

## Some interesting functions for scouting

Here some user defined function can be found. They are suitable for statistical players data in order to analyze the individual performance. They can be also integrated in the scouting context analysis, in order to evaluate the best performances.

you'll find:

- 1 - Expected goals |
- 2 - final third |
- 3 - won duels |
- 4 - goal assist key passes |
- 5 - key passes |

```
In [9]: import pandas as pd
import numpy as np
import json
# plotting
import matplotlib.pyplot as plt
# statistical fitting of models
import statsmodels.api as sm
import statsmodels.formula.api as smf
#opening data
import os
import pathlib
import warnings
#used for plots
from scipy import stats
from mplsoccer import PyPizza, FontManager

pd.options.mode.chained_assignment = None
warnings.filterwarnings('ignore')
```

```
In [5]: path = 'C:/Users/kecco/Desktop/MyScout/events/events_England.json'

with open(path) as f:
    data = json.load(f)
    train = pd.DataFrame(data)

players_path = 'C:/Users/kecco/Desktop/MyScout/players.json'

with open(players_path) as f:
    players = json.load(f)
    player_df = pd.DataFrame(players)
```

```
In [6]: train = pd.concat([train, pd.DataFrame(data)], ignore_index = True)
#potential data collection error handling
train = train.loc[train.apply (lambda x: len(x.positions) == 2, axis = 1)]
```

```
In [20]: def calculatexG(df, npxG):

    #very basic xG model based on
    shots = df.loc[df["eventName"] == "Shot"].copy()
    shots["X"] = shots.positions.apply(lambda cell: (100 - cell[0]['x']) * 105/100)
    shots["Y"] = shots.positions.apply(lambda cell: cell[0]['y'] * 68/100)
    shots["C"] = shots.positions.apply(lambda cell: abs(cell[0]['y'] - 50) * 68/100)
    #calculate distance and angle
    shots["Distance"] = np.sqrt(shots["X"]**2 + shots["C"]**2)
    shots["Angle"] = np.where(np.arctan(7.32 * shots["X"] / (shots["X"]**2 + shots["C"]**2 - (7.32/2)**2)) > 0, np.arctan(7.32 * shots["X"] / (shots["X"]**2 + shots["C"]**2 - (7.32/2)**2)), np.arctan(7.32 * shots["X"] / (shots["X"]**2 + shots["C"]**2)))
    #if you ever encounter problems (like you have seen that model treats 0 as 1 and 1 as 0) while modelling - change the dependant variable to object
    shots["Goal"] = shots.tags.apply(lambda x: 1 if ('id':101) in x else 0).astype(object)
    #headers have id = 403
    headers = shots.loc[shots.apply (lambda x: ('id':403) in x.tags, axis = 1)]
    non_headers = shots.drop(headers.index)

    headers_model = smf.glm(formula="Goal ~ Distance + Angle" , data=headers,
                           family=sm.families.Binomial()).fit()
    #non-headers
    nonheaders_model = smf.glm(formula="Goal ~ Distance + Angle" , data=non_headers,
                              family=sm.families.Binomial()).fit()

    #assigning xG
    #headers
    b_head = headers_model.params
    xG = 1/(1+np.exp(b_head[0]+b_head[1]*headers['Distance'] + b_head[2]*headers['Angle']))
    headers = headers.assign(xG = xG)

    #non-headers
    b_nhead = nonheaders_model.params
    xG = 1/(1+np.exp(b_nhead[0]+b_nhead[1]*non_headers['Distance'] + b_nhead[2]*non_headers['Angle']))
    non_headers = non_headers.assign(xG = xG)

    if npxG == False:
        #find pens
        penalties = df.loc[df["subEventName"] == "Penalty"]
        #assign 0.8
        penalties = penalties.assign(xG = 0.8)
        #concat, group and sum
        all_shots_xg = pd.concat([non_headers[["playerId", "xG"]], headers[["playerId", "xG"]], penalties[["playerId", "xG"]]])
        xG_sum = all_shots_xg.groupby(["playerId"])["xG"].sum().sort_values(ascending = False).reset_index()
    else:
        #concat, group and sum
        all_shots_xg = pd.concat([non_headers[["playerId", "xG"]], headers[["playerId", "xG"]]])
        all_shots_xg.rename(columns = ("xG": "npxG"), inplace = True)
        xG_sum = all_shots_xg.groupby(["playerId"])["npxG"].sum().sort_values(ascending = False).reset_index()
    #group by player and sum

    return xG_sum

#making function
npxg = calculatexG(train, npxG = True)
#investigate structure
npxg.head(3)
```

```
Out[20]:   playerId    npxG
0         8717  44.028360
1        120353  34.431637
2         11066  28.288968
```

The function returns a DataFrame (xG\_sum) containing the total xG (or npxG) for each player, aggregated from all their shots. It provides a way to estimate a player's goal-scoring threat by combining shot positions, logistic regression modeling, and the concept of expected goals. The handling of penalties allows for flexibility in considering different scenarios.

```
In [16]: def FinalThird(df):

    df = df.copy()
    #need player who had received the ball
    df["nextPlayerId"] = df["playerId"].shift(-1)
    passes = df.loc[train["eventName"] == "Pass"].copy()
    #changing coordinates
    passes["X"] = passes.positions.apply(lambda cell: (cell[0]['x']) * 105/100)
    passes["Y"] = passes.positions.apply(lambda cell: (100 - cell[0]['y']) * 68/100)
    passes["end_x"] = passes.positions.apply(lambda cell: (cell[1]['x']) * 105/100)
    passes["end_y"] = passes.positions.apply(lambda cell: (100 - cell[1]['y']) * 68/100)

    #get accurate passes
    accurate_passes = passes.loc[passes.apply (lambda x: ('id':1801) in x.tags, axis = 1)]
    #get passes into final third
    final_third_passes = accurate_passes.loc[accurate_passes["end_x"] > 2*105/3]

    #passes into final third by player
    ftp_player = final_third_passes.groupby(["playerId"]).end_x.count().reset_index()
    ftp_player.rename(columns = {'end_x': 'final_third_passes'}, inplace=True)

    #receptions of accurate passes in the final third
    rtp_player = final_third_passes.groupby(["nextPlayerId"]).end_x.count().reset_index()
    rtp_player.rename(columns = {'end_x': 'final_third_receptions', "nextPlayerId": "playerId"}, inplace=True)

    #outer join not to lose values
    final_third = ftp_player.merge(rtp_player, how = "outer", on = ["playerId"])
    return final_third

final_third = FinalThird(train)
#investigate structure
final_third.head(3)
```

```
Out[16]:   playerId  final_third_passes  final_third_receptions
0         36.0                 372.0                 166.0
1         38.0                 124.0                 132.0
2         48.0                 784.0                 376.0
```

It isn't enough for a striker to be a good passer of the ball he or she should be able to perform well in the final third. The function merges the counts of final third passes and receptions into a DataFrame, ensuring that players with no passes or receptions are not lost. It is designed to summarize a player's involvement in passing within the final third of the pitch, providing insights into their playmaking abilities in key attacking areas.

```
In [18]: def wonDuels(df):
    #find air duels
    air_duels = df.loc[df["subEventName"] == "Air duel"]
    #703 is the id of a won duel
    won_air_duels = air_duels.loc[air_duels.apply (lambda x: ('id':703) in x.tags, axis = 1)]

    #group and sum air duels
    wad_player = won_air_duels.groupby(["playerId"]).eventId.count().reset_index()
    wad_player.rename(columns = {'eventId': 'air_duels_won'}, inplace=True)

    #find ground duels won
    ground_duels = df.loc[df["subEventName"].isin(["Ground attacking duel"])]
    won_ground_duels = ground_duels.loc[ground_duels.apply (lambda x: ('id':703) in x.tags, axis = 1)]

    wgd_player = won_ground_duels.groupby(["playerId"]).eventId.count().reset_index()
    wgd_player.rename(columns = {'eventId': 'ground_duels_won'}, inplace=True)

    #outer join
    duels_won = wgd_player.merge(wad_player, how = "outer", on = ["playerId"])
    return duels_won

duels = wonDuels(train)
#investigate structure
duels.head(3)
```

```
Out[18]:   playerId  ground_duels_won  air_duels_won
0         0                 4488.0             2122.0
1         36                  26.0              46.0
2         38                  14.0              22.0
```

The counts of ground duels and air duels won are merged into a DataFrame ("duels\_won") using an outer join on the player ID. It summarizes a player's success in both ground and air duels, providing insights into their physical prowess and effectiveness in challenging situations during a match.

```
In [19]: def GoalsAssistsKeyPasses(df):

    #get goals
    shots = df.loc[df["subEventName"] == "Shot"]
    goals = shots.loc[shots.apply (lambda x: ('id':101) in x.tags, axis = 1)]
    #get assists
    passes = df.loc[df["eventName"] == "Pass"]
    assists = passes.loc[passes.apply (lambda x: ('id':301) in x.tags, axis = 1)]
    #get key passes
    key_passes = passes.loc[passes.apply (lambda x: ('id':302) in x.tags, axis = 1)]

    #goals by player
    g_player = goals.groupby(["playerId"]).eventId.count().reset_index()
    g_player.rename(columns = {'eventId': 'goals'}, inplace=True)

    #assists by player
    a_player = assists.groupby(["playerId"]).eventId.count().reset_index()
    a_player.rename(columns = {'eventId': 'assists'}, inplace=True)

    #key passes by player
    kp_player = key_passes.groupby(["playerId"]).eventId.count().reset_index()
    kp_player.rename(columns = {'eventId': 'key_passes'}, inplace=True)

    data = g_player.merge(a_player, how = "outer", on = ["playerId"]).merge(kp_player, how = "outer", on = ["playerId"])
    return data

gakp = GoalsAssistsKeyPasses(train)
#investigate structure
gakp.head(3)
```

```
Out[19]:   playerId  goals  assists  key_passes
0         54      20.0      10.0      50.0
1         74       2.0      NaN       2.0
2         93       4.0      10.0      28.0
```

The counts of goals, assists, and key passes are merged into a DataFrame ("data") using outer joins on the player ID. It provides a consolidated view of a player's attacking contributions, including their goal-scoring ability, assisting prowess, and the creation of key passes. It offers valuable insights into a player's overall offensive impact on the game.

```
In [7]: def smartPasses(df):
    #get smart passes
    smart_passes = df.loc[df["subEventName"] == "Smart pass"]
    #find accurate
    smart_passes_made = smart_passes.loc[smart_passes.apply (lambda x: ('id':1801) in x.tags, axis = 1)]

    #sum by player
    sp_player = smart_passes_made.groupby(["playerId"]).eventId.count().reset_index()
    sp_player.rename(columns = {'eventId': 'smart_passes'}, inplace=True)

    return sp_player

smart_passes = smartPasses(train)
#investigate structure
smart_passes.head(3)
```

```
Out[7]:   playerId  smart_passes
0         36             2
1         38             2
2         48             6
```

Another statistic that we want to add are accurate smart passes. Those are the passes that break the opponent defensive line. It provides a clear summary of the number of smart passes made by each player.