# Milestone 4 Project Results

Riley Taylor

# Objectives

- In this presentation, I will address the following items:
  - Client Information, Initial Hypotheses, and Initial Approach
- Analysis Results
  - Correlations and Metrics
  - Relation to Initial Hypotheses
- Graphics and Visualizations

# Client Description

- I've decided to work with SportStats to analyze data about previous Olympic medal winners
  - SportsStats is a sports analysis firm that works to provide insights to their partners
- I'm looking to provide an analysis that will develop a news story or discover key health insights based on geography

# Preliminary Questions

- Is there a geographic pattern that correlates with the events that each country succeeds in?
  - How would climate affect the number of medals won in the Summer vs Winter Games?
- Is there a geographic pattern that correlates with the number of medals received by each country?

# Initial Hypothesis

- There will be a correlation between geography and performance
  - Countries with colder climates will perform best in the Winter Games
  - Countries with warmer climates will perform best in the Summer Games
- Countries with higher populations will have higher medal counts
  - A higher population will be correlated with a higher number of competitive athletes to choose for the national team for each event

# Approach

- I'll primarily be looking at the frequency of medal wins and will separate by Summer vs Winter Games
  - From there, I'll analyze by country, sport, and event
- Columns I expect to primarily analyze:
  - Team, Games, Year, Season, Sport, Event, and Medal
- Target Metric:
  - Count of medals by season and country

# Importing the Athlete Dataset

- I imported the Athlete Dataset using Python's Pandas Library

- The data was imported as athlete_data and the info is displayed

# Importing Additional Data

- I decided to import another Dataset from Kaggle to be able to analyze data from each country

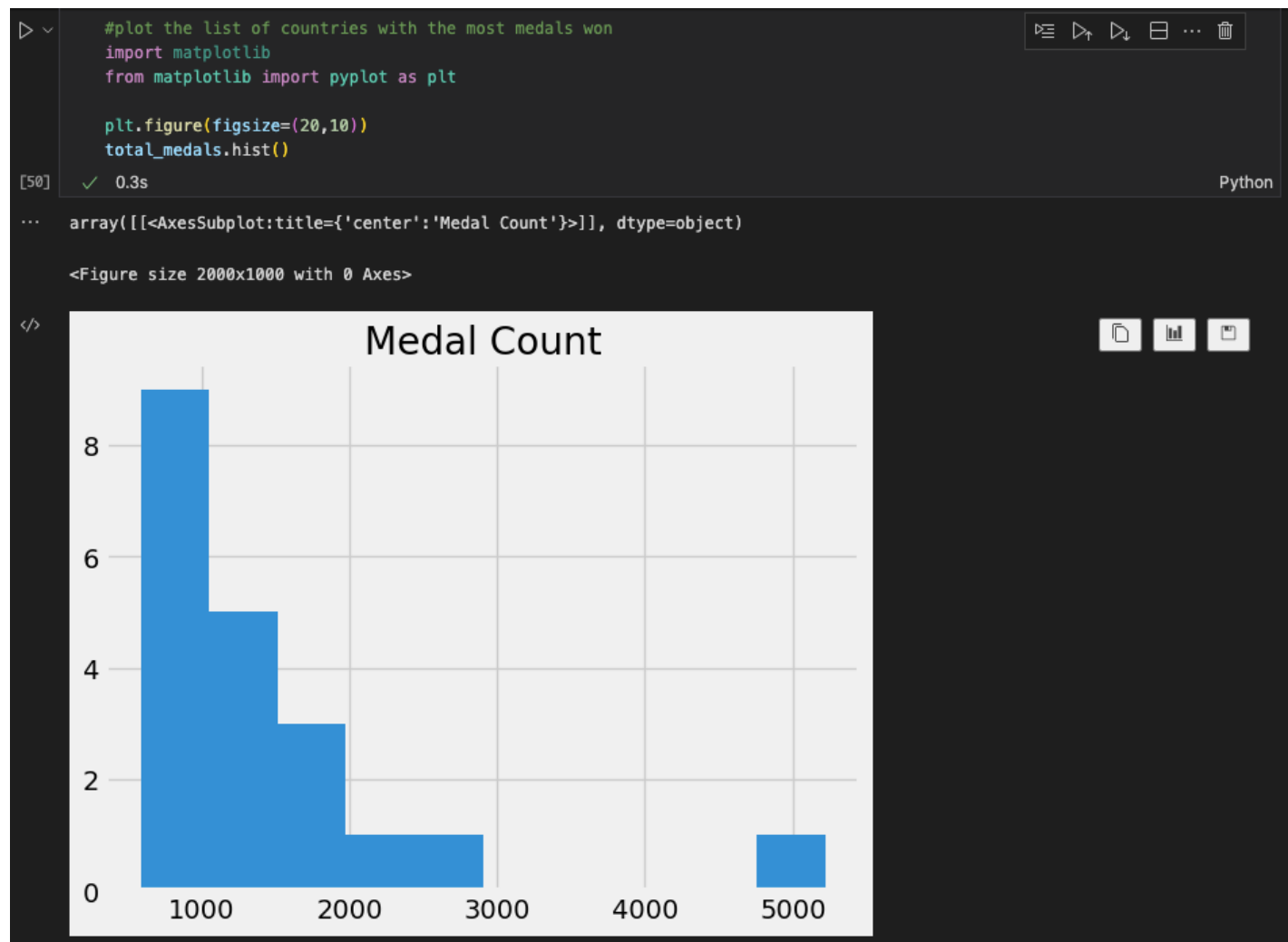- The data was imported as country_data and the info is displayed

# Initial Exploration

```
!select the 20 countries with the most medals from athlete_data set
sqlit("SELECT Team, COUNT(Medal) AS 'Medal Count' FROM athlete_data GROUP BY Team ORDER BY COUNT(Medal) DESC LIMIT 20")
```

[22]  ✓ 10.7s                                                                                    Python  Python

|    | Team          | Medal Count |
|----|---------------|-------------|
| 0  | United States | 5219        |
| 1  | Soviet Union  | 2451        |
| 2  | Germany       | 1984        |
| 3  | Great Britain | 1673        |
| 4  | France        | 1550        |
| 5  | Italy         | 1527        |
| 6  | Sweden        | 1434        |
| 7  | Australia     | 1306        |
| 8  | Canada        | 1243        |
| 9  | Hungary       | 1127        |
| 10 | Russia        | 1110        |
| 11 | Netherlands   | 988         |
| 12 | East Germany  | 941         |
| 13 | Japan         | 911         |
| 14 | Norway        | 910         |
| 15 | China         | 901         |
| 16 | Finland       | 876         |
| 17 | Romania       | 651         |
| 18 | South Korea   | 592         |
| 19 | Switzerland   | 588         |

# Technical Challenges

- Since SQL queries were run within Python using SQLite, there were a few limitations
    - However, these were later resolved once the data was cleaned accordingly, and the formatting issues were addressed within my queries
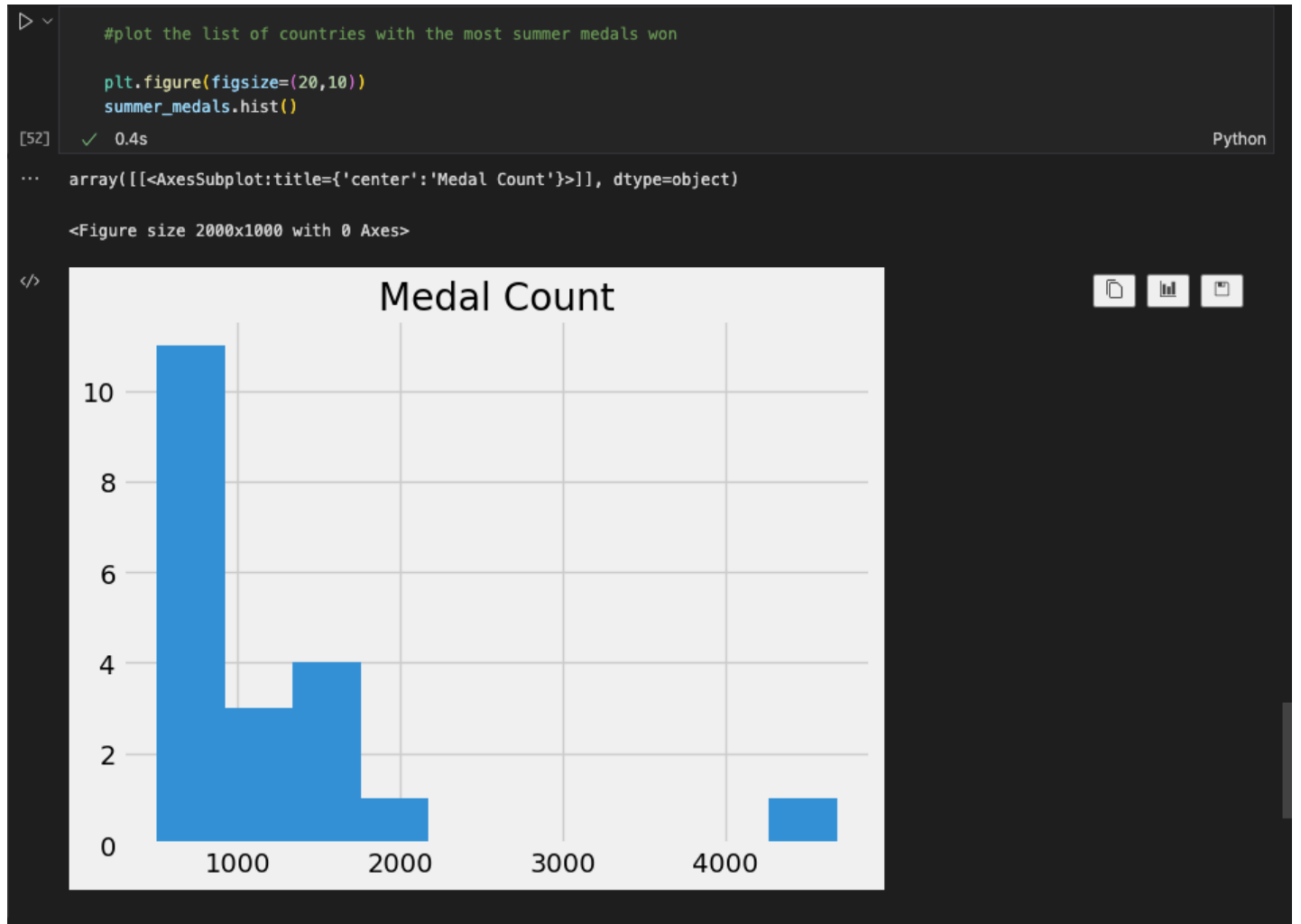
# Initial Findings

# Initial Findings (cont.)

- The lists of the countries with the most medals by Winter vs Summer Games were incredibly similar
  - The Winter Games list suggests a slight correlation of climate vs number of medals won, but wasn't statistically significant upon further analysis
- These results led me toward further analysis based on population and the percentage of the population of each country living in urban areas

# Visualizations (cont.)

# Deeper Analysis and Final Findings

- Deeper analysis showed a more significant correlation between population size and the percentage of the population living in urban areas
  - Additional analysis may be needed, but I suspect that pulling in additional data would provide more insight
  - For example, a more urbanized population may suggest a higher GDP, which I suspect may correlate with an increase in the number of medals won