

[1]: # Import the libraries import pandas as pd import numpy as np import matplotlib.pyplot as plt import seaborn as sns matplotlib inline

Matplotlib is building the font cache; this may take a moment.

In [4]: #Loading Dataset epl_df = pd.read_csv('E:\Data analysis training\ep1_analysis\EPL_20_21.csv') epl_df.head()

Out[4]:

	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	2690	6	5	1981	82.3	1	0	0.41	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	2602	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	2373	1	2	1987	85.0	0	0	0.06	0.12	3	0

In [5]: # columns and their names epl_df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 532 entries, 0 to 531
Data columns (total 18 columns):
 #   Column                Non-Null Count  Dtype
--  --
 0   Name                  532 non-null    object
 1   Club                  532 non-null    object
 2   Nationality           532 non-null    object
 3   Position              532 non-null    object
 4   Age                   532 non-null    int64
 5   Matches               532 non-null    int64
 6   Starts                532 non-null    int64
 7   Mins                  532 non-null    int64
 8   Goals                 532 non-null    int64
 9   Assists               532 non-null    int64
10  Passes_Attempted      532 non-null    int64
11  Perc_Passes_Completed  532 non-null    float64
12  Penalty_Goals          532 non-null    int64
13  Penalty_Attempted     532 non-null    int64
14  xG                     532 non-null    float64
15  xA                     532 non-null    float64
16  Yellow_Cards           532 non-null    int64
17  Red_Cards             532 non-null    int64
dtypes: float64(3), int64(11), object(4)
memory usage: 74.9+ KB
```

In [6]: # summary statistics about dataset, works only on numeric columns epl_df.describe()

Out[6]:

	Age	Matches	Starts	Mins	Goals	Assists	Passes_Attempted	Perc_Passes_Completed	Penalty_Goals	Penalty_Attempted	xG	xA	Yellow_Cards	Red_Cards
count	532.000000	532.000000	532.000000	532.000000	532.000000	532.000000	532.000000	532.000000	532.000000	532.000000	532.000000	532.000000	532.000000	532.000000
mean	25.500000	19.535714	15.714286	1411.443609	1.853883	1.287594	717.750000	77.828172	0.191729	0.234662	0.113289	0.074650	2.114662	0.080226
std	4.319404	11.844549	11.921161	1043.171856	3.338009	2.095191	631.372522	13.011631	0.850881	0.975818	0.184174	0.090072	2.269054	0.293068
min	16.000000	1.000000	0.000000	1.000000	0.000000	0.000000	0.000000	-1.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	22.000000	9.000000	4.000000	426.000000	0.000000	0.000000	171.500000	73.500000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
50%	26.000000	21.000000	15.000000	1345.000000	1.000000	0.000000	573.500000	79.200000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
75%	29.000000	30.000000	27.000000	2303.500000	2.000000	2.000000	1129.500000	84.625000	0.000000	0.000000	0.150000	0.115000	3.000000	0.000000
max	38.000000	38.000000	38.000000	3420.000000	23.000000	14.000000	3214.000000	100.000000	9.000000	10.000000	1.850000	0.800000	12.000000	2.000000

In [7]: # finding total number of null values epl_df.isna().sum()

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
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
5	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
6	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
7	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
8	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
9	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
17	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
dtypes:	int64(3)	int64(11)	object(4)															

In [11]: # Create 2 new columns epl_df['MinperMatch'] = epl_df['Mins'] / epl_df['Matches'].astype(int) epl_df['GoalsPerMatch'] = epl_df['Goals'] / epl_df['Matches'].astype(int) epl_df.head()

Out[11]:

	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	MinperMatch	Goals
0	Mason Mount	Chelsea	ENG	MF/FW	21	36	32	2690	6	5	1981	82.3	1	0	0.41	0.24	2	0	74.342857	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	88.543807	0.000000
2	Timo Werner	Chelsea	GER	FW	24	35	29	2602	6	8	826	77.2	0	0	0.41	0.21	2	0	74.342857	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	84.666667	0.111111
4	Reece James	Chelsea	ENG	DF	20	32	25	2373	1	2	1987	85.0	0	0	0.06	0.12	3	0	74.156250	0.062500

In [12]: # Total Goals Total_Goals = epl_df['Goals'].sum() print(Total_Goals)

986

In [13]: # Total penalty goals Total_Penalty = epl_df['Penalty_Goals'].sum() print(Total_Penalty)

182

In [15]: # Total penalty attempts Total_PenaltyAttempts = epl_df['Penalty_Attempted'].sum() print(Total_PenaltyAttempts)

325

In [21]: # Pie chart for penalties missed vs scored plt.figure(figsize=(12,6)) plt_not_scored = epl_df['Penalty_Attempted'].sum() - Total_Penalty data = [plt_not_scored, Total_Penalty] labels = ['Penalties missed', 'Penalties Scored'] color = sns.color_palette('Set2') plt.pie(data, labels=labels, colors= color, autopct= '%.0f%%')

Penalties missed
18%

Penalties Scored
82%

In [22]: #Unique positions epl_df['Position'].unique()

array(['MF', 'FW', 'GK', 'FW', 'DF', 'MF', 'FW', 'MF', 'FW', 'DF', 'DF', 'MF', 'MF', 'DF', 'FW'], dtype=object)

Out[23]:

Total FW players epl_df[epl_df['Position'] == 'FW']

#subsetting the position

Out[23]:

	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	MinperMatch
2	Timo Werner	Chelsea	GER	FW	24	35	29	2602	6	8	826	77.2	0	0	0.41	0.21	2	0	74.342857
16	Abraham	Chelsea	ENG	FW	22	22	12	1040	6	1	218	68.3	0	0	0.56	0.07	0	0	47.272727
19	Oliver	Chelsea	FRA	FW	33	17	8	748	4	0	217	74.2	0	0	0.58	0.09	1	0	44.000000
23	Ruben Loftus-Cheek	Chelsea	ENG	FW	24	1	1	60	0	0	16	68.8	0	0	0.00	0.00	0	0	60.000000
30	Raheem Sterling	Manchester City	ENG	FW	25	31	28	2536	10	7	1127	85.4	0	0	1.43	0.17	4	0	81.894562
516	Orkun Kökçü	Sheffied United	SCO	FW	23	25	14	1269	1	1	262	70.6	0	0	0.17	0.13	2	0	50.760000
518	Oliver McBurnie	Sheffied United	SCO	FW	24	23	12	1324	1	0	426	62.9	0	0	0.21	0.07	2	0	57.560217
519	Rhian Brewster	Sheffied United	ENG	FW	20	27	12	1128	0	0	225	69.3	0	0	0.14	0.13	1	0	41.777778
523	Billy Sharp	Sheffied United	ENG	FW	34	16	7	735	3	0	123	69.9	2	2	0.33	0.07	1	0	45.937500
526	Daniel Jebbison	Sheffied United	ENG	FW	17	4	3	284	1	0	34	70.6	0	0	0.50	0.01	0	0	71.000000

81 rows x 20 columns

In [24]: #players from different Nations np.size(epl_df['Nationality'].unique())

59

Out[24]:

In [25]: #what the nations they are representing epl_df['nationality'].unique()

array(['ENG', 'SEN', 'GER', 'ESP', 'ITA', 'BRA', 'CRO', 'USA', 'SEN', 'MAR', 'SCO', 'ARG', 'POR', 'BEL', 'ALG', 'UKR', 'NED', 'SWE', 'URU', 'SRB', 'WAL', 'CIV', 'NGA', 'EGY', 'TUR', 'CHN', 'GUA', 'SUI', 'JPN', 'IRN', 'IRI', 'ISR', 'IRQ', 'AUT', 'JAM', 'RSA', 'CZE', 'POL', 'PAR', 'COD', 'KOR', 'COL', 'GAB', 'MDR', 'RUS', 'ETH', 'ISL', 'PRK', 'GHA', 'ZIM', 'SWE', 'HKG', 'CAN', 'MLI', 'IRN', 'NZL', 'MTN', 'SKN'], dtype=object)

In [27]: # Which country is providing more players nationality = epl_df.groupby('Nationality').size().sort_values(ascending = False) nationality.head(10).plot(kind = 'bar', figsize=(12,6), color = sns.color_palette("magma"))

Out[27]:

<AxesSubplot: xlabel= 'Nationality'>

In [30]: #clubs with maximum players in their squad epl_df['Club'].value_counts().nlargest(5).plot(kind = 'bar', color = sns.color_palette("viridis"))

Out[30]:

<AxesSubplot>

In [31]: #clubs with least player in the squad epl_df['Club'].value_counts().nsmallest(5).plot(kind = 'bar', color = sns.color_palette("viridis"))

Out[31]:

<AxesSubplot>

In [39]: #players based on age group Under20 = epl_df[epl_df['Age']<= 20] age20_25 = epl_df[(epl_df['Age'] >= 20) & (epl_df['Age'] <= 25)] age25_30 = epl_df[(epl_df['Age'] >= 25) & (epl_df['Age'] <= 30)] Above30 = epl_df[epl_df['Age'] >= 30]

In [40]: x = np.array(Under20['Name'].count(), age20_25['Name'].count(), age25_30['Name'].count(), Above30['Name'].count()) labels = epl_df.groupby('Age').size().sort_values(ascending = False) nationality.head(10).plot(kind = 'bar', figsize=(12,6), color = sns.color_palette("magma"))

Out[40]:

<AxesSubplot: xlabel= 'Nationality'>

In [44]: # Total under 20 players in each club players_under_20 = epl_df[epl_df['Age'] <= 20] players_under_20['Club'].value_counts().plot(kind = 'bar', color = sns.color_palette("cubehelix"))

Out[44]:

<AxesSubplot>

In [46]: # under 20 players in Manchester united players_under_20[players_under_20['Club']=='Manchester United']

Out[46]:

	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	MinperMatch
61	Mason Mount	Manchester United	ENG	FW	18	31	21	1822	7	2	732	83.1	0	0	0.37	0.09	2	0	58.741934
72	Brandon Williams	Manchester United	ENG	DF	19	4	2	188	0	0	140	85.7	0	0	0.05	0.01	0	0	47.000000
73	Anthony Elanga	Manchester United	CIV	FW	18	3	2	166	0	1	64	84.4	0	0	0.02	0.26	0	0	55.333333
74	Anthony Elanga	Manchester United	SWE	FW	18	2	2	155	1	0	53	81.1	0	0	0.16	0.02	0	0	77.500000
76	Sheep	Manchester United	ENG	FW	16	2	0	11	0	0	8	75.0	0	0	0.00	0.00	0	0	5.500000
78	Hanshal Bhatnagar	Manchester United	FRA	MF	17	1	0	9	0	0	3	100.0	0	0	0.00	0.00	0	0	9.000000
79	William Thomas Fish	Manchester United	ENG	DF	17	1	0	1	0	0	1	0.0	0	0	0.00	0.00	0	0	1.000000

In [48]: #under 20 players in Liverpool players_under_20[players_under_20['Club']=='Liverpool FC']

Out[48]:

	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	MinperMatch	Go
91	Curtis Jones	Liverpool FC	ENG	MF	19	24	13	1179	1	2	976	91.2	0	0	0.11	0.12	2	0	49.125000	
96	Rhys Williams	Liverpool FC	ENG	DF	19	9	7	661	0	0	451	92.0	0	0	0.07	0.05	0	0	73.444444	
102	Neco Williams	Liverpool FC	WAL	DF	19	6	3	249	0	0	213	73.2	0	0	0.03	0.10	1	0	41.500000	

In [50]: #Average age of players in each club plt.figure(figsize=(12,6)) sns.boxplot(x = 'Club', y = 'Age', data = epl_df) plt.xticks(rotation = 90)

Out[50]:

array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19])

Text(0, 0, 'Chelsea'), Text(1, 0, 'Manchester United'), Text(2, 0, 'Liverpool FC'), Text(3, 0, 'Manchester United'), Text(4, 0, 'Leicester City'), Text(5, 0, 'West Ham United'), Text(6, 0, 'Tottenham Hotspur'), Text(7, 0, 'Arsenal'), Text(8, 0, 'Leeds United'), Text(9, 0, 'Everton'), Text(10, 0, 'Aston Villa'), Text(11, 0, 'Newcastle United'), Text(12, 0, 'Wolverhampton Wanderers'), Text(13, 0, 'Crystal Palace'), Text(14, 0, 'Southampton'), Text(15, 0, 'Brighton'), Text(16, 0, 'Burnley'), Text(17, 0, 'Fulham'), Text(18, 0, 'Sheffild United'), Text(19, 0, 'Sheffild United')]

In [52]: # Average age of players in each club in numerical manner num_player = epl_df.groupby('Club').size() data = epl_df.groupby('Club')['Age'].sum() num_player data.sort_values(ascending = False)

Out[52]:

Club	Sum
Crystal Palace	28.333333
West Ham United	27.580000
Burnley	27.400000
West Bromwich Albion	26.766667
Newcastle United	26.074074
Manchester City	25.769230
Tottenham Hotspur	25.625000
Chelsea	25.582500
Leicester City	25.582500
Liverpool FC	25.574285
Everton	25.537500
Leeds United	25.347826
Fulham	25.057143
Arsenal	24.985517
Sheffild United	24.814815
Brighton	24.505000
Wolverhampton Wanderers	24.444444
Aston Villa	24.350000
Southampton	24.179331
Manchester United	