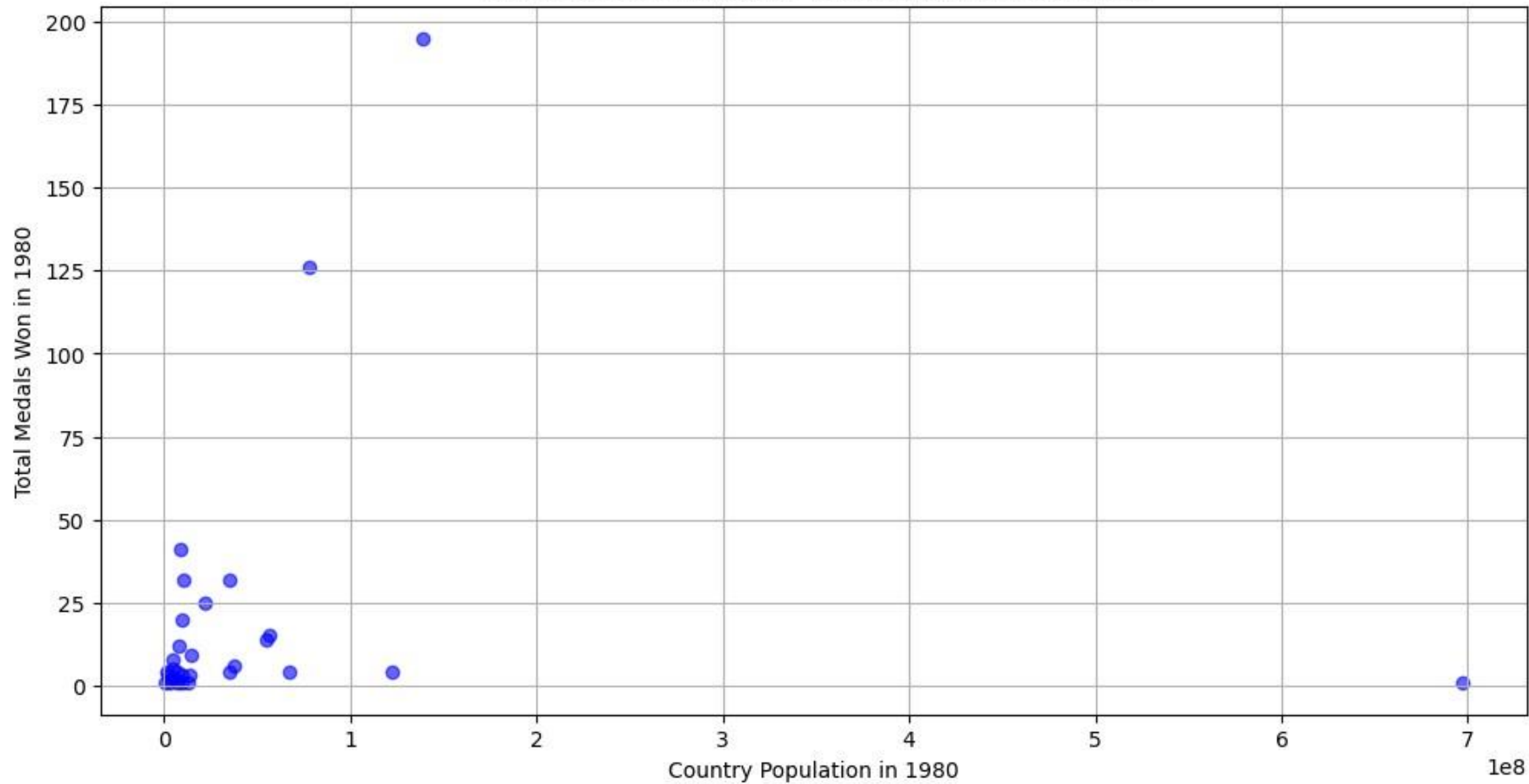Top Ten Highest Populated Countries in 2020 (Pop. vs Medals)

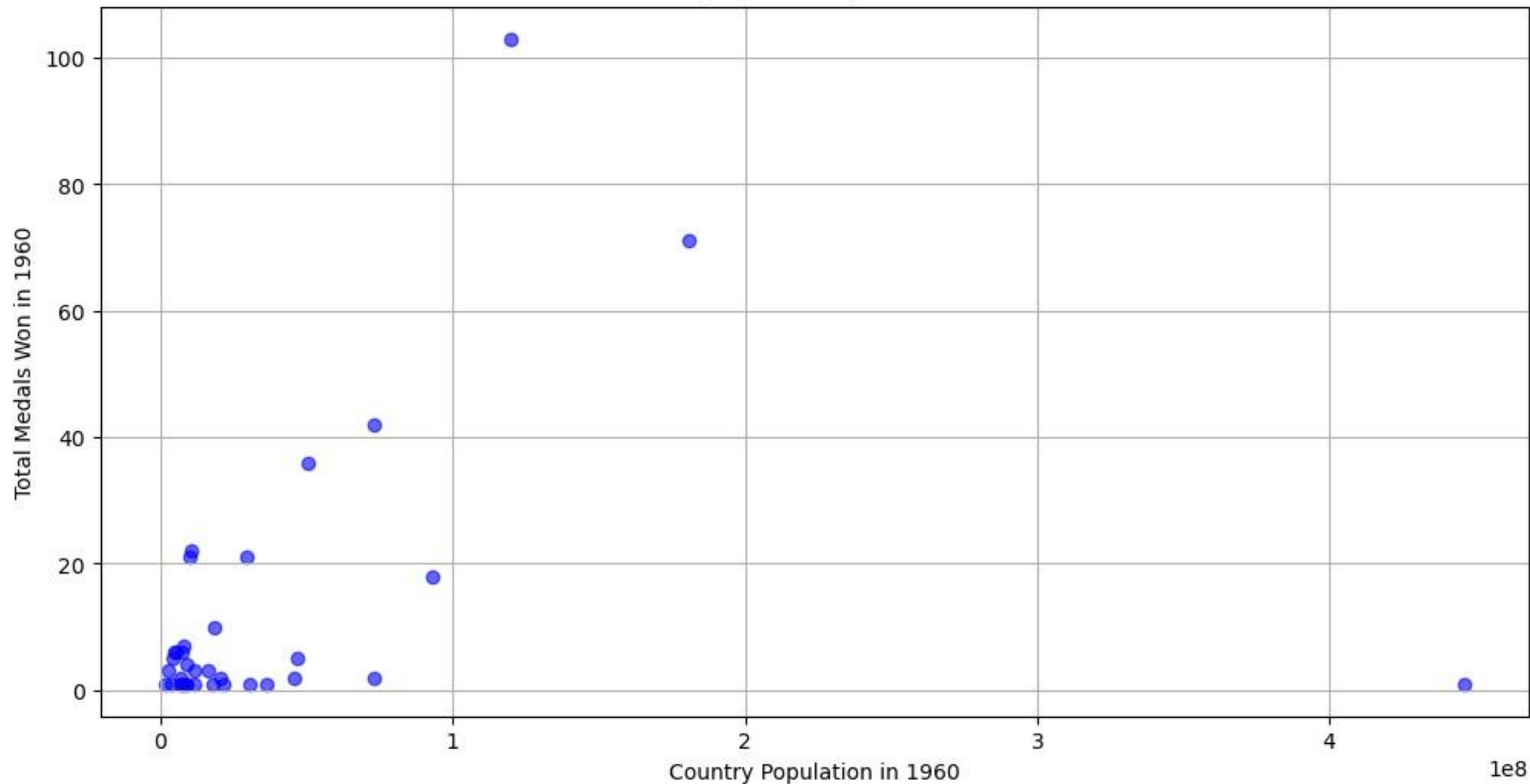60 Lowest Populated Countries in 2020 (Pop. vs Medals)

2000 Summer Olympics - Population vs Medals Won

1980 Summer Olympics - Population vs Medals Won

1960 Summer Olympics - Population vs Medals Won

# Limitations

- Assuming the GDP and Population numbers are correct in these CSV

- No Data charted for countries that did not win medals

- No Consideration for the Athletes who made it to the Finals

- How much a country invests into olympic training is probably a better indicator (India)

- Outside Factors of Cheating and Favoritism (Drugs, Biased Judges, Bribed Judges)

fppt.com

# Conclusions

Based on this Data, GDP and Population are one of the positive factors for how many total medals a country wins at the Summer Olympics.