



UNIFIED MENTOR
YOUR SKILL. SUCCESS & JOURNEY

A Project Report

By

Alladi Nikhila Sarea

Project Title	Olympics Data Analysis
Domain	Data Analyst
Contribution	Individual

Dataset: [Dataset](#)

To access the code, here is the link for Google Colab: [Colab](#)

GitHub Link: [Github](#)

Introduction

The Summer Olympics: A Global Celebration of Sport and Unity

The Summer Olympics, officially known as the Games of the Olympiad, stand as the pinnacle of international sports competition, bringing together athletes from across the globe every four years. Rooted in ancient Greek traditions, the modern Olympic Games have grown to symbolize not just athletic excellence, but also unity, cultural exchange, and peace.

The Summer Olympics feature a wide array of events that span across over 30 sports, including athletics, swimming, gymnastics, basketball, boxing, tennis, and more recently, sports like skateboarding and surfing. This diversity allows athletes of various talents and backgrounds to showcase their skills on a global stage. The competition is not just about medals—it celebrates perseverance, discipline, and sportsmanship.

Beyond the medals and world records, the Summer Olympics carry profound social and cultural importance. They provide a platform for emerging nations to shine, foster international camaraderie, and often highlight broader global issues.

The Summer Olympics continue to captivate the world, not just for the competition, but for the shared human spirit they represent. As athletes push the boundaries of physical achievement, they also build bridges between cultures, foster peace, and inspire generations. In a world often divided by conflict and difference, the Olympics serve as a powerful reminder of what humanity can achieve when it comes together in the spirit of fair play and mutual respect.

Project Description

The primary focus of this project will be on exploring and understanding the dataset, performing exploratory data analysis (EDA), and uncovering trends and insights related to athletes, countries, and sports over the years.

This dataset contains information on all medal winners from the Summer Olympics held between the 1976 Montreal Games and the 2008 Beijing Games, covering every medal awarded during that time span.

The following columns has been used for the analysis:

- **City:** The city where the Olympics took place.
- **Year:** The year of the Olympics.
- **Sport:** The sport the event is categorized under.
- **Discipline:** A subcategory of the sport.
- **Event:** The specific event within a discipline.

- **Athlete:** The name of the athlete who participated.
- **Gender:** The gender of the athlete.
- **Country_Code:** The country code (abbreviation).
- **Country:** The full name of the country.
- **Event_gender:** The gender category of the event.
- **Medal:** The medal won (Gold, Silver, Bronze).

Project Objective

The main objective of Olympic Data Analysis is to examine and derive insights from the Olympic history and data. The goal is to identify patterns, trends, and connections within the data, and to reveal interesting discoveries about the Olympics, its athletes, and its broader societal impact. This project involves applying a range of data analysis and machine learning methods to preprocess, clean, and interpret the data in order to uncover valuable and meaningful insights.

The project focuses on analyzing the dataset to uncover trends in medal distribution over time. It aims to identify the most successful countries and standout athletes, highlighting their achievements across different Olympic editions. Additionally, it explores the gender distribution of events and medals to understand how participation and success have evolved for male and female athletes. To support these insights, various data visualization techniques in Python are employed, making the analysis more intuitive and insightful.

Project steps:

- i. **Data Preparation:** The first step involves setting up the environment by importing the necessary Python libraries required for data analysis and visualization. Once the libraries are in place, the Olympic dataset is loaded into the workspace. Following this, the dataset is cleaned to address any inconsistencies or missing values, ensuring that the data is ready for analysis.
- ii. **Exploratory Data Analysis (EDA):** With the dataset cleaned, the next phase is to perform exploratory data analysis, starting with generating summary statistics to understand the structure and basic characteristics of the data. The analysis then moves on to examining medal trends across different years, and pinpointing the most successful athletes and countries throughout Olympic history.
- iii. **Visualizing Key Insights:** To make the findings more interpretable, the third step focuses on visualizing key insights using Python's plotting libraries. This includes illustrating medal distributions by country, year, and sport, as well as conducting a detailed analysis of gender representation in various sports and events to understand participation and achievement trends among male and female athletes.

- iv. **Predictive Analysis:** The final step involves the use of machine learning techniques to build a predictive model. The goal is to train the model to predict the likelihood of an athlete winning a medal based on features such as their country, sport, and other relevant attributes. This step adds a forward-looking element to the analysis by transforming historical patterns into predictive insights.

Task 1: Data Preparation

1.1 Import the necessary libraries

```
import pandas as pd
import numpy as np

import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import OneHotEncoder
from sklearn.compose import ColumnTransformer
from sklearn.metrics import classification_report, accuracy_score, confusion_matrix
```

The libraries imported are explained below:

- i. **NumPy (np):** Fundamental package for numerical computing in Python. Used for arrays, mathematical functions, linear algebra, etc.
- ii. **Pandas (pd):** Essential for data manipulation and analysis. Provides Data frames and Series for structured data.
- iii. **Matplotlib (plt):** The most commonly used library for 2D plotting in Python.
- iv. **Seaborn (sns):** Built on top of matplotlib, used for making statistical graphics with better aesthetics.
- v. **train_test_split:** Splits dataset into training and testing subsets.
- vi. **Logistic Regression:** A classification algorithm used when the target variable is categorical. Useful for binary classification problems.
- vii. **Pipeline:** The Pipeline class helps in assembling several steps (like preprocessing and modeling) into one object. This ensures that all the steps are applied consistently to both training and testing data.
- viii. **OneHotEncoder:** Converts categorical features into a format that works better with ML algorithms by creating binary columns (one-hot encoding).
- ix. **ColumnTransformer:** Helps apply different transformations to different columns in the dataset.
- x. **Classification report:** Displays precision, recall, F1-score, and support.

- xi. **Accuracy score:** Measures the ratio of correctly predicted samples.
- xii. **Confusion matrix:** Summarizes true positives, false positives, true negatives, and false negatives.

1.2 Load the dataset

```
data = pd.read_csv(r"olympic.csv", encoding='latin1')
```

The dataset has been imported for the model.

1.2.1 Cleaning the dataset

After importing the dataset, it is cleaned to make sure there are no duplicate values and null values.

Task 2: Exploratory Data Analysis (EDA)

The summary statistics of the dataset has been observed. The dataset was verified to be clean, with no duplicate entries or missing values. It consists of 15,316 rows and 11 columns. A total of 127 countries took part in the events. The participation spanned across 28 main sports categories, which included sub-categories, amounting to a total of 41 distinct sports.

2.1 Trends of medals across years

The trends of medals of the countries across years has been observed here.

The dataset contains a total of 5,041 gold medals, 5,016 silver medals, and 5,258 bronze medals. The years in which the Olympics were conducted have been identified. It was observed that the highest number of gold and bronze medals were awarded in 2008, while the maximum number of silver medals were distributed in 2000. The lowest overall medal distribution occurred in the year 1976. India appears to have achieved its highest medal count in the 1980 Olympics.

```
df = pd.concat([df, pd.get_dummies(df['Medal'])], axis=1)
```

Binary (dummy) columns were created for the 'Medal' category to simplify the analysis, where each column represents a specific medal type and contains values in binary format (0 or 1).

```
print("Years the events are held: ")
year = np.unique(df['Year'].dropna().values).tolist()
year.sort()
year
```

```
Years the events are held:
[1976.0, 1980.0, 1984.0, 1988.0, 1992.0, 1996.0, 2000.0, 2004.0, 2008.0]
```

The years during which the events were held were found.

2.1.1 Function year_list():

```
# A plot and list of years and their respective medal counts, arranged in ascending order according to the number of gold medals won.

def year_list(df):
    total_medals = df.groupby('Year')[['Gold', 'Silver', 'Bronze']].sum()
    reset = total_medals.reset_index()

    plt.figure(figsize=(14, 6))
    bar_width = 0.25
    x = range(len(reset))

    plt.bar([i - bar_width for i in x], reset['Gold'], width=bar_width, label='Gold', color='gold')
    plt.bar(x, reset['Silver'], width=bar_width, label='Silver', color='silver')
    plt.bar([i + bar_width for i in x], reset['Bronze'], width=bar_width, label='Bronze', color='peru')

    plt.xticks(ticks=x, labels=reset['Year'], rotation=45)
    plt.xlabel('Year', fontsize=12)
    plt.ylabel('Number of Medals', fontsize=12)
    plt.title(f'Medals acheived over the years', fontsize=16)
    plt.legend()
    plt.grid(axis='y', linestyle='--', alpha=0.6)
    plt.tight_layout()
    plt.show()

    return reset

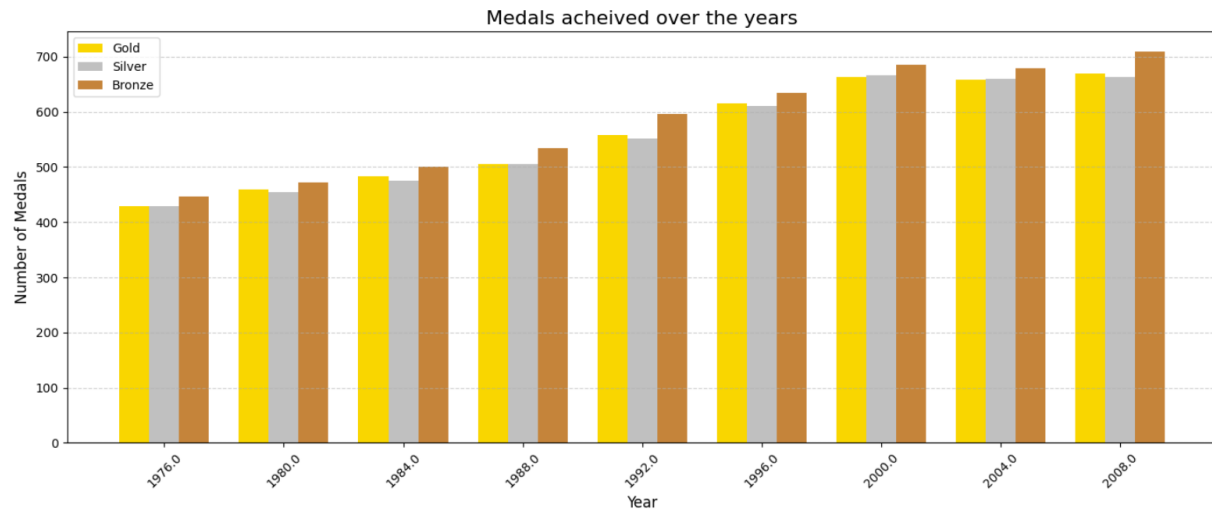
year_list(df)
```

This function generates a grouped bar chart showing how many gold, silver, and bronze medals were won each year.

- ✓ It takes a DataFrame 'df' as input.
- ✓ Groups the dataset by the 'Year' column.
- ✓ Sums up the number of medals (binary dummy columns) for each year.
- ✓ The resulting total_medals DataFrame will have years as index and total counts of gold, silver, and bronze medals as columns.
- ✓ Converts the index (which is currently the Year) into a normal column for easier plotting.

The plot is created:

- Initializes a new figure for the plot.
 - Sets the figure size to 14 inches wide and 6 inches tall.
 - bar_width: Sets the width of each bar.
 - x: A sequence of integers representing the positions of each year on the x-axis.
 - Plots gold medal bars slightly left of center (i - bar_width).
 - Uses gold color for gold medals. Plots silver medal bars in the center. lots bronze medal bars slightly right of center using a bronze-like color (peru). These three lines create grouped bars for each year.
 - Renders the final plot.
- ✓ Returns the reset DataFrame (summary of medals per year).



	Year	Gold	Silver	Bronze
0	1976.0	429	429	447
1	1980.0	459	455	472
2	1984.0	483	476	500
3	1988.0	506	505	535
4	1992.0	558	551	596
5	1996.0	615	610	634
6	2000.0	663	667	685
7	2004.0	659	660	679
8	2008.0	669	663	710

This plot shows consistent upward trend in the total number of medals awarded over the years. Likely reflecting the increase in the number of events and participating countries in each Olympic edition. The plot effectively illustrates the expansion of the Olympic Games in terms of event offerings and global participation. The increase in bronze medals highlights structural rules in certain sports that lead to more than one bronze being awarded.

2.1.2 Function country_MedalsOverYears():

```
# A function that generates a plot and a list of years along with the number of medals the country won, based on the specified country and type of medal
# provided as input.

def country_MedalsOverYears(df, COUNTRY):
    # dataframe for the specified country
    country_df = df[df['Country'] == COUNTRY]

    # Group by year and sum all medal types
    medals = country_df.groupby('Year')[['Gold', 'Silver', 'Bronze']].sum().reset_index()

    # Plot
    plt.figure(figsize=(10, 5))
    width = 0.25
    x = medals['Year']
    x_indexes = range(len(x))

    plt.bar([i - width for i in x_indexes], medals['Gold'], width=width, label='Gold', color='gold')
    plt.bar(x_indexes, medals['Silver'], width=width, label='Silver', color='silver')
    plt.bar([i + width for i in x_indexes], medals['Bronze'], width=width, label='Bronze', color='peru')

    plt.title(f'{COUNTRY} - Medal Counts Over the Years', fontsize=16)
    plt.xlabel('Year', fontsize=12)
    plt.ylabel('Number of Medals', fontsize=12)
    plt.xticks(ticks=x_indexes, labels=medals['Year'], rotation=90)
    plt.legend()
    plt.grid(axis='y', linestyle='--', alpha=0.6)
    plt.tight_layout()
    plt.show()

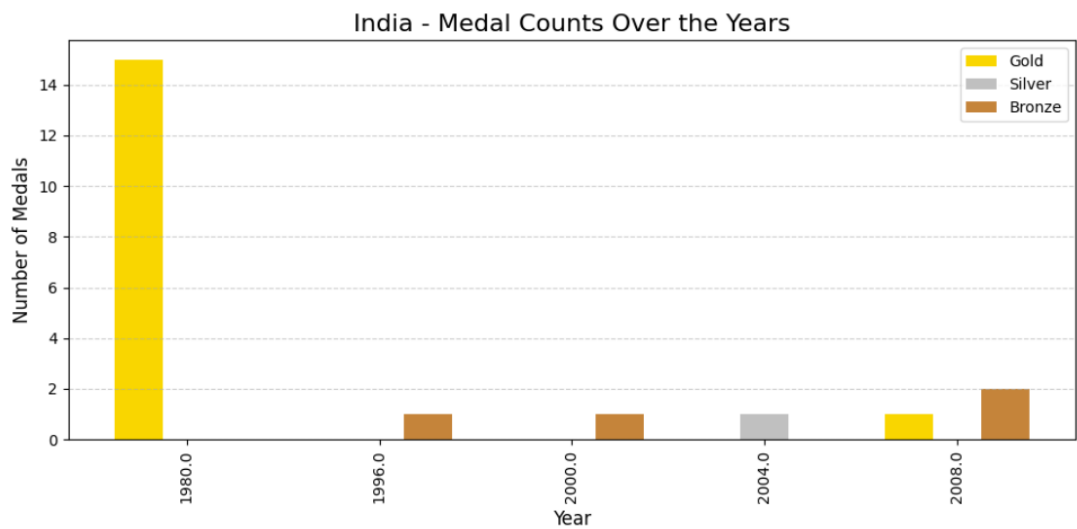
    # Return the result table
    return medals
```

The function plots the gold, silver, and bronze medal counts over time for a specific country using a grouped bar chart. It also returns a summary table of medals by year for that country. This function is ideal for performing a country-wise temporal analysis of Olympic performance. It's useful for visual comparisons across years, identifying peak performance years, detecting trends like growth or decline in medal counts.

- ✓ This function accepts a dataset and a string 'COUNTRY' indicating the country you want to analyze (e.g., 'USA', 'India')
- ✓ Filters the dataframe to only include rows where the 'Country' matches the specified 'COUNTRY'.
- ✓ Groups the country-specific data by Year.
- ✓ Sums up the gold, silver, and bronze medal columns for each year.
- ✓ Resets the index to convert Year from index to a column for plotting.
- ✓ Displays the final plot.
- ✓ Returns a DataFrame showing total medals by year for the selected country.

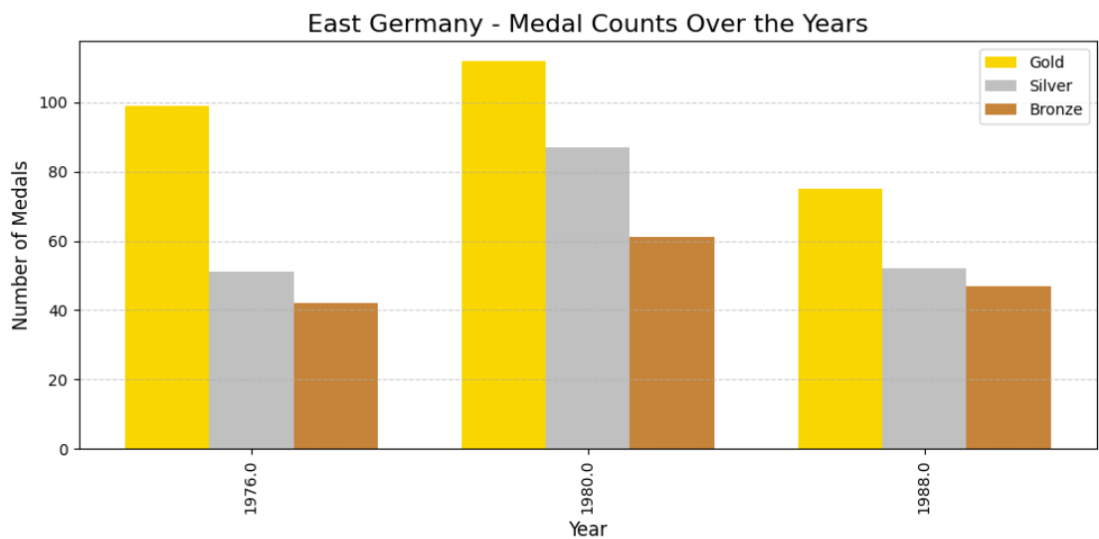
Examples:

```
country_MedalsOverYears(df, 'India')
```



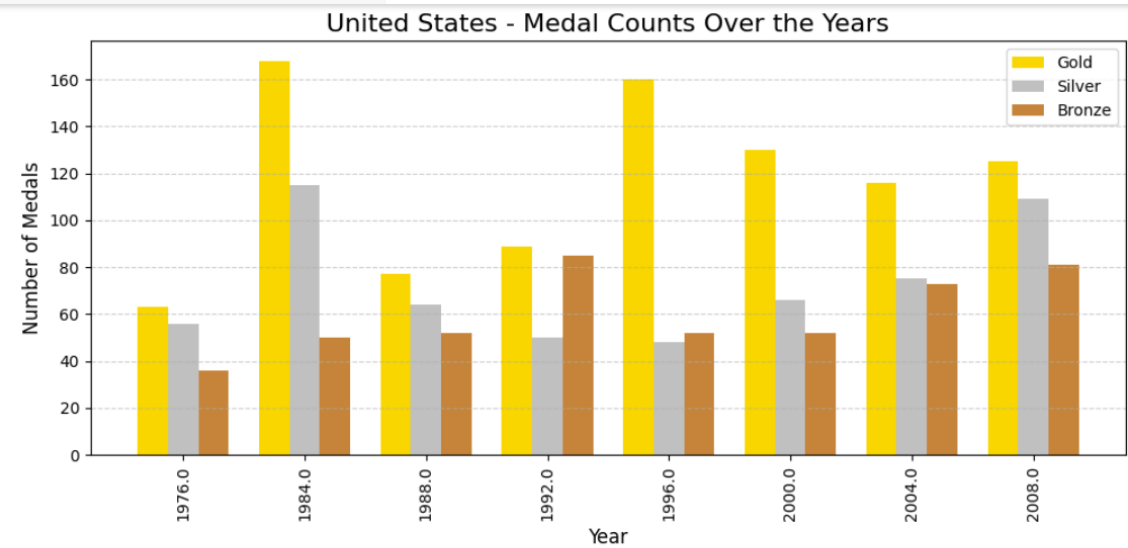
	Year	Gold	Silver	Bronze
0	1980.0	15	0	0
1	1996.0	0	0	1
2	2000.0	0	0	1
3	2004.0	0	1	0
4	2008.0	1	0	2

```
country_MedalsOverYears(df, 'East Germany')
```



	Year	Gold	Silver	Bronze
0	1976.0	99	51	42
1	1980.0	112	87	61
2	1988.0	75	52	47

```
country_MedalsOverYears(df, 'United States')
```



	Year	Gold	Silver	Bronze
0	1976.0	63	56	36
1	1984.0	168	115	50
2	1988.0	77	64	52
3	1992.0	89	50	85
4	1996.0	160	48	52
5	2000.0	130	66	52
6	2004.0	116	75	73
7	2008.0	125	109	81

2.2 Top-performing athletes and countries

The **United States** emerged as the leading country, securing 928 gold, 583 silver, and 481 bronze medals, followed by the **Soviet Union** and **East Germany**.

India secured 16 gold, 1 silver, and 4 bronze medals through the entire Olympics and the data of its achievements has been extracted.

2.2.1 Function years_Top10Countries():

```
# A function that generates a plot and a list of countries along with the number of medals they won in a specific year

def years_Top10Countries(df, YEAR):
    # Validate input year
    valid_years = [1976.0, 1980.0, 1984.0, 1988.0, 1992.0, 1996.0, 2000.0, 2004.0, 2008.0]
    if YEAR not in valid_years:
        raise ValueError(f"YEAR must be one of: {valid_years}")

    # data for the given year
    year_df = df[df['Year'] == YEAR]

    # Group by country and sum medals
    medals = year_df.groupby('Country')[['Gold', 'Silver', 'Bronze']].sum()

    # Add total column to get top 10
    medals['Total'] = medals['Gold'] + medals['Silver'] + medals['Bronze']
    top_countries = medals.sort_values('Total', ascending=False).head(10).reset_index()

    # Plot
    plt.figure(figsize=(14, 6))
    bar_width = 0.25
    x = range(len(top_countries))

    plt.bar([i - bar_width for i in x], top_countries['Gold'], width=bar_width, label='Gold', color='gold')
    plt.bar(x, top_countries['Silver'], width=bar_width, label='Silver', color='silver')
    plt.bar([i + bar_width for i in x], top_countries['Bronze'], width=bar_width, label='Bronze', color='peru')

    plt.xticks(ticks=x, labels=top_countries['Country'], rotation=45)
    plt.xlabel('Country', fontsize=12)
    plt.ylabel('Number of Medals', fontsize=12)
    plt.title(f'Top 10 Countries by Total Medals in {int(YEAR)}', fontsize=16)
    plt.legend()
    plt.grid(axis='y', linestyle='--', alpha=0.6)
    plt.tight_layout()
    plt.show()

    return top_countries[['Country', 'Gold', 'Silver', 'Bronze', 'Total']]
```

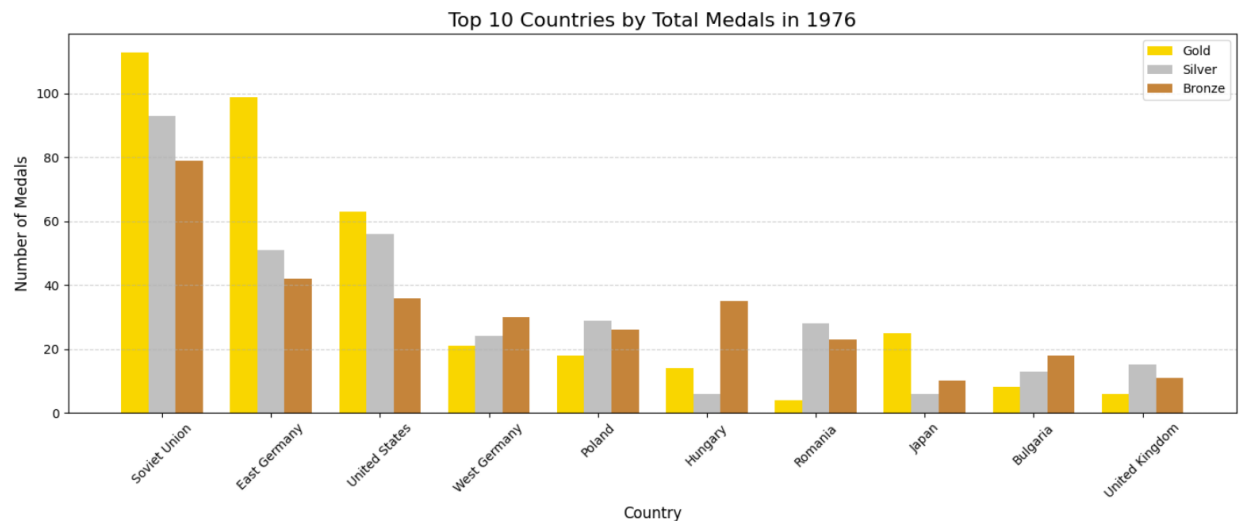
This function displays the top 10 countries based on total medals in a given Olympic year.

- ✓ The function accepts a DataFrame 'df' containing Olympic data and a specific 'YEAR' for which to analyze the top-performing countries.
- ✓ Defines a list of years that are valid for analysis (presumably based on dataset coverage).
- ✓ The values are floats because years in the dataset might be stored as floats (e.g., 1976.0).

- ✓ Ensures the input YEAR is in the valid_years list. If not, raises an error with a clear message.
- ✓ Filters the main DataFrame to include **only the rows for the specified year**.
- ✓ Groups the filtered data by country. Sums the number of gold, silver, and bronze medals won by each country that year.
- ✓ Adds a new column 'Total' representing the sum of all three medal types.
- ✓ Sorts the countries in **descending order** by total medals. Selects the **top 10 countries**. Resets the index so 'Country' becomes a normal column again.
- ✓ Displays the plot.
- ✓ Returns a table with the top 10 countries, showing their medal breakdown and total.

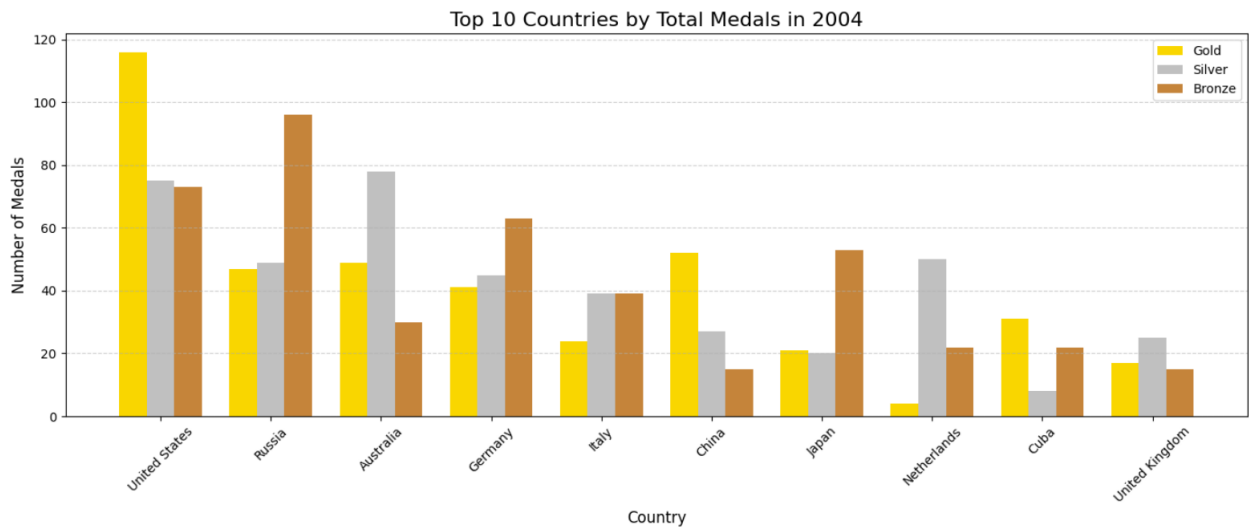
Examples:

```
years_Top10Countries(df, 1976)
```



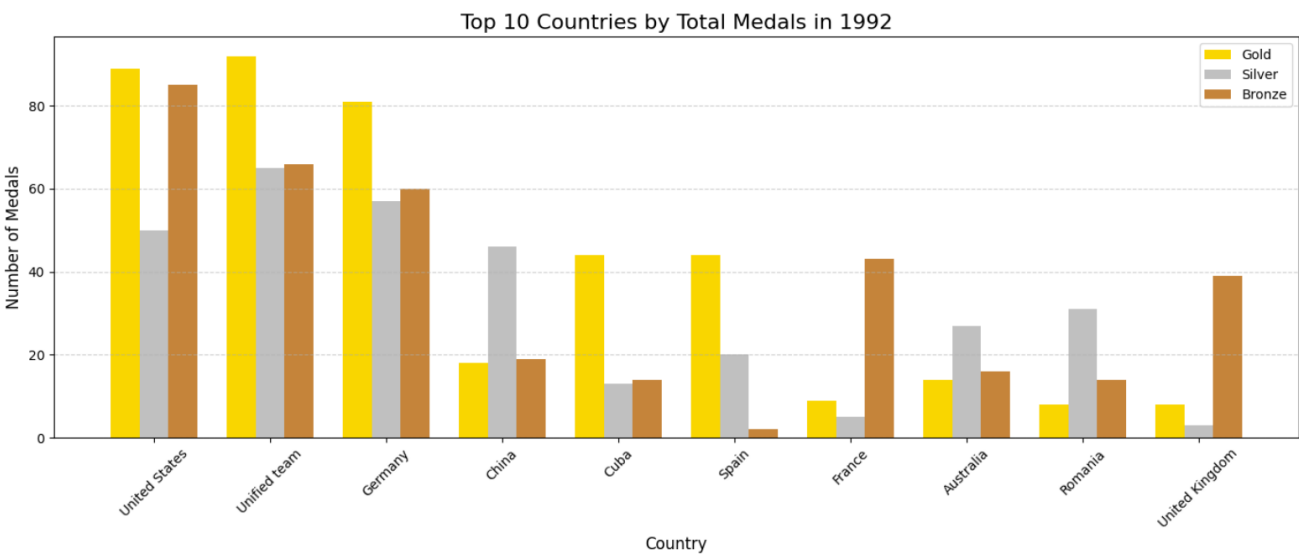
	Country	Gold	Silver	Bronze	Total
0	Soviet Union	113	93	79	285
1	East Germany	99	51	42	192
2	United States	63	56	36	155
3	West Germany	21	24	30	75
4	Poland	18	29	26	73
5	Hungary	14	6	35	55
6	Romania	4	28	23	55
7	Japan	25	6	10	41
8	Bulgaria	8	13	18	39
9	United Kingdom	6	15	11	32

years_Top10Countries(df, 2004)



	Country	Gold	Silver	Bronze	Total
0	United States	116	75	73	264
1	Russia	47	49	96	192
2	Australia	49	78	30	157
3	Germany	41	45	63	149
4	Italy	24	39	39	102
5	China	52	27	15	94
6	Japan	21	20	53	94
7	Netherlands	4	50	22	76
8	Cuba	31	8	22	61
9	United Kingdom	17	25	15	57

years_Top10Countries(df, 1992)



	Country	Gold	Silver	Bronze	Total
0	United States	89	50	85	224
1	Unified team	92	65	66	223
2	Germany	81	57	60	198
3	China	18	46	19	83
4	Cuba	44	13	14	71
5	Spain	44	20	2	66
6	France	9	5	43	57
7	Australia	14	27	16	57
8	Romania	8	31	14	53
9	United Kingdom	8	3	39	50

2.2.2 Function medals_InYearCountry():

```
# A function that generates a list of medals they won, based on the specified year and country provided as input.
# If the input: 'overall' is given, then the data across the column is taken.

def medals_InYearCountry(df, YEAR, COUNTRY):
    # make sure the input is valid
    if YEAR not in [1976.0, 1980.0, 1984.0, 1988.0, 1992.0, 1996.0, 2000.0, 2004.0, 2008.0, 'overall']:
        raise ValueError("YEAR must be one of: [1976.0, 1980.0, 1984.0, 1988.0, 1992.0, 1996.0, 2000.0, 2004.0, 2008.0]")

    # dataframe with the required year and country
    if (YEAR == 'overall'):
        medal_list = df[(df['Country'] == COUNTRY)]
    elif (COUNTRY == 'overall'):
        medal_list = df[(df['Year'] == YEAR)]
    elif ((YEAR == 'overall') & (COUNTRY == 'overall')):
        medal_list = df
    else:
        medal_list = df[(df['Year'] == YEAR) & (df['Country'] == COUNTRY)]

    # Count medals by type
    medal_counts = medal_list['Medal'].value_counts().reset_index()
    medal_counts.columns = ['Medal Type', 'Count']

    return medal_counts
```

This function returns a count of medals by type (Gold, Silver, Bronze) for a specific year, a specific country, or both, depending on the input. It also allows for overall medal counts across all years or all countries. This function is very useful for comparing performance across different countries, tracking a country's performance over all Olympic Games, viewing overall medal trends year-wise or country-wise.

- ✓ The function accepts a DataFrame 'df' containing Olympic data, a specific year and a country or 'overall' for either.
- ✓ Checks whether the input YEAR is valid — either one of the specified Olympic years or the string 'overall'. If it's not valid, it raises an error and exits the function.
- ✓ In the if-else condition:

- If 'overall' is chosen for YEAR, it selects all data for the specified COUNTRY, across **all years**.
 - If 'overall' is chosen for COUNTRY, it selects all countries' data for the **specific year**.
 - If **both** YEAR and COUNTRY are 'overall', it selects the **entire dataset**.
 - In all other cases (specific year and country), it filters the data for the **intersection of both conditions**.
- ✓ Counts the number of occurrences of each medal type (Gold, Silver, Bronze) in the filtered data and converts the result into a DataFrame and resets the index for easier handling.
 - ✓ Renames the columns to make them more descriptive.
 - ✓ Returns the final DataFrame showing how many gold, silver, and bronze medals were won based on the input criteria.

Examples:

```
print("Medals INDIA acheived over the years:")
medals_InYearCountry(df, 'overall', 'India')
```

Medals INDIA acheived over the years:

	Medal Type	Count
0	Gold	16
1	Bronze	4
2	Silver	1

```
print("Medals the US acheived in 2000:")
medals_InYearCountry(df, 2000, 'United States')
```

Medals the US acheived in 2000:

	Medal Type	Count
0	Gold	130
1	Silver	66
2	Bronze	52

```
print("Medals Distributed in 2008:")
medals_InYearCountry(df, 2008, 'overall')
```

Medals Distributed in 2008:

	Medal Type	Count
0	Bronze	710
1	Gold	669
2	Silver	663

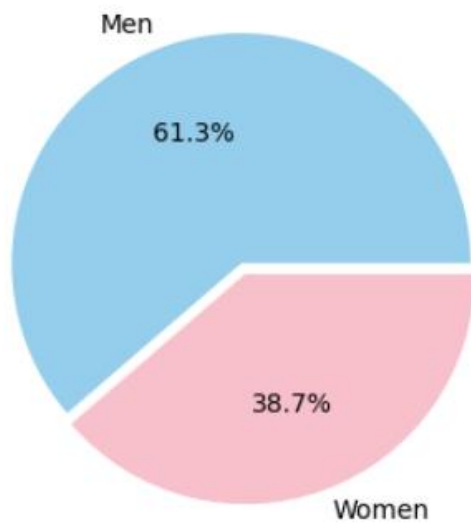
2.2.3 Top-performing Athletes

The top 5 athletes with the best performance in overall Olympics are found to be:

	Country	Sport	Athlete	Gender	Gold	Silver	Bronze
0	United States	Aquatics	PHELPS, Michael	Men	14	0	2
1	United States	Athletics	LEWIS, Carl	Men	9	1	0
2	United States	Aquatics	THOMPSON, Jenny	Women	8	3	1
3	United States	Aquatics	BIONDI, Matthew	Men	8	2	1
4	Soviet Union	Gymnastics	ANDRIANOV, Nikolay	Men	6	4	2

2.2.4 Gender Distribution

Gender Distribution in Olympics Events



Approx. 61% of athletes who participated in the Olympics are found to be men and only 39% of athletes are women.

The list of the best performing male and female athlete in the whole summer Olympics is given below:

It has both the name and the number of medals the athlete achieved in the games.

Sport	Male Athlete (G/S/B)	Female Athlete (G/S/B)
Aquatics	PHELPS, Michael (14/0/2)	THOMPSON, Jenny (8/3/1)
Archery	JANG, Yong-Ho (2/1/0)	KIM, Soo-Nyung (4/1/1)
Athletics	LEWIS, Carl (9/1/0)	ASHFORD, Evelyn (4/1/0)
Badminton	KIM, Dong Moon (2/0/1)	GAO, Ling (2/1/1)
Baseball	LAZO, Pedro Luis (2/2/0)	N/A
Basketball	ROBINSON, David Maurice (2/0/1)	EDWARDS, Teresa (4/0/1)
Boxing	SAVON, Felix (3/0/0)	N/A
Canoe / Kayak	FERGUSON, Ian Gordon (4/1/0)	FISCHER, Birgit (8/4/0)
Cycling	HOY, Chris (4/1/0)	ZIJLAARD-VAN MOORSEL, Leontien (4/1/1)
Equestrian	KLIMKE, Reiner (4/0/1)	WERTH, Isabell (5/3/0)
Fencing	POZDNYAKOV, Stanislav (4/0/1)	VEZZALI, Valentina (5/1/1)
Football	MASCHERANO, Javier (2/0/0)	CHASTAIN, Brandi (2/1/0)
Gymnastics	ANDRIANOV, Nikolay (6/4/2)	COMANECI, Nadia (5/3/1)
Handball	LAVROV, Andrey (3/0/1)	KARLOVA, Larisa (2/0/1)
Hockey	DE NOOIJER, Teun (2/1/0)	HAWKES, Rechelle (3/0/0)
Judo	NOMURA, Tadahiro (3/0/0)	TANI, Ryoko (2/2/1)
Modern Pentathlon	MASALA, Daniele (2/1/0)	COOK, Stephanie (1/0/0)
Rowing	REDGRAVE, Steven (5/0/1)	LIPA, Elisabeta (5/2/1)
Sailing	AINSLIE, Ben (3/1/0)	AYTON, Sarah (2/0/0)
Shooting	SCHUMANN, Ralf (3/2/0)	LOGVINENKO, Marina (2/1/2)
Softball	N/A	BERG, Laura (3/1/0)
Table Tennis	MA, Lin (3/0/0)	WANG, Nan (4/1/0)
Taekwondo	LOPEZ, Steven (2/0/1)	CHEN, Zhong (2/0/0)
Tennis	MASSU, Nicolas (2/0/0)	WILLIAMS, Venus (3/0/0)
Triathlon	WHITFIELD, Simon (1/1/0)	ALLEN, Kate (1/0/0)
Volleyball	KIRALY, Charles (3/0/0)	FERNANDEZ, Ana Ibis (3/0/1)
Weightlifting	DIMAS, Pyrros (3/0/1)	CHEN, Yanqing (2/0/0)
Wrestling	KARELIN, Aleksandr (3/1/0)	ICHO, Kaori (2/0/0)

2.2.5 Function athletes_Event():

```
import difflib

def athletes_Event(df, Sport, Discipline=None):
    # Check if the Sport is valid
    available_sports = df['Sport'].unique()
    if Sport not in available_sports:
        close_matches = difflib.get_close_matches(Sport, available_sports, n=3, cutoff=0.5)
        print(f"✗ Sport '{Sport}' not found.")
        if close_matches:
            print("Did you mean:")
            for match in close_matches:
                print("-", match)
        else:
            print("No similar sport names found.")
        return None

    # If Discipline is not provided, list all disciplines under the Sport
    if Discipline is None:
        disciplines = df[df['Sport'] == Sport]['Discipline'].unique()
        print(f"Available disciplines under '{Sport}':")
        for d in disciplines:
            print("-", d)
        return None
```

```

# Check if the Discipline is valid for the given Sport
valid_disciplines = df[df['Sport'] == Sport]['Discipline'].unique()
if Discipline not in valid_disciplines:
    close_matches = difflib.get_close_matches(Discipline, valid_disciplines, n=3, cutoff=0.5)
    print(f"✗ Discipline '{Discipline}' not found under sport '{Sport}'.")
    if close_matches:
        print("Did you mean:")
        for match in close_matches:
            print("-", match)
    else:
        print("No similar discipline names found.")
    return None

# Filter the dataset for valid Sport and Discipline
list_ath = df[(df['Sport'] == Sport) & (df['Discipline'] == Discipline)]

# Group by Event and Athlete, then sum medal counts
ath1 = list_ath.groupby(['Event', 'Athlete'])[['Gold', 'Silver', 'Bronze']].sum().reset_index()
ath2 = ath1.sort_values(by=['Gold', 'Silver', 'Bronze'], ascending=[False, False, False])
ath3 = ath2.groupby('Event').first().reset_index()

return ath3

```

This function helps you retrieve top-performing athletes (based on medal count) for each event under a specific sport and discipline from a given Olympic dataset. It also provides suggestions if the input is misspelled or incomplete. This function is ideal for finding top athletes per event in a specific sport and exploring event-level Olympic performances.

- ✓ Imports the “difflib” module, which helps find close string matches (used for suggesting alternatives if input sport/discipline is incorrect).
- ✓ The function accepts a DataFrame ‘df’ containing Olympic data, the sport and the discipline for filtering.
- ✓ Validates the sport name
- ✓ If the user doesn’t provide a discipline, lists all available disciplines under the specified sport and exits the function. Gets valid disciplines for the given sport and checks whether the provided Discipline is in the list. If it's not found, tries to suggest similar discipline names. If no valid match is found, exits the function.
- ✓ Filters the dataset to include only rows where both sport and discipline match the user's input.
- ✓ Groups by Event and Athlete, Sum Medals
- ✓ Sorts Athletes by Medal Ranking
- ✓ Selects Top Athlete per Event
- ✓ Returns a DataFrame with top athletes per event for the given sport and discipline.

Examples:

```
# top performing athletes in 'aquatics' - 'diving'
```

```
athletes_Event(df, 'Aquatics', 'Diving')
```

	Event	Athlete	Gold	Silver	Bronze
0	10m platform	LOUGANIS, Gregory	2	1	0
1	3m springboard	GUO, Jingjing	2	1	0
2	synchronized diving 10m platform	TIAN, Liang	1	1	0
3	synchronized diving 3m springboard	GUO, Jingjing	2	1	0

```
# top performing athletes in 'Cycling' - 'Cycling Track'
```

```
athletes_Event(df, 'Cycling', 'Cycling Track')
```

	Event	Athlete	Gold	Silver	Bronze
0	1km time trial	GRÜNKE, Klaus-Jürgen	1	0	0
1	500m time trial	BALLANGER, Felicia	1	0	0
2	Individual Pursuit	WIGGINS, Bradley	2	0	0
3	Keirin	BAYLEY, Ryan	1	0	0
4	Madison	AITKEN, Brett	1	0	0
5	Olympic Sprint	GANE, Laurent	1	0	0
6	Points Race	LLANERAS, Joan	2	1	0
7	Sprint individual	FIEDLER, Jens	2	0	1
8	Team Pursuit (4000m)	FULST, Guido	2	0	0
9	Team Sprint	NIMKE, Stefan	1	0	1
10	individual pursuit	ZIJLAARD-VAN MOORSEL, Leontien	1	0	1
11	points race	BELLUTTI, Antonella	1	0	0
12	sprint	BALLANGER, Felicia	2	0	0

2.2.6 Summary

- **The United States** emerged as the top-performing country.
- India secured 16 gold, 1 silver, and 4 bronze medals through the entire Olympics.

- The top 10 countries with the highest number of medals were identified over several years.
- A function has been created to find the medal list of countries over the years or a particular year.
- **PHELPS, Michael** of United States is the Top- Performing Athlete in the Overall Summer Olympics.
- A complete list of Top performing male and female athletes in each sport has been found.
- A function has been created to find the top performing athlete in each particular event.

Task 3: Visualizing Key Insights

3.1 Function `top_countries_by_sport()`:

```
# Find the top 3 countries for a specified year and sport.

def top_countries_by_sport(df, year, sport):
    # Filter for the given year and sport
    df_year = df[(df['Year'] == year) & (df['Sport'] == sport)]

    # Count medals by Country
    medal_counts = df_year.groupby('Country').size().reset_index(name='Medal Count')

    # Sort and get top 3 countries
    top3 = medal_counts.sort_values('Medal Count', ascending=False).head(3)

    # Plotting
    plt.figure(figsize=(8, 5))
    ax = sns.barplot(data=top3, x='Country', y='Medal Count', hue='Country', palette='Set2', dodge=False, width=0.5)

    # Add value labels
    for container in ax.containers:
        ax.bar_label(container, fmt='%.0f', label_type='edge', padding=3)

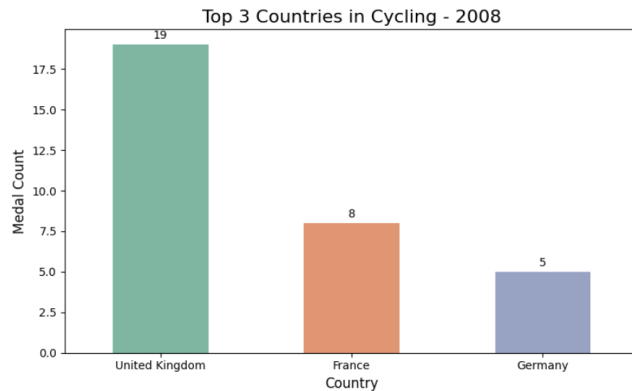
    plt.title(f"Top 3 Countries in {sport} - {year}", fontsize=16)
    plt.xlabel("Country", fontsize=12)
    plt.ylabel("Medal Count", fontsize=12)
    plt.xticks(fontsize=10)
    plt.yticks(fontsize=10)
    plt.tight_layout()
    plt.show()
```

This function filters the dataset for a specific year and sport, counts how many medals each country won, identifies the top 3 performing countries and visualizes the result using a bar plot.

The function accepts a DataFrame 'df' containing Olympic data, the year and the sport for filtering. It filters, counts the medals, gets the top 3 countries and plots a bar graph to represent them.

Example:

```
# Top 3 countries in cycling
top_countries_by_sport(df, 2008, 'Cycling')
```



3.2 Function top_country_per_sport():

```
# Find the top performing country for every sport in a given year or all the years.

def top_country_per_sport(df, year):
    # make sure the year is valid
    if year not in [1976.0, 1980.0, 1984.0, 1988.0, 1992.0, 1996.0, 2000.0, 2004.0, 2008.0, 'overall']:
        raise ValueError("YEAR must be one of: [1976.0, 1980.0, 1984.0, 1988.0, 1992.0, 1996.0, 2000.0, 2004.0, 2008.0]")

    # Filter the dataset for the given year
    if (year == 'overall'):
        df_year = df
    else:
        df_year = df[df['Year'] == year]

    # Count medals by sport and country
    medal_counts = df_year.groupby(['Sport', 'Country']).size().reset_index(name='Medal Count')

    # Find top 1 country per sport
    top_countries = medal_counts.sort_values(['Sport', 'Medal Count'], ascending=[True, False])
    top1 = top_countries.groupby('Sport').first().reset_index()

    # Plotting
    plt.figure(figsize=(12, len(top1) * 0.4))
    ax = sns.barplot(data=top1, y='Sport', x='Medal Count', color='skyblue', width=0.6)

    # Add country name labels on the bars
    for index, row in top1.iterrows():
        ax.text(row['Medal Count'] + 0.1, index, row['Country'], va='center', fontsize=9)

    plt.title(f"Top Country in Each Sport - {year}", fontsize=16)
    plt.xlabel("Medal Count")
    plt.ylabel("Sport")
    plt.grid(axis='x', linestyle='--')
    plt.tight_layout()
    plt.show()

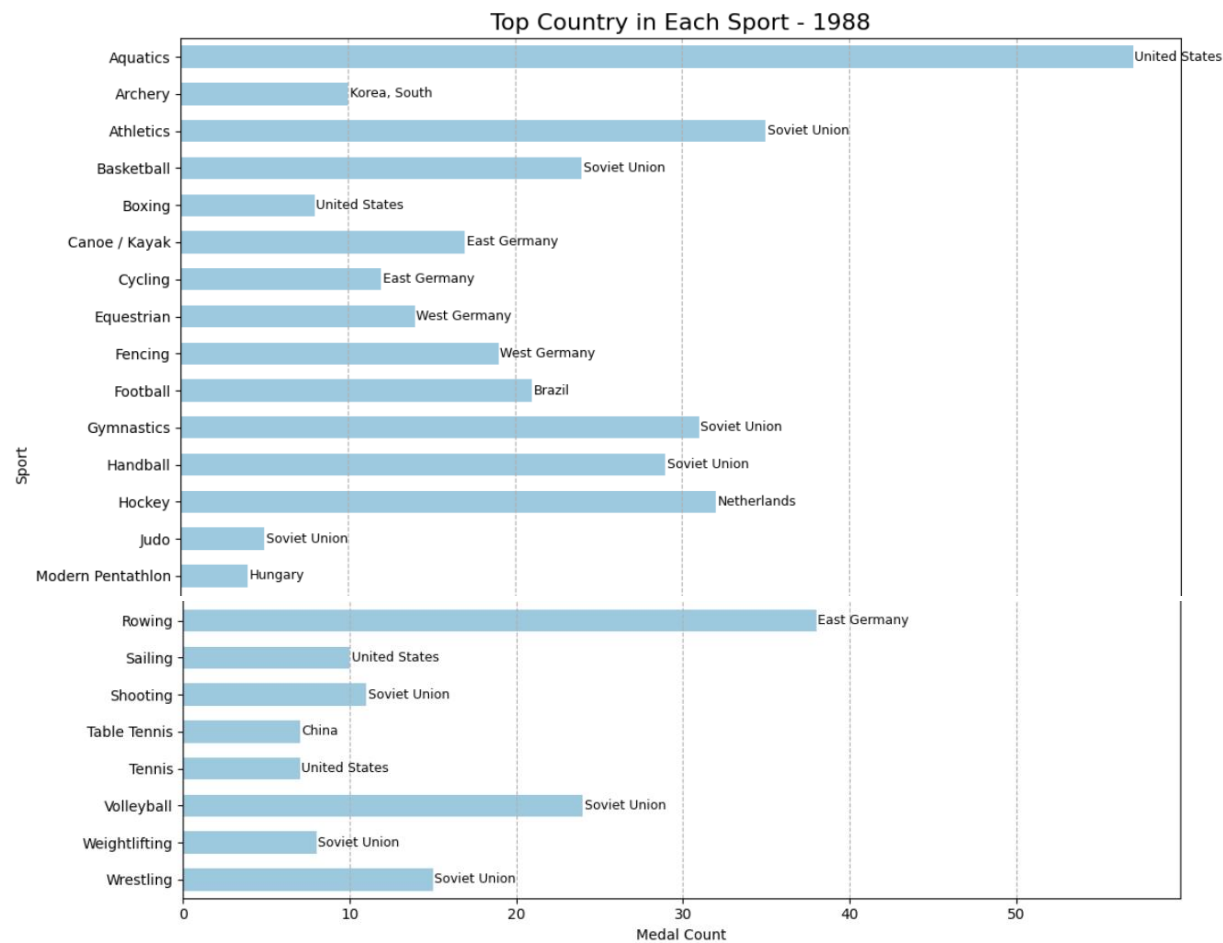
    return top1
```

This function determines the top country (based on medal count) in each sport for a specified year or overall and visualizes the result using a horizontal bar plot, where each bar shows the sport name, the number of medals and the top-performing country (labeled on the bar).

The function accepts a DataFrame 'df' containing Olympic data and the year or 'overall' for filtering. It validates the input first. Then filters and groups the data to find the top performing countries. It finally plots the data and returns the dataframe.

Example:

```
top_country_per_sport(df, 1988)
```



	Sport	Country	Medal Count
0	Aquatics	United States	57
1	Archery	Korea, South	10
2	Athletics	Soviet Union	35
3	Basketball	Soviet Union	24
4	Boxing	United States	8
5	Canoe / Kayak	East Germany	17
6	Cycling	East Germany	12
7	Equestrian	West Germany	14
8	Fencing	West Germany	19
9	Football	Brazil	21
10	Gymnastics	Soviet Union	31

(the list continues for all the sports)

3.3 Function country_acheivements():

```
def country_acheivements(df, country):
    # Filter data for the given country
    country_df = df[df['Country'] == country].copy()

    # Create a new column for total medals
    country_df.loc[:, 'Total'] = country_df[['Gold', 'Silver', 'Bronze']].sum(axis=1)

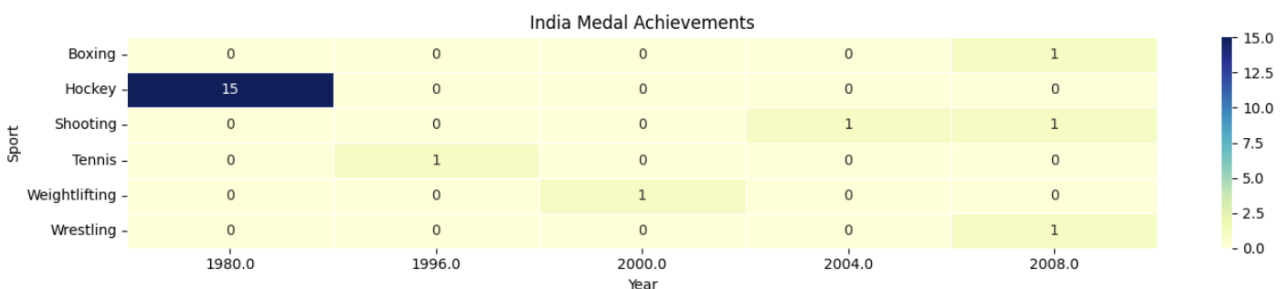
    # Group by Sport and Year, summing total medals
    pivot_data = country_df.groupby(['Sport', 'Year'])['Total'].sum().unstack(fill_value=0)

    # Plot heatmap
    plt.figure(figsize=(14, len(pivot_data) * 0.5))
    sns.heatmap(pivot_data, cmap='YlGnBu', linewidths=0.5, annot=True, fmt='d')
    plt.title(f"{country} Medal Achievements")
    plt.xlabel("Year")
    plt.ylabel("Sport")
    plt.tight_layout()
    plt.show()
```

This function filters Olympic data for a specific country, aggregates the total number of medals by sport and year and displays the result in a **heatmap**, where the intensity of color represents how successful the country was in each sport over the years. It takes the dataframe 'df' and the country name as input.

Example:

```
country_acheivements(df, 'India')
```



3.4 Gender distribution

3.4.1 Function gender_distribution_by_sport():

```
def gender_distribution_by_sport(df):
    # Group by Sport and Sex to count participants
    gender_counts = df.groupby(['Sport', 'Gender']).size().reset_index(name='Count')

    # Plot
    plt.figure(figsize=(14, 7))
    ax = sns.barplot(data=gender_counts, x='Sport', y='Count', hue='Gender', palette={'Men': 'skyblue', 'Women': 'pink'})

    plt.title("Gender Distribution Across Different Sports", fontsize=16)
    plt.xlabel("Sport", fontsize=12)
    plt.ylabel("Number of Participants", fontsize=12)
    plt.xticks(rotation=90)
    plt.tight_layout()
    plt.show()
```

3.4.2 Function gender_distribution_by_discipline():

```
def gender_distribution_by_discipline(df, sport):
    df_sport = df[df['Sport'] == sport]
    # Group by Sport and Sex to count participants
    gender_counts = df_sport.groupby(['Discipline', 'Gender']).size().reset_index(name='Count')

    # Plot
    plt.figure(figsize=(14, 7))
    ax = sns.barplot(data=gender_counts, x='Discipline', y='Count', hue='Gender', palette={'Men': 'skyblue', 'Women': 'pink'})

    plt.title(f"Gender Distribution Across Different Disciplines in {sport}", fontsize=16)
    plt.xlabel("Discipline", fontsize=12)
    plt.ylabel("Number of Participants", fontsize=12)
    plt.xticks(rotation=90)
    plt.tight_layout()
    plt.show()
```

3.4.3 Function gender_distribution_by_event():

```
# gender distribution accross events for a specific sport and discipline.

def gender_distribution_by_event(df, Sport, Discipline = None):
    # Check if the Sport is valid
    available_sports = df['Sport'].unique()
    if Sport not in available_sports:
        close_matches = difflib.get_close_matches(Sport, available_sports, n=3, cutoff=0.5)
        print(f"✗ Sport '{Sport}' not found.")
        if close_matches:
            print("Did you mean:")
            for match in close_matches:
                print("-", match)
        else:
            print("No similar sport names found.")
        return None

    # If Discipline is not provided, list all disciplines under the Sport
    if Discipline is None:
        disciplines = df[df['Sport'] == Sport]['Discipline'].unique()
        print(f"Available disciplines under '{Sport}':")
        for d in disciplines:
            print("-", d)
        return None
```



```

# Check if the Discipline is valid for the given Sport
valid_disciplines = df[df['Sport'] == Sport]['Discipline'].unique()
if Discipline not in valid_disciplines:
    close_matches = difflib.get_close_matches(Discipline, valid_disciplines, n=3, cutoff=0.5)
    print(f"✗ Discipline '{Discipline}' not found under sport '{Sport}'.")
    if close_matches:
        print("Did you mean:")
        for match in close_matches:
            print("-", match)
    else:
        print("No similar discipline names found.")
    return None

df_sport = df[(df['Sport'] == Sport) & (df['Discipline'] == Discipline)]
# Group by Sport and Sex to count participants
gender_counts = df_sport.groupby(['Event', 'Gender']).size().reset_index(name='Count')

# Plot
plt.figure(figsize=(14, 7))
ax = sns.barplot(data=gender_counts, x='Event', y='Count', hue='Gender', palette={'Men': 'skyblue', 'Women': 'pink'})

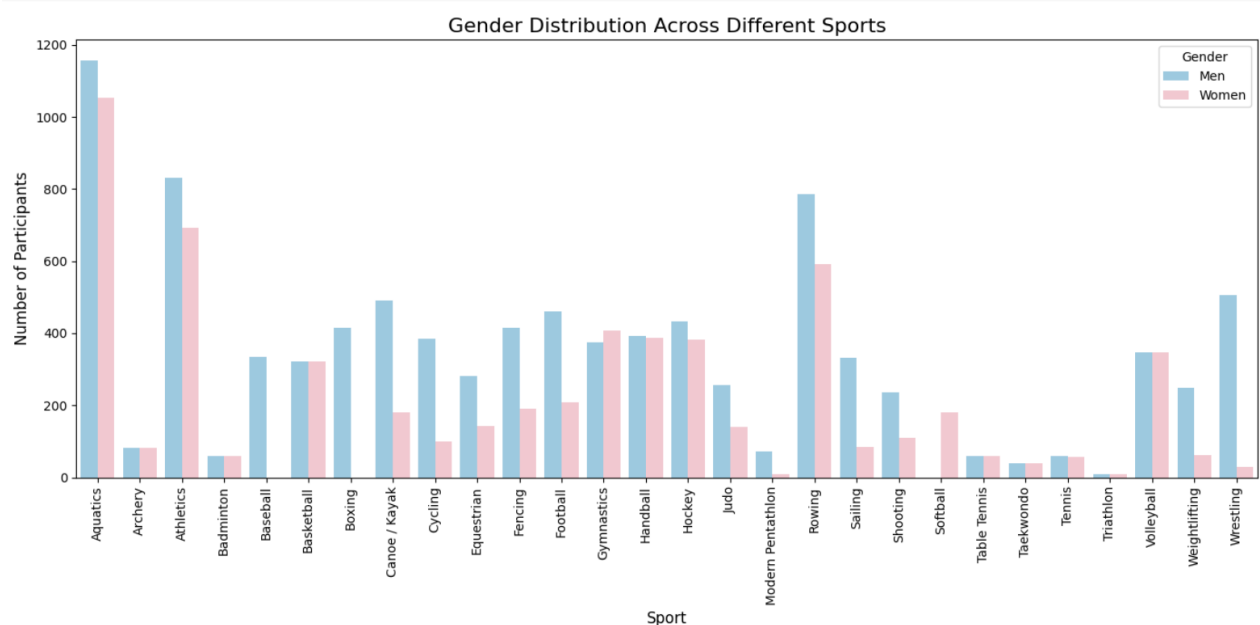
plt.title(f"Gender Distribution across events in {Sport} - {Discipline}", fontsize=16)
plt.xlabel("Event", fontsize=12)
plt.ylabel("Number of Participants", fontsize=12)
plt.xticks(rotation=90)
plt.tight_layout()
plt.show()

```

The above three functions groups Olympic data by Sport/Discipline/Event and Gender, counts the number of participants in each category and finally plots the gender distribution in each category using a grouped bar chart. They take the dataframe 'df' and the respective category as input. These functions are perfect for highlighting gender balance or disparity in Olympic sports and detecting which sports or their respective categories are male-dominated or female-dominated.

Example:

```
gender_distribution_by_sport(df)
```



Task 4: Predictive Analysis

Train a machine learning model to predict whether an athlete will win a medal based on their country, sport, and other attributes.

```
df_model = data.drop(columns = ['City', 'Discipline', 'Event', 'Year', 'Country', 'Athlete'])
df_model['Medal'] = df_model['Medal'].map({'Gold': 1, 'Silver': 1, 'Bronze': 1, np.nan: 0})
df_model
```

- ✓ Created a new dataframe 'df_model' by dropping all the unnecessary columns to make sure overfitting is absent.
- ✓ The data in column 'Medal' has been replaced with Boolean values where if the medal is won it becomes '1' else '0'.

```
x = df_model.drop('Medal', axis=1)
y = df_model['Medal']
```

- ✓ Set x contains all features except the target 'Medal'.
- ✓ Set y is the target variable — the column we want to predict.

```
categorical = ['Sport', 'Gender', 'Country_Code', 'Event_gender']
preprocessor = ColumnTransformer(
    transformers=[('cat', OneHotEncoder(handle_unknown='ignore'), categorical)]
)
```

- ✓ A list "categorical" of all categorical columns that need to be encoded is made.
- ✓ ColumnTransformer applies transformations only to specified columns.
- ✓ It uses OneHotEncoder to convert the categorical features into binary vectors.
- ✓ handle_unknown='ignore': avoids errors from categories that appear in test data but not in training.

```
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.3, random_state=42)
```

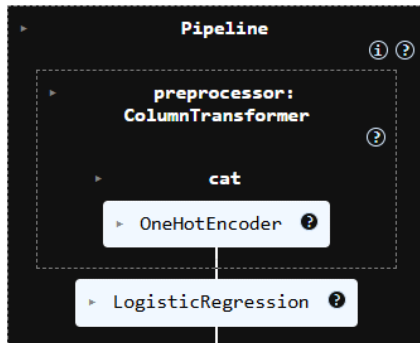
- ✓ Splits data into training (70%) and testing (30%) sets.
- ✓ random_state=42 ensures reproducibility.

```
pipeline = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('classifier', LogisticRegression(max_iter=1000))
])
```

- ✓ Builds a Pipeline with two steps:
 - 'preprocessor': One-hot encodes categorical features.

- 'classifier': Trains a logistic regression model.
- ✓ max_iter=1000 ensures the logistic regression solver has enough iterations to converge (fixes warnings).

```
pipeline.fit(x_train, y_train)
```



- ✓ Trains the entire pipeline on the training data:
 - Automatically applies encoding to the categorical features.
 - Fits the logistic regression model on the transformed features.

```
y_pred = pipeline.predict(x_test)
```

- ✓ Predicts the 'Medal' class labels for the test set using the trained model.

```
--- Model Evaluation ---
```

```
Accuracy Score: 1.0
```

```
Confusion Matrix:
```

```
[[ 28  0]
 [  0 4602]]
```

```
Classification Report:
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	28
1	1.00	1.00	1.00	4602
accuracy			1.00	4630
macro avg	1.00	1.00	1.00	4630
weighted avg	1.00	1.00	1.00	4630

The evaluation metrics of the model suggests that the model has been created successfully.

Conclusion

This project aimed to explore, analyze, and predict Olympic medal outcomes using historical data of athletes, events, countries, and their performances. Through a comprehensive blend of exploratory data analysis (EDA), data visualization, and machine learning techniques, we were able to derive meaningful insights into the dynamics of success at the Olympic Games. The initial phase involved cleaning and preparing the dataset, followed by rich exploratory analysis. We investigated trends in medal distribution across years, disciplines, and countries.

Some of the major insights include:

1. The United States was consistently the top-performing country, securing the highest number of medals over the years.
2. India's performance was specifically analyzed, revealing a tally of 16 gold, 1 silver, and 4 bronze medals across the years considered, highlighting the sports and events where the country has had the most impact.
3. Using group-by and aggregation techniques, we identified the top athletes in each sport and discipline for both men and women, providing an athlete-level perspective on excellence.
4. The use of heatmaps gave a visual representation of a country's performance by year and sport, making it easier to spot strengths, gaps, and growth areas over time.

The project incorporated a variety of insightful visualizations, including:

1. Bar plots to highlight top athletes and medal counts by sport and various relations between them.
2. Stacked bar charts showing distribution of medal types.
3. Heatmaps to explore medal wins over time and across sports for individual countries. These visualizations served not only as analytical tools but also as effective communication aids for stakeholders and audiences.

Beyond analysis, the project ventured into predictive modeling using logistic regression to classify whether an athlete would win a medal based on categorical inputs such as Country code, Sport, Gender, Event gender.

To ensure robust model performance and avoid data leakage or bias:

1. We used OneHotEncoding for categorical features instead of label encoding, avoiding unintended ordinal relationships.
2. Data was split into training and test sets using stratified sampling, preserving class balance.
3. The model pipeline streamlined preprocessing and training, with performance evaluated using accuracy score, confusion matrix, and a classification report.

Despite the simplicity of logistic regression, the model achieved a fair accuracy and demonstrated the feasibility of using basic demographic and event features for medal prediction. This approach can serve as a foundation for more advanced classification models.