Out[3]: Click here to display/hide the code.

Introduction:

In anticipation of the 2022 World Cup in Qatar, this is a summary analysis of the 2018 World Cup. This analysis will be used to predict performance in this upcoming tournament, and so readers will have a better understanding of this years tournament.

Selection of Data:

This summary will use the raw data sources provided by FBREF, and these data is available at: https://fbref.com/en/comps/1/stats/FIFA-World-Cup-Stats#stats_standard and https://github.com/fghaddar/data-science-final-project/blob/main/scores_fixtures.txt. This project consists of 4 sources of data, all of which add up to 668 rows and 45 features. The 3 source files are called player.csv, team_for, and team_against.

The majority of the data is clean in the sense that most 3 of the data sources do not have null values in any of the rows or columns. However, one data source did have null values in it's rows, and so these rows have been dropped. Additionally, randomly generated letters have been appended to each name of every player in the player.csv file. These random letters are added after a backslash, and so I plan to extract the name before the backslash and used that extracted name to update that row. In addition, the data has some

irrelevant features that are not needed for this discussion, and so they have been dropped. They are found within the player.csv file and these features are: Rk, Pos, Born, G-PK, PK, PKatt, and CrdY.

Within the team_against and team_for files, all of the features are numeric except for the Squad name, which is a string. Within the player.csv file all features are numeric except for player name.

The features used as indicators for performance will be several fold. These features and indicators include Goals Scored, Goals Conceded, Expected Goals generated by each team, Expected Assists generated by each team, Expected Goals faced by each team, Expected Assists faced by each team, the average age of the squad at the time, and how many notable performers a subset of teams have. Exoected Goals (A.K.A. xG) is the likelihood of a given shot becoming a goal whilst Expected Assists (A.K.A. xA) is the likelihood of a given pass becoming an assist. Both xG and xA are measured out of a scale 1, where 1 means that the likelihood is a certainity and 0 means that the likelihood is never happening.

Methods:

Tools:

Pandas & Matplotlib for data analysis and visualization

Scikit-learn for inference

Github for version control

VS Code as IDE

Juypter Notebook within Anaconda

```
In [313... import pandas as pd import matplotlib.pyplot as plt
```

Data Prep/Cleaning

```
# Read the data, check for NA values, drop irrelevant features

stats_against = pd.read_csv("https://raw.githubusercontent.com/fghaddar/data-science-final-project/main/stats_against.t

stats_for = pd.read_csv("https://raw.githubusercontent.com/fghaddar/data-science-final-project/main/stats_for.txt")

player_stats = pd.read_csv("https://raw.githubusercontent.com/fghaddar/data-science-final-project/main/player.txt")

scores = pd.read_csv("https://raw.githubusercontent.com/fghaddar/data-science-final-project/main/scores_fixtures.txt")
```

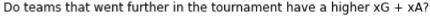
```
# Some non-null values but all headers describe the columns well
# stats against.info()
# stats for.info()
# player stats.info()
# scores.info
# Get rid of non-null values and irrelevant features
stats against = stats against.drop(columns=['playing time 90s', 'performance G-PK', 'performance PK', 'per 90 G-PK'])
stats for = stats for drop(columns=['Playing Time 90s', 'Performance G-PK', 'Performance PK', 'Per 90 Minutes G-PK'])
player stats = player stats.drop(columns=['Performance CrdY', 'Performance CrdR'])
scores = scores.drop(columns=['Wk', 'Attendance', 'Venue', 'Match Report', 'Notes', 'Referee'])
# Update scores table so that the Home and Away team names are by ID and not country names
for i in scores.index:
    scores.at[i, "Home"] = scores.at[i, "Home"][-2:]
   scores.at[i, "Away"] = scores.at[i, "Away"][:2]
# Update stats against and stats for so that the Home and Away team names are by ID and not country names
for i in stats_against.index:
    stats against.at[i, "Squad"] = stats against.at[i, "Squad"][:2]
    stats_for.at[i, "Squad"] = stats_for.at[i, "Squad"][:2]
```

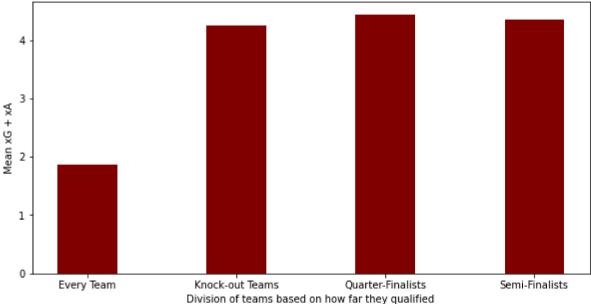
lets start by examining the performance of teams that went further into the tournament. Do these teams have higher xG's and xA's than average?

```
# Get the semi-finalists, quarter-finalists, and teams that qualified from the group stage (r16)
In [317...
         semi finalists = []
         qrtr_finalists = []
         r 16 = []
         for i in range(4):
             index = -(i + 1)
             semi finalists += [scores.iloc[index]['Home'], scores.iloc[index]['Away']]
         for i in range(8):
             index = -(i + 1)
             qrtr finalists += [scores.iloc[index]['Home'], scores.iloc[index]['Away']]
         for i in range(16):
             index = -(i + 1)
             r 16 += [scores.iloc[index]['Home'], scores.iloc[index]['Away']]
         # Get the xG+xA performance of every team
         avg xG xA performance = stats for['Per 90 xG+xA'].mean()
         # Get the xG+xA performance of teams that qualified from their groups, quarter finalists, and semi finalists
         sf mean = 0
```

```
qf mean = 0
r16 mean = 0
for semi finalist in semi finalists:
    index = stats for.index[stats for['Squad'] == semi finalist]
    sf mean += stats for.iloc[index[0]]['Per 90 xG+xA']
for qrtr finalist in qrtr finalists:
    index = stats for.index[stats for['Squad'] == qrtr finalist]
    qf mean += stats for.iloc[index[0]]['Per 90 xG+xA']
for r in r 16:
   index = stats for.index[stats for['Squad'] == r]
    r16 mean += stats for.iloc[index[0]]['Per 90 xG+xA']
sf mean = sf mean / 4
qf mean = qf mean / 8
r16 mean = r16 mean / 16
data = {'Every Team': avg xG xA performance, 'Knock-out Teams': r16 mean, 'Quarter-Finalists': qf mean, 'Semi-Finalists'
fig = plt.figure(figsize = (10, 5))
plt.bar(list(data.keys()), list(data.values()), color ='maroon', width = 0.4)
plt.ylabel('Mean xG + xA')
plt.xlabel('Division of teams based on how far they qualified')
plt.title('Do teams that went further in the tournament have a higher xG + xA?')
```

Out[317]: Text(0.5, 1.0, 'Do teams that went further in the tournament have a higher xG + xA?')





In []: