

```

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import plotlylib inline

# Insight 1 : Brazil has won the tournament most number of times
data=pd.read_csv(r"C:\Users\Debansh\Desktop\Fifa project\WorldcupMatches.csv")
data=pd.read_csv(r"C:\Users\Debansh\Desktop\Fifa project\WorldCupPlayers.csv")
data=pd.read_csv(r"C:\Users\Debansh\Desktop\Fifa project\WorldCupsGoals.csv")

In [50]:
data.head()

Out[50]:
   Year  Datetime Stage Stadium City Home Team Home Team Goals Away Team Goals Away Team Name Win conditions Attendance Half-time Home Goals Half-time Away Goals Referee Assistant 1 Assistant 2 RoundID MatchHD Home Team Initials Away Team Initials
0 1930.0 12-Jul-1930 Group 1 Potos Montevideo France 4.0 1.0 Mexico LOMBARDO Domingo CRISTOPHE Henry (BEL) REGO Gilberto (BRA) 201.0 1096.0 FRA MEX
1 1934.0 12-Jun-1934 Group 4 Parque Central Montevideo USA 3.0 0.0 Belgium MACIAS Jose (ARG) MATEUCCI Francisco (URU) WARWEN Alberto (CHI) 201.0 1090.0 USA BEL
2 1930.0 1930-12-25 Group 2 Parque Central Montevideo Yugoslavia 2.0 1.0 Brazil TELADA Anibal (URU) VALLARINO Ricardo (FRA) BALWAY Thomas (FRA) 201.0 1093.0 YUG BRA
3 1930.0 14-Jul-1930 Group 3 Potos Montevideo Romania 3.0 1.0 Peru 2549.0 1.0 0.0 WARWEN Alberto (CHI) LANGENIUS Juan (BEL) MATEUCCI Francisco (URU) 201.0 1098.0 ROU PER
4 1930.0 15-Jul-1930 Group 1 Parque Central Montevideo Argentina 1.0 0.0 France 23409.0 0.0 0.0 REGO Gilberto (BRA) SAUCEDO Ulises (BOL) RADULESCU Constantin (ROU) 201.0 1095.0 ARG FRA

In [51]:
data2=data.head()

Out[51]:
RoundID MatchHD Team Initials Coach Name Line-up Shirt Number Player Name Position Event
0 201 1096 FRA CAUDRON Raul(FRA) S 0 Alex THEPOT GK NAN
1 201 1096 MEX LUQUE Juan(MEX) S 0 Oscar BONFIGLIO GK NAN
2 201 1096 FRA CAUDRON Raul(FRA) S 0 Marcel LANGILLER NAN G40
3 201 1096 MEX LUQUE Juan(MEX) S 0 Juan CARRENO NAN G70
4 201 1096 FRA CAUDRON Raul(FRA) S 0 Ernest LIBERATI NAN NAN

In [52]:
data3=data.head()

Out[52]:
Year Country Winner Runners-Up Third Fourth GoalsScored QualifiedTeams MatchesPlayed Attendance
0 1930 Uruguay Uruguay Argentina USA Yugoslavia 70 13 18 590549
1 1934 Italy Czechoslovakia Germany Austria 70 16 17 363000
2 1938 France Italy Hungary Brazil Sweden 84 15 18 275700
3 1950 Brazil Uruguay Brazil Sweden Spain 88 13 22 1045246
4 1954 Switzerland Germany FR Hungary Austria Uruguay 140 16 26 769067

In [53]:
import plotly as py
import cufflinks as cf
from plotly.offline import plot
py.figure_init_notebook_mode(connected=True)
cf.go_offline()

In [54]:
#----DATASET LOADED NOW WE WILL GET SOME INSIGHTS ABOUT THE DATA-----

In [55]:
# COUNTRY WISE ANALYSIS

In [56]:
data_countries = pd.DataFrame(data3['winner'].value_counts())
data_countries

Out[56]:
Winner
Brazil      5
Italy       4
Germany FR  3
Uruguay     2
Argentina   2
England     1
France      1
Spain       1
Germany     1

In [57]:
# Insight 1 : Brazil has won the tournament most number of times

In [58]:
# Tabular Representation of above mentioned data
data_countries.plot(kind='bar',y='winner',title='Countries who have won worldcups',colors='blue')

Out[58]:


In [59]:
# Now collaborating for first three positions we are having
data_winner_up=pd.DataFrame(data3['runner-up'].value_counts())
data_runner_up=pd.DataFrame(data3['runners-up'].value_counts())
data_third=pd.DataFrame(data3['third'].value_counts())

In [60]:
data_winner.head()

Out[60]:
Winner
Brazil      5
Italy       4
Germany FR  3
Argentina   2

In [61]:
data_runner_up.head()

Out[61]:
Runners-Up
Argentina    3
Germany FR   3
Netherlands  3
Czechoslovakia 2
Hungary      2

In [62]:
data_third.head()

Out[62]:
Third
Germany      3
Brazil       2
Sweden       2
France       2
Poland       2

In [63]:
teams = pd.concat([data_winner, data_runner_up, data_third], axis = 1)

Out[63]:
Winner Runners-Up Third
Brazil      5.0      2.0      2.0
Italy       4.0      2.0      1.0
Germany FR  3.0      3.0      1.0
Uruguay     2.0      NaN      NaN
Argentina   2.0      3.0      NaN
England     1.0      NaN      NaN
France      1.0      1.0      2.0
Spain       1.0      NaN      NaN
Germany     1.0      1.0      1.0
Netherlands NaN      3.0      1.0
Czechoslovakia NaN      2.0      NaN
Hungary     NaN      2.0      NaN
Sweden      NaN      1.0      2.0
Poland      NaN      NaN      2.0
USA         NaN      NaN      1.0
Austria     NaN      NaN      1.0
Chile       NaN      NaN      1.0
Paraguay   NaN      NaN      1.0
Costa Rica  NaN      NaN      1.0
Turkey     NaN      NaN      1.0

In [64]:
# Dealing with NaN values
teams.fillna(0,inplace=True)

In [65]:
teams=teams.astype(int)

In [66]:
# Insight 2 : A complete depiction of number of world cups won, first runner-up, and second runner-up positions by various participating teams
teams.plot(kind='bar',yTitle='Count',title='Country wise analysis',xTitle='Country')

Out[66]:


In [67]:
# NUMBER OF GOALS PER COUNTRY

In [68]:
data1=data.head()

Out[68]:
   Year  Datetime Stage Stadium City Home Team Home Team Goals Away Team Goals Away Team Name Win conditions Attendance Half-time Home Goals Half-time Away Goals Referee Assistant 1 Assistant 2 RoundID MatchHD Home Team Initials Away Team Initials
0 1930.0 12-Jul-1930 Group 1 Potos Montevideo France 4.0 1.0 Mexico LOMBARDO Domingo CRISTOPHE Henry (BEL) REGO Gilberto (BRA) 201.0 1096.0 FRA MEX
1 1930.0 12-Jun-1934 Group 4 Parque Central Montevideo USA 3.0 0.0 Belgium MACIAS Jose (ARG) MATEUCCI Francisco (URU) WARWEN Alberto (CHI) 201.0 1090.0 USA BEL
2 1930.0 14-Jul-1930 Group 2 Parque Central Montevideo Yugoslavia 2.0 1.0 Brazil 24059.0 2.0 0.0 TELADA Anibal (URU) VALLARINO Ricardo (FRA) BALWAY Thomas (FRA) 201.0 1093.0 YUG BRA
3 1930.0 14-Jul-1930 Group 3 Potos Montevideo Romania 3.0 1.0 Peru 2549.0 1.0 0.0 WARWEN Alberto (CHI) LANGENIUS Juan (BEL) MATEUCCI Francisco (URU) 201.0 1098.0 ROU PER
4 1930.0 15-Jul-1930 Group 1 Parque Central Montevideo Argentina 1.0 0.0 France 23409.0 0.0 0.0 REGO Gilberto (BRA) SAUCEDO Ulises (BOL) RADULESCU Constantin (ROU) 201.0 1095.0 ARG FRA

In [69]:
# Separating data based on goals scored by teams
data_home=data[[('Home Team Name','Home Team Goals')].dropna()]
data_away=data[[('Away Team Name','Away Team Goals')].dropna()]

In [70]:
data_home.head()

Out[70]:
Home Team Name Home Team Goals
0 France 4.0
1 USA 3.0
2 Yugoslavia 2.0
3 Romania 3.0
4 Argentina 1.0

In [71]:
data_away.head()

Out[71]:
Away Team Name Away Team Goals
0 Mexico 1.0
1 Belgium 0.0
2 Brazil 1.0
3 Peru 1.0
4 France 0.0

In [72]:
# Setting up the columns in both the tables
data_home.columns=['Country','Goals']
data_away.columns=['Country','Goals']

In [73]:
data_country_goals = pd.concat([data_home, data_away], ignore_index=True)

In [74]:
data_country_goals

Out[74]:
Country Goals
0 France 4.0
1 USA 3.0
2 Yugoslavia 2.0
3 Romania 3.0
4 Argentina 1.0
... ..
1699 Costa Rica 0.0
1700 Germany 7.0
1701 Argentina 0.0
1702 Netherlands 3.0
1703 Argentina 0.0
1704 rows x 2 columns

In [75]:
# The above table do contain all the goals both home and away but can have different values for same countries, so...
data_final_country_goal=data_country_goals.groupby('Country').sum()

In [76]:
# Arranging by number of goals
final_data=data_final_country_goal.sort_values(by='Goals',ascending=False)

In [77]:
final_data[final_data[:10]]

Out[77]:
Country Goals
0 Brazil 225.0
1 Argentina 133.0
Germany FR 131.0
Italy 128.0
France 108.0
Germany 104.0
Spain 92.0
Netherlands 87.0
Hungary 87.0
Uruguay 80.0

In [78]:
# Insight 3 : Brazil scored the most number of goals throughout the history of worldcup followed by Argentina and Germany.
final_data.plot(kind='bar',yTitle='No of Goals',title='Countries with maximum number of goals',colors='red',xTitle='Country')

Out[78]:


In [79]:
# Comparing half time home goals scored and half time away goals scored
half_time_home=pd.DataFrame(data1[['Home Team Name','Half-time Home Goals']])
half_time_away=pd.DataFrame(data1[['Away Team Name','Half-time Away Goals']])

half_time_home = half_time_home.groupby('Home Team Name').sum()
half_time_away = half_time_away.sort_values(by='Half-time Away Goals',ascending=False)
half_time_away

Out[79]:
Home Team Name Half-time Home Goals
Brazil 66.0
Argentina 48.0
Germany FR 38.0
Italy 36.0
Hungary 33.0
-- --
Norway 0.0
Iran 0.0
New Zealand 0.0
Iraq 0.0
United Arab Emirates 0.0
78 rows x 1 columns

In [80]:
half_time_away = half_time_away.groupby('Away Team Name').sum()
half_time_away = half_time_away.sort_values(by='Half-time Away Goals',ascending=False)
half_time_away

Out[80]:
Away Team Name Half-time Away Goals
Spain 20.0
Germany 18.0
Netherlands 18.0
France 17.0
Brazil 17.0
-- --
China PR 0.0
Dutch East Indies 0.0
Slovenia 0.0
Haiti 0.0
83 rows x 2 columns

In [81]:
# Concatinating both the tables on team name
total = pd.concat([half_time_home, half_time_away], axis = 1)
total

Out[81]:
Half-time Home Goals Half-time Away Goals
Brazil 66.0 17.0
Argentina 48.0 8.0
Germany FR 38.0 12.0
Italy 36.0 13.0
Hungary 33.0 7.0
-- --
Egypt NaN 2.0
Israel NaN 0.0
Kuwait NaN 0.0
El Salvador NaN 0.0
Dutch East Indies NaN 0.0
83 rows x 2 columns

In [82]:
# Creating total goals columns to order the table based on total number of goals scored by a team
total['total_goals'] = total['Half-time Home Goals'] + total['Half-time Away Goals']
total = total.sort_values(by='total_goals',ascending=False)
total[total[:18]]

Out[82]:
Team Half-time Home Goals Half-time Away Goals total_goals
Brazil 66.0 17.0 83.0
Argentina 48.0 8.0 56.0
Germany FR 38.0 12.0 50.0
Italy 36.0 13.0 49.0
Germany 28.0 18.0 46.0
France 25.0 17.0 42.0
Hungary 33.0 7.0 40.0
Netherlands 19.0 18.0 37.0
Spain 16.0 20.0 36.0
Uruguay 27.0 6.0 33.0
-- --
Korea Rep 0.0 0.0 0.0
South Korea 0.0 0.0 0.0
North Korea 0.0 0.0 0.0
Vietnam 0.0 0.0 0.0
Laos 0.0 0.0 0.0
Macao 0
```