# Bigdata Assignment 3.2

This Data set is about Olympics. You can download the data set from the below link:

https://drive.google.com/open?id=0ByJLBTmJojjzV1czX3Nha0R3bTQ

DATE SET DESCRIPTION

The data set consists of the following fields.

Athlete: This field consists of the athlete name

Age: This field consists of athlete ages

Country: This fields consists of the country names which participated in Olympics

Year: This field consists of the year

Closing Date: This field consists of the closing date of ceremony

Sport: Consists of the sports name

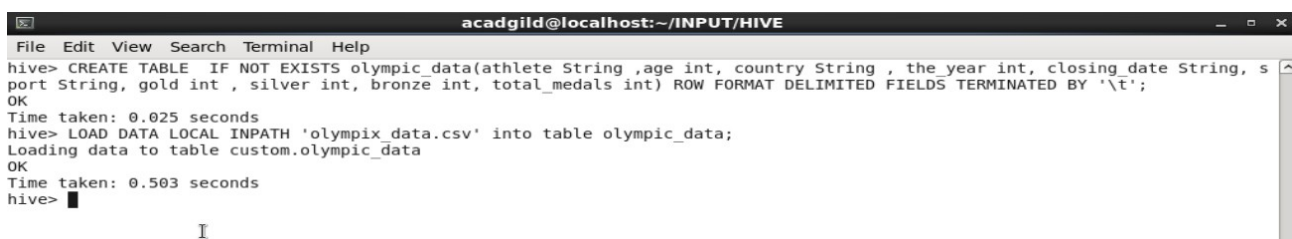Gold Medals: No. of Gold medals

Silver Medals: No. of Silver medals

Bronze Medals: No. of Bronze medals

Total Medals: Consists of total no. of medals

<u>Table Creation:-</u>  Dataset olympix.csv was loaded into a table called olmpic_data with the given columns.

**CREATE TABLE IF NOT EXISTS olympic_data(athlete String , age int , country String, the_year int , closing_date String , sport String , gold int , silver int , bronze int , total_medals int) ROW FORAT DELIMITED FIELDS TERMINATED BY '\t' ;**
**LOAD DATA LOCAL INPATH 'olmpix.csv' into table olympic_data ;**



1. Write a Hive program to find the number of medals won by each country in swimming.

Ans –  We selected each country and sum of total medals in Swimming using GROUP BY.

**select country , SUM(total_medals) from olmpic_data where sport =**

**'Swimming'  GROUP BY country;**



Output- We got the list of each country and its corresponding no of total medals
 in Swimming.



2. Write a Hive program to find the number of medals that India won year wise.

Ans –  We calculated in each year no of medals won by India using GROUP
BY.

**select the_year , SUM(total_medals) from olympic_data where country =
'India' GROUP BY the_year;**

In the above screenshot we got the output in which year India got medals.

3. Write a Hive Program to find the total number of medals each country won.



Ans – We calculated no of medals won by each country using GROUP By.
**select country , SUM(total_medals) from olmpic_data GROUP BY country;**

OUTPUT -
We got the list of country and its corresponding no of total medals won.



```
                                                              acadgild
 File   Edit   View   Search   Terminal   Help
Total MapReduce CPU Time Spent: 4 seconds 810 msec
OK
Afghanistan       2
Algeria 8
Argentina         141
Armenia 10
Australia         609
Austria 91
Azerbaijan        25
Bahamas 24
Bahrain 1
Barbados          1
Belarus 97
Belgium 18
Botswana          1
Brazil  221
Bulgaria          41
Cameroon          20
Canada  370
Chile   22
China   530
Chinese Taipei  20
Colombia          13
Costa Rica        2
Croatia 81
Cuba    188
Cyprus  1
Czech Republic  81
Denmark 89
Dominican Republic      5
Ecuador 1
Egypt   8
Eritrea 1
Estonia 18
Ethiopia          29
Finland 118
France  318
Gabon   1
Georgia 23
```

4. Write a Hive program to find the number of gold medals each country won.

Ans -   We calculated no of gold medals won by each country using GROUP By.
**select country , SUM(gold) from olmpic_data GROUP BY country;**

OUTPUT -  We got the list of country and its corresponding no of total gold medals won.

```
                                                    acadgild

 File  Edit  View  Search  Terminal  Help
Total MapReduce CPU Time Spent: 5 seconds 780 msec
OK
Afghanistan     0
Algeria 2
Argentina       49
Armenia 0
Australia       163
Austria 36
Azerbaijan      6
Bahamas 11
Bahrain 0
Barbados        0
Belarus 17
Belgium 2
Botswana        0
Brazil  46
Bulgaria        8
Cameroon        20
Canada  168
Chile   3
China   234
Chinese Taipei  2
Colombia        2
Costa Rica      0
Croatia 35
Cuba    57
Cyprus  0
Czech Republic  14
Denmark 46
Dominican Republic      3
Ecuador 0
Egypt   1
Eritrea 0
Estonia 6
Ethiopia        13
Finland 11
France  108
Gabon   0
Georgia 6
```

```
Panama    1
Paraguay            0
Poland    20
Portugal            1
Puerto Rico         0
Qatar     0
Romania 57
Russia    234
Saudi Arabia        0
Serbia    1
Serbia and Montenegro    11
Singapore           0
Slovakia            10
Slovenia            5
South Africa        10
South Korea         110
Spain     19
Sri Lanka           0
Sudan     0
Sweden    57
Switzerland         21
Syria     0
Tajikistan          0
Thailand            6
Togo      0
Trinidad and Tobago      1
Tunisia 2
Turkey    9
Uganda    1
Ukraine 31
United Arab Emirates     1
United States       552
Uruguay 0
Uzbekistan          5
Venezuela           1
Vietnam 0
Zimbabwe            2
Time taken: 25.775 seconds, Fetched: 110 row(s)
hive> █
```