**INTRODUCTION TO DATA MANAGEMENT**

**PROJECT REPORT**

(Project Semester August-December 2021)

# *FIFA FOOTBALL WORLDCUP(1930-2014)*
# *DATA ANALYSIS*
# *IN EXCEL*

Submitted by

**Rajeev Ranjan Pan**

Registration No:11902929

Computer Science and Engineering

Section:KM009

Roll:RKM009A23

Course Code INT217

Under the Guidance of

**Sandeep Kaur**

**(UID:23614)**

**Discipline of CSE/IT**

**Lovely School of Computer Science and Engineering**

**Lovely Professional University, Phagwara**

# CERTIFICATE

This is to certify that **Rajeev Ranjan Pan** bearing Registration no.**11902929** has completed **INT217** project titled, **"FIFA Football Data Analysis in Excel(1930-2014)"** under my guidance and supervision. To the best of my knowledge, the present work is the result of his/her original development, effort and study.

**Signature and Name of the Supervisor**
**Designation of the Supervisor**
**School of Computer Science**
Lovely Professional University
Phagwara, Punjab.

Date: 08/12/2021

# DECLARATION

I, **Rajeev Ranjan Pan**, student of **Data Science** under CSE/IT Discipline at, Lovely Professional University, Punjab, hereby declare that all the information furnished in this project report is based on my own intensive work and is genuine.

Name of the student: Rajeev Ranjan Pan

Date:08/12/2021

Signature:

Registration No.11902929

# ACKNOWLEDGEMENT

A project work is a combination of views, ideas, suggestions and contribution of many people. Thus, one of the pleasant parts of writing the report is to thank those who have contributed towards its fulfilment.

I consider it as great privilege to have esteemed Lecturer **Ms. Sandeep Kaur** as my project guide. I take this opportunity to express my sincere gratitude to her through constant advice and constructive criticism nourished my interest in the subject and provided a free and pleasant atmosphere to work against all odd situations. I avail this opportunity to extend my heart full thanks and deep respect to faculty member for their able guidance during this project.

My gratitude to all those, who responded to my questionnaire in a well-defined manner and helped me acquiring knowledge.

I would like to communicate a deep sense of gratitude to all these people without whom my project would not have been such a great learning experience.

Rajeev Ranjan Pan

KM009

Reg no: 11902929

Lovely Professional University
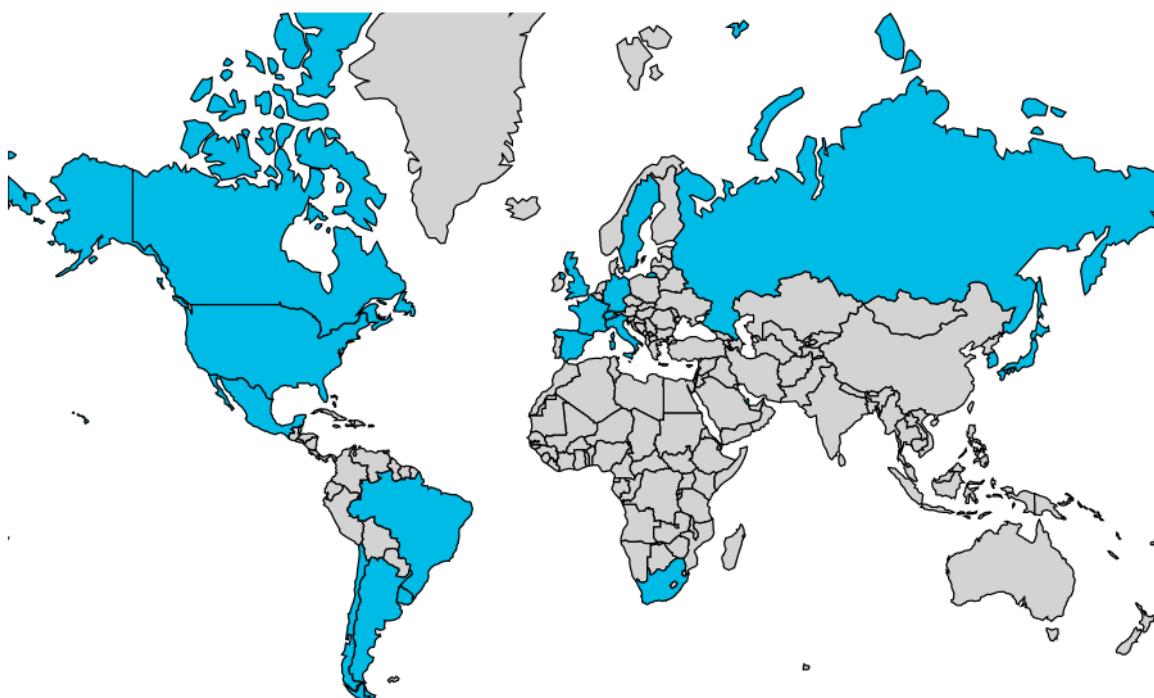
# **Table of Contents**

# **Introduction**

The FIFA World Cup, often simply called the World Cup, is an international association football competition contested by the senior men's national teams that belong to the Fédération Internationale de Football Association (FIFA), the sport's global governing body. The championship has been awarded every four years since the inaugural tournament in 1930, except in 1942 and 1946 when it was not held because of the Second World War. The current champion is France, which won its second title at the 2018 tournament in Russia and the next world cup will happen in 2022 in Qatar. As of 2019, there are 211 teams eligible to qualify for the championship. Only 32 teams make it to the finals to determine one winner. IFA's Council chooses the host countries. A balloting system is used to determine which bidding nations will become the host country. This system was put into place to avoid boycotts and controversies – problems that had plagued the tournament in its early years.

The first host country in 1930 was Uruguay. In 1934, Italy was the host country. The next host country selected was France. All host countries of the FIFA World Cup are as follows:

- Canada, United States, and Mexico: 2026
- Qatar: 2022
- Russia: 2018
- Brazil: 2014
- South Africa: 2010
- Germany: 2006
- Japan/South Korea: 2002
- United States: 1994
- Mexico: 1986
- Spain: 1982
- Argentina: 1978
- West Germany: 1974
- Mexico: 1970
- England: 1966
- Chile: 1962
- Sweden: 1958
- Switzerland: 1954
- Brazil: 1950, 2014
- France: 1938, 1998
- Italy: 1934, 1990
- Uruguay: 1930

These datasets includes the data about FIFA worldcup teams, venues, winner, attendance and matches for worldcup years between 1930 to 2014(20 years).

FIFA World Cup Host Countries



These datasets includes the worldcup results and team information for 20 years in which worldcup was held between 1930 to 2014. The First Dataset(Winning Title Dataset) contain year in which worldcup is played, hosting countries, Winner, Runner-Up,Third,Fourth,Goals     Scored,Qualified     Teams,Matches played,Attendence of the audience.The second Dataset(Matches Dataset) contains Year,Datetime,Stage,stadium,city,home team and away team name,Home team and away team goal,Referee name.

- **Datasets**

  The World Cups dataset show all information about all the World Cups in the history, while the World Cup Matches dataset shows all the results from the matches contested as part of the cups.

| S.no | Dataset Name | Definition |
|------|--------------|------------|
| 1. | **Winning Titles** | Contains data of teams according to titles won per year. |

| 2. | Matches | Contains data according to goals and home/away team name. |
|----|---------|----------------------------------------------------------|

## 1.Winning Titles dataset

| S.No | Column name | Definition |
|------|-------------|------------|
| 1. | **Year** | Year of the worldcup |
| 2. | **Country** | Hosting country of the worldcup of the particular year |
| 3. | **Winner** | Team who won the worldcup |
| 4. | **Runner-Up** | Team who was the second place |
| 5. | **Third** | Team who was the third place |
| 6. | **Fourth** | Team who was the fourth place |
| 7. | **GoalsScored** | Total goals scored in the worldcup |
| 8. | **QualifiedTeams** | Total participating teams |
| 9. | **MatchesPlayed** | Total matches played in the cup |
| 10. | **Attendance** | Total attendance of the audience in the worldcup |

## 2.Matches dataset

| S.No | Column name | Definition |
|---|---|---|
| 1. | **Year** | The year in which the match was played |
| 2. | **Datetime** | The Date on which the match was played along with a 24 hour format time |
| 3. | **Stage** | The stage at which the match was played |
| 4. | **Stadium** | Stadium name where the match was held |
| 5. | **City** | The city name, where the match was played |
| 6. | **Home Team Name** | Home team country name |
| 7. | **Home Team Goals** | Total goals scored by the home team by the end of the match |
| 8. | **Away Team Goals** | Total goals scored by the away team by the end of the match |
| 9. | **Away Team Name** | Away team country name |

# <u>Scope of The Analysis</u>

To analyze a data set related to FIFA World Cup using a suitable method. In this study we have taken up the data sets of the FIFA World Cup(1930-2014) and analyzed them using Excel.

The analysis focused on:

 a) which team got the titles of winner,runner-up,third for a desired year

b) Which team won the most titles of winner,runner-up and third

c) country hitting most number of goals per country and number of goals per country.

d) Per Year:

- Attendance of audience
- No.of qualified teams
- Goals scored
- Matches played

e) According to attendance:

- Matches with highest of attendance
- Stadium with highest average attendance

f) Matches outcome according to home and away team

.

# Source of The Dataset

- The datasets are taken from the Kaggle with the name 'FIFA World Cup'.

   https://www.kaggle.com/abecklas/fifa-world-cup
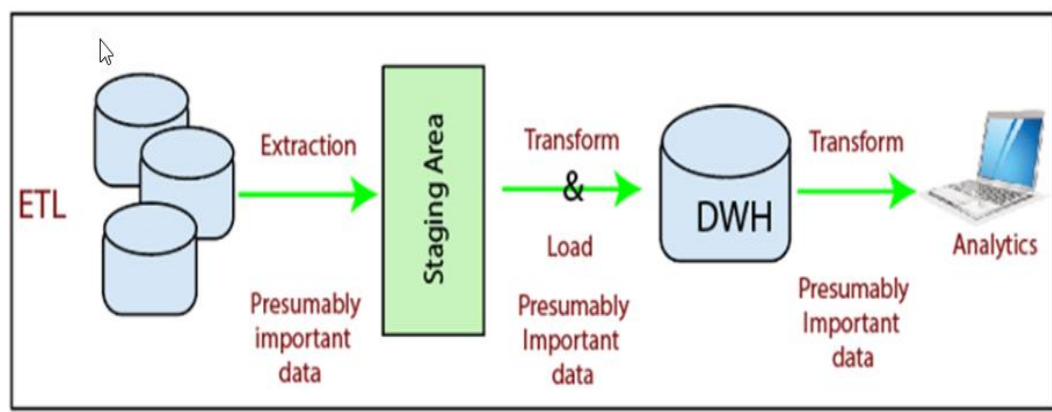
- Author of the Datasets

   Andre Becklas

- Data last updated

   2017

# ETL Process

ETL, which stands for Extraction, Transformation and Loading, is a data integration process that combines data from multiple data sources into a single, consistent data store that is loaded into a data warehouse or other target system. It is the foundation of data analytics and a process for integrating and loading data for computation and analysis, eventually becoming the primary method to process data for data warehousing projects.



- **Extraction**: Extraction is the operation of extracting information from a source system for further use in a data warehouse environment. This is the first stage of the ETL process. It is often one of the most time-consuming tasks in the ETL. The source systems might be complicated and poorly documented, and thus determining which data needs to be extracted can be difficult. The data has to be extracted

several times in a periodic manner to supply all changed data to the warehouse and keep it up to date.

- **Cleansing:** The cleansing stage is crucial in a data warehouse technique because it is supposed to improve data quality. The primary data cleansing features found in ETL tools are rectification and homogenization. They use specific dictionaries to rectify typing mistakes and to recognize synonyms, as well as rule-based cleansing to enforce domain-specific rules and defines appropriate associations between values.
- **Transformation**: Transformation is the core of the reconciliation phase. It converts records from its operational source format into a particular data warehouse format. If we implement a three-layer architecture, this phase outputs our reconciled data layer.
- **Load:** The **Load** is the process of writing the data into the target database. During the load step, it is necessary to ensure that the load is performed correctly and with as little resources as possible. Loading can be carried in two ways:

    o **Refresh:** Data Warehouse data is completely rewritten. This means that older file is replaced. Refresh is usually used in combination with static extraction to populate a data warehouse initially.

    o **Update:** Only those changes applied to source information are added to the Data Warehouse. An update is typically carried out without deleting or modifying pre-existing data. This method is used in combination with incremental extraction to update data warehouses regularly.

The whole ETL process for this project is done in Microsoft Excel.

**Cleaning in this project:**

In this project Cleaning is an important part of the analysis as some of the initial data is been cleaned to get desired data to perform the analysis.

1)Cleaning 1:Some of the Home Team Name in Matches Data set are having "rn">" in their name.So we need to remove this by replacing them.

```
rn">Republic of Ireland
rn">Trinidad and Tobago
rn">Bosnia and Herzegovina
rn">Serbia and Montenegro
rn">United Arab Emirates
```

2)Cleaning 2:The Stadium Maracana was written as "Maracanï¿½ - Estï¿½dio Jornalista Mï¿½rio Filho".So we replaced it with Maracana

| | | | |
|---|---|---|---|
| Victor Boucquey | | | |
| Fort Carree | | | |
| Maracanï¿½ - Estï¿½dio Jornalista Mï¿½rio Filho | | | |
| Durival de Brito | | | |
| Pacaembu | | | |
| Independencia | | | |
| Eucaliptos | | | |

3)Cleaning 3: The Home team Name and Away team name have Germany's other mistaken name as Germany FR.So we should replace Germany FR to Germany.

| | |
|---|---|
| Germany | |
| Spain | |
| Italy | |
| Czechoslovakia | |
| Cuba | |
| England | |
| Germany FR | |

# Analysis of Dataset

## 1.Objective 1:Most Number of Winning Titles(Winner,Runner-up,Third)

**a) Introduction:** The analysis shows the countries having winners, runner-ups, third titles in the FIFA football worldcup.For this analysis we used **Winning Title Datasheet**.

**b) Specific Requirements/Functions and Formulas:**

1) Pivot table of Winning Title Data containing country name and count of winner titles per country

| Row Labels | Count of Winner |
|---|---|
| Argentina | 2 |
| Brazil | 5 |
| England | 1 |
| France | 1 |
| Germany | 4 |
| Italy | 4 |
| Spain | 1 |
| Uruguay | 2 |
| **Grand Total** | **20** |

2) Pivot table of Winning Title Data containing country name and count of Runner-Up titles per country

| Row Labels | Count of Runners-Up |
|---|---|
| Argentina | 3 |
| Brazil | 2 |
| Czechoslovakia | 2 |
| France | 1 |
| Germany | 4 |
| Hungary | 2 |
| Italy | 2 |
| Netherlands | 3 |
| Sweden | 1 |
| **Grand Total** | **20** |

3) Pivot table of Winning Title Data containing country name and count of Runner-Up titles per country

| Row Labels | Count of Third |
|---|---|
| Austria | 1 |
| Brazil | 2 |
| Chile | 1 |
| Croatia | 1 |
| France | 2 |
| Germany | 4 |
| Italy | 1 |
| Netherlands | 1 |
| Poland | 2 |
| Portugal | 1 |
| Sweden | 2 |
| Turkey | 1 |
| USA | 1 |
| **Grand Total** | **20** |

4) COUNTIF Function

| Country | Total winning title | Total Ru |
|---|---|---|
| Uruguay | =COUNTIF(C2:C21,L2) | |
| Italy | COUNTIF(range, criteria) | |
| | 1 | |

5) Addition formula

| Total Top 3 Titles |
|---|
| =M2+N2+O2 |

6) Slicer for year, country column
7) Clustered Column chart for winners, Runner-up and Third
8) Hyperlink

**c) Analysis Results:**

- Brazil has most number of winner title of 5 worldcups.
- Germany has most number of runner-up titles of 4 worldcups
- Germany also have most number of third(second runner-up) titles of 4 worldcups
- Germany and Brazil are most consistent teams with total sum of all top 3 titles of 8 each.
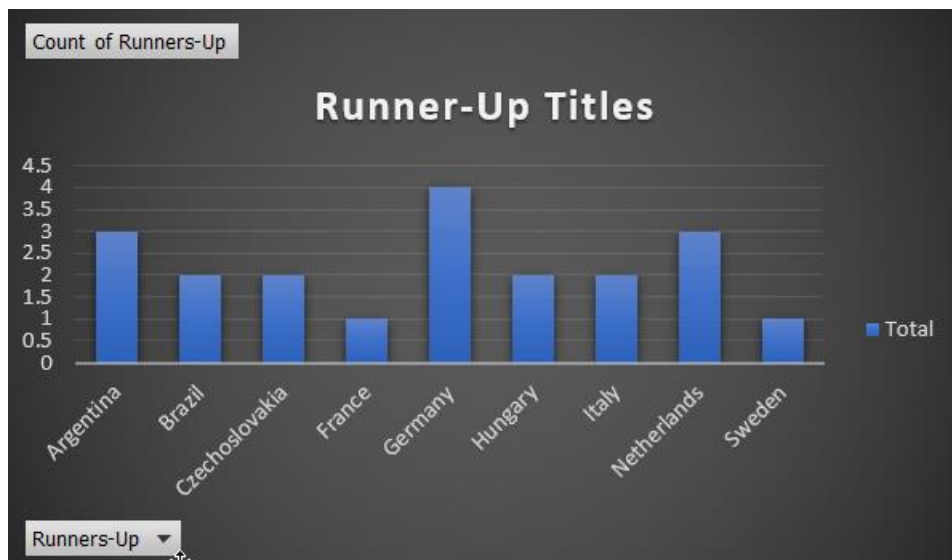
**d)Visualization:**

Runner-Up Titles



Second Runner-Up Titles

**Table for sum of top 3 titles of every country:**

| Country | Total winning title | Total Runner-up title | Total Third Titles | Total Fourth Title | Total Top 3 Titles |
|---|---|---|---|---|---|
| Uruguay | 2 | 0 | 0 | 3 | 2 |
| Italy | 4 | 2 | 1 | 1 | 7 |
| France | 1 | 1 | 2 | 1 | 4 |
| Brazil | 5 | 2 | 1 | 2 | 8 |
| Switzerland | 0 | 0 | 0 | 0 | 0 |
| Sweden | 0 | 1 | 1 | 0 | 2 |
| Chile | 0 | 0 | 1 | 0 | 1 |
| England | 1 | 0 | 0 | 1 | 1 |
| Mexico | 0 | 0 | 0 | 0 | 0 |
| Germany | 3 | 3 | 2 | 0 | 8 |
| Argentina | 2 | 2 | 0 | 0 | 4 |
| Spain | 1 | 0 | 0 | 0 | 1 |
| Mexico | 0 | 0 | 0 | 0 | 0 |
| Italy | 1 | 1 | 1 | 0 | 3 |
| USA | 0 | 0 | 0 | 0 | 0 |
| France | 1 | 1 | 0 | 0 | 2 |
| Korea/Japan | 0 | 0 | 0 | 0 | 0 |
| Germany | 1 | 0 | 2 | 0 | 3 |
| South Africa | 0 | 0 | 0 | 0 | 0 |
| Brazil | 0 | 0 | 0 | 1 | 0 |

**Dashboard of Objective 1:**



# 2.Objective 2: Number of goals per country and country scoring most number of goals

**a) Introduction:** The analysis shows the countries having total goals in the FIFA football worldcup and the country scoring most number of goals.For this analysis we used **Matches Dataset.**

**b) Specific Requirements/Functions and Formulas:**

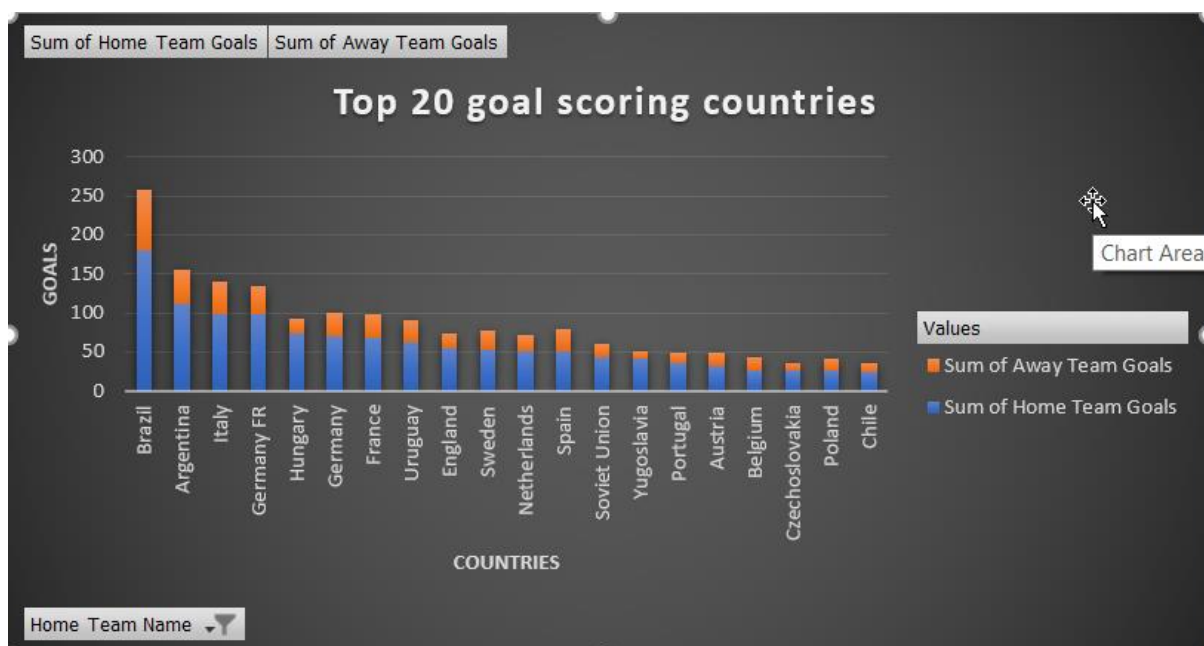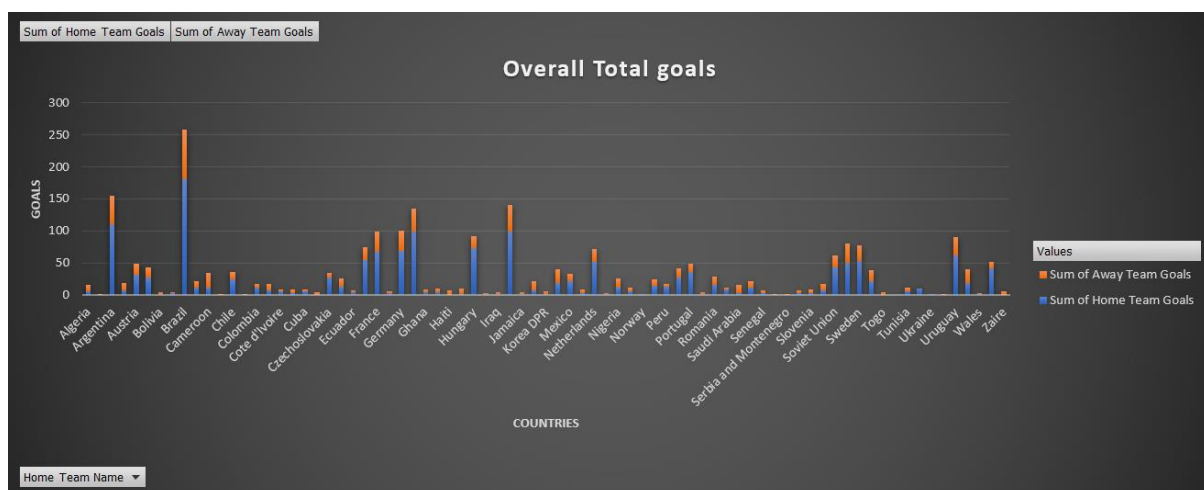1) Pivot table of Matches Dataset containing sum of Home Team Goals and Sum of Away Team Goals

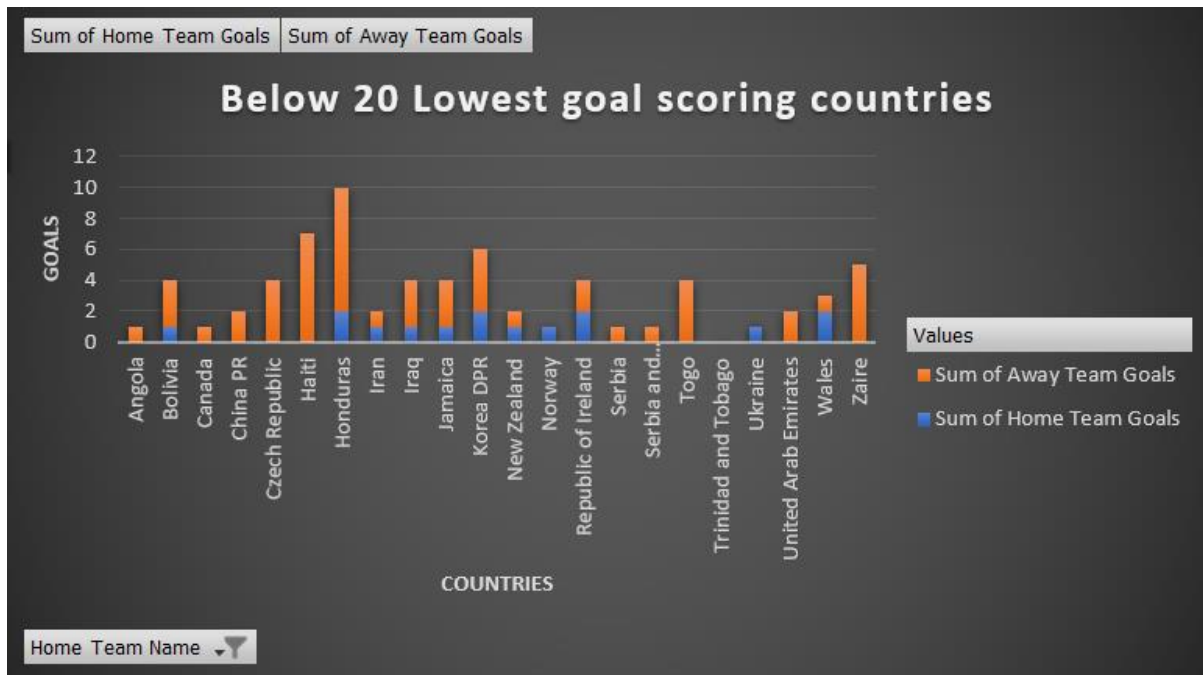| Home Team Name | Sum of Home Team Goals | Sum of Away Team Goals |
|---|---|---|
| Brazil | 180 | 78 |
| Argentina | 111 | 44 |
| Italy | 99 | 41 |
| Germany FR | 99 | 36 |
| Hungary | 73 | 19 |
| Germany | 69 | 32 |
| France | 68 | 31 |
| Uruguay | 62 | 29 |
| England | 54 | 20 |
| Sweden | 53 | 25 |
| Netherlands | 51 | 21 |
| Spain | 50 | 30 |
| Soviet Union | 43 | 18 |
| Yugoslavia | 42 | 9 |
| Portugal | 36 | 13 |
| Austria | 31 | 17 |
| Belgium | 27 | 16 |
| Czechoslovakia | 27 | 8 |
| Poland | 27 | 14 |
| Chile | 25 | 11 |
| **Grand Total** | **1227** | **512** |

2) Applying custom filter in the pivot to get below 20 and top 20 total goals scoring countries
3) Stacked Column Chart for Overall Total goals for every country, below 20 countries with lowest total goals and top 20 countries with highest total goals
4) Hyperlink
5) Slicer of Home and Away team name

c) Analysis Results:

- Brazil has most number of overall goals with 258 goals(180 home goals and 78 away goals)
- Brazil, Argentina and Italy are the top 3 countries with most overall goals 258,155,140 goals respectively.
- Trinidad and Tobago have lowest overall goals with 0 goals

d)Visualization:

**Dashboard of Objective 2:**



## 3.Objective 3: Per Year:i)Attendance of audience ii)No.of Qualified Teams iv)Goals scored v)Matches Played

**a) Introduction:** The analysis shows the Year highest attendance of audience,No.of team qualified teams, total goals scored and matches played per year comparison in the FIFA.For this analysis we used **Winning Title Dataset.**

**b) Specific Requirements/Functions and Formulas:**

1) Pivot table of Winning Title Dataset containing Year of world cup and Sum of Attendance.

| Row Labels | Sum of Attendance |
|---|---|
| 1930 | 590549 |
| 1934 | 363000 |
| 1938 | 375700 |
| 1950 | 1045246 |
| 1954 | 768607 |
| 1958 | 819810 |
| 1962 | 893172 |
| 1966 | |
| 1970 | |
| 1974 | |
| 1978 | 1545791 |
| 1982 | 2109723 |
| 1986 | 2394031 |
| 1990 | 2516215 |
| 1994 | 3587538 |
| 1998 | 2785100 |
| 2002 | 2705197 |
| 2006 | 3359439 |
| 2010 | 3178856 |
| 2014 | 3386810 |
| Grand Total | 37457647 |

Sum of Atten
Value: 893172
Row: 1962

2) Pivot table of Winning Title Dataset containing Year of world cup and Sum of QualifiedTeams

| Row Labels | Sum of QualifiedTeams |
|---|---|
| 1930 | 13 |
| 1934 | 16 |
| 1938 | 15 |
| 1950 | 13 |
| 1954 | 16 |
| 1958 | 16 |
| 1962 | 16 |
| 1966 | 16 |
| 1970 | 16 |
| 1974 | 16 |
| 1978 | 16 |
| 1982 | 24 |
| 1986 | 24 |
| 1990 | 24 |
| 1994 | 24 |
| 1998 | 32 |
| 2002 | 32 |
| 2006 | 32 |
| 2010 | 32 |
| 2014 | 32 |
| Grand Total | 425 |

3) Pivot table of Winning Title Dataset containing Year of world cup and Sum of MatchesPlayed and sum of GoalsScored.

| Row Labels | Sum of MatchesPlayed | Sum of GoalsScored |
|---|---|---|
| 1930 | 18 | 70 |
| 1934 | 17 | 70 |
| 1938 | 18 | 84 |
| 1950 | 22 | 88 |
| 1954 | 26 | 140 |
| 1958 | 35 | 126 |
| 1962 | 32 | 89 |
| 1966 | 32 | 89 |
| 1970 | 32 | 95 |
| 1974 | 38 | 97 |
| 1978 | 38 | 102 |
| 1982 | 52 | 146 |
| 1986 | 52 | 132 |
| 1990 | 52 | 115 |
| 1994 | 52 | 141 |
| 1998 | 64 | 171 |
| 2002 | 64 | 161 |
| 2006 | 64 | 147 |
| 2010 | 64 | 145 |
| 2014 | 64 | 171 |
| **Grand Total** | **836** | **2379** |

4)Slicer for year

5)Line graph with trendline
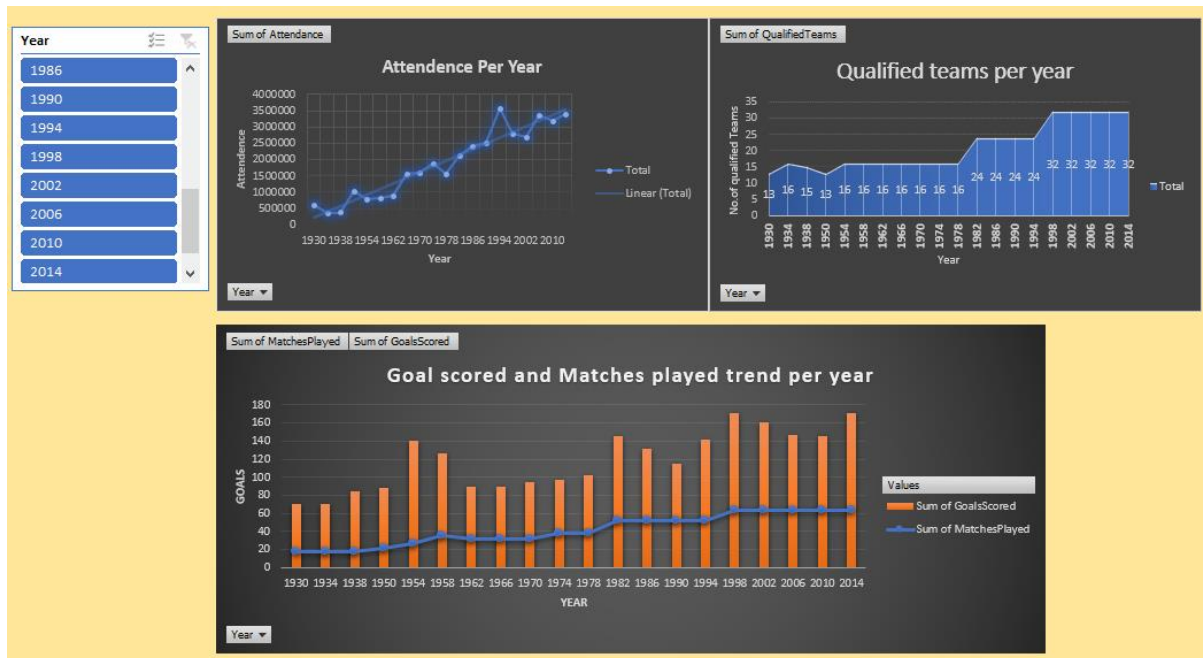
6)Area graph

7)Clustered -Line Combo Chart

**c) Analysis Results:**

- There is a rapid increase of the attendance in year 1994 of 3587538 people.
- There is an increase of popularity trend of football worldcup among the people within the years.
- In 1930 and 1950 lowest number of teams have qualified with 13 teams each.
- In 1998 the trend of highest number of teams qualification increase to 32 teams.
- At an average of every 5 year there is an increase of qualified teams.
- In 1998 and 2014 there were most number of goals scored 171 goals and number of matches played is same with 64 matches which are record breaking years.
- With increase in number of matches there is uneven increase in goals per year.

**d)Visualization:**

**Dashboard for Objective 3:**



# 4.Objective 4:According to attendance: i)Matches with highest number of attendance ii)Stadium with highest average attendence

**a) Introduction:** The analysis shows the match between two different countries having highest number of audience attendance and Stadium in which average attendance of audience in the FIFA football worldcup.For this analysis we used **Matches Datasheet**.

**b) Specific Requirements/Functions and Formulas:**

1)  **CONCAT** Function for joining two teams(Home and away) name to create new column Team1vsTeam2.



2)  Pivot table for Top 10 Match names and average of attendance

| Row Labels | Average of Attendance |
|---|---|
| England Vs Germany FR | 96924 |
| Iraq | 103763 |
| Italy Vs Germany FR | 96222 |
| Mexico Vs Belgium | 108192 |
| Mexico Vs El Salvador | 103058 |
| Mexico Vs Paraguay | 114600 |
| Mexico Vs Soviet Union | 107160 |
| Uruguay Vs Brazil | 173850 |
| USA Vs Colombia | 93869 |
| USA Vs Romania | 93869 |
| **Grand Total** | **107975.3636** |

3) Pivot Table for top 10 Stadium with their average attendance

| Row Labels | Average of Attendance |
|---|---|
| Estadio Azteca | 100924 |
| Estadio do Maracana | 74197 |
| Giants Stadium | 73690 |
| Rose Bowl | 92601 |
| Santiago Bernabeu | 82522 |
| Soccer City Stadium | 83857 |
| Stade de France | 78222 |
| Stanford Stadium | 81737 |
| Wembley Stadium | 86448 |
| Maracana | 101693 |
| **Grand Total** | **87535** |

4) Slicer for Year,Stadium,Home team name and Away team name
5) Clustered Bar Chart
6) Hyperlink
7) Filter

**c) Analysis Results:**

- The match between **Uruguay** and **Brazil** has highest average attendance of audience of **173850** people. It is the most popular match, which means that Uruguay and Brazil are most popular teams in FIFA worldcup.
- **Maracana** and **Estadio Azteca stadium** has the highest number of average attendance of audience of **101696** and **100924**. Which means that they must have high audience capacity in the stadium. I did research and found that **Estadio Azteca stadium** is the **4th** largest stadium in the world with a capacity of **81044**.
- As **Estadio Azteca stadium** is in **Spain** and **Uruguay** is a **Spanish** count ry so that's conclude that the reason behind **Uruguay and Brazil match** to be most popular match.

## d)Visualization:





## Dashboard for Objective 4:

# 5.Objective 5:Matches outcome with Home and Away team

**a) Introduction:** The analysis shows the matches outcome of the matches according to home and away team and can be used to analyse about individual team winning percentage playing as Home team and Playing as Away team in the FIFA football worldcup.For this analysis we used **Matches Datasheet**.

**b) Specific Requirements/Functions and Formulas:**

1) IFS Function to find whether Home teams wins or Away team wins.

| Match Outcome(home/away) | | | |
|---|---|---|---|
| =IFS(G2>H2,"Home team Wins",H2>G2,"Away team Wins",G2=H2,"Draw") | | | |
| IFS(**logical_test1**, value_if_true1, [logical_test2, value_if_true2], [logical_test3, value_if_t |

2) Pivot table of Match outcome and Count of Match Outcome(Home/Away)

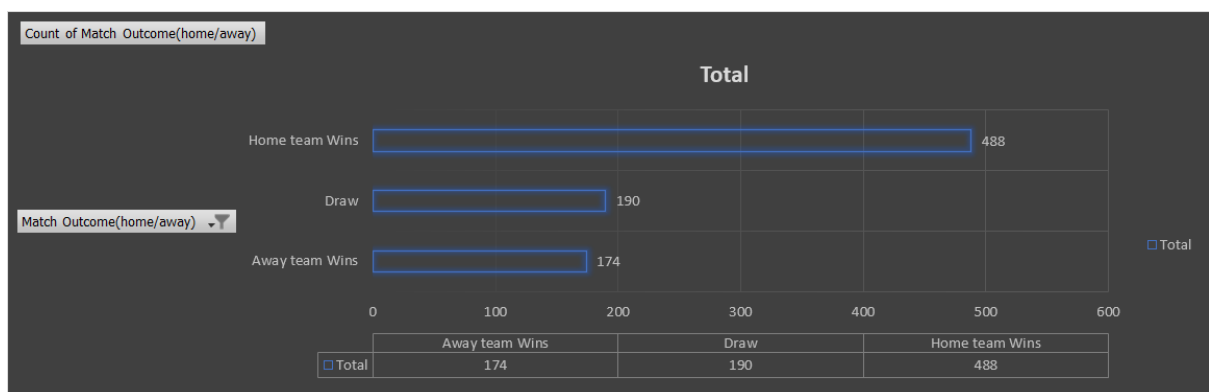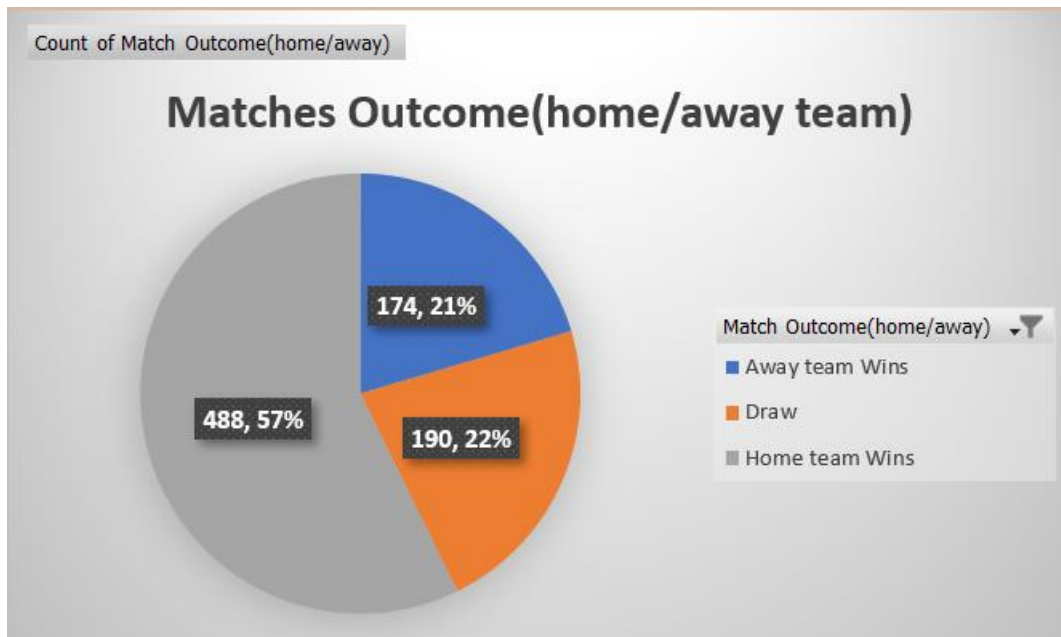| Row Labels | Count of Match Outcome(home/away) |
|---|---|
| Away team Wins | 174 |
| Draw | 190 |
| Home team Wins | 488 |
| **Grand Total** | **852** |

3) Pie chart
4) Clustered Bar Chart
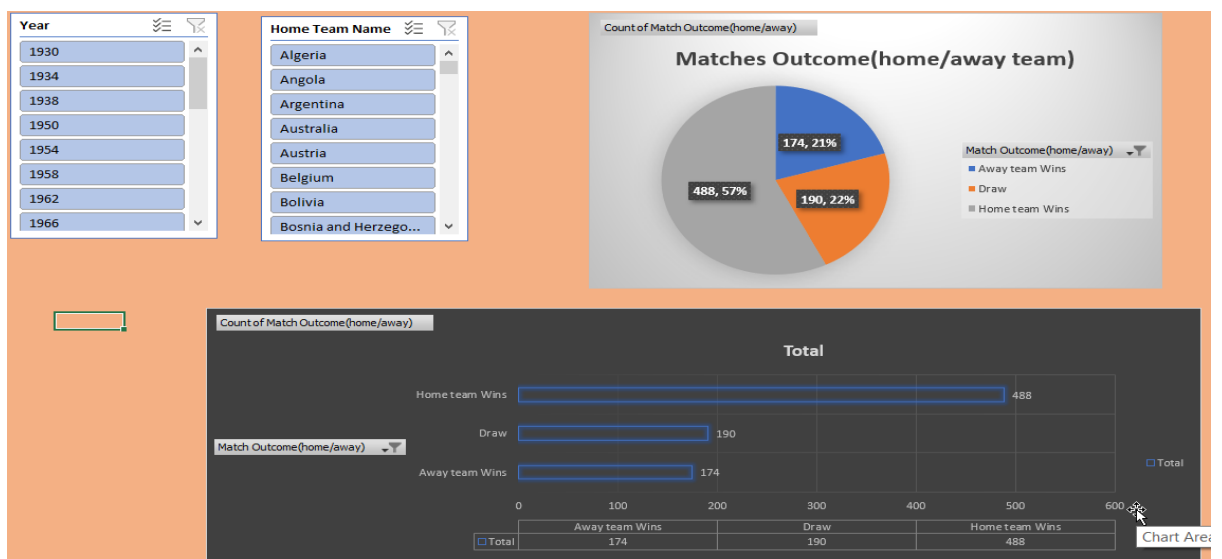5) Slicer of Year and Home Team Name

**c) Analysis Results:**

- Teams playing at their home country have the highest chances to win the worldcup as the Home Team has highest percentage of winning percentage of 57% and away team winning percentage is 21%.
- The draw percentage which is 22% is greater than away team wins percentage which is 21%,which means that teams playing in away country have very less chances to win.

**d)Visualization:**





## Dashboard for Objective 5:

# 5.Objective 5:Top 5 referees with most number of matches

**a) Introduction:** The analysis shows Top 5 referees with most number of matches

in the FIFA football worldcup. For this analysis we used **Matches Datasheet**.

**b) Specific Requirements/Functions and Formulas:**

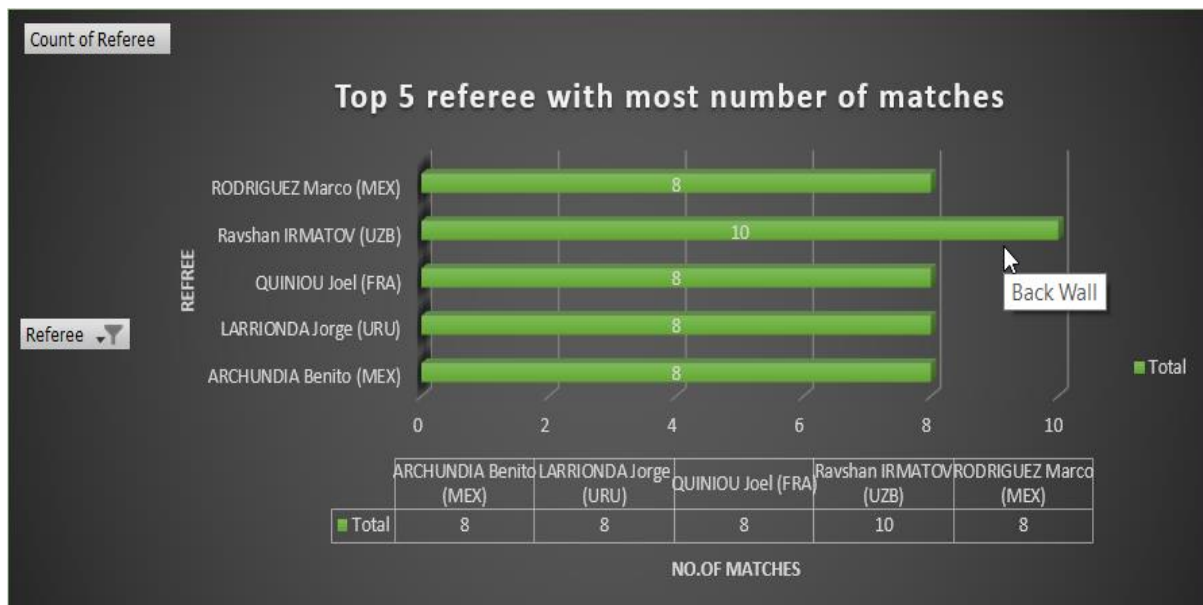1) pivot of Referee Name and count of Referee Name

| Row Labels | Count of Referee |
|---|---|
| ARCHUNDIA Benito (MEX) | 8 |
| LARRIONDA Jorge (URU) | 8 |
| QUINIOU Jo Chart Area | 8 |
| Ravshan IRMATOV (UZB) | 10 |
| RODRIGUEZ Marco (MEX) | 8 |
| **Grand Total** | **42** |

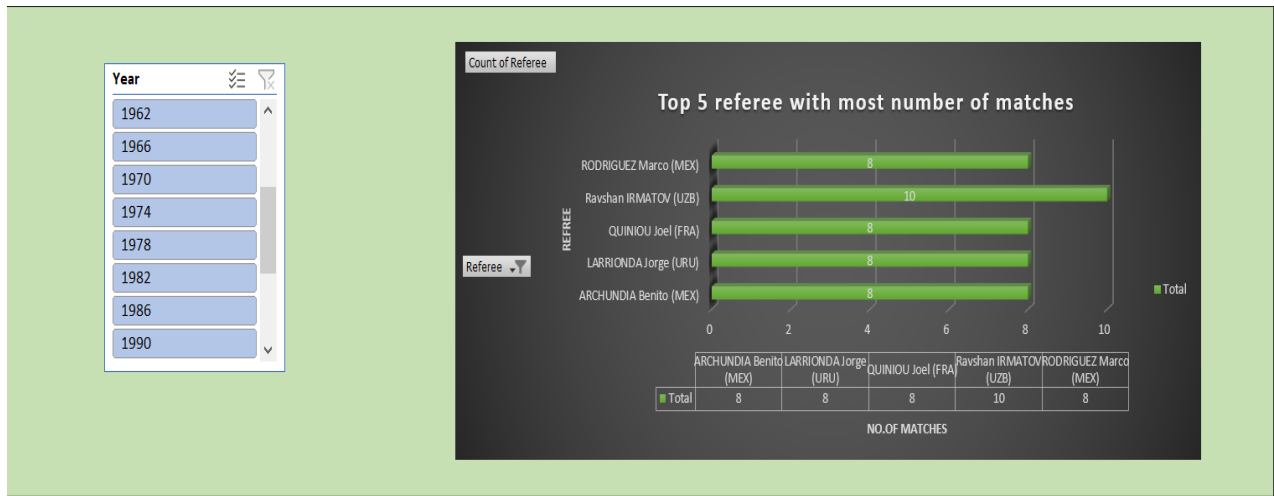2) Filter for top 5
3) Clustered Bar Chart

**c) Analysis Results:**

- **Ravshan Irmatov** was the referee of most number of matches of 10

**d)Visualization:**

**Dashboard for Objective 6:**



# In this Project

- Formulas used:7
- Hyperlink used:32
- Pivot used:14
- Slicers used:12
- Graphs:14
- Icon:15
- Shapes:32
- Excel sheets:13
- Data Sets:2
- Dataset 1(winning title):10 Columns,21 Rows
- Dataset 2(Matches):22 Columns,853 Rows

# List of Analysis with Results

- Brazil has most number of overall goals with 258 goals(180 home goals and 78 away goals)
- Brazil, Argentina and Italy are the top 3 countries with most overall goals 258,155,140 goals respectively.
- Trinidad and Tobago have lowest overall goals with 0 goals
- Brazil has most number of overall goals with 258 goals(180 home goals and 78 away goals)
- Brazil, Argentina and Italy are the top 3 countries with most overall goals 258,155,140 goals respectively.
- Trinidad and Tobago have lowest overall goals with 0 goals
- There is a rapid increase of the attendance in year 1994 of 3587538 people.
- There is an increase of popularity trend of football worldcup among the people within the years.
- In 1930 and 1950 lowest number of teams have qualified with 13 teams each.
- In 1998 the trend of highest number of teams qualification increase to 32 teams.
- At an average of every 5 year there is an increase of qualified teams.
- In 1998 and 2014 there were most number of goals scored 171 goals and number of matches played is same with 64 matches which are record breaking years.
- With increase in number of matches there is uneven increase in goals per year.
- The match between **Uruguay** and **Brazil** has highest average attendance of audience of **173850** people. It is the most popular match, which means that Uruguay and Brazil are most popular teams in FIFA worldcup.
- **Maracana** and **Estadio Azteca stadium** has the highest number of average attendance of audience of **101696** and **100924**. Which means that they must have high audience capacity in the stadium. I did research and found that **Estadio Azteca stadium** is the **4th** largest stadium in the world with a capacity of **81044**.
- As **Estadio Azteca stadium** is in **Spain** and **Uruguay** is a **Spanish** country so that's conclude that the reason behind **Uruguay and Brazil match** to be most popular match.
- Teams playing at their home country have the highest chances to win the worldcup as the Home Team has highest percentage of winning percentage of 57% and away team winning percentage is 21%.
- The draw percentage which is 22% is greater than away team wins percentage which is 21%,which means that teams playing in away country have very less chances to win.

# References

- https://www.kaggle.com/abecklas/fifa-world-cup

- http://www.youtube.com/

- http://www.google.com/
- http://www.stackoverflow.com/
- http://www.github.com/

# Bibliography

1. Microsoft Excel 2016 Bible: The Comprehensive Tutorial Resource by John Walkenbach, Wiley

2. Fundamentals of Business Analytics by R.N. Prasad, Seema Acharya, Wiley