| In [2]: In [6]: Out[6]:             | import numpy as np  ### import |
|-------------------------------------|--|
| In [7]: In [8]:                     | <pre>Nationality  #importing orginal dataset nat=pd.read_csv('/Users/alexandergandji/Desktop/Course Folder/Capstone Project SQL Data Science /Sports/Sports?  #creating a dataframe for all participating nations nat_dfl=pd.DataFrame(nat, columns=['NOC']) nat_dfl=nat_dfl.append({'NOC':'SGP'}, ignore_index=True)  #alocating a unique number to every value nat_dfl['NOC_ID']=range(1,232)</pre>  |
| In [11]:<br>Out[11]:                | #setting a new index nat_dfl=nat_dfl.set_index('NOC_ID')  nat_dfl  NOC  NOC_ID  1 AFG 2 AHO 3 ALB 4 ALG 5 AND 227 YMD 228 YUG 229 ZAM 230 ZIM  |
| In [12]:                            | 231 SGP  231 rows x 1 columns  #saving it on a csv file nat_dfl.to_csv('/Users/alexandergandji/Desktop/Course Folder/Capstone Project SQL Data Science /Sports/national Medal  sports_events    ID   Name   Sex   Age   Height   Weight   Team   NOC   Games   Year   Season   City   Sport   Events   |
|                                     | Aaby   |
| In [18]: In [19]: In [20]: Out[20]: | Tomasz James Winter Nagano Bobsleigh Men's F  Bobsleigh Men's F  Tomasz James Winter Nagano Bobsleigh Men's F  Bobsleigh Men's F  Tomasz James Winter Nagano Bobsleigh Men's F  Tomasz James Winter Nagano Bobsleigh Men's F  Bobsleigh Men's F  Winter Nagano Bobsleigh Men's F  Salt Lake City Bobsleigh Men's F  Wen's F  Tomasz James Winter Nagano Bobsleigh Men's F  Tomasz James Winter Nagano Bobsleigh Men's F  Bobsleigh Men's F  Winter Nagano Bobsleigh Men's F  Salt Lake City Bobsleigh Men's F  Wen's F  Tomasz James Winter Nagano Bobsleigh Men's F  Wen's F  Tomasz James Winter Nagano Bobsleigh Men's F  Bobsleigh Men's F  Winter Nagano Bobsleigh Men's F  Salt Lake City Bobsleigh Men's F  Wen's F  Tomasz James Winter Nagano Bobsleigh Men's F  Wen's F  Tomasz James Winter Nagano Bobsleigh Men's F  Wen's F  Tomasz James Winter Nagano Bobsleigh Men's F  Salt Lake Bobsleigh Men's F  Wen's F  Tomasz James Winter Nagano Bobsleigh Men's F  Wen's F  Tomasz James Winter Nagano Bobsleigh Men's F  Wen's F  Tomasz James Winter Nagano Bobsleigh Men's F  Tomasz James Winter Nagano Bob |
| In [22]: Out[22]: In [23]:          | <pre>#using lambda to allocate a number to each medal type medal_df['Medal_ID']=medal_df['Medal'].apply(lambda x: 1 if x=='Gold' else (2 if x=='Silver' else</pre>   |
| In [25]:                            | <pre>#setting a new index medal_df.set_index('Medal_ID', inplace=True)  #saving it on a csv file medal_df.to_csv('/Users/alexandergandji/Desktop/Course Folder/Capstone Project SQL Data Science /Sports/olympic  Olympic games  sports_events.head()  ID Name Sex Age Height Weight Team NOC Games Year Season City Sport Event Medal</pre>   |
| In [27]:                            | 0       1       A Dijiang       M       24.0       180.0       80.0       China       CHN       1992 Summer       1992 Summer       Barcelona       Basketball Basketball Men's Basketball Men's Basketball Men's Basketball       Nan Men's Basketball Men's Basketball       Nan Men's Basketball       Nan Men's Basketball Men's Basketball       Nan Men's Basket   |
| In [29]: In [30]: Out[30]:          | <pre>#sorting values based on Games games.sort_values('Games', inplace=True)  #using lambda to allocate to summer and winter games an unique number/ID games['Season_ID']=games['Season'].apply(lambda x: 1 if x=='Summer' else 2)  len(games.Games)  52  #allocating every game an unique number/ID</pre>   |
| In [32]: In [33]: Out[33]:          | #droping column 'Season', this could have also been done with reindex() games.drop('Season',1, inplace=True)  /var/folders/06/r72gs45101ddlg4t75zm78vc0000gn/T/ipykernel_4326/2062314925.py:1: FutureWarning: In a future version of pandas all arguments of DataFrame.drop except for the argument 'labels' will be keyword-only games.drop('Season',1, inplace=True)  games.head()  Games Year City Season_ID Games_ID  3079 1896 Summer 1896 Athina 1 1  3 1900 Summer 1900 Paris 1 2  711 1904 Summer 1904 St. Louis 1 3  268 1906 Summer 1906 Athina 1 4  1149 1908 Summer 1908 London 1 5  |
| In [35]: In [36]: In [37]: Out[37]: | <pre>games_df=games.reindex(columns=['Games', 'Games_ID', 'Year', 'Season_ID', 'City'])  #setting a new index games_df.set_index('Games_ID', inplace=True)  #saving it on a csv file games_df.to_csv('/Users/alexandergandji/Desktop/Course Folder/Capstone Project SQL Data Science /Sports/olympic  athletes  sports_events.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 Barcelona Basketball Men's NaN Basketball Judo Men's  Judo Men's NaN Basketball Judo Men's</pre>  |
| In [39]: In [40]: In [42]:          | Gunnar  Gunnar  Gunnar  Gunnar  Aaby  Aaby |
| In [105<br>In [106                  | <pre>inplace=True) athlete.replace(' Eleonora Margarida Josephina Scmitt', 'Eleonora Margarida Josephina Scmitt', inplace=True) athlete.replace(' Jean Hauptmanns', 'Jean Hauptmanns', inplace=True) athlete.replace(' Luis ngel Fernando de los Santos Grossi', 'Luis ngel Fernando de los Santos Grossi', inplace=athlete.replace(' Th Anh', 'Th Anh', inplace=True) athlete.replace(' Th Ngn Thng', 'Th Ngn Thng', inplace=True) athlete.replace(' Tin Tun', 'Tin Tun', inplace=True) athlete.sort_values('Name', inplace=True)  #allocating an unique number/ID to every athlete and merging the dataframe with the orginal data to get the get athlete['Athlete_ID']=range(1,134788) athlete=athlete.merge(sports_events[['Name', 'Sex']], how='left').drop_duplicates()  #using lambda to allocate an unique number to every gender athlete['Sex_ID']=athlete.Sex.apply(lambda x: 1 if x=='M' else 2)  #setting a new index athlete.set_index('Athlete_ID', inplace=True)  athl=athlete.dropna()</pre>   |
| In [48]:<br>Out[48]:                | #saving it on a csv file athl.to_csv('/Users/alexandergandji/Desktop/Course Folder/Capstone Project SQL Data Science /Sports/athlete.csv  Sports  sports_events.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 Barcelona Basketball Men's NaN Basketball  |
| In [50]: In [51]: Out[51]:          | #creating a new dataframe for every sport event sports=pd.DataFrame(sports_events.Event.unique().tolist(), columns=['Sport'])  #sorting data by sport events sports.sort_values('Sport', inplace=True)  Sport  Alpine Skiing Men's Combined Alpine Skiing Men's Giant Slalom Alpine Skiing Men's Slalom  Minter 1968 Willter 1968 Willter 2968 Willter Calgary Skating Women's Skating Women's Skating Women's National |
| In [52]: In [53]: In [54]: Out[54]: | 340 Wrestling Women's Heavyweight, Freestyle  292 Wrestling Women's Light-Heavyweight, Freestyle  356 Wrestling Women's Lightweight, Freestyle  609 Wrestling Women's Middleweight, Freestyle  765 rows × 1 columns  #allocating an unique ID to every sport event sport=sports['Sport_ID']=range(1,766)  #setting a new index sports.set_index('Sport_ID', inplace=True)  Sport  Sport_ID  Aeronautics Mixed Aeronautics  |
| In [55]:                            | Alpine Skiing Men's Combined  Alpine Skiing Men's Downhill  Alpine Skiing Men's Giant Slatom  Alpine Skiing Men's Slatom  Alpine Skiing Men's Slatom  Mrestling Women's Flyweight, Freestyle  Westling Women's Heavyweight, Freestyle  Westling Women's Light-Heavyweight, Freestyle  Westling Women's Lightweight, Freestyle  Westling Women's Lightweight, Freestyle  Westling Women's Middleweight, Freestyle   |
| In [248 In [249 In [250             | <pre>#merging orginal data with the nations dataframe results1=sports_events.merge(nat_df1.reset_index(), how='left')  #adding a new column, which contains the medal id; using lambda results1['Medal_ID']=results1.Medal.apply(lambda x: 1 if x=='Gold' else</pre>   |
| In [253 In [254                     | <pre>#merging results dataframe with the athlete dataframe results3=results2.merge(athlete.reset_index(), how='left')  #dropping 'Sport' column=&gt; this column just gave an general overview of the events results3.drop('Sport',1,inplace=True)  /var/folders/06/r72gs4510lddlg4t75zm78vc0000gn/T/ipykernel_4326/93642161.py:1: FutureWarning: In a future version of pandas all arguments of DataFrame.drop except for the argument 'labels' will be keyword-only results3.drop('Sport',1,inplace=True)  #renaming Event to Sport=&gt; the new Sport column has a more detailed overview of the different events results3.rename(columns={'Event':'Sport'}, inplace=True)  #merging results dataframe with sports dataframe results4=results3.merge(sports.reset_index(), how='left')</pre>  |
| In [256                             | #sort result dataframe by Year results4.sort_values('Year', inplace=True)  results4    ID   Name   Sex   Age   Height   Weight   Team   NOC   Games   Year     City   Sport   Medal   NOC        |
|                                     | 181753         91353         Paraskevopoulos         M         Nan         Nan         Nan         Greece         GRE         Summer         1896          Athlina         Cycling Men's 12-Hours Race         Nan           204736         102806         Richard Rstel         M         Nan         Nan         Nan         Germany         GER         1896 Summer         1896          Athlina         Gymnastics Men's Pommelled Horse                Athlina         Men's Pommelled Horse         Nan           40051         20599         Phillip Hung Chew         M         22.0         173.0         91.0         United States         USA         2016 Summer         2016         Rio de Janeiro         Badminton Men's Doubles   |
| In [258                             | 194190         97481         Bence Pulai         M         24.0         193.0         89.0         Hungary         HUN         2016 Summer         2016 Summer         Rio de Janeiro         Swimming Men's 100 metres Butterfly         NaN           194191         97481         Bence Pulai         M         24.0         193.0         89.0         Hungary         HUN         2016 Summer         2016 Summer         Swimming Men's 100 metres Butterfly         NaN           40081         20616         Danny Chia         M         43.0         170.0         75.0         Malaysia         MAS         2016 Summer         2016 Summer         Rio de Janeiro         Golf Men's Individual Individual NaN         NaN           245289         122828         Nigora Tursunkulova         F         17.0         181.0         67.0         Uzbekistan         UZB         2016 Summer         Rio de Janeiro         Taekwondo Women's Welterweight         NaN           271119 rows × 21 columns         4due to the fact that we have data points which did not have any value in the Age, Height and Weight column we #used ffill()         Fermion Tursunkulova         Height and Weight column we #used ffill()           4due to the fact that we have data points which did not have any value in the Age, Height and Weight column         Papplying ffill() in the first row of the Age column  |
| In [259                             | results4.Age.replace(results4.Age.iloc[0],21, inplace=True) #transforming Age from float to integer results4['Age']=results4.Age.apply(np.int64)  results4=results4.copy()  #applying ffill() in the first row of the Height column results4['Height']=results4.Height.replace(results4.Height.iloc[0],176) #transforming Height from float to integer results4['Height']=results4.Height.apply(np.int64)  ### 176 was chosen, because the height was found at other athletes in the same time period and disceplines  results4  |
|                                     | 204735         102806         Richard Rstel         M         21         176         NaN         Germany         GER         1896 Summer         1896 Summer          Athina         Men's Horizontal Bar, Teams         Gold           181753         91353         Georgios Paraskevopoulos         M         21         176         NaN         Greece         GRE         1896 Summer         1896 Summer          Athina         Cycling Men's Road Race, Individual           181754         91353         Georgios Paraskevopoulos         M         21         176         NaN         Greece         GRE         1896 Summer         1896 Summer          Athina         Cycling Men's Road Race, Individual           204736         102806         Richard Rstel         M         21         176         NaN         Germany         GER         1896 Summer         1896 Summer          Athina         Men's Pommelled Horse         NaN           204736         102806         Richard Rstel         M         21         176         NaN         Germany         GER         1896 Summer         1896 Summer          Athina         Men's Pommelled Horse         NaN           40051         20599         Phillip Hung Chew  |
|                                     | 194190         97481         Bence Pulai         M         24         193         89.0         Hungary         HUN         2016 Summer         2016 Summer          Rio de Janeiro         Swimming Men's 100 metres Butterfly         NaN           194191         97481         Bence Pulai         M         24         193         89.0         Hungary         HUN         2016 Summer         2016         Rio de Janeiro         Swimming Men's 4 x 100 metres Medley Relay           40081         20616         Danny Chia         M         43         170         75.0         Malaysia         MAS         2016 Summer         2016         Rio de Janeiro         Golf Men's Individual         NaN           245289         122828         Nigora Nigora F 17         181         67.0         Uzbekistan         UZB Summer         2016         Rio de Janeiro         Taekwondo Women's NaN   |
| In [262<br>In [263                  | NIGOTA F 17 191 670 Limbelistan LIZE 2010 2016 RIO de Wemenle Nell   |
|                                     | results4['Weight']=results4.Weight.apply(np.int64)  results4[results4.Athlete_ID.isnull()]  ID Name Sex Age Height Weight Team NOC Games Year City Sport Medal NOC_I  92100 46661  |
|                                     | Sind Santos Grossi  Eleonora Margarida Josephina Scmitt  Eleonora Margarida Josephina Josephin |
|                                     | 51335 26388 Fernando de los Santos Grossi  M 27 176 67 Uruguay URU 1952 Helsinki Men's Road Race, Individual  51336 26388 Luis ngel Fernando de los Santos Grossi  |
|                                     | 51338         26388         Fernando de los Santos Grossi         M         27         176         67         Uruguay         URU         1952 Summer         1952         Helsinki         Men's Team Pursuit, 4,000 metres         NaN         22           58302         29843         Th Ngn Thng         F         19         147         47         Vietnam         VIE         2008 Summer         2008         Beijing         Gymnastics Women's Individual All-Around           58303         29843         Th Ngn Thng         F         19         147         47         Vietnam         VIE         2008 Summer         2008         Beijing         Gymnastics Women's Floor Exercise           Th Ngn         Th Ngn         F         19         147         47         Vietnam         VIE         2008 Summer         2008         Beijing         MaN         2008 Summer  |
|                                     | 58304         29843         Thong         F         19         147         47         Vietnam         VIE         2008 Summer         2008         Beijing         Women's Horse Vault         NaN         23           58305         29843         Th Ngn Thng         F         19         147         47         Vietnam         VIE         2008 Summer         2008         Beijing         Gymnastics Women's Uneven Bars         NaN         23           58306         29843         Th Ngn Thng         F         19         147         47         Vietnam         VIE         2008 Summer         2008         Beijing         Gymnastics Women's Balance Beam         NaN         23           58306         29843         Th Ngn Thng         F         19         147         47         Vietnam         VIE         2008 Summer         2008         Beijing         Gymnastics           58306         29843         Th Ngn Thng         F         19         147         47         Vietnam         VIE         2008 Summer         2008         Beijing         Gymnastics           69mnastics         Gymnastics         Beijing         Women's NaN         23         2008         2008         2008         2008  |
|                                     | 58308         29843         Th Ngn Thng         F         23         147         47         Vietnam         VIE         2012 Summer         2012         London         Women's Balance Beam         NaN         23           58307         29843         Th Ngn Thng         F         23         147         47         Vietnam         VIE         2012 Summer         2012         London         Gymnastics Women's Uneven Bars         NaN         23           1560         869         "Gabrielle Marie "Gabby" Adcock (White-)         F         25         167         67         Great Britain         GBR Summer         2016         Rio de Janeiro         Badminton Mixed Doubles         NaN         7           3016         Women's Summer         Pio de Women's Summer         Women's Summer         Pio de Women's Summer         Women's Summer   |
| In [265                             | 58301 29842 Th Anh F 20 165 58 Vietnam VIE 2016 Summer 2016 Rio de Janeiro Foil, Individual  17 rows × 21 columns  #these data points were dropped, because even after the use of ffill(), the data points were too ambigous results4.drop([92100, 58301, 1560, 58307, 58308, 51338, 58306, 58305, 58304, 58303, 58302, 51337, 51336, 51335, 51334, 215055, 215056], 0, inplace=True)  /var/folders/06/r72gs4510lddlg4t75zm78vc0000gn/T/ipykernel_4326/564590031.py:1: FutureWarning: In a future version of pandas all arguments of DataFrame.drop except for the argument 'labels' will be keyword-only results4.drop([92100, 58301, 1560, 58307, 58308, 51338, 58306, 58305, 58304, 58303, 58302, 51337, 51336, 51336, 51338, 58306, 58305, 58304, 58303, 58302, 51337, 51336, 51338, 58306, 58305, 58304,      |
|                                     | ID   Name   Sex   Age   Height   Weight   Team   NOC   Games   Year     City   Sport   Medal   NOC   N       |
|                                     | 181753         91353         Georgios Paraskevopoulos         M         21         176         67         Greece         GRE         1896 Summer   |
|                                     |  |
|                                     | 271102 rows × 21 columns   |
| In [267 In [268                     | <pre>#transforming Athlete_ID from float to integer results4['Athlete_ID']=results4.Athlete_ID.apply(np.int64)  #reindexing the columns results5=results4.reindex(columns=['Athlete_ID', 'Sex', 'Age', 'Height', 'Weight', 'NOC_ID',</pre>   |
| In [267 In [268                     | <pre>#transforming Athlete_ID from float to integer results4['Athlete_ID']=results4.Athlete_ID.apply(np.int64)  #reindexing the columns results5=results4.reindex(columns=['Athlete_ID', 'Sex', 'Age', 'Height', 'Weight', 'NOC_ID',</pre>   |