# 1. Step 1: Preparing for Your Proposal

## 1. Which client/dataset did you select and why?

Client 3: SportsStats (Olympics Dataset - 120 years of data)

SportsStats is a sports analysis firm partnering with local news and elite personal trainers to provide "interesting" insights to help their partners.  Insights could be patterns/trends highlighting certain groups/events/countries, etc. for the purpose of developing a news story or discovering key health insights.

Using this dataset to get some insights on games.

## 2. Describe the steps you took to import and clean the data.

Downloaded the data from

https://www.dropbox.com/sh/0wqw8fmiwrzr8ef/AABQijjQM522INXX1FCdamzma?dl=0

The source provides two tables.

1. athlete_events.csv and

2. noc_regions.csv

```
Number of records athlete_events: 271116
Number of records noc_regions: 230
```

Columns and data type for athlete_events table.

```
 1  df_athlete_events.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 271116 entries, 0 to 271115
Data columns (total 15 columns):
 #   Column  Non-Null Count   Dtype
---  ------  --------------   -----
 0   ID      271116 non-null  int64
 1   Name    271116 non-null  object
 2   Sex     271116 non-null  object
 3   Age     261642 non-null  float64
 4   Height  210945 non-null  float64
 5   Weight  208241 non-null  float64
 6   Team    271116 non-null  object
 7   NOC     271116 non-null  object
 8   Games   271116 non-null  object
 9   Year    271116 non-null  int64
 10  Season  271116 non-null  object
 11  City    271116 non-null  object
 12  Sport   271116 non-null  object
 13  Event   271116 non-null  object
 14  Medal   39783 non-null   object
dtypes: float64(3), int64(2), object(10)
memory usage: 31.0+ MB
```

Columns and data type for noc_region.

```
 1  df_noc_regions.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 230 entries, 0 to 229
Data columns (total 3 columns):
 #   Column  Non-Null Count  Dtype
---  ------  --------------  -----
 0   NOC     230 non-null    object
 1   region  227 non-null    object
 2   notes   21 non-null     object
dtypes: object(3)
memory usage: 5.5+ KB
```

## 3. Perform initial exploration of data and provide some screenshots or display some stats of the data you are looking at.

```
1  # moments of numerical data
2  df_athlete_events.describe()
```
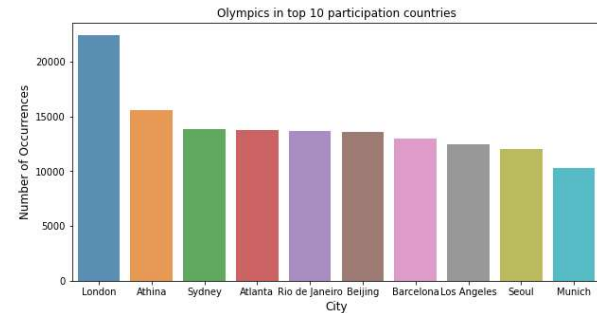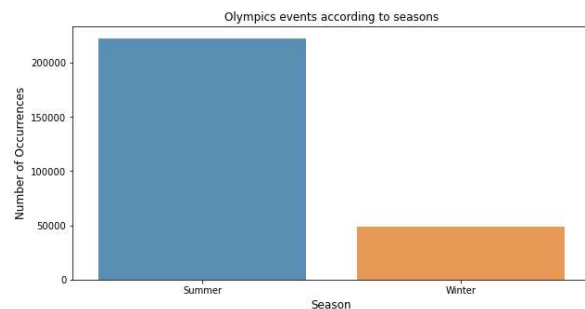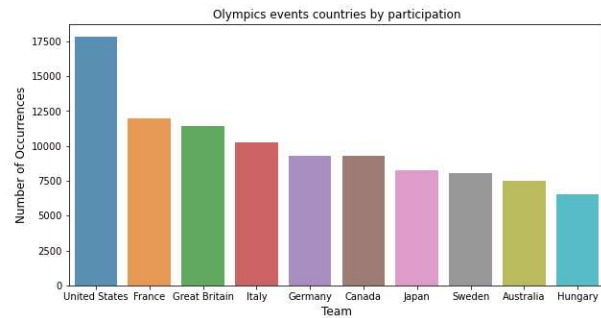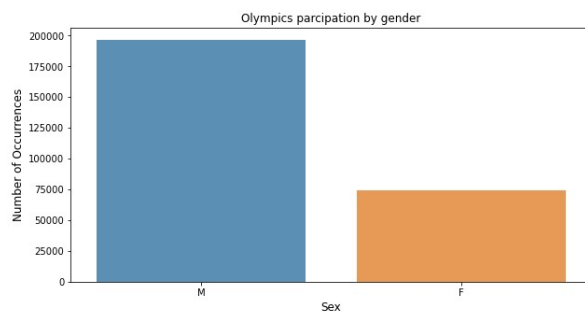
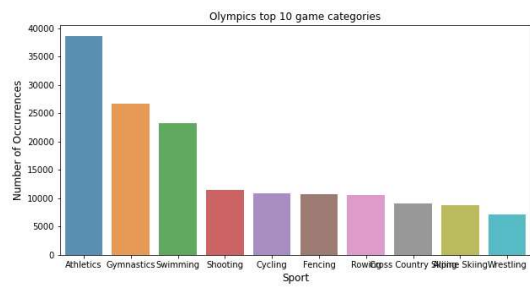|       | ID            | Age           | Height        | Weight        | Year          |
|-------|---------------|---------------|---------------|---------------|---------------|
| count | 271116.000000 | 261642.000000 | 210945.000000 | 208241.000000 | 271116.000000 |
| mean  | 68248.954396  | 25.556898     | 175.338970    | 70.702393     | 1978.378480   |
| std   | 39022.286345  | 6.393561      | 10.518462     | 14.348020     | 29.877632     |
| min   | 1.000000      | 10.000000     | 127.000000    | 25.000000     | 1896.000000   |
| 25%   | 34643.000000  | 21.000000     | 168.000000    | 60.000000     | 1960.000000   |
| 50%   | 68205.000000  | 24.000000     | 175.000000    | 70.000000     | 1988.000000   |
| 75%   | 102097.250000 | 28.000000     | 183.000000    | 79.000000     | 2002.000000   |
| max   | 135571.000000 | 97.000000     | 226.000000    | 214.000000    | 2016.000000   |

## Missing values in athlete_events table

```
1  # as we see above the counts a
2  # lets see the null values in
3  df_athlete_events.isna().sum()
```

```
ID            0
Name          0
Sex           0
Age        9474
Height    60171
Weight    62875
Team          0
NOC           0
Games         0
Year          0
Season        0
City          0
Sport         0
Event         0
Medal    231333
dtype: int64
```

## Visualizing categorical values in athlete_events table

Olympics top 10 game categories

Number of Occurrences (y-axis): 0, 5000, 10000, 15000, 20000, 25000, 30000, 35000, 40000

Sport (x-axis): Athletics, Gymnastics, Swimming, Shooting, Cycling, Fencing, Rowing, Cross Country Skiing, Alpine Skiing, Wrestling

## 4. Entity diagrams

### Initial ER diagram for athlete_events table

```
#    Column   Non-Null Count     Dtype
---  ------   --------------     -----
 0   ID        271116 non-null   int64
 1   Name      271116 non-null   object
 2   Sex       271116 non-null   object
 3   Age       261642 non-null   float64
 4   Height    210945 non-null   float64
 5   Weight    208241 non-null   float64
 6   Team      271116 non-null   object
 7   NOC       271116 non-null   object
 8   Games     271116 non-null   object
 9   Year      271116 non-null   int64
10   Season    271116 non-null   object
11   City      271116 non-null   object
12   Sport     271116 non-null   object
13   Event     271116 non-null   object
14   Medal      39783 non-null   object
```

### Initial ER diagram for noc_regions table

```
#    Column   Non-Null Count    Dtype
---  ------   --------------    -----
 0   NOC       230 non-null     object
 1   region    227 non-null     object
 2   notes      21 non-null     object
```

## Step 2: Develop Project Proposal

There are two tables provided. Athlete events table contains most of the information's which later further split into multiple tables like

1. countries,
2. sports,
3. medals,
4. categories.

**Questions**

1. are the winning countries make more medals in future games. Did they find the secret of winning?
2. Are there new countries emerging with more than 10 medals?
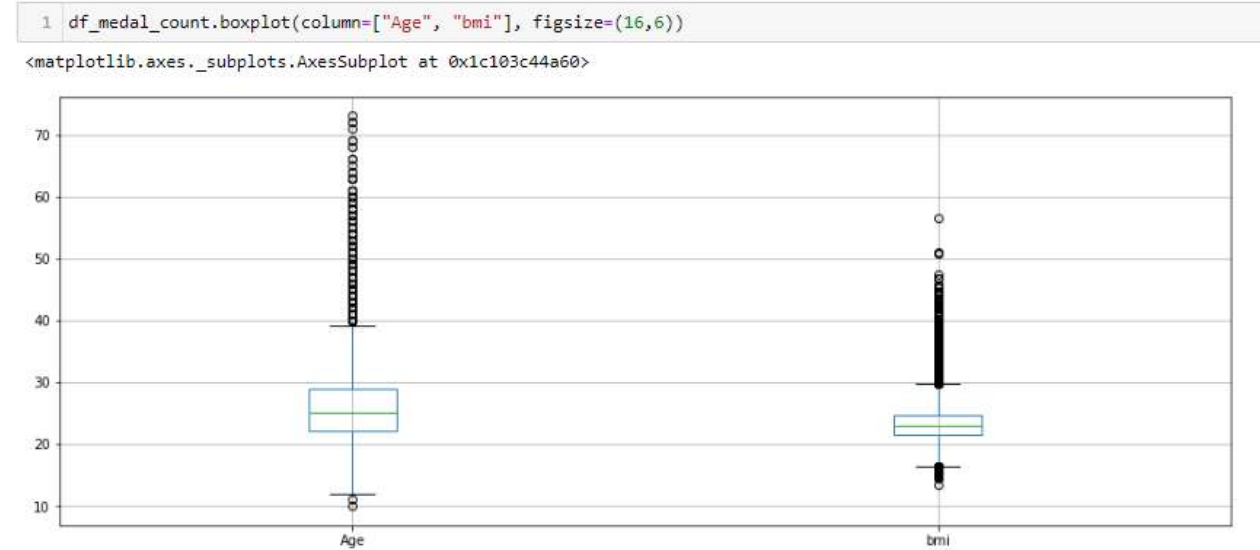
**Hypothesis**

1. The athletes who participate more Olympic games wins more medals in later games.

2. The motivation outweighs physical characteristics to win the game is the hypothesis.

**Approach**

1. All numerical and categorical fields will be used.

2. Find the relation between physical characteristics with respect to winning the game

3. physical characteristics for top 5 games across years to see how it varies?
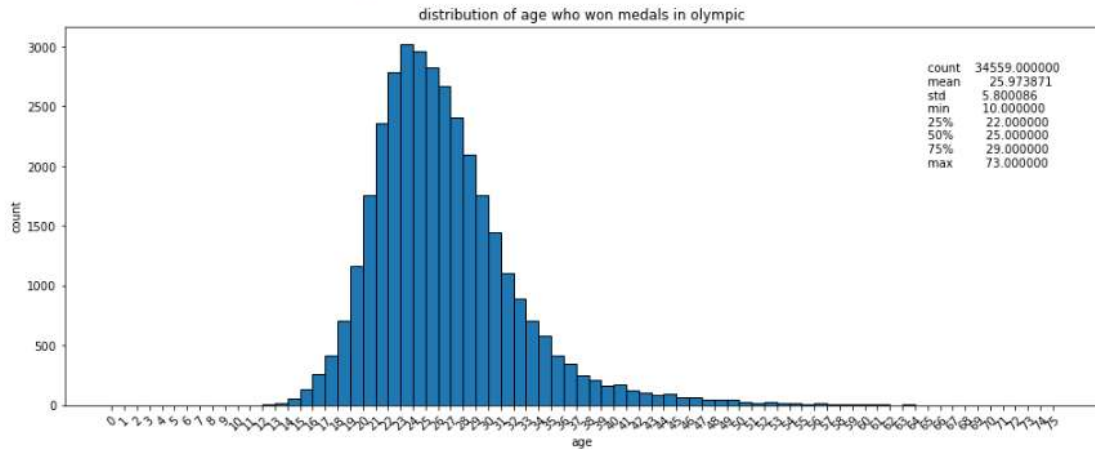
The bmi is calculated using a standard formula (weight/height**2)

Following shows the box plot for age and bmi. And we can see the both distributions are right tailed.

```
1  df_medal_count.boxplot(column=["Age", "bmi"], figsize=(16,6))
```

<matplotlib.axes._subplots.AxesSubplot at 0x1c103c44a60>

Distribution of Age and Winning a medal, The distribution is right tailed.
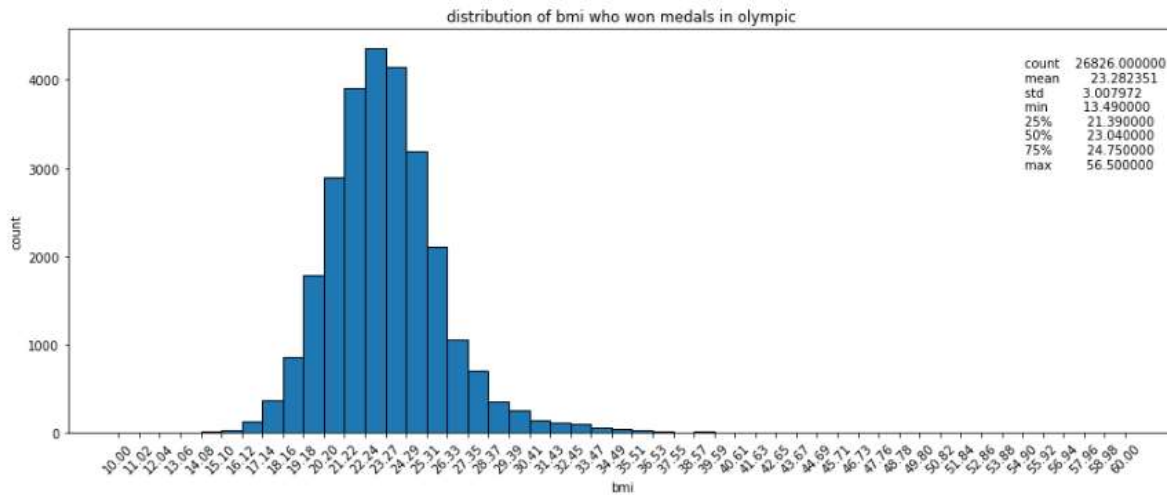
```
In [26]:   1  # distribution of age who won medals in olympic
           2  import matplotlib.pyplot as plt
           3  plt.figure(figsize=(16,6))
           4  plt.hist(df_medal_count['Age'].values, bins=age_bins, edgecolor="k")
           5  plt.title("distribution of age who won medals in olympic")
           6  plt.xticks(rotation=45)
           7  plt.xlabel("age")
           8  plt.ylabel("count")
           9  plt.xticks(age_bins);
          10  plt.text(65, 2000, df_medal_count['Age'].describe().to_string());
```



distribution of age who won medals in olympic

| | |
|---|---|
| count | 34559.000000 |
| mean | 25.973871 |
| std | 5.800086 |
| min | 10.000000 |
| 25% | 22.000000 |
| 50% | 25.000000 |
| 75% | 29.000000 |
| max | 73.000000 |

Interesting to note that the youngest winner is 10 year old and oldest winner is 73 year old.

Distribution of BMI and Winning a Medal, this distribution also right tailed.

```
1   # distribution of age who won medals in olympic
2   import matplotlib.pyplot as plt
3   plt.figure(figsize=(16,6))
4   plt.hist(df_medal_count['bmi'].values, bins=bmi_bins, edgecolor="k")
5   plt.title("distribution of bmi who won medals in olympic")
6   plt.xticks(rotation=45)
7   plt.xlabel("bmi")
8   plt.ylabel("count")
9   plt.xticks(bmi_bins)
10  plt.text(55, 3000, df_medal_count['bmi'].describe().to_string());
```



distribution of bmi who won medals in olympic

| | |
|---|---|
| count | 26826.000000 |
| mean | 23.282351 |
| std | 3.007972 |
| min | 13.490000 |
| 25% | 21.390000 |
| 50% | 23.040000 |
| 75% | 24.750000 |
| max | 56.500000 |

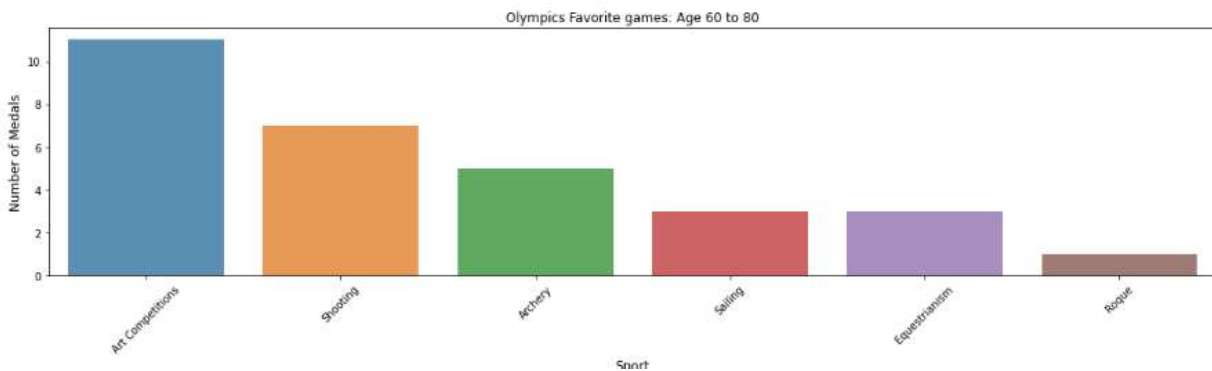There is a game for everyone i.e, underweight, normal, overweight, and obese.

From BMI 13 underweight to BMI 56 obese has games to win.

From age 60-80 years old participated in games and won the medals.

```
1  # The Sports senior citizens playing
2  query = " \
3      select Year,Team, Sport,Name,Age,Height,Weight,bmi,count(Name) as medal_count \
4      from df_medal_count \
5      where Age between 60 and 80 \
6      group by Year,Sport,Name,Age,Height,Weight,bmi \
7      order by Age desc"
8
9  df_above_60_80 = pysqldf(query)
10
11 city_count = df_above_60_80["Sport"].value_counts()
12 city_count = city_count[:20,]
13 plt.figure(figsize=(16,5))
14 sns.barplot( city_count.index, city_count.values, alpha=0.8)
15 plt.title('Olympics Favorite games: Age 60 to 80')
16 plt.xticks(rotation=45)
17 plt.ylabel('Number of Medals', fontsize=12)
18 plt.xlabel('Sport', fontsize=12)
19 plt.tight_layout()
20 plt.show()
```
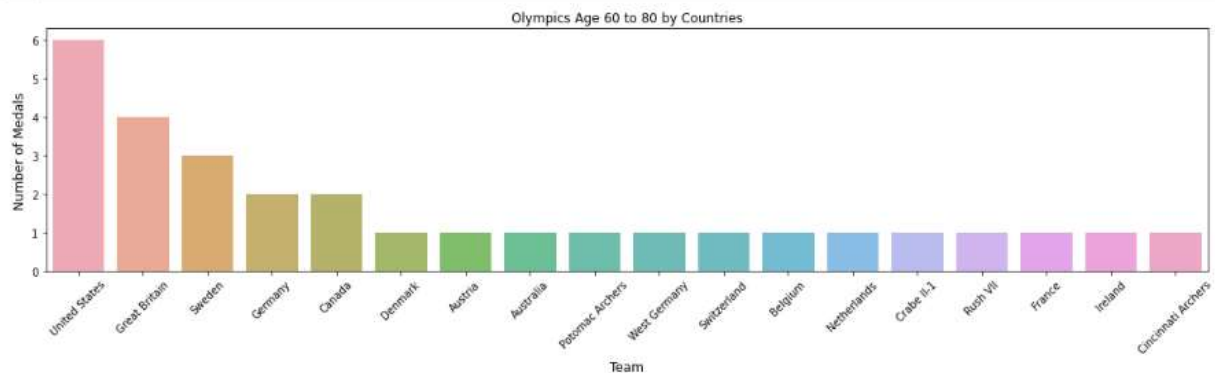


From age 60-80 years represented following teams and won gold.
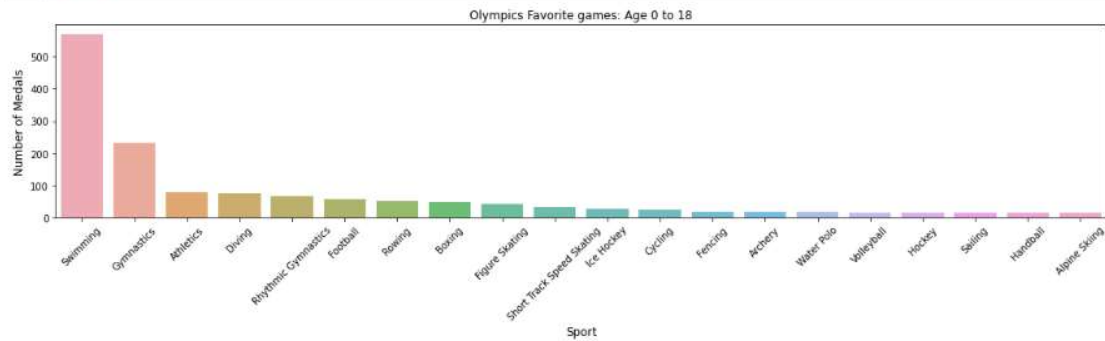
```
1  city_count = df_above_60_80["Team"].value_counts()
2  city_count = city_count[:20,]
3  plt.figure(figsize=(16,5))
4  sns.barplot( city_count.index, city_count.values, alpha=0.8)
5  plt.title('Olympics Age 60 to 80 by Countries')
6  plt.xticks(rotation=45)
7  plt.ylabel('Number of Medals', fontsize=12)
8  plt.xlabel('Team', fontsize=12)
9  plt.tight_layout()
10 plt.show()
```
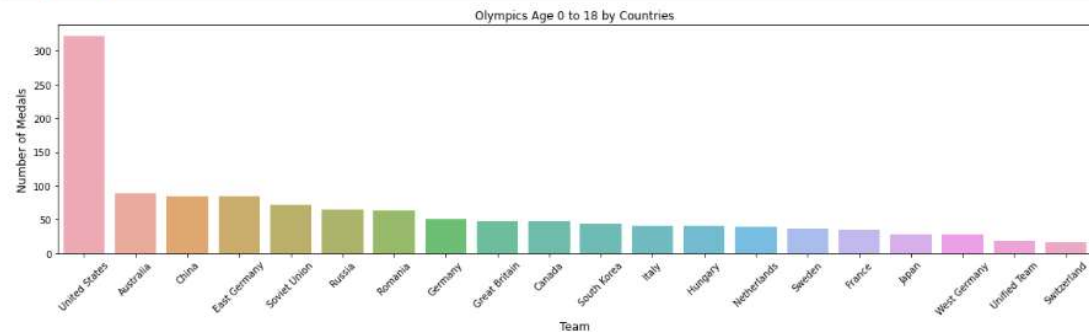
## Teenager's also played in following games

```
1  # The Sports senior citizens playing
2  query = " \
3      select Year,Team, Sport,Name,Age,Height,Weight,count(Name) as medal_count \
4      from df_medal_count \
5      where Age between 0 and 18 \
6      group by Year,Sport,Name,Age,Height,Weight \
7      order by Age desc"
8
9  df_above_0_18 = pysqldf(query)
10
11 city_count = df_above_0_18["Sport"].value_counts()
12 city_count = city_count[:20,]
13 plt.figure(figsize=(16,5))
14 sns.barplot( city_count.index, city_count.values, alpha=0.8)
15 plt.title('Olympics Favorite games: Age 0 to 18')
16 plt.xticks(rotation=45)
17 plt.ylabel('Number of Medals', fontsize=12)
18 plt.xlabel('Sport', fontsize=12)
19 plt.tight_layout()
20 plt.show()
```



Olympics Favorite games: Age 0 to 18

## Teenager's represented following Teams

```
1  # The Sports senior citizens playing
2  query = " \
3      select Year,Team, Sport,Name,Age,Height,Weight,count(Name) as medal_count \
4      from df_medal_count \
5      where Age between 0 and 18 \
6      group by Year,Sport,Name,Age,Height,Weight \
7      order by Age desc"
8
9  df_above_0_18 = pysqldf(query)
10
11 city_count = df_above_0_18["Team"].value_counts()
12 city_count = city_count[:20,]
13 plt.figure(figsize=(16,5))
14 sns.barplot( city_count.index, city_count.values, alpha=0.8)
15 plt.title('Olympics Age 0 to 18 by Countries')
16 plt.xticks(rotation=45)
17 plt.ylabel('Number of Medals', fontsize=12)
18 plt.xlabel('Team', fontsize=12)
19 plt.tight_layout()
20 plt.show()
```
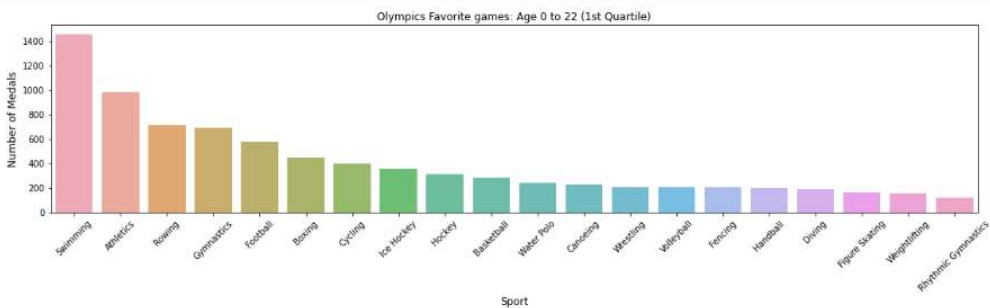


Olympics Age 0 to 18 by Countries

The 1st Quartile of Age group won in following games.

```
1  # The Sports senior citizens playing
2  query = " \
3      select Year,Team, Sport,Name,Age,Height,Weight,count(Name) as medal_count \
4      from df_medal_count \
5      where Age between 0 and 22 \
6      group by Year,Sport,Name,Age,Height,Weight \
7      order by Age desc"
8
9  df_above_0_22 = pysqldf(query)
10
11 city_count = df_above_0_22["Sport"].value_counts()
12 city_count = city_count[:20,]
13 plt.figure(figsize=(16,5))
14 s=sns.barplot(city_count.index, city_count.values, alpha=0.8)
15
16 plt.title('Olympics Favorite games: Age 0 to 22 (1st Quartile)')
17 plt.xticks(rotation=45)
18 plt.ylabel('Number of Medals', fontsize=12)
19 plt.xlabel('Sport', fontsize=12)
20 plt.tight_layout()
21 plt.show()
```



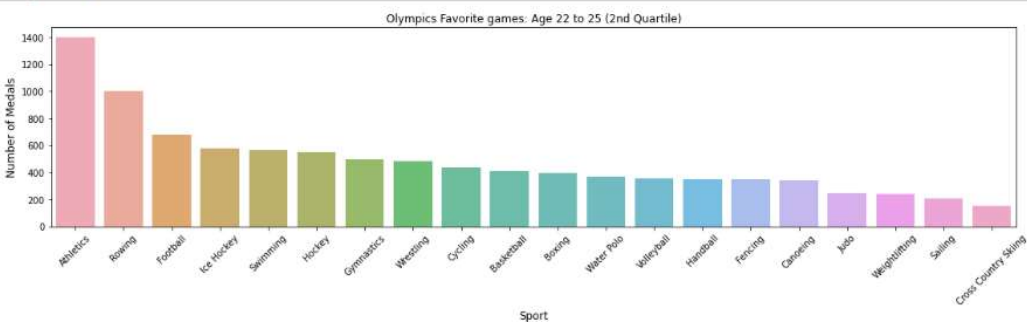Olympics Favorite games: Age 0 to 22 (1st Quartile)

2nd Quartile of Age group won following games

```
1  # The Sports senior citizens playing
2  query = " \
3      select Year,Team, Sport,Name,Age,Height,Weight,count(Name) as medal_count \
4      from df_medal_count \
5      where Age between 22 and 25 \
6      group by Year,Sport,Name,Age,Height,Weight \
7      order by Age desc"
8  df_above_22_25 = pysqldf(query)
9
10 city_count = df_above_22_25["Sport"].value_counts()
11 city_count = city_count[:20,]
12 plt.figure(figsize=(16,5))
13 sns.barplot( city_count.index, city_count.values, alpha=0.8)
14 plt.title('Olympics Favorite games: Age 22 to 25 (2nd Quartile)')
15 plt.xticks(rotation=45)
16 plt.ylabel('Number of Medals', fontsize=12)
17 plt.xlabel('Sport', fontsize=12)
18 plt.tight_layout()
19 plt.show()
```



Olympics Favorite games: Age 22 to 25 (2nd Quartile)

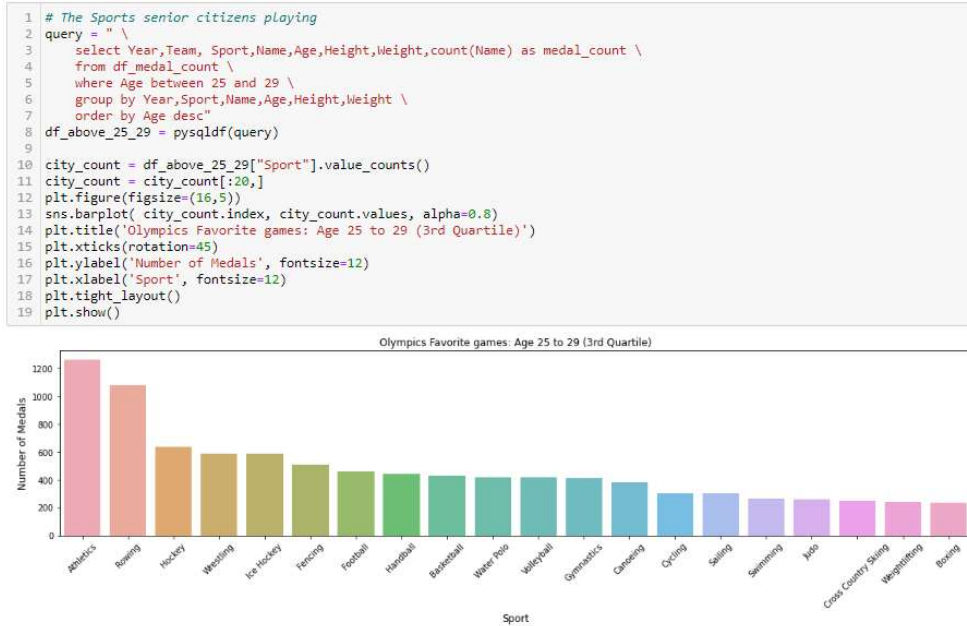# 3rd Quartile of Age group won following games.

```
1  # The Sports senior citizens playing
2  query = " \
3      select Year,Team, Sport,Name,Age,Height,Weight,count(Name) as medal_count \
4      from df_medal_count \
5      where Age between 25 and 29 \
6      group by Year,Sport,Name,Age,Height,Weight \
7      order by Age desc"
8  df_above_25_29 = pysqldf(query)
9
10 city_count = df_above_25_29["Sport"].value_counts()
11 city_count = city_count[:20,]
12 plt.figure(figsize=(16,5))
13 sns.barplot( city_count.index, city_count.values, alpha=0.8)
14 plt.title('Olympics Favorite games: Age 25 to 29 (3rd Quartile)')
15 plt.xticks(rotation=45)
16 plt.ylabel('Number of Medals', fontsize=12)
17 plt.xlabel('Sport', fontsize=12)
18 plt.tight_layout()
19 plt.show()
```
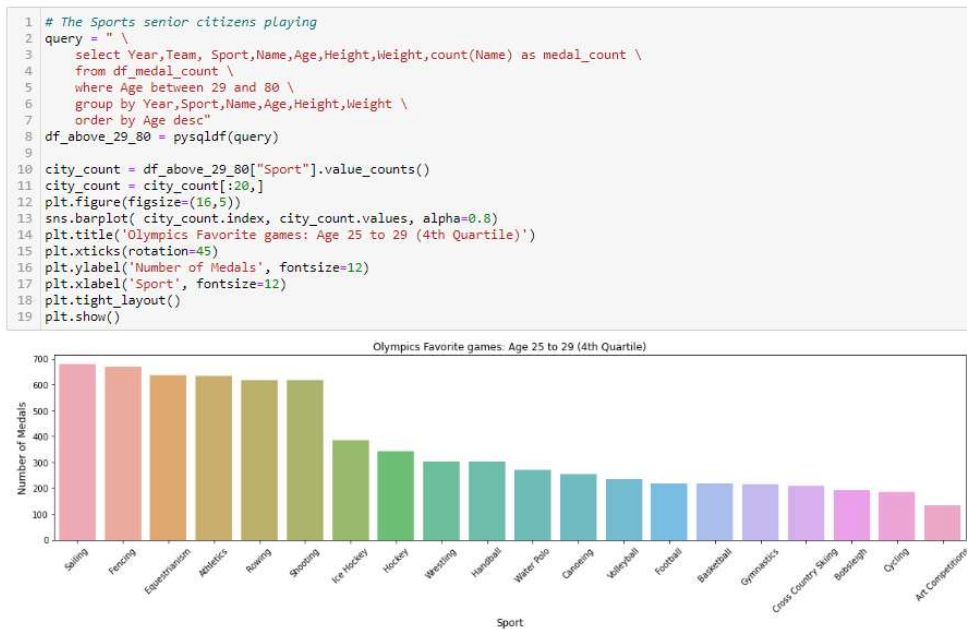


Olympics Favorite games: Age 25 to 29 (3rd Quartile)

## The 4th Quartile won in following games.

```
1  # The Sports senior citizens playing
2  query = " \
3      select Year,Team, Sport,Name,Age,Height,Weight,count(Name) as medal_count \
4      from df_medal_count \
5      where Age between 29 and 80 \
6      group by Year,Sport,Name,Age,Height,Weight \
7      order by Age desc"
8  df_above_29_80 = pysqldf(query)
9
10 city_count = df_above_29_80["Sport"].value_counts()
11 city_count = city_count[:20,]
12 plt.figure(figsize=(16,5))
13 sns.barplot( city_count.index, city_count.values, alpha=0.8)
14 plt.title('Olympics Favorite games: Age 25 to 29 (4th Quartile)')
15 plt.xticks(rotation=45)
16 plt.ylabel('Number of Medals', fontsize=12)
17 plt.xlabel('Sport', fontsize=12)
18 plt.tight_layout()
19 plt.show()
```



Olympics Favorite games: Age 25 to 29 (4th Quartile)
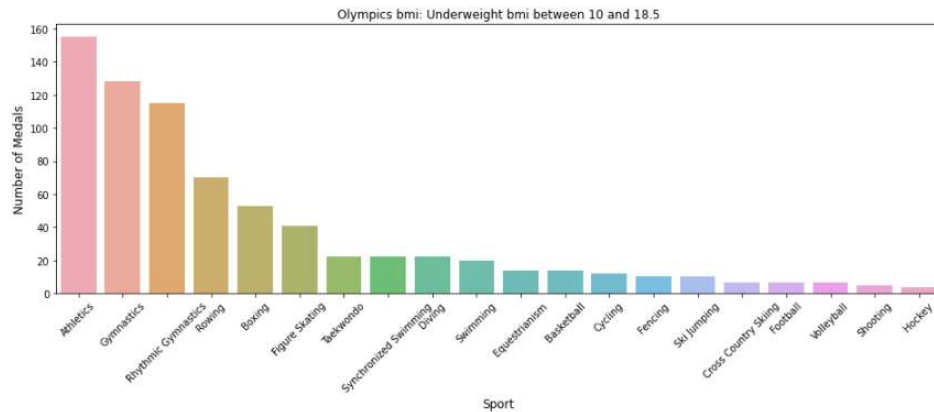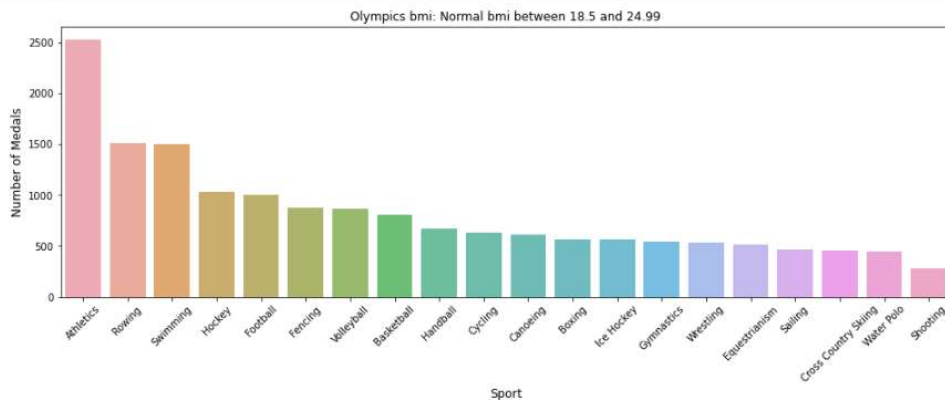
The 1ˢᵗ Quartile of bmi group won in following games.

```
1  query = " \
2      select Year,Team, Sport,Name,Age,Height,Weight,bmi,count(Name) as medal_count \
3      from df_medal_count \
4      where bmi between 10 and 18.5 \
5      group by Year,Sport,Name,Age,Height,Weight \
6      order by Age desc"
7  df_bmi_1stQ = pysqldf(query)
8
9  city_count = df_bmi_1stQ["Sport"].value_counts()
10 city_count = city_count[:20,]
11 plt.figure(figsize=(16,5))
12 sns.barplot( city_count.index, city_count.values, alpha=0.8)
13 plt.title('Olympics bmi: Underweight bmi between 10 and 18.5')
14 plt.xticks(rotation=45)
15 plt.ylabel('Number of Medals', fontsize=12)
16 plt.xlabel('Sport', fontsize=12)
17 plt.show()
```



Olympics bmi: Underweight bmi between 10 and 18.5

The 2nd Quartile of bmi group won in following games

```
1  # The Sports senior citizens playing
2  query = " \
3      select Year,Team, Sport,Name,Age,Height,Weight,bmi,count(Name) as medal_count \
4      from df_medal_count \
5      where bmi between 18.5 and 24.99 \
6      group by Year,Sport,Name,Age,Height,Weight \
7      order by Age desc"
8  df_bmi_1stQ = pysqldf(query)
9
10 city_count = df_bmi_1stQ["Sport"].value_counts()
11 city_count = city_count[:20,]
12 plt.figure(figsize=(16,5))
13 sns.barplot( city_count.index, city_count.values, alpha=0.8)
14 plt.title('Olympics bmi: Normal bmi between 18.5 and 24.99')
15 plt.xticks(rotation=45)
16 plt.ylabel('Number of Medals', fontsize=12)
17 plt.xlabel('Sport', fontsize=12)
18 plt.show()
```



Olympics bmi: Normal bmi between 18.5 and 24.99
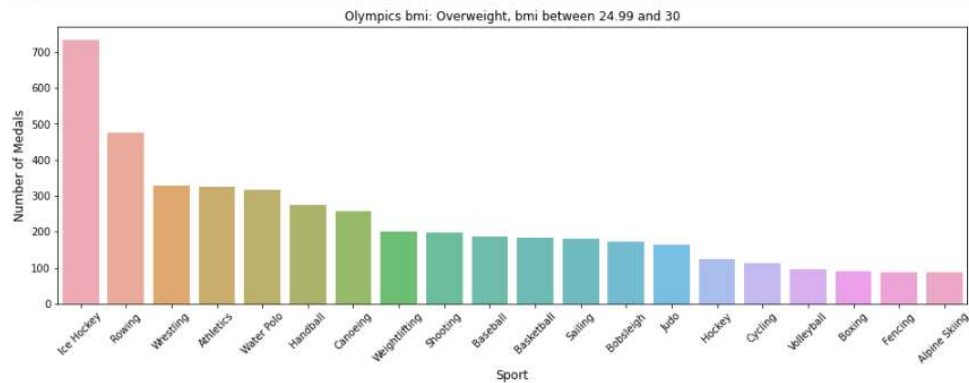
# The 3rd Quartile of bmi group won in following games

```
1  query = " \
2      select Year,Team, Sport,Name,Age,Height,Weight,bmi,count(Name) as medal_count \
3      from df_medal_count \
4      where bmi between 24.99 and 30 \
5      group by Year,Sport,Name,Age,Height,Weight \
6      order by Age desc"
7  df_bmi_1stQ = pysqldf(query)
8
9  city_count = df_bmi_1stQ["Sport"].value_counts()
10 city_count = city_count[:20,]
11 plt.figure(figsize=(16,5))
12 sns.barplot( city_count.index, city_count.values, alpha=0.8)
13 plt.title('Olympics bmi: Overweight, bmi between 24.99 and 30')
14 plt.xticks(rotation=45)
15 plt.ylabel('Number of Medals', fontsize=12)
16 plt.xlabel('Sport', fontsize=12)
17 plt.show()
```
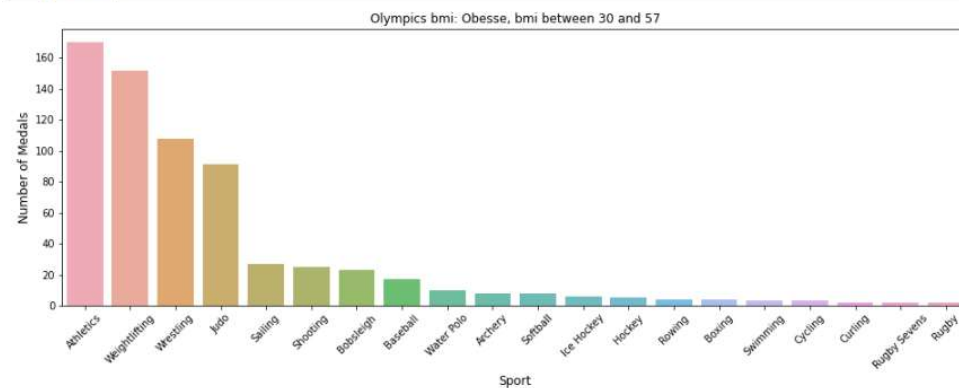


Olympics bmi: Overweight, bmi between 24.99 and 30

# The 4th Quartile of bmi group won in following games

```
1  query = " \
2      select Year,Team, Sport,Name,Age,Height,Weight,bmi,count(Name) as medal_count \
3      from df_medal_count \
4      where bmi between 30 and 57 \
5      group by Year,Sport,Name,Age,Height,Weight \
6      order by Age desc"
7  df_bmi_1stQ = pysqldf(query)
8
9  city_count = df_bmi_1stQ["Sport"].value_counts()
10 city_count = city_count[:20,]
11 plt.figure(figsize=(16,5))
12 sns.barplot( city_count.index, city_count.values, alpha=0.8)
13 plt.title('Olympics bmi: Obesse, bmi between 30 and 57')
14 plt.xticks(rotation=45)
15 plt.ylabel('Number of Medals', fontsize=12)
16 plt.xlabel('Sport', fontsize=12)
17 plt.show()
```
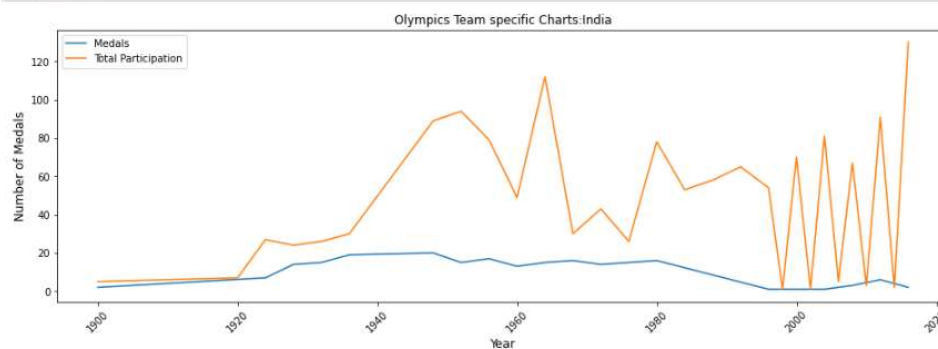


Olympics bmi: Obesse, bmi between 30 and 57

Few games shows more participation brings more winning. But for few other countries, the relation not showing. So the hypothesis is partially true.
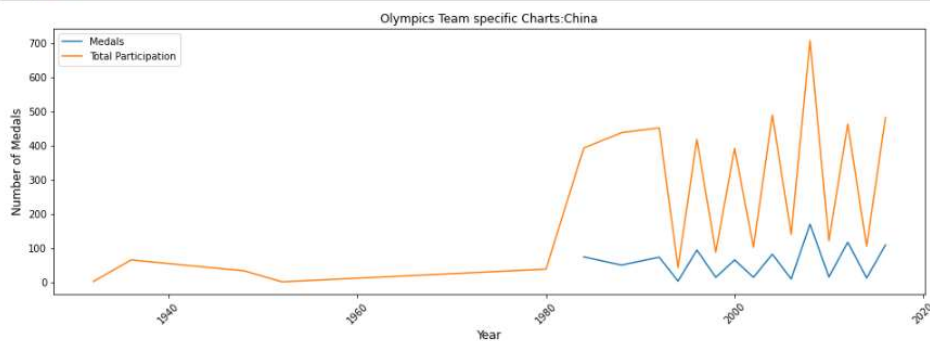
The participation Vs winning for India

```
1  country = 'India'
2  df_by_team = get_data_by_team(df_few_columns, country)
3  df_by_year = df_by_team['Year'].value_counts().sort_index()
4
5  df1 = df_athlete_events[df_athlete_events.Team == country]
6  df_by_parcipation = df1['Year'].value_counts().sort_index()
7
8  city_count = df_by_year[:20,]
9  plt.figure(figsize=(16,5))
10 plt.plot(city_count, label='Total Medals')
11 plt.plot(df_by_parcipation, label="Total Participation")
12 plt.legend()
13 plt.title('Olympics Team specific Charts:{0}'.format(country))
14 plt.xticks(rotation=45)
15 plt.ylabel('Number of Medals', fontsize=12)
16 plt.xlabel('Year', fontsize=12)
17 plt.show()
```
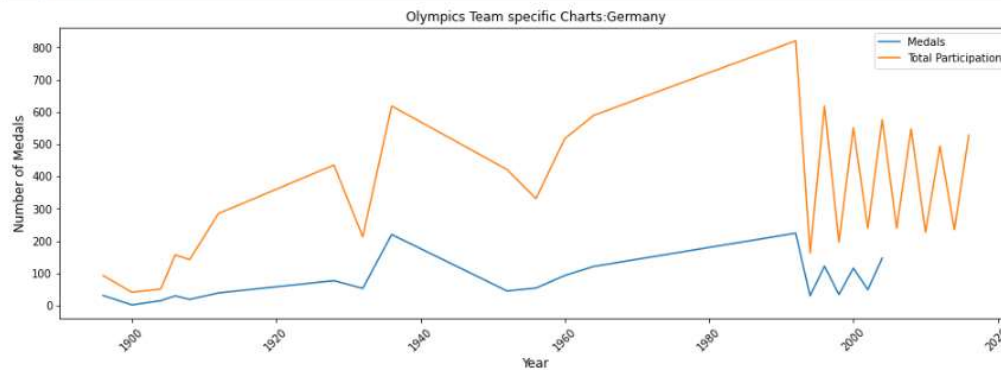

Olympics Team specific Charts:India

The participation Vs winning for China

```
1  country = 'China'
2  df_by_team = get_data_by_team(df_few_columns, country)
3  df_by_year = df_by_team['Year'].value_counts().sort_index()
4
5  df1 = df_athlete_events[df_athlete_events.Team == country]
6  df_by_parcipation = df1['Year'].value_counts().sort_index()
7
8  city_count = df_by_year[:20,]
9  plt.figure(figsize=(16,5))
10 plt.plot(city_count, label='Total Medals')
11 plt.plot(df_by_parcipation, label="Total Participation")
12 plt.legend()
13 plt.title('Olympics Team specific Charts:{0}'.format(country))
14 plt.xticks(rotation=45)
15 plt.ylabel('Number of Medals', fontsize=12)
16 plt.xlabel('Year', fontsize=12)
17 plt.show()
```


Olympics Team specific Charts:China

## The participation Vs winning for Germany

```
1  country = 'Germany'
2  df_by_team = get_data_by_team(df_few_columns, country)
3  df_by_year = df_by_team['Year'].value_counts().sort_index()
4
5  df1 = df_athlete_events[df_athlete_events.Team == country]
6  df_by_parcipation = df1['Year'].value_counts().sort_index()
7
8  city_count = df_by_year[:20,]
9  plt.figure(figsize=(16,5))
10 plt.plot(city_count, label='Total Medals')
11 plt.plot(df_by_parcipation, label="Total Participation")
12 plt.legend()
13 plt.title('Olympics Team specific Charts:{0}'.format(country))
14 plt.xticks(rotation=45)
15 plt.ylabel('Number of Medals', fontsize=12)
16 plt.xlabel('Year', fontsize=12)
17 plt.show()
```



## The participation Vs winning for Japan

```
1  country = 'Japan'
2  df_by_team = get_data_by_team(df_few_columns, country)
3  df_by_year = df_by_team['Year'].value_counts().sort_index()
4
5  df1 = df_athlete_events[df_athlete_events.Team == country]
6  df_by_parcipation = df1['Year'].value_counts().sort_index()
7
8  city_count = df_by_year[:20,]
9  plt.figure(figsize=(16,5))
10 plt.plot(city_count, label='Total Medals')
11 plt.plot(df_by_parcipation, label="Total Participation")
12 plt.legend()
13 plt.title('Olympics Team specific Charts:{0}'.format(country))
14 plt.xticks(rotation=45)
15 plt.ylabel('Number of Medals', fontsize=12)
16 plt.xlabel('Year', fontsize=12)
17 plt.show()
```