

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline

In [ ]:
```

```
In [8]: tokyo = pd.read_excel("athlete_excel.xlsx")
tokyo.head()
```

ID	Name	Sex	Age	Height	Weight	Team	NOC	Games	Year	Season	City	Sport	Event	Medal	
0	1	A Dijiang	M	24.0	180.0	80.0	China	CHN	1992 Summer	1992	Summer	Barcelona	Basketball	Basketball Men's Basketball	NaN
1	2	A Lamusi	M	23.0	170.0	60.0	China	CHN	2012 Summer	2012	Summer	London	Judo	Judo Men's Extra-Lightweight	NaN
2	3	Gunnar Nielsen Aaby	M	24.0	NaN	NaN	Denmark	DEN	1920 Summer	1920	Summer	Antwerpen	Football	Football Men's Football	NaN
3	4	Edgar Lindenu Aabye	M	34.0	NaN	NaN	Denmark/Sweden	DEN	1900 Summer	1900	Summer	Paris	Tug-Of-War	Tug-Of-War Men's Tug-Of-War	Gold
4	5	Christine Jacoba Aaftink	F	21.0	185.0	82.0	Netherlands	NED	1988 Winter	1988	Winter	Calgary	Speed Skating	Speed Skating Women's 500 metres	NaN

```
In [9]: olympic = pd.read_excel("region_excel.xlsx")
olympic.head()
```

NOC	region	notes	
0	AFG	Afghanistan	NaN
1	AHO	Curacao	Netherlands Antilles
2	ALB	Albania	NaN
3	ALG	Algeria	NaN
4	AND	Andorra	NaN

```
In [11]: data = tokyo.merge(olympic,how = 'left', on = 'NOC')
data.head()
```

```
Out[11]:
```

ID	Name	Sex	Age	Height	Weight	Team	NOC	Games	Year	Season	City	Sport	Event	Medal	region	notes	
0	1	A Dijiang	M	24.0	180.0	80.0	China	CHN	1992 Summer	1992	Summer	Barcelona	Basketball	Basketball Men's Basketball	NaN	China	NaN
1	2	A Lamusi	M	23.0	170.0	60.0	China	CHN	2012 Summer	2012	Summer	London	Judo	Judo Men's Extra-Lightweight	NaN	China	NaN
2	3	Gunnar Nielsen Aaby	M	24.0	NaN	NaN	Denmark	DEN	1920 Summer	1920	Summer	Antwerpen	Football	Football Men's Football	NaN	Denmark	NaN
3	4	Edgar Lindenuau Aabye	M	34.0	NaN	NaN	Denmark/Sweden	DEN	1900 Summer	1900	Summer	Paris	Tug-Of-War	Tug-Of-War Men's Tug-Of-War	Gold	Denmark	NaN
4	5	Christine Jacoba Aaftink	F	21.0	185.0	82.0	Netherlands	NED	1988 Winter	1988	Winter	Calgary	Speed Skating	Speed Skating Women's 500 metres	NaN	Netherlands	NaN

```
In [12]: data.shape

Out[12]: (271116, 17)
```

```
In [13]: data.rename(columns={'region':'Region','notes':'Notes'},inplace=True)
data.head()
```

```
Out[13]:
```

ID	Name	Sex	Age	Height	Weight	Team	NOC	Games	Year	Season	City	Sport	Event	Medal	Region	Notes	
0	1	A Dijiang	M	24.0	180.0	80.0	China	CHN	1992 Summer	1992	Summer	Barcelona	Basketball	Basketball Men's Basketball	NaN	China	NaN
1	2	A Lamusi	M	23.0	170.0	60.0	China	CHN	2012 Summer	2012	Summer	London	Judo	Judo Men's Extra-Lightweight	NaN	China	NaN
2	3	Gunnar Nielsen Aaby	M	24.0	NaN	NaN	Denmark	DEN	1920 Summer	1920	Summer	Antwerpen	Football	Football Men's Football	NaN	Denmark	NaN
3	4	Edgar Lindenuau Aabye	M	34.0	NaN	NaN	Denmark/Sweden	DEN	1900 Summer	1900	Summer	Paris	Tug-Of-War	Tug-Of-War Men's Tug-Of-War	Gold	Denmark	NaN
4	5	Christine Jacoba Aaftink	F	21.0	185.0	82.0	Netherlands	NED	1988 Winter	1988	Winter	Calgary	Speed Skating	Speed Skating Women's 500 metres	NaN	Netherlands	NaN

```
In [14]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 271116 entries, 0 to 271115
Data columns (total 17 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  218945 non-null      float64
5   Weight  288241 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
15  Region  270746 non-null      object
16  Notes   5839 non-null        object
dtypes: float64(3), int64(2), object(12)
memory usage: 37.2+ MB
```

```
In [15]: data.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

```
In [16]: data.isnull().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
Region	370
Notes	266077
dtype:	int64

```
In [24]: data.query('Sport == "Athletics").head()
```

```
Out[24]:
```

ID	Name	Sex	Age	Height	Weight	Team	NOC	Games	Year	Season	City	Sport	Event	Medal	Region	Notes	
26	8	Cornelia "Cor" Aalten (- Strannood)	F	18.0	168.0	NaN	Netherlands	NED	1932 Summer	1932	Summer	Los Angeles	Athletics	Athletics Women's 100 metres	NaN	Netherlands	NaN
27	8	Cornelia "Cor" Aalten (- Strannood)	F	18.0	168.0	NaN	Netherlands	NED	1932 Summer	1932	Summer	Los Angeles	Athletics	Athletics Women's 4 x 100 metres Relay	NaN	Netherlands	NaN
57	18	Timo Antero Aaltonen	M	31.0	189.0	130.0	Finland	FIN	2000 Summer	2000	Summer	Sydney	Athletics	Athletics Men's Shot Put	NaN	Finland	NaN
94	31	Evald rma (rman-)	M	24.0	174.0	70.0	Estonia	EST	1936 Summer	1936	Summer	Berlin	Athletics	Athletics Men's Pole Vault	NaN	Estonia	NaN
95	32	Olav Augunson Aarnes	M	23.0	NaN	NaN	Norway	NOR	1912 Summer	1912	Summer	Stockholm	Athletics	Athletics Men's High Jump	NaN	Norway	NaN

```
In [25]: data.query('Team == "Japan"').head()
```

```
Out[25]:
```

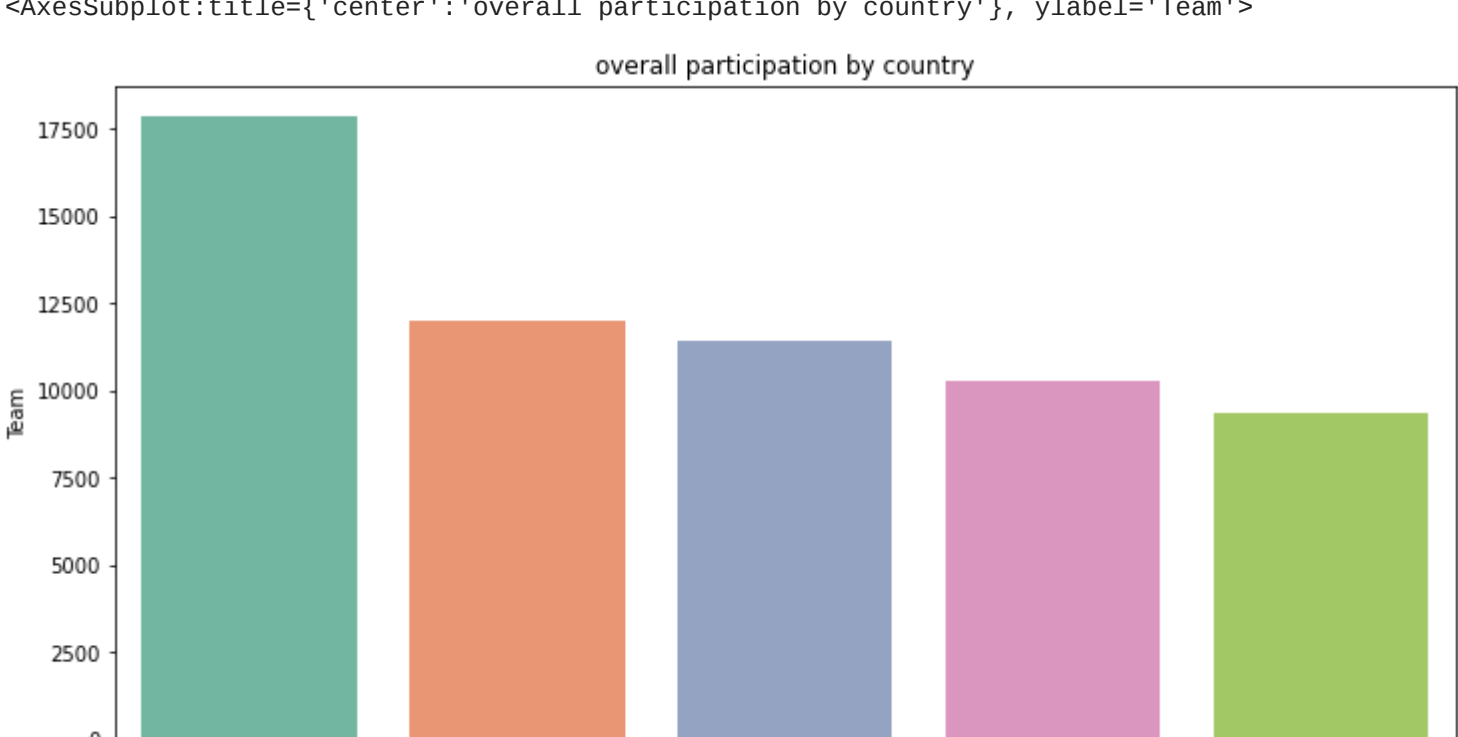
ID	Name	Sex	Age	Height	Weight	Team	NOC	Games	Year	Season	City	Sport	Event	Medal	Region	Notes	
625	362	Isao Ko Abe	M	24.0	177.0	75.0	Japan	JPN	1936 Summer	1936	Summer	Berlin	Athletics	Athletics Men's Hammer Throw	NaN	Japan	NaN
629	363	Kazumi Abe	M	28.0	178.0	67.0	Japan	JPN	1976 Winter	1976	Winter	Innsbruck	Bobsleigh	Bobsleigh Men's Four	NaN	Japan	NaN
630	364	Kazuo Abe	M	25.0	166.0	69.0	Japan	JPN	1960 Summer	1960	Summer	Roma	Wrestling	Wrestling Men's Lightweight, Freestyle	NaN	Japan	NaN
631	365	Kinya Abe	M	23.0	168.0	68.0	Japan	JPN	1992 Summer	1992	Summer	Barcelona	Fencing	Fencing Men's Foil, Individual	NaN	Japan	NaN
632	366	Kiyoshi Abe	M	25.0	167.0	62.0	Japan	JPN	1972 Summer	1972	Summer	Munich	Wrestling	Wrestling Men's Featherweight, Freestyle	NaN	Japan	NaN

```
In [32]: top_10 = data.Team.value_counts().sort_values(ascending=False).head()
top_10
```

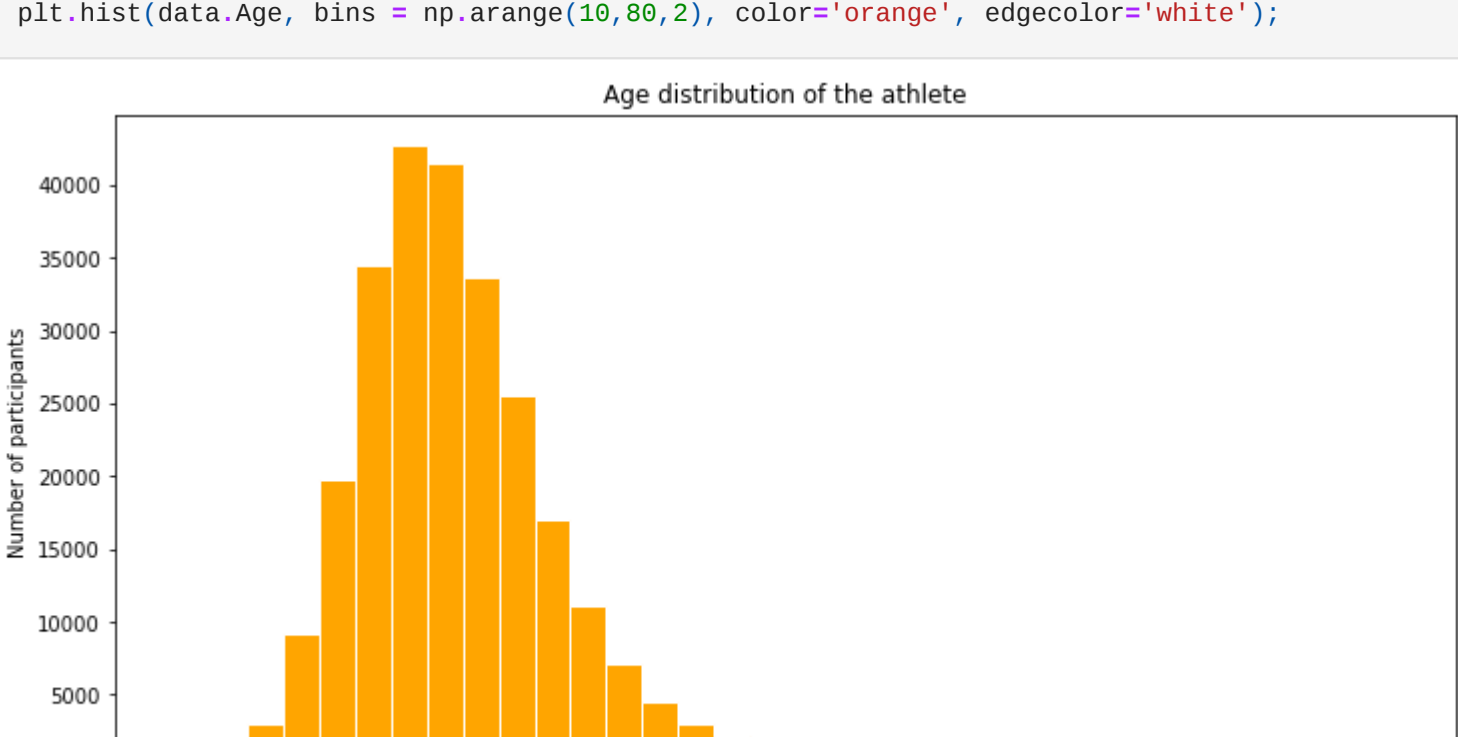
United States	17847
France	11988
Great Britain	11404
Italy	10260
Germany	9326
Name:	Team, dtype: int64

```
In [33]: plt.figure(figsize=(12,6))
plt.title("overall participation by country")
sns.barplot(x=top_10.index , y=top_10 ,palette='Set2')
```

```
Out[33]: <AxesSubplot:title='center':'overall participation by country', ylabel='Team'>
```



```
In [37]: plt.figure(figsize=(12,6))
plt.title("Age distribution of the athlete")
plt.xlabel('Age')
plt.ylabel("Number of participants")
plt.hist(data.Age, bins = np.arange(10,80,2), color='orange', edgecolor='white');
```



```
In [42]: winter_sports = data[data.Season == 'Winter'].Sport.unique()
winter_sports
```

```
Out[42]: array(['Speed Skating', 'Cross Country Skiing', 'Ice Hockey', 'Biathlon',
        'Alpine Skiing', 'Luge', 'Bobsleigh', 'Figure Skating',
        'Nordic Combined', 'Freestyle Skiing', 'Ski Jumping', 'Curling',
        'Snowboarding', 'Short Track Speed Skating', 'Skeleton',
        'Military Ski Patrol', 'Alpinism'], dtype=object)
```

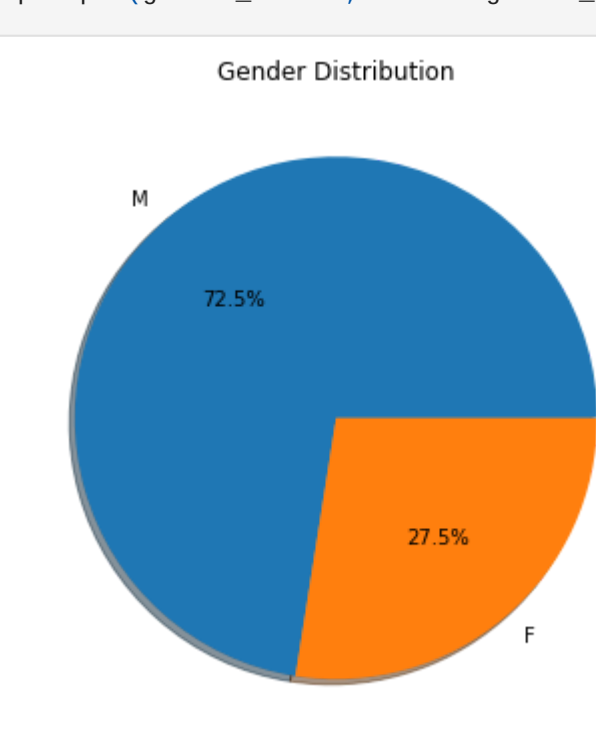
```
In [45]: summer_sports = data[data.Season == 'Summer'].Sport.unique()
summer_sports
```

```
Out[45]: array(['Basketball', 'Judo', 'Football', 'Tug-Of-War', 'Athletics',
        'Swimming', 'Badminton', 'Sailing', 'Gymnastics',
        'Art Competitions', 'Handball', 'Weightlifting', 'Wrestling',
        'Water Polo', 'Hockey', 'Rowing', 'Fencing', 'Equestrianism',
        'Shooting', 'Boxing', 'Taekwondo', 'Cycling', 'Diving', 'Canoeing',
        'Tennis', 'Modern Pentathlon', 'Golf', 'Softball', 'Archery',
        'Volleyball', 'Synchronized Swimming', 'Table Tennis', 'Baseball',
        'Rhythmic Gymnastics', 'Rugby Sevens', 'Trampolineing',
        'Beach Volleyball', 'Triathlon', 'Rugby', 'Lacrosse', 'Polo',
        'Cricket', 'Ice Hockey', 'Racquets', 'Motorboating', 'Croquet',
        'Figure Skating', 'Jeu De Paume', 'Roque', 'Basque Pelota',
        'Alpinism', 'Aeronautics'], dtype=object)
```

```
In [46]: gender_counts = data.Sex.value_counts()
gender_counts
```

M	196594
F	74522
Name:	Sex, dtype: int64

```
In [52]: plt.figure(figsize=(12,6))
plt.title('Gender Distribution')
plt.pie(gender_counts, labels=gender_counts.index, autopct= '%1.1f%', startangle=360, shadow=True);
```



```
In [53]: data.Medal.value_counts()
```

Gold	13372
Bronze	13295
Silver	13116
Name:	Medal, dtype: int64