Analyzing 120 years (1896-2016) Olympic Dataset

# Team Members:

- 1. Rajath John Bosco
- 2. Manshi Shah
- 3. Chaitya Mehta



#### Data Cleaning

Does team enjoy Home Advantage?

Top 5 performing events in a Country?

Analyzing number of athletes in Summer and Winter Olympics over the years

Analyzing the number of events in Summer and Winter Olympics over the years

Calculating top performing countries in summer vs winter

BMI analysis over the years

Distribution of Age vs Height from 1896 to 2016

Does GDP play a part to decide if a country would win medal or not

Calculating the medals won by country in summer and winter olympics



Olympics Dataset -

https://www.kaggle.com/heesoo37/120-years-of-olympic-history-athletes-and-results

GDP Dataset - <a href="https://datahub.io/core/gdp">https://datahub.io/core/gdp</a>

Population Dataset - <a href="https://vizhub.com/celtic660/datasets/world\_pop">https://vizhub.com/celtic660/datasets/world\_pop</a>

Cleaning data of the medal column by replacing NA with DNW(Did not Win)

Data Cleaning:

Categorizing events into Individual and Team events

We awarded 1 medal per team instead of awarding medal to each player in the team.

If a Country has won medal in soccer, it should count one single medal instead of 11 medals (since 11 players in a team)

Removing events where age of athlete is more than 80 years

After verifying data from internet, there were players who were dead, but their age was still updating. So, we removed events whose name started from "Arts"

Does Team enjoy Home Advantage?

22	Year	Host Country	Team	Season	Medal_Won_Prev_Year	Medal_Won_Host_Year	Medal_Won_Next_Year
0	1900	France	France	Summer	NaN	66.0	1.0
1	1904	United States	United States	Summer	45.0	193.0	46.0
2	1912	Sweden	Sweden	Summer	24.0	57.0	NaN
3	1924	France	France	Summer	NaN	39.0	24.0
4	1924	France	France	Winter	NaN	39.0	24.0
5	1928	Netherlands	Netherlands	Summer	10.0	22.0	8.0
6	1932	United States	United States	Summer	56.0	114.0	59.0
7	1932	United States	United States	Winter	56.0	114.0	59.0
8	1936	Germany	Germany	Winter	25.0	103.0	NaN
9	1936	Germany	Germany	Summer	25.0	103.0	NaN
10	1952	Norway	Norway	Winter	6.0	18.0	3.0
11	1952	Finland	Finland	Summer	22.0	30.0	15.0
12	1956	Sweden	Sweden	Summer	37.0	17.0	7.0
13	1956	Australia	Australia	Summer	11.0	33.0	0.0

Top 5 performances by Germany, Russia, USA

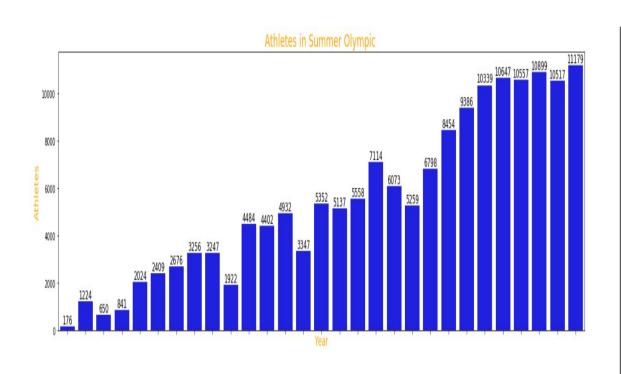
Team	Event	Gold_Medal_Count
Germany	Luge Men's Singles	6
Germany	Equestrianism Mixed Dressage, Team	5
Germany	Luge Women's Singles	5
Germany	Alpine Skiing Women's Combined	4
Germany	Canoeing Men's Canadian Doubles, 1,000 metres	4
Russia	Gymnastics Women's Uneven Bars	4
Russia	Rhythmic Gymnastics Women's Group	4
Russia	Rhythmic Gymnastics Women's Individual	4
Russia	Synchronized Swimming Women's Duet	4
Russia	Synchronized Swimming Women's Team	4
United States	Athletics Men's Long Jump	19
United States	Athletics Men's 110 metres Hurdles	17
United States	Athletics Men's 400 metres Hurdles	17
United States	Athletics Men's 400 metres	16

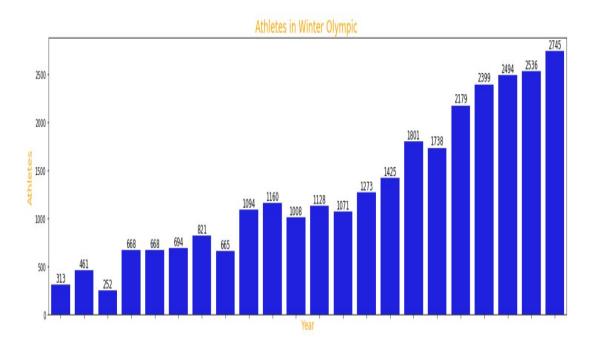
Changing
Trends in
women
participation
in Olympics

	Year	Women Athlete Count
0	1900	23
1	1904	6
2	1906	6
3	1908	44
4	1912	53
5	1920	78
6	1924	169
7	1928	340
8	1932	222
9	1936	441
10	1948	523
11	1952	629
12	1956	516
13	1960	757
14	1964	880
15	1968	994
16	1972	1266

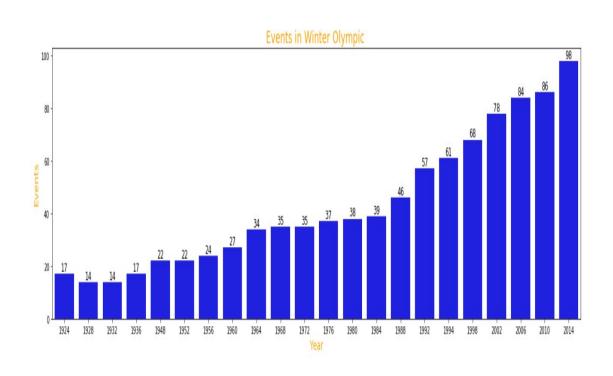
16	1972	1266
17	1976	1490
18	1980	1358
19	1984	1845
20	1988	2511
21	1992	3205
22	1994	522
23	1996	3511
24	1998	788
25	2000	4065
26	2002	885
27	2004	4288
28	2006	955
29	2008	4598
30	2010	1032
31	2012	4645
32	2014	1102
33	2016	5031

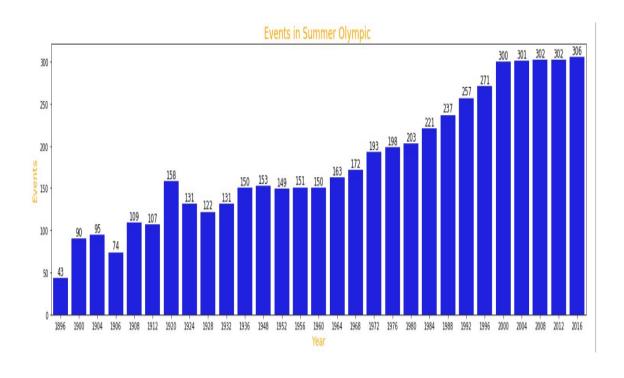
# Analyzing number of athletes in summer vs winter over the years



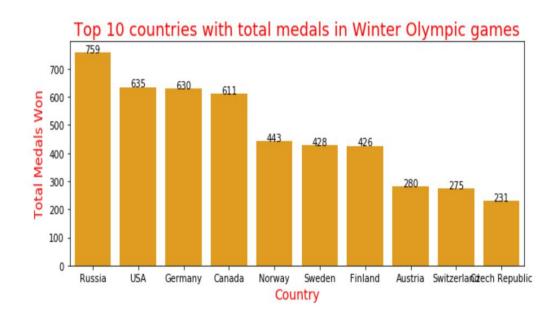


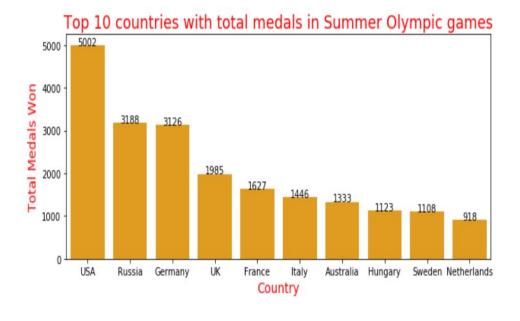
### Analyzing number of events in summer vs winter over the years



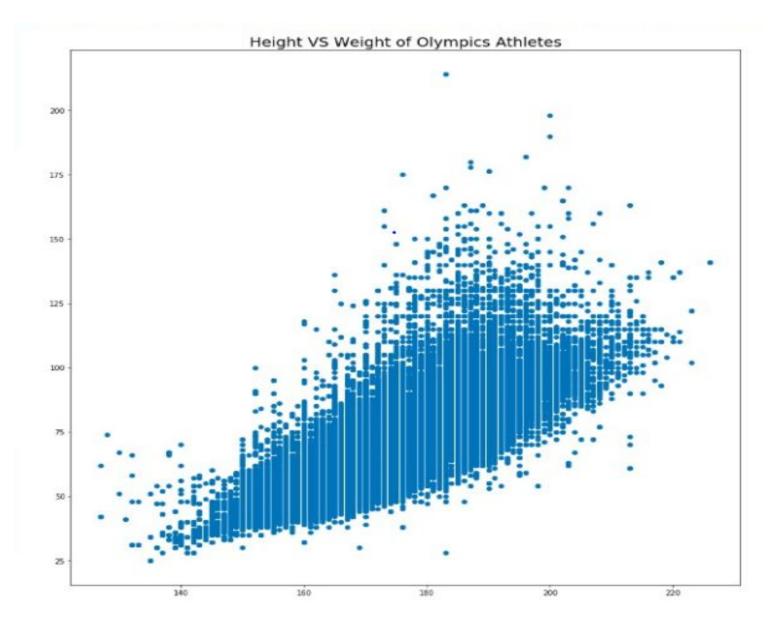


### Calculating top performing countries in summer vs winter

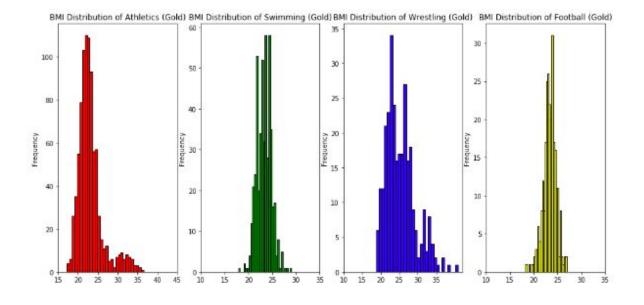




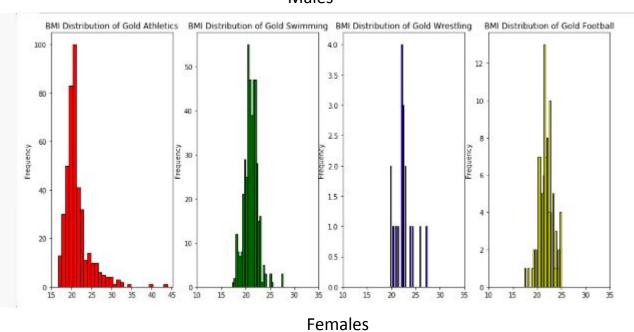
Distribution of Age vs Height from 1896 to 2016



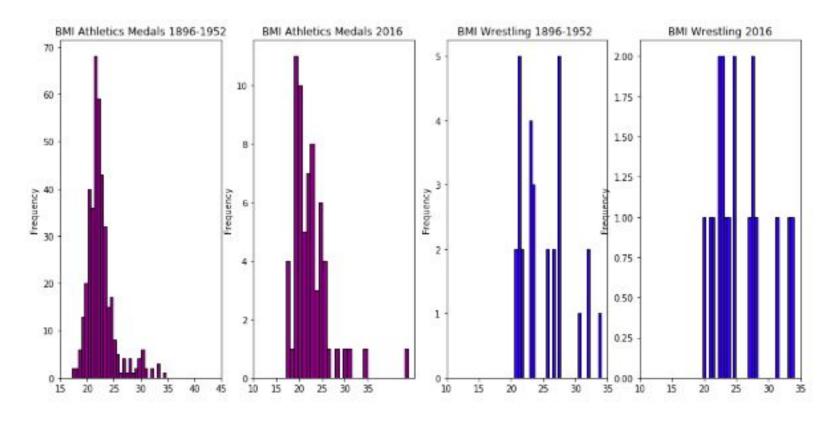




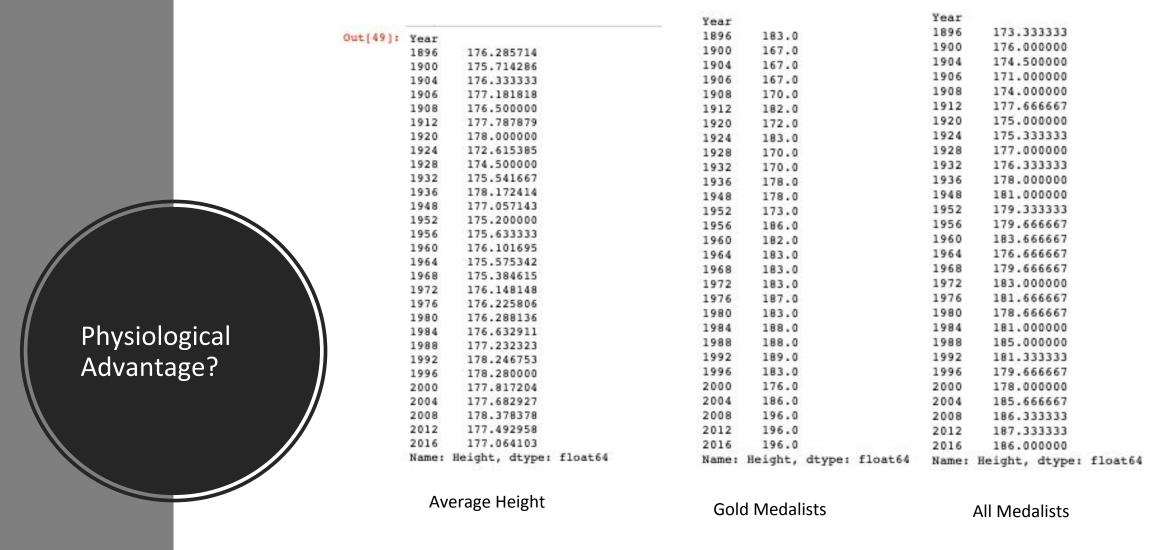
#### Males



BMI analysis over the years:



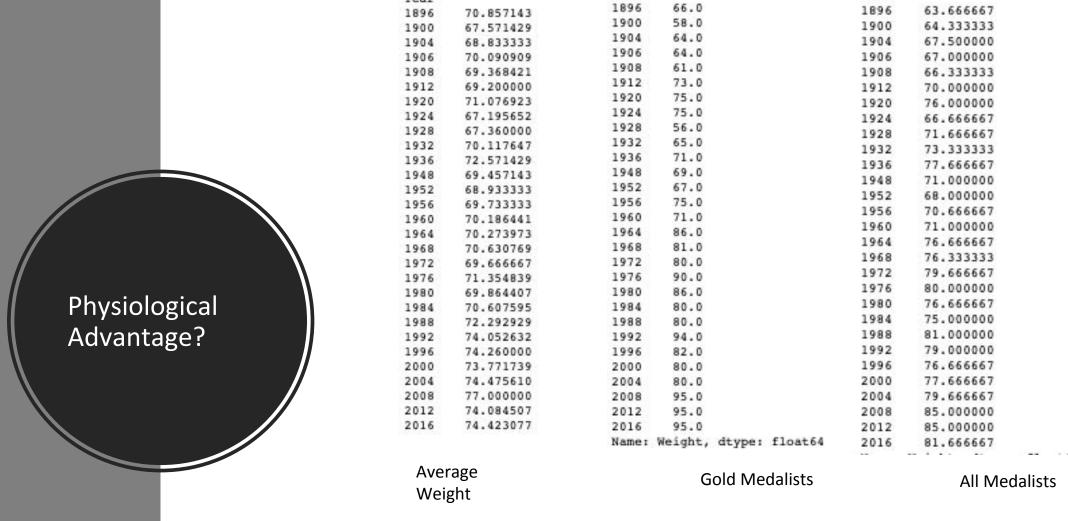
- 1. BMI of Athletes ( 1896 1928)
- 2. BMI of Athletes Gold Medalists (2016)
- 3. BMI of Wrestling Events ( 1896 1928)
- 4. BMI of Wrestling Events Gold Medalists 2016



Height Comparison of Men's 100m Runners

Standards on age, height and weight in Olympic running events for men.

https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1859632/?page=2



Year

Year

Weight Comparison of Men's 100m Runners

Year

Standards on age, height and weight in Olympic running events for men.

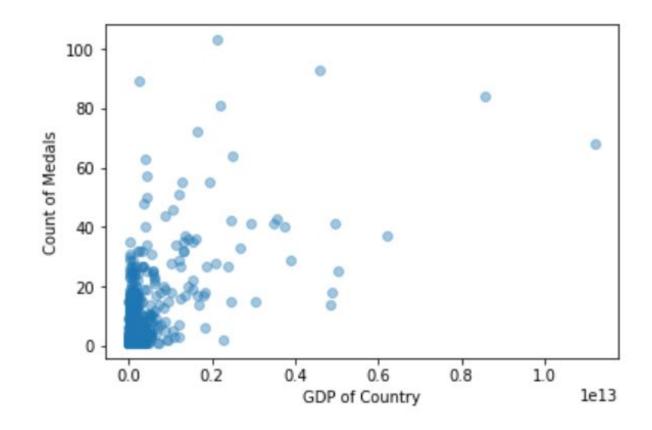


Dimension	Depende nt Variable:	Indicator	Hypot hesized Sign
Olympic Performance	Medal Count (M <sub>t</sub> )	The number of medal won by a country in a particular Olympics	N/A
	Independ ent Variables		
Demographic Environment	Populatio n (N <sub>t</sub> )	The population size of a country at a particular Olympic year	+
Economic Development	GDP per capita (Y <sub>t</sub> /N <sub>t</sub> )	The per capita GDP (measured in PPP current international dollars) of a country at a particular Olympic year	+
Social Development	Human Developm ent Index HDI <sub>t</sub>	The Human Deployment Index of a country at a particular Olympic year	(±)
Political Environment	Political System (Pol)	1 if the country is or used to be a socialist country or 0 otherwise	+
Geographic Environment	Hosting country (H <sub>t</sub> )	1 if the country is the hosting country of the year or 0 otherwise	+

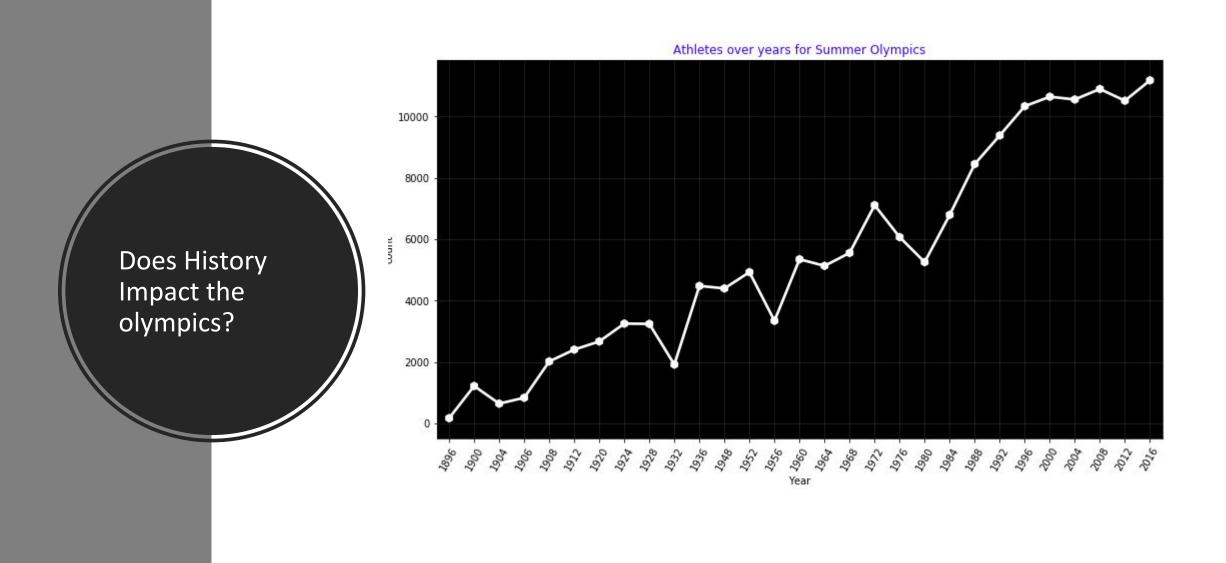
#### Model - Linear Regression

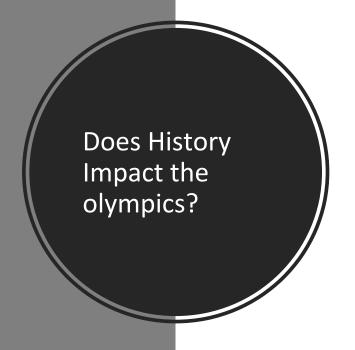
$$\begin{array}{ll} \mathbf{\hat{M}_{t}} &= \\ + \, \beta_{5} \, H_{t} + E \end{array} \qquad \beta_{0} + \beta_{1} \left( \mathbf{N}_{t} \right) + \, \beta_{2} \left( \mathbf{Y}_{t} \! / \! \mathbf{N}_{t} \right) + \, \beta_{3} H D I_{t} \, + \beta_{4} \left( Pol \right) \end{array}$$

Predictions -R-Squared ~ 0.75 RMSE ~ 7 How does GDP affect one of the most popular country ie USA, Germany, China?



Dataset - <a href="https://datahub.io/core/gdp">https://datahub.io/core/gdp</a>





#### 1936 Summer Olympics (Berlin, Germany) -

Hitler's rise to power, less countries were invited, Racism.

1940 and 1944 Summer Olympics (not held due to World War II)

1956 Summer Olympics (Melbourne, Australia) -

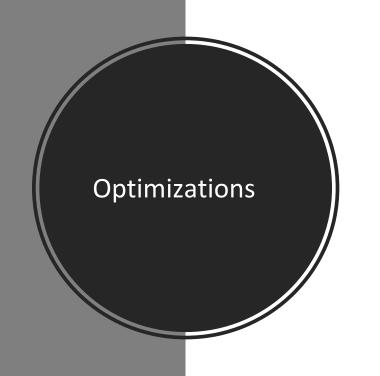
Suez Crisis (Egypt, Iraq, and Lebanon), Hungarian revolution (Netherlands, Spain, and Switzerland), Republic of China (Formosa).

1980 Summer Olympics (Moscow, Soviet Union) -

Soviet invasion of Afghanistan (USA and its allies boycotted)

1984 Summer Olympics (Los Angeles, United States) -

Russia Boycotted (Security Reasons)



- 1. Index Optimalization (merge tables by index)
- 2. Use vectorized operations (.iloc , .apply , lambda expressions)
- 3. Numba (Optimize Numpy)
- 4. Modin (read csv)

#### Reference -

- 1. <a href="https://towardsdatascience.com/understanding-the-need-for-optimization-when-using-pandas-8ce23b83330c">https://towardsdatascience.com/understanding-the-need-for-optimization-when-using-pandas-8ce23b83330c</a>
- 2. <a href="https://towardsdatascience.com/understanding-the-need-for-optimization-when-using-pandas-8ce23b83330c">https://towardsdatascience.com/understanding-the-need-for-optimization-when-using-pandas-8ce23b83330c</a>

## THANK YOU!