

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings("ignore")
```

```
In [2]: football_data = pd.read_csv("C:/Users/mtis/Desktop/EPLDATA/FootballDataset.csv", encoding="unicode_escape")
```

```
Out[3]: football_data.head()
```

	Name	Club	Nationality	Position	Age	Matches	Starts	Mins	Goals	Assists	Passes_Attempted	Perc_Passes_Completed	Penalty_Goals	Penalty_Attempted	xG	xA	Yellow_Cards	Red_Cards
0	Mason Mount	Chelsea	ENG	MF,FW	21	36	32	2890	6	5	1881	82.3	1	1	0.21	0.24	2	0
1	Edouard Mendy	Chelsea	SEN	GK	28	31	31	2745	0	0	1007	84.6	0	0	0.00	0.00	2	0
2	Timo Werner	Chelsea	GER	FW	24	35	29	2002	6	8	826	77.2	0	0	0.41	0.21	2	0
3	Ben Chilwell	Chelsea	ENG	DF	23	27	27	2286	3	5	1806	78.6	0	0	0.10	0.11	3	0
4	Reece James	Chelsea	ENG	DF	20	32	25	2273	1	2	1807	85.0	0	0	0.06	0.12	3	0

```
In [4]: football_data.describe
```

```
Out[5]:
```

```
<bound method NDFrame.describe of
0      Mason Mount      Chelsea      ENG      MF,FW      21      36
1      Edouard Mendy      Chelsea      SEN      GK      28      31
2      Timo Werner      Chelsea      GER      FW      24      35
3      Ben Chilwell      Chelsea      ENG      DF      23      27
4      Reece James      Chelsea      ENG      DF      20      32
...
527      Lys Mousset      Sheffield United      FRA      MF,MF      24      11
528      Jack O'Connell      Sheffield United      ENG      DF      26      2
529      Ainslie Hockford      Sheffield United      ENG      GF,FW      16      1
531      Ferni Seriki      Sheffield United      ENG      DF      17      1
...
8      Start's      Mins      Goals      Assists      Passes_Attempted      Perc_Passes_Completed
1      0      0      0      0      0      0
2      19      0      0      0      0      0
3      31      2745      0      6      1887      84.6
4      25      2373      0      5      1807      77.2
1      27      2286      3      5      1806      78.6
4      25      2373      1      2      1807      85.0
...
527      2      296      0      8      59      89.8
528      2      180      0      8      77      77.8
529      8      12      0      8      3      109.8
530      8      11      0      8      1      180.8
531      0      1      0      0      0      -1.0
...
Penalty_Goals      Penalty_Attempted      xG      xA      Yellow_Cards      Red_Cards
0      0      0      0.21      0.24      2      0
1      0      0      0.88      0.09      2      0
2      0      0      0.42      0.21      2      0
3      0      0      0.18      0.11      3      0
4      0      0      0.06      0.12      3      0
...
527      0      0      0.22      0.16      8      0
528      0      0      0.88      0.00      8      0
529      0      0      0.88      0.00      8      0
530      0      0      1.16      0.00      8      0
531      0      0      0.88      0.00      8      0
[532 rows x 18 columns]>
```

Cleaning the Dataset

```
In [5]: football_data.isnull()
```

```
Out[5]:
```

	Name	Club	Nationality	Position	Age	Matches	Starts	Mins	Goals	Assists	Passes_Attempted	Perc_Passes_Completed	Penalty_Goals	Penalty_Attempted	xG	xA	Yellow_Cards	Red_Cards
0	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse
1	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse
2	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse
3	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse
4	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse
...
527	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse
528	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse
529	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse
530	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse
531	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse	Fabse
532 rows x 18 columns																		

```
In [6]: football_data.isna().sum()
```

```
Out[6]:
```

```
Name      0
Club      0
Nationality      0
Position      0
Age      0
Matches      0
Starts      0
Mins      0
Goals      0
Assists      0
Passes_Attempted      0
Perc_Passes_Completed      0
Penalty_Goals      0
Penalty_Attempted      0
xG      0
xA      0
Yellow_Cards      0
Red_Cards      0
dtypes: int64
```

Data Exploration

```
In [7]: #adding 3 new columns to the dataset to further explore
football_data['MinsPerMatch'] = (football_data['Mins']/ football_data['Matches']).astype(int)
football_data['GoalsPerMatch'] = (football_data['Goals']/ football_data['Matches']).astype(float)
football_data.head()
```

```
Out[7]:
```

	Name	Club	Nationality	Position	Age	Matches	Starts	Mins	Goals	Assists	Passes_Attempted	Perc_Passes_Completed	Penalty_Goals	Penalty_Attempted	xG	xA	Yellow_Cards	Red_Cards	MinsPerMatch	GoalsPerMatch
0	Mason Mount	Chelsea	ENG	MF,FW	21	36	32	2890	6	5	1881	82.3	1	1	0.21	0.24	2	0	80	0.166667
1	Edouard Mendy	Chelsea	SEN	GK	28	31	31	2745	0	0	1007	84.6	0	0	0.00	0.00	2	0	88	0.000000
2	Timo Werner	Chelsea	GER	FW	24	35	29	2002	6	8	826	77.2	0	0	0.41	0.21	2	0	74	0.171429
3	Ben Chilwell	Chelsea	ENG	DF	23	27	27	2286	3	5	1806	78.6	0	0	0.10	0.11	3	0	84	0.111111
4	Reece James	Chelsea	ENG	DF	20	32	25	2273	1	2	1807	85.0	0	0	0.06	0.12	3	0	74	0.032250

```
In [8]: Total_Goals = football_data['Goals'].sum()
print(Total_Goals)
```

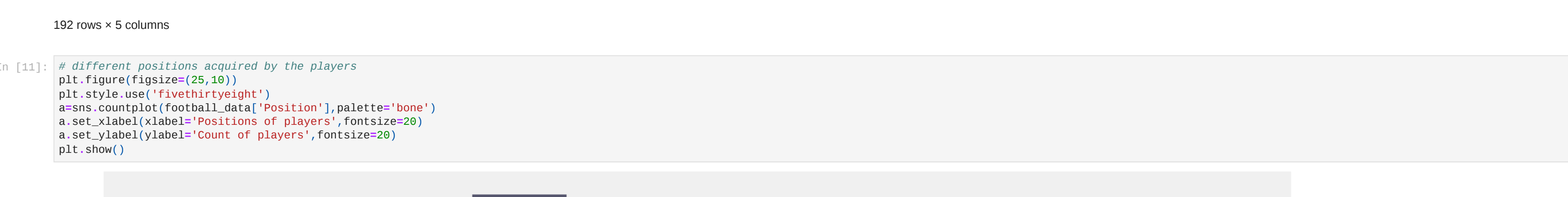
```
In [9]: def country(x):
    return football_data[football_data['Nationality'] == x][['Name','Position','Goals','MinsPerMatch','GoalsPerMatch']]
```

```
In [10]: country('ENG')
```

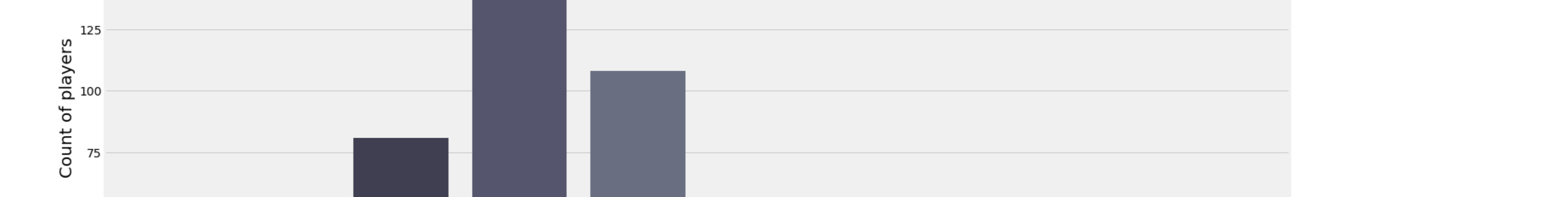
```
Out[10]:
```

	Name	Position	Goals	MinsPerMatch	GoalsPerMatch
0	Mason Mount	MF,FW	6	80	0.166667
3	Ben Chilwell	DF	3	84	0.111111
4	Reece James	DF	1	74	0.031250
16	Tammy Abraham	FW	6	47	0.272727
18	Callum Hudson-Odoi	FW,DF	2	46	0.086957
...
525	Phil Jagieka	DF	0	52	0.000000
528	Daniel Johnson	FW	1	71	0.030000
528	Jack O'Connell	DF	0	90	0.000000
530	Ainslie Hockford	DF,FW	0	11	0.000000
531	Ferni Seriki	DF	0	1	0.000000
152 rows x 5 columns					

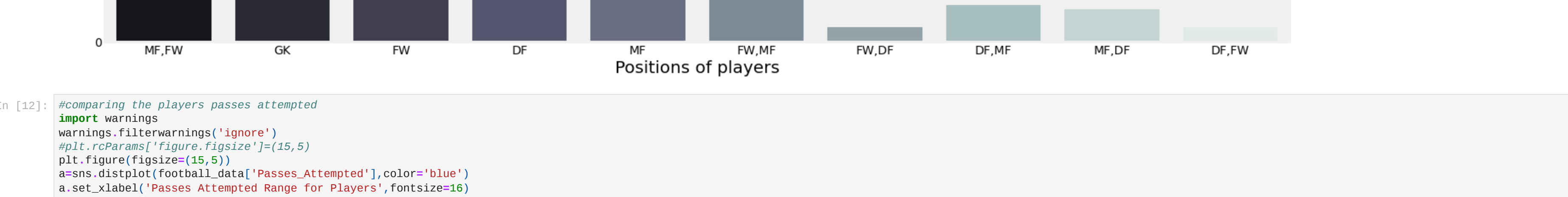
```
In [11]: # different positions acquired by the players
plt.figure(figsize=(25,8))
plt.style.use('fivethirtyeight')
sns.countplot(football_data['Position'],palette='bone')
a.set_xlabel('Positions of players',fontsize=28)
a.set_ylabel('Count of players',fontsize=28)
plt.show()
```



```
In [12]: #comparing the players passes attempted
sapsors_warnings
warnings.filterwarnings("ignore")
fig=plt.figure(figsize=(15,5))
plt.rcParams['figure.figsize']=(15,5)
plt.figure(figsize=(15,5))
sns.distplot(football_data['Passes_Attempted'],color='blue')
a.set_xlabel('Passes Attempted Range for Players',fontsize=26)
a.set_ylabel('Count of the Players',fontsize=26)
a.set_title('Distribution of Passes Attempted for players',fontsize=28)
plt.show()
```



```
In [13]: plt.figure(figsize=(10,8))
sns.countplot(football_data['Yellow_Cards'],palette='dark')
p.set_title('Count of players on basis of yellow cards',fontsize=28)
p.set_xlabel('Cards',fontsize=28)
p.set_ylabel('Count of Players',fontsize=28)
plt.show()
```



```
In [14]: # best players per each position with their age, club, and nationality based on their xG
football_data.iloc[football_data.groupby(football_data['Position'])['xG'].idxmax()][['Position','Name','Age','Club','Nationality']].style.background_gradient('Blues')
```

The histogram displays the frequency of players within various age brackets. The highest frequency is observed in the 20-25 age group, with a count of approximately 0.00075. The frequency decreases as age increases beyond 25, with a slight secondary peak around the 30-35 age group.

Age Group	Count of the Players
15-20	0.00015
20-25	0.00075
25-30	0.00065
30-35	0.00055
35-40	0.00045
40-45	0.00040
45-50	0.00035
50-55	0.00025
55-60	0.00020

```
In [15]: # worst players per each position with their age, club, and nationality based on their xG
football_data.iloc[football_data.groupby(football_data['Position'])['xG'].idxmin()][['Position','Name','Age','Club','Nationality']].style.background_gradient('Blues')
```

```
p.set_title('Count of players on basis of yellow cards',fontsize=28)
p.set_xlabel(xlabel='Cards',fontsize=16)
p.set_ylabel(ylabel='Count of Players',fontsize=16)
plt.show()
```

Count of players on basis of yellow cards

Cards	Count of Players
1	180

```
In [16]: #Countries with most number of players
football_data['Nationality'].value_counts().head(10)
```

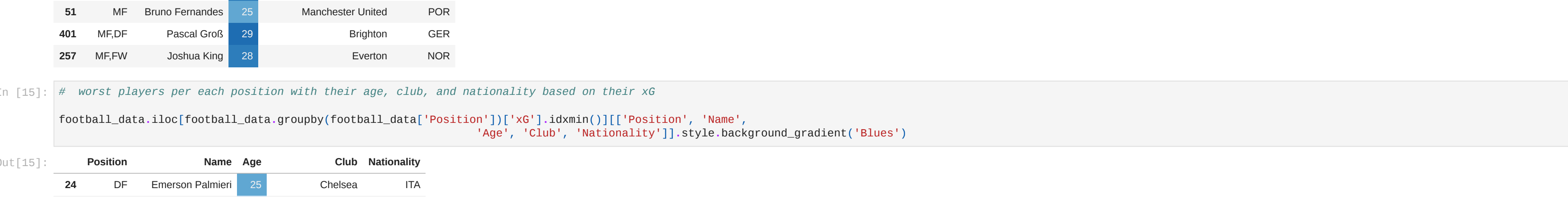
ENG	192
FRA	31
BRA	27
ESP	26
IRL	21
POR	21
SCO	20
NED	16
WAL	12
BEL	11

```
In [17]: # Every Nations' Player and their xG
plt.rcParams['figure.figsize']=(15,7)
countries=['ENG','ESP','BRA','POR','FRA','NED','SCO','IRL']
football_data.countries=football_data.loc[football_data['Nationality'].isin(countries)] & football_data['xG']
sns.barplot(football_data.countries['Nationality'],y=football_data.countries['xG'],color='red')
a.set_xlabel('Nationality',fontsize=16)
a.set_ylabel('xG',fontsize=16)
a.set_title('Violin Plot',fontsize=28)
plt.show()
```

```
In [17]: Text(0.5,1.8,"Violin Plot")
```



```
In [18]: # Every Nations' Player and their Perc_Passes_Completed
plt.rcParams['figure.figsize']=(15,7)
countries=['ENG','ESP','BRA','POR','FRA','NED','SCO','IRL']
football_data.countries=football_data.loc[football_data['Nationality'].isin(countries)] & football_data['Perc_Passes_Completed']
sns.barplot(football_data.countries['Nationality'],y=football_data.countries['Perc_Passes_Completed'],palette='Purple')
a.set_xlabel('Nationality',fontsize=16)
a.set_ylabel('Perc_Passes_Completed',fontsize=16)
a.set_title('Bar Plot',fontsize=20)
plt.show()
```



I now want to find out whether age can influence a player's expected goals. i therefore posed the question and found the most suitable algorithm to use.

Can Age influence player performance through expected Goals?

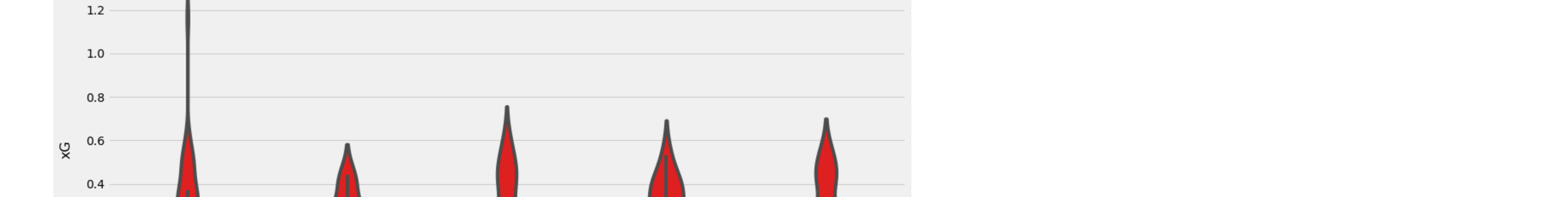
```
In [19]: #I decided to explore age and the other columns
# comparing the performance of age and matches of footballers
sns.leplot(x='Age',y='Matches', data = football_data)
```

```
Out[19]: <seaborn.axisgrid.FacetGrid at 8x2c7f78dfca>
```



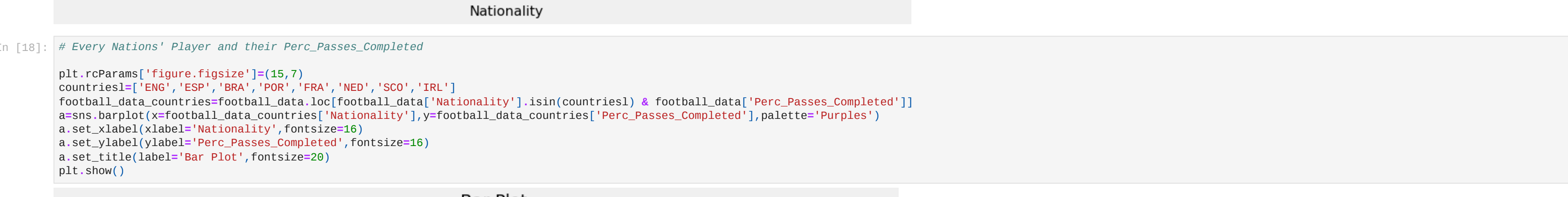
```
In [20]: # comparing the performance of age and xG of footballers
sns.leplot(x='Age',y='xG', data = football_data)
```

```
Out[20]: <seaborn.axisgrid.FacetGrid at 8x2c7f78dc886>
```



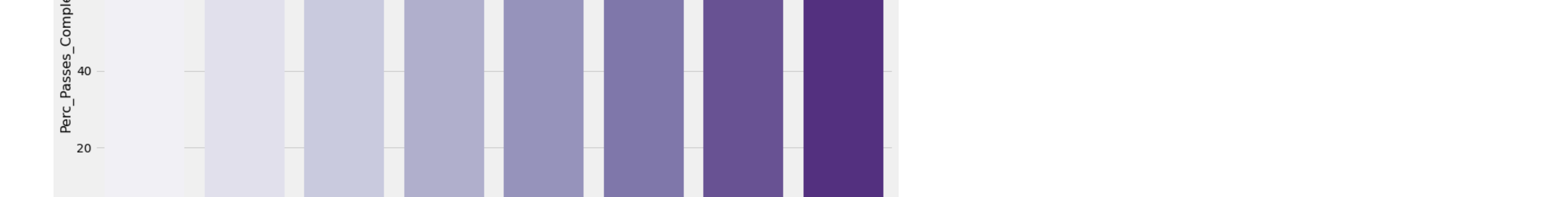
```
In [21]: # comparing the performance of age and xA of footballers
sns.leplot(x='Age',y='xA', data = football_data)
```

```
Out[21]: <seaborn.axisgrid.FacetGrid at 8x2c7f78916d8>
```



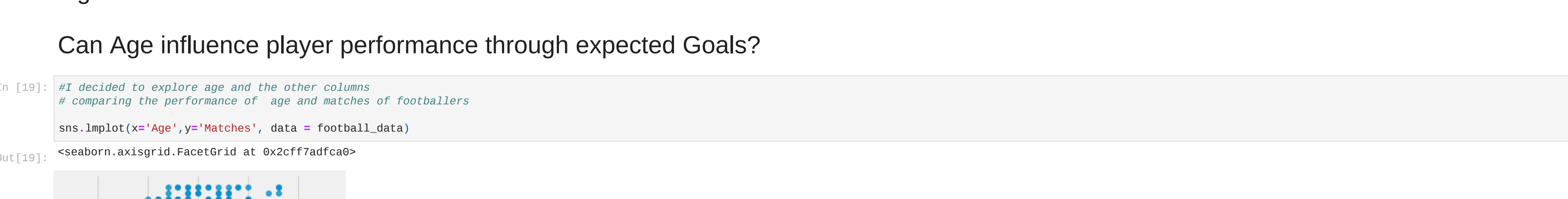
```
In [22]: # comparing the performance of age and assists per match of footballers
sns.leplot(x='Age',y='MinsPerMatch', data = football_data)
```

```
Out[22]: <seaborn.axisgrid.FacetGrid at 8x2c7f78a4886>
```



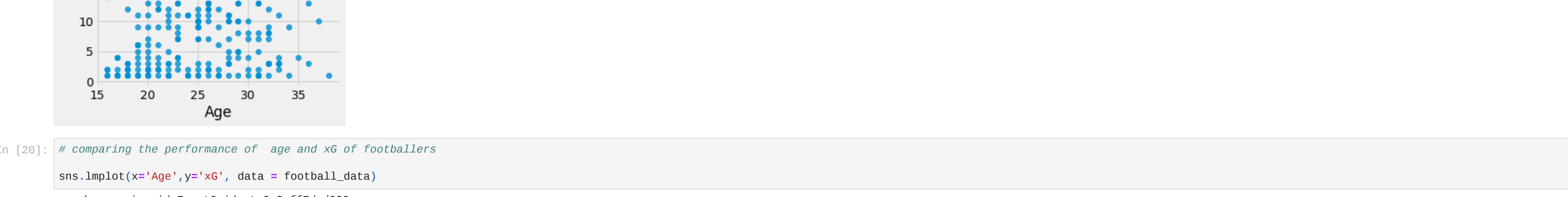
```
In [23]: # comparing the performance of age and goals per match of footballers
sns.leplot(x='Age',y='GoalsPerMatch', data = football_data)
```

```
Out[23]: <seaborn.axisgrid.FacetGrid at 8x2c7f78b6a86>
```



```
In [24]: # comparing the performance of age and goals of footballers
sns.leplot(x='Age',y='Goals', data = football_data)
```

```
Out[24]: <seaborn.axisgrid.FacetGrid at 8x2c7f78c1376>
```



Linear Regression algorithm to find the probability that age can influence expected goals outcome of a player

```
In [25]: #step 1 Import the required modules
from sklearn.datasets import make_classification
from sklearn.linear_model import LinearRegression
```

```
In [26]: # we need to set our independent and dependent variables
y = football_data.xG
x = football_data.Age.values.reshape(-1,1)
```

```
In [27]: print(x.shape , y.shape)
```

```
(532, 1) (532,)
```

```
In [28]: model = LinearRegression().fit(x,y)
```

```
In [29]: r_sq = model.score(x,y)
intercept = model.intercept_
slope = model.coef_
```

```
In [30]: print(r_sq)
```

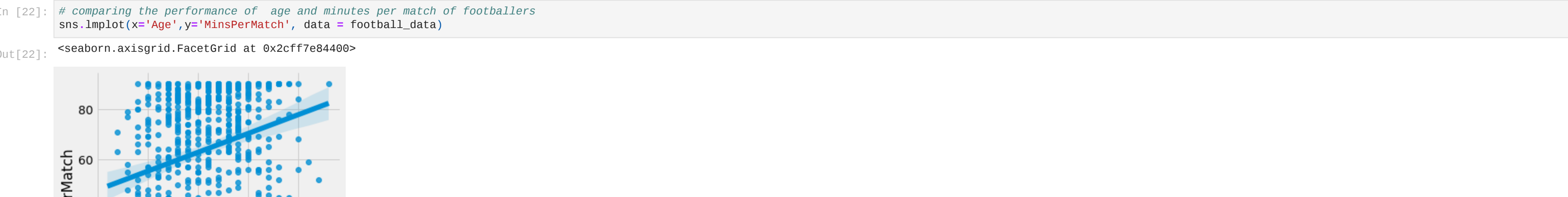
```
0.0031469840809148996
```

```
In [31]: y_pred = intercept + slope * x
```

```
In [32]: fig, ax = plt.subplots(figsize = (10,10))
# here we are going to create a scatterplot
plt.scatter(x,y)
plt.plot(x,y_pred, color = 'red')
plt.xlabel('Age')
plt.ylabel('xG')
plt.title('Evaluating the relationship between age and expected goals of premier league players')
```

```
Out[32]: Text(0.5,1.8,"Evaluating the relationship between age and expected goals of premier league players")
```

Evaluating the relationship between age and expected goals of premier league players



In conclusion, the study indicates R² of the study to be 0.00315 after using Linear Regression. This means that age doesn't influence the expected Goals of a premier league player, this is indicated in the scatterplot as players who are regarded to be in their prime years (18-30) can have a lower expected goal, and players on their 30s have also shown they can have a higher expected goal. However, the study included all premier league players, not regarding their playing position as a limiting factor. Therefore, the study can be improved further by considering player position for example assessing whether age can influence the expected goals of Premier League Forwards(CF).