# Using Statistics for Data Analysis

David Carbajal

# The Data

Dataset:

"big_5_players_stats_2023_2024
.csv" – Kaggle

- Player Statistics from the Top 5 Leagues of Soccer in Europe from the 2023-2024 league season.

- Statistics include player name, league, position, team, goals, assists, G+A, minutes played, expected goals/assist, age, nationality, progression, etc.
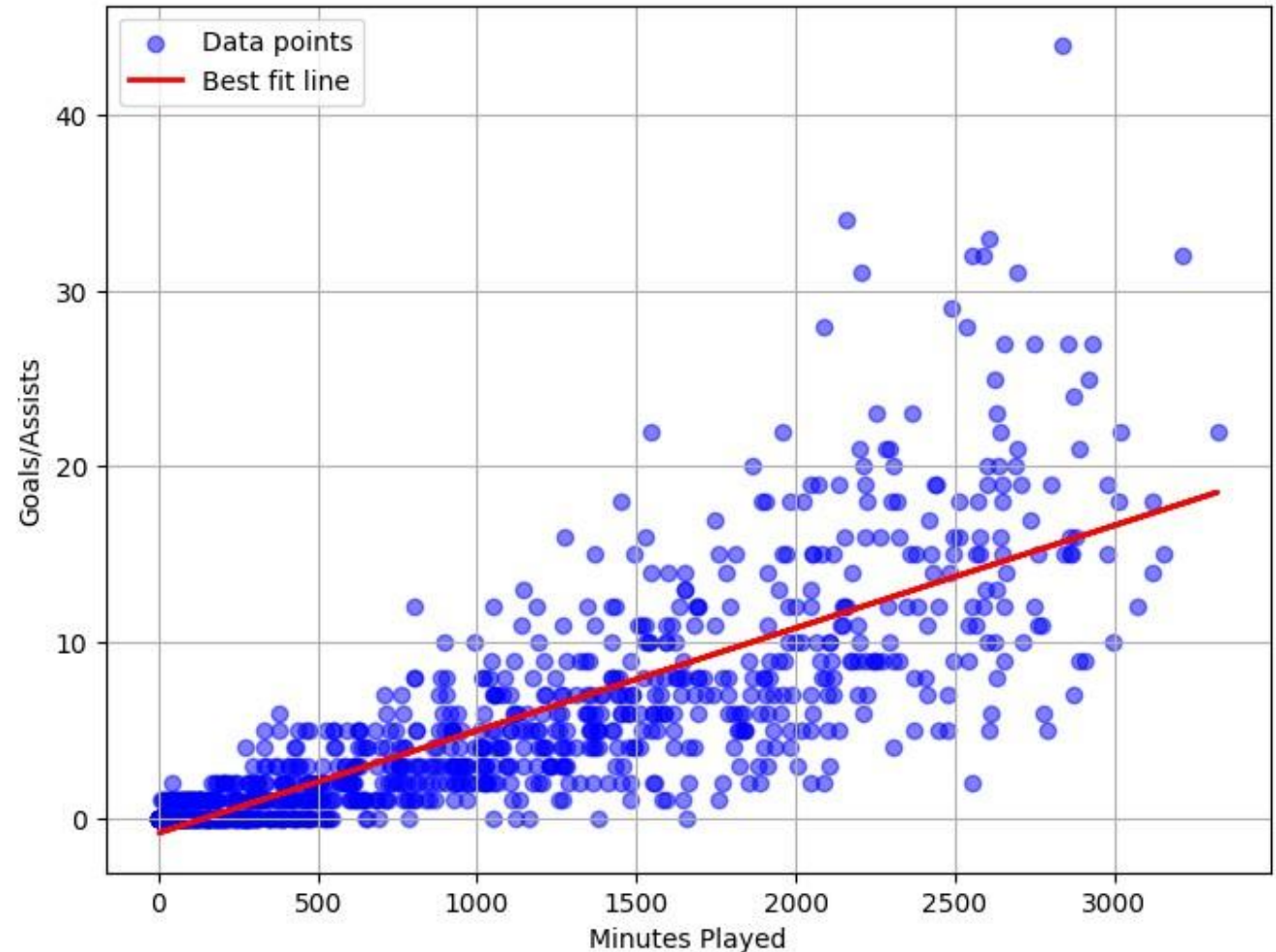
- Across 2,959 players

# Case #1:

- Is there a correlation between total amount of minutes played and goals + assists (Strictly considering players playing in the forward positions).

- Fields Relevant: Player Name, Position, Minutes Played, Goals, and Assists

# Model

- Calculated the correlation coefficient between the minutes players and goals + assists by forwards.

- Included a scatter Plot with a graphed best fit line.

- Demonstrates the relationship between minutes played and goals/assists for forwards, where the regression line indicates the general trend.



Correlation coefficient: 0.8101015524402119
Strong correlation

Total Minutes Played vs. Goals + Assists for Forwards
2023-2024 Season

# Implementation

- Python tools and libraries such as numpy, matplotlib, and csv were utilized.

- Read and uploaded data from the kaggle csv file.

- In order to compute the correlation calculation, np.corrcoef() was used.

- Visualization created through the use of matplotlib.

```python
countGA = np.array(countGA)
mintuesPlayed = np.array(minutesPlayed)

#using numpy to calculate the correlation coefficient
correlationMatrix = np.corrcoef(countGA, mintuesPlayed)
correlation = correlationMatrix[0, 1]

print("Correlation coefficient:", correlation)
if(correlation > 0.75):
    print("Strong correlation")
elif(correlation > 0.50):
    print("Moderate correlation")
else:
    print("Weak correlation")
```

```python
# Calculate the best fit line using polyfit (degree 1 for a linear line)
slope, intercept = np.polyfit(minutesPlayed, countGA, 1)

# Create a polynomial function for the best fit line
best_fit_line = np.poly1d([slope, intercept])

# Plotting the data
plt.figure(figsize=(8,6))
plt.scatter(minutesPlayed, countGA, color='blue', alpha=0.5, label="Data points")
plt.plot(minutesPlayed, best_fit_line(minutesPlayed), color='red', label="Best fit line", linewidth=2)
plt.title("Total Minutes Played vs. Goals + Assists for Forwards\n2023-2024 Season")
plt.xlabel("Minutes Played")
plt.ylabel("Goals/Assists")
plt.legend()
plt.grid(True)
plt.show()
```

# Case 2:

- Finding the average goals + assists and age of the players by position. (Forward, Midfielder, and Defense)


- Finding whether or not there is a correlation between players ages, and their goals and assists. This was done by positions (Forward, Midfielder, and Defense)

# Implementation of mean

```python
import csv
import numpy as np
import matplotlib.pyplot as plt

playerNames = [];
countGAForward = [];
ageForward = [];
countGAMidfielder = [];
ageMidfielder = [];
countGADefender = [];
ageDefender = [];

filePath = "/content/big_5_players_stats_2023_2024.csv"

#opening the file and scanning it
with open(filePath, "r") as file:
  scanner = csv.reader(file)
  next(scanner)  # Skip the header line

  for values in scanner:
    if len(values) > 1:
      if "FW" in values[3].split(","): #if the player is a forward / plays forward
        goals_assists = int(values[14])  # Convert goals/assists to integer
        age = int(values[6])  # Convert age played to integer
        playerNames.append(values[1])  # Append player name for reference
        countGAForward.append(goals_assists)
        ageForward.append(age)

      elif "MF" in values[3].split(","): #if the player is a midfield / plays midfield
        goals_assists = int(values[14])  # Convert goals/assists to integer
        age = int(values[6])  # Convert age played to integer
        playerNames.append(values[1])  # Append player name for reference
        countGAMidfielder.append(goals_assists)
        ageMidfielder.append(age)

    elif "DF" in values[3].split(","): #if the player is a defender / plays defense
      goals_assists = int(values[14])  # Convert goals/assists to integer
      age = int(values[6])  # Convert age played to integer
      playerNames.append(values[1])  # Append player name for reference
      countGADefender.append(goals_assists)
      ageDefender.append(age)


FwGAMean = np.mean(countGAForward)
MfGAMean = np.mean(countGAMidfielder)
DfGAMean = np.mean(countGADefender)
FwAgeMean = np.mean(ageForward)
MfAgeMean = np.mean(ageMidfielder)
DfAgeMean = np.mean(ageDefender)
```

- Imported python tools and libraries such as numpy, matplotlib, and csv.

- Read data and created lists using data from the csv file.

- Calculated mean using np.mean(list).

# Implementation cont.

Using python libraries and previously created lists, I found the correlation between ages and goals + assists using np.corrcoef

Depending on the result, I stated whether correlation was visible or not.

If the correlation was near –1, this meant that younger players proved to score + assist more

If the correlation was near 1, this proved older players proved to score + assist more

```python
print("The following data is from a data set that provides information on every player playing in one of the Big 5 european soccer leagues")
print()


print("Forward Mean goals + assists: ", FwGAMean)
print("Midfielder Mean goals + assists: ", MfGAMean)
print("Defender Mean goals + assists: ", DfGAMean)
print("Forward Mean age: ", FwAgeMean)
print("Midfielder Mean age: ", MfAgeMean)
print("Defender Mean age: ", DfAgeMean)
print()

#using numpy to calculate the correlation coefficient
FWcorrelationMatrix = np.corrcoef(ageForward, countGAForward)
FWcorrelation = FWcorrelationMatrix[0, 1]
MFcorrelationMatrix = np.corrcoef(ageMidfielder, countGAMidfielder)
MFcorrelation = MFcorrelationMatrix[0, 1]
DFcorrelationMatrix = np.corrcoef(ageDefender, countGADefender)
DFcorrelation = DFcorrelationMatrix[0, 1]


print("Correlation between age and goals + assists based on position:")
print("A posoitve correlation which is close to 1, will suggest older players have higher goals + assists")
print("A weak correlation which is close to -1, will suggest younger players have higher goals + assists")
print("And a correlation near 0 suggests there is no relationship between age and goals + assists")
print()

print("FW Correlation coefficient:", FWcorrelation)
if(FWcorrelation > 0.75):
    print("Older players playing forward tend to score and assist more")
elif(FWcorrelation > -.75 and FWcorrelation < .25):
    print("There is no linear relationship between older and younger players playing forward and their goals + assists")
elif(FWcorrelation < -0.75):
    print("Younger players playing forward tend to score and assist more")
else:
    print("Modertae correlation for players playing forward")
```
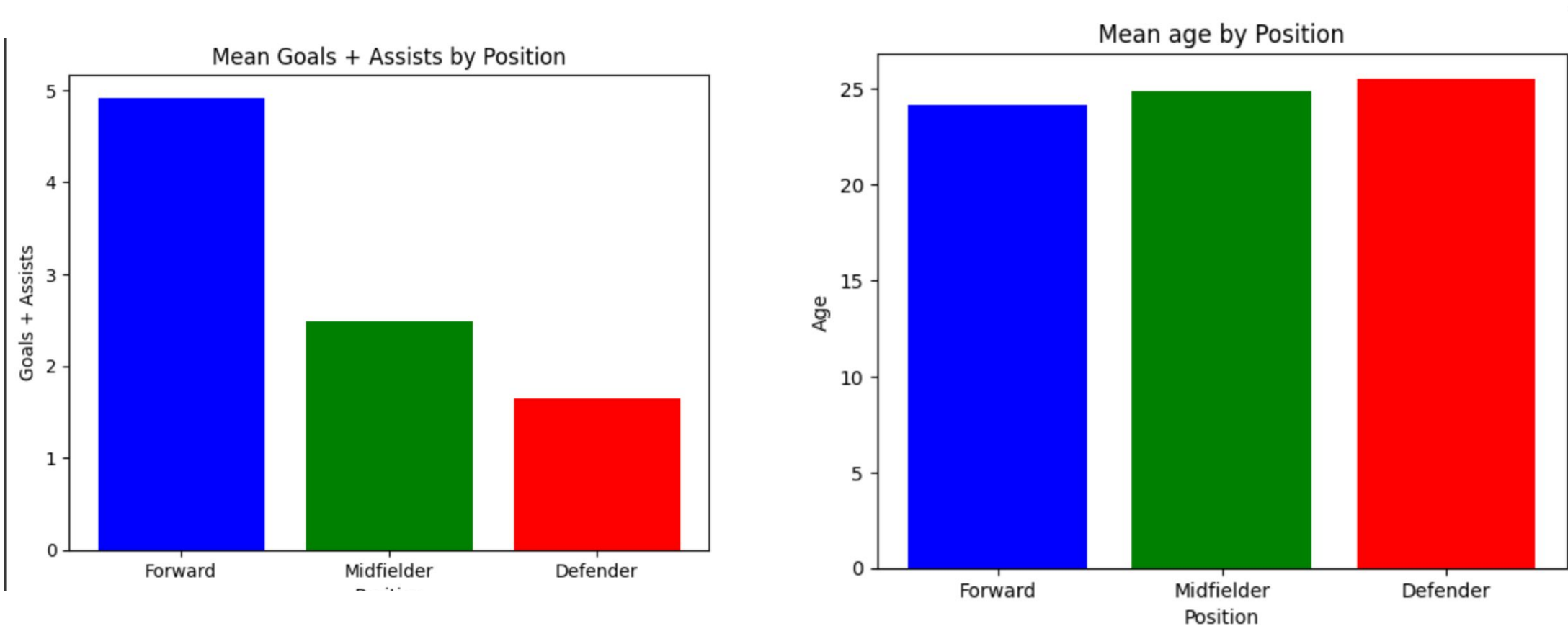
# Results:

- Forward Mean goals + assists: 4.917307692307692
- Midfielder Mean goals + assists: 2.4788359788359786
- Defender Mean goals + assists: 1.6477541371158393
- Forward Mean age: 24.139423076923077
- Midfielder Mean age: 24.8505291005291
- Defender Mean age: 25.541371158392433

# Models

- Mean goals + assist and age models by position

# Results Cont.

- No linear relationship was found between players age and their goals + Assists.

- A correlation near 0 suggests there is no relationship between age and goals + assists

- Forward Correlation coefficient: 0.22912766798821843
- There is no linear relationship between older and younger players playing forward and their goals + assists.

- Midfielder Correlation coefficient: 0.1756854272595472
- There is no linear relationship between older and younger players playing midfield and their goals + assists.

- Defender Correlation coefficient: 0.07301425197406965
- There is no linear relationship between older and younger players playing defense and their goals + assists.

# Model

Demonstrates no correlation between ages and goals + assists



Correlation of Age and Goals + Assists by Position