

```

1 #%%
2 ### Contributions
3 # Tasks were distributed as follows:
4 # [Wenhan Zhan]: [Part 1 in EDA]
5 # [Ruoyao Yan]: [Part 1 in Data Cleaning, Part 4,7
6 # in EDA, Part 5 in Advanced Analysis, Part 1,2,3 in
7 # Prediction]
8 # [Jun Sun]: [Part 2 in Data Cleaning, Part 1,2,3,5
9 ,6 in EDA, Part 1,2,3,4 in Advanced Analysis, Part
10 3 in Prediction]
11 #%% md
12 # 1. Data Cleaning
13 #%%
14 # Module: [Data Cleaning1]
15 # Contributor: [Ruoyao Yan]
16 import pandas as pd
17
18 df = pd.read_csv('Topic1_english_pl_dataset.csv')
19
20 numeric_cols = ['FullTimeHomeGoals', 'FullTimeAwayGoals', 'HalfTimeHomeGoals',
21                 'HalfTimeAwayGoals', 'HomeShots', 'AwayShots', 'HomeShotsOnTarget',
22                 'AwayShotsOnTarget', 'HomeCorners', 'AwayCorners', 'HomeFouls',
23                 'AwayFouls', 'HomeYellowCards', 'AwayYellowCards', 'HomeRedCards',
24                 'AwayRedCards']
25
26 df[numeric_cols] = df[numeric_cols].fillna(df[
27     numeric_cols].mean())
28 categorical_cols = ['FullTimeResult', 'HalfTimeResult', 'HomeTeam', 'AwayTeam']
29 for col in categorical_cols:
30     df[col] = df[col].fillna(df[col].mode()[0])
31

```

```

32 #check check check check list
33 df.loc[df['FullTimeHomeGoals'] < 0, 'FullTimeHomeGoals'] = 0
34 df.loc[df['FullTimeAwayGoals'] < 0, 'FullTimeAwayGoals'] = 0
35 df.loc[df['HalfTimeHomeGoals'] < 0, 'HalfTimeHomeGoals'] = 0
36 df.loc[df['HalfTimeAwayGoals'] < 0, 'HalfTimeAwayGoals'] = 0
37
38
39 mask = df['HalfTimeHomeGoals'] > df['FullTimeHomeGoals']
40 df.loc[mask, 'HalfTimeHomeGoals'] = df.loc[mask, 'FullTimeHomeGoals']
41
42 mask = df['HalfTimeAwayGoals'] > df['FullTimeAwayGoals']
43 df.loc[mask, 'HalfTimeAwayGoals'] = df.loc[mask, 'FullTimeAwayGoals']
44
45 df['FullTimeHomeGoals'] = df['FullTimeHomeGoals'].astype(int)
46 df['FullTimeAwayGoals'] = df['FullTimeAwayGoals'].astype(int)
47 df['HomeShotsOnTarget'] = df['HomeShotsOnTarget'].astype(int)
48 df['HomeFouls'] = df['HomeFouls'].astype(int)
49 df['HalfTimeAwayGoals'] = df['HalfTimeAwayGoals'].astype(int)
50 df['AwayRedCards'] = df['AwayRedCards'].astype(int)
51
52 valid_results = ['H', 'D', 'A']
53 df['FullTimeResult'] = df['FullTimeResult'].apply(
    lambda x: x if x in valid_results else df['FullTimeResult'].mode()[0])
54
55
56 def check_result_consistency(row):
57     if row['FullTimeHomeGoals'] > row['FullTimeAwayGoals'] and row['FullTimeResult'] != 'H'

```

```

57  '':
58      return False
59      elif row['FullTimeHomeGoals'] < row['
60          FullTimeAwayGoals'] and row['FullTimeResult'] != 'A
61      '':
62          return False
63      elif row['FullTimeHomeGoals'] == row['
64          FullTimeAwayGoals'] and row['FullTimeResult'] != 'D
65      '':
66          return False
67      return True
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88

```

57 '':

58 return False

59 elif row['FullTimeHomeGoals'] < row['
 FullTimeAwayGoals'] and row['FullTimeResult'] != 'A
 '':

60 return False

61 elif row['FullTimeHomeGoals'] == row['
 FullTimeAwayGoals'] and row['FullTimeResult'] != 'D
 '':

62 return False

63 return True

64

65

66 df['Consistent'] = df.apply(
 check_result_consistency, axis=1)

67 inconsistent = df[~df['Consistent']]

68

69 for idx, row in inconsistent.iterrows():

70 if row['FullTimeHomeGoals'] > row['
 FullTimeAwayGoals']:
 df.at[idx, 'FullTimeResult'] = 'H'
 elif row['FullTimeHomeGoals'] < row['
 FullTimeAwayGoals']:
 df.at[idx, 'FullTimeResult'] = 'A'
 else:
 df.at[idx, 'FullTimeResult'] = 'D'

71

72

73

74

75

76

77 df['HomeShots'] = df[['HomeShots', '
 HomeShotsOnTarget']].max(axis=1)

78 df['AwayShots'] = df[['AwayShots', '
 AwayShotsOnTarget']].max(axis=1)

79

80 df.drop('Consistent', axis=1, inplace=True)

81

82 df.to_csv('cleaned_football_data.csv', index=False)

83 print("data has been cleaned and saved! ")

84 print("\nMissing value statistics after cleaning:")

85 print(df.isnull().sum())

86

87 #for ppt

88 print("\ncleaned data sample:")

```
89 print(df.head())
90 time_df = pd.DataFrame({
91     'MatchDate': ['2000-08-19', '2001-09-20',
92     '#####', '2003-07-01']
93 })
94 time_df['MatchDate'] = pd.to_datetime(time_df['MatchDate'].replace('#####', pd.NaT), errors='coerce')
95
96 print(time_df)
97
98 def determine_outcome(home_goals, away_goals):
99     if home_goals > away_goals:
100         return 'Home Win'
101     elif home_goals < away_goals:
102         return 'Away Win'
103     else:
104         return 'Draw'
105
106
107 df['MatchOutcome'] = df.apply(
108     lambda row: determine_outcome(row['FullTimeHomeGoals'], row['FullTimeAwayGoals']),
109     axis=1
110 )
111
112 print(df[['MatchDate', 'HomeTeam', 'AwayTeam', 'FullTimeHomeGoals', 'FullTimeAwayGoals', 'MatchOutcome']].head())
113
114 count_home = (df['FullTimeResult'] == 'H').sum()
115 print(f"Total home win: {count_home}")
116 count_away = (df['FullTimeResult'] == 'A').sum()
117 print(f"Total away win: {count_away}")
118 print("This shows that they have a higher winning rate play at home compare to play away")
119 num_teams = df['HomeTeam'].nunique()
120 print(f"total numbers of football teams: {num_teams}")
121 #%%
```

```

122 # Module: [Data Cleaning2]
123 # Contributor: [Jun Sun]
124 import pandas as pd
125 df = pd.read_csv('cleaned_football_data.csv')
126 df['HomeGoalDifference'] = df['FullTimeHomeGoals']
127     ] - df['FullTimeAwayGoals']
127 df['AwayGoalDifference'] = df['FullTimeAwayGoals']
128     ] - df['FullTimeHomeGoals']
128 df['GoalDifference'] = df['HomeGoalDifference']
129 df.to_csv('cleaned_football_data.csv', index=False)
130
130 import pandas as pd
131 df = pd.read_csv('cleaned_football_data.csv')
132 print(df.head())
133 #%% md
134 # 2. Basic Exploratory Data Analysis (EDA)
135 2.1 Trends in number of matches played (by season)
136 #%%
137 # Module: [EDA1]
138 # Contributor: [Wenhan Zhan] and [Junxiang Wang]
139 import pandas as pd
140 import matplotlib.pyplot as plt
141 import seaborn as sns
142
143 df = pd.read_csv("cleaned_football_data.csv")
144
145 df['MatchDate'] = pd.to_datetime(df['MatchDate'])
146
147 def get_season(date):
148     if date.month >= 8:
149         return f"{date.year}-{date.year+1}"
150     else:
151         return f"{date.year-1}-{date.year}"
152
153 df['Season'] = df['MatchDate'].apply(get_season)
154
155 season_counts = df['Season'].value_counts().
156     sort_index()
157 plt.figure(figsize=(12,6))
158 sns.lineplot(x=season_counts.index,y=season_counts

```

```
158 .values, markers='o')
159 plt.xticks(rotation = 30)
160 plt.title('Number of Matches per Season')
161 plt.xlabel('Season')
162 plt.ylabel('Number of Matches')
163 plt.tight_layout()
164 plt.show()
165 #%% md
166 # Trends in number of matches played (by Months)
167 #%%
168 # Module: [EDA2]
169 # Contributor: [Jun Sun]
170 import pandas as pd
171 import matplotlib.pyplot as plt
172 from datetime import datetime
173
174 df = pd.read_csv(
175     'cleaned_football_data.csv')
176
177 def is_valid_date(date_str):
178     try:
179         datetime.strptime(date_str, "%Y-%m-%d")
180     # Change format if needed
181     return True
182     except:
183         return False
184
185 valid_dates = df['MatchDate'].apply(is_valid_date)
186 df = df[valid_dates]
187
188 df['MatchDate'] = pd.to_datetime(df['MatchDate'])
189
190 df['Month'] = df['MatchDate'].dt.month_name()
191
192 month_order = ['August', 'September', 'October', 'November',
193                 'December',
194                 'January', 'February', 'March', 'April', 'May']
195 monthly_matches = df['Month'].value_counts().
```

```
195 reindex(month_order)
196
197 plt.figure(figsize=(12, 6))
198 monthly_matches.plot(kind='bar', color='blue')
199 plt.title('Total Number of Matches by Month (All
    Seasons Combined)')
200 plt.xlabel('Month')
201 plt.ylabel('Number of Matches')
202 plt.grid(axis='y', linestyle='--', alpha=0.7)
203 plt.xticks(rotation=45)
204 plt.tight_layout()
205 plt.show()
206


---


207 #%% md
208 ## Distribution of home and away team goals (bar
chart)


---


209 #%%
210 # Module: [EDA3]
211 # Contributor: [Jun Sun]
212 import pandas as pd
213 import matplotlib.pyplot as plt
214
215 df = pd.read_csv(
216     'cleaned_football_data.csv')
217
218 home_goals = df['FullTimeHomeGoals'].value_counts
().sort_index()
219 away_goals = df['FullTimeAwayGoals'].value_counts
().sort_index()
220
221 max_goals = max(home_goals.index.max(), away_goals
.index.max())
222 x_range = range(0, max_goals + 1)
223
224 home_goals = home_goals.reindex(x_range,
fill_value=0)
225 away_goals = away_goals.reindex(x_range,
fill_value=0)
226
227 fig, ax = plt.subplots(1, 2, figsize=(14, 5),
sharey=True)
```

```

228
229 ax[0].bar(home_goals.index, home_goals.values,
   color='skyblue', edgecolor='black')
230 ax[0].set_title('Home Team Goals Distribution',
   fontsize=14)
231 ax[0].set_xlabel('Goals Scored', fontsize=12)
232 ax[0].set_ylabel('Number of Matches', fontsize=12)
233 ax[0].set_xticks(x_range)
234 ax[0].grid(axis='y', linestyle='--', alpha=0.5)
235
236 ax[1].bar(away_goals.index, away_goals.values,
   color='salmon', edgecolor='black')
237 ax[1].set_title('Away Team Goals Distribution',
   fontsize=14)
238 ax[1].set_xlabel('Goals Scored', fontsize=12)
239 ax[1].set_xticks(x_range)
240 ax[1].grid(axis='y', linestyle='--', alpha=0.5)
241
242 plt.tight_layout()
243 plt.show()
244


---


245 #%%
246 # Module: [EDA4]
247 # Contributor: [Ruoyao Yan]
248 import pandas as pd
249 import matplotlib.pyplot as plt
250 import seaborn as sns
251
252
253 df = pd.read_csv("cleaned_football_data.csv")
254 df['MatchDate'] = pd.to_datetime(df['MatchDate'])
255
256 def get_season(date):
257     if date.month >= 8:
258         return f"{date.year}-{date.year+1}"
259     else:
260         return f"{date.year-1}-{date.year}"
261
262 df['Season'] = df['MatchDate'].apply(get_season)
263 season_goals = df.groupby('Season')[['
   FullTimeHomeGoals', 'FullTimeAwayGoals']].sum().

```

```

263 sum(axis=1)
264 first_15_seasons = [s for s in season_goals.index
265     if int(s.split('-')[0]) <= 2014]
266 last_10_seasons = [s for s in season_goals.index
267     if 2015 <= int(s.split('-')[0]) <= 2024]
268
269
270 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(16,
271                                 6))
272 ax1.bar(range(len(first_15_seasons)),
273         season_goals[first_15_seasons].values,
274         color='skyblue', edgecolor='black')
275 ax1.set_title('Total Goals (2000-2015)')
276 ax1.set_xticks(range(len(first_15_seasons)))
277 ax1.set_xticklabels(first_15_seasons, rotation=45
278 , ha='right')
279 ax1.set_xlim(0, 1300)
280 ax2.bar(range(len(last_10_seasons)),
281         season_goals[last_10_seasons].values,
282         color='salmon', edgecolor='black')
283 ax2.set_title('Total Goals (2015-2025)')
284 ax2.set_xticks(range(len(last_10_seasons)))
285 ax2.set_xticklabels(last_10_seasons, rotation=45,
286 ha='right')
287 ax2.set_xlim(0, 1300)
288 plt.tight_layout()
289 plt.show()
290 #%% md
291 ## Full-time/half-time result ratio (pie chart)
292 #%%
293 # Module: [EDA5]
294 # Contributor: [Jun Sun]
295 import matplotlib.pyplot as plt
296
297 result_order = ['H', 'D', 'A']
298
299 fulltime_counts = df['FullTimeResult'].
300     value_counts().reindex(result_order, fill_value=0)
301 fulltime_labels = ['Home Win (H)', 'Draw (D)', '
302     Away Win (A)']

```

```
297 fulltime_colors = ['lightgreen', 'lightgray', 'tomato']
298
299
300 halftime_counts = df['HalfTimeResult'].
    value_counts().reindex(result_order, fill_value=0)
301 halftime_labels = ['Home Lead (H)', 'Draw (D)', 'Away Lead (A)']
302 halftime_colors = ['lightblue', 'lightgray', 'orange']
303
304 fig, ax = plt.subplots(1, 2, figsize=(14, 6))
305
306
307 ax[0].pie(fulltime_counts, labels=fulltime_labels,
    , autopct='%.1f%%',
308             colors=fulltime_colors, startangle=140,
    wedgeprops={'edgecolor': 'black'})
309 ax[0].set_title('Full-Time Match Result Proportion',
    , fontsize=14)
310
311 # Half-Time Pie Chart
312 ax[1].pie(halftime_counts, labels=halftime_labels,
    , autopct='%.1f%%',
313             colors=halftime_colors, startangle=140,
    wedgeprops={'edgecolor': 'black'})
314 ax[1].set_title('Half-Time Match Result Proportion',
    , fontsize=14)
315
316 plt.tight_layout()
317 plt.show()
318
319 #%% md
320 ## Home vs Away Win Percentage (bar graph)
321 #%%
322 # Module: [EDA6]
323 # Contributor: [Jun Sun]
324 import matplotlib.pyplot as plt
325
326 total_matches = len(df)
327 total_home_matches = df['HomeTeam'].count()
```

```

328 total_away_matches = df['AwayTeam'].count()
329
330 num_home_wins = (df['FullTimeResult'] == 'H').sum()
331 num_away_wins = (df['FullTimeResult'] == 'A').sum()
332
333 home_win_rate = num_home_wins / total_home_matches
334 away_win_rate = num_away_wins / total_away_matches
335
336 win_rate_data = {
337     'Home Win Rate': home_win_rate,
338     'Away Win Rate': away_win_rate
339 }
340
341
342 fig, ax = plt.subplots(figsize=(8, 5))
343 bars = ax.bar(win_rate_data.keys(), win_rate_data.
    values(), color=['royalblue', 'darkorange'])
344
345 for bar in bars:
346     height = bar.get_height()
347     ax.text(bar.get_x() + bar.get_width()/2,
348             height + 0.02,
349                 f'{height:.1%}', ha='center', va='bottom',
350                 fontsize=12)
351 ax.set_title('Home vs Away Win Rate', fontsize=14)
352 ax.set_ylabel('Win Rate', fontsize=12)
353 ax.set_ylim(0, 1)
354 ax.grid(axis='y', linestyle='--', alpha=0.7)
355 plt.tight_layout()
356 plt.show()
357 #%%
358 # Module: [EDA7]
359 # Contributor: [Ruoyao Yan]
360 import matplotlib.pyplot as plt
361 import pandas as pd
362
363 df = pd.read_csv('cleaned_football_data.csv')

```

```

364
365 # Filter Arsenal vs Brighton matches
366 AvsB = []
367 for index, row in df.iterrows():
368     home = row['HomeTeam']
369     away = row['AwayTeam']
370     if (home == 'Arsenal' and away == 'Brighton') or (home == 'Brighton' and away == 'Arsenal'):
371         AvsB.append(row)
372
373 arsenal_vs_brighton = pd.DataFrame(AvsB)
374
375 total_arsenal_wins = len(arsenal_vs_brighton[
376     ((arsenal_vs_brighton['HomeTeam'] == 'Arsenal') & (arsenal_vs_brighton['FullTimeResult'] == 'H')) |
377     ((arsenal_vs_brighton['AwayTeam'] == 'Arsenal') & (arsenal_vs_brighton['FullTimeResult'] == 'A'))
378 ])
379 total_brighton_wins = len(arsenal_vs_brighton[
380     ((arsenal_vs_brighton['HomeTeam'] == 'Brighton') & (arsenal_vs_brighton['FullTimeResult'] == 'H')) |
381     ((arsenal_vs_brighton['AwayTeam'] == 'Brighton') & (arsenal_vs_brighton['FullTimeResult'] == 'A'))
382 ])
383 total_draws = len(arsenal_vs_brighton[
384     arsenal_vs_brighton['FullTimeResult'] == 'D'])
385 labels = ['Arsenal Wins', 'Brighton Wins', 'Draws']
386 counts = [total_arsenal_wins, total_brighton_wins, total_draws]
387 colors = ['red', 'blue', 'grey']
388
389 plt.figure(figsize=(6, 4))
390 plt.pie(
391     counts,
392     labels=labels,

```

```
393     autopct='%1.2f%',  
394     startangle=90,  
395     colors=colors,  
396 )  
397  
398 plt.title('Arsenal vs Brighton Match Results',  
            fontsize=14)  
399 plt.tight_layout()  
400 plt.show()  
401 #%% md  
402  
403 # 3. Advanced Statistical Analysis  
404 #%% md  
405 Average shots per match, shots on target, and goal  
      conversion rate  
406 #%%  
407 # Module: [Advanced Analysis1]  
408 # Contributor: [Jun Sun]  
409 df.fillna({  
410     'HomeShots': 0, 'AwayShots': 0,  
411     'HomeShotsOnTarget': 0, 'AwayShotsOnTarget': 0  
        ,  
412     'FullTimeHomeGoals': 0, 'FullTimeAwayGoals': 0  
413 }, inplace=True)  
414  
415 df['TotalShots'] = df['HomeShots'] + df['AwayShots']  
        [  
416 df['TotalShotsOnTarget'] = df['HomeShotsOnTarget']  
        + df['AwayShotsOnTarget']  
417 df['TotalGoals'] = df['FullTimeHomeGoals'] + df['  
        FullTimeAwayGoals']  
418  
419 average_shots = df['TotalShots'].mean()  
420 average_shots_on_target = df['TotalShotsOnTarget']  
        .mean()  
421 average_goals = df['TotalGoals'].mean()  
422  
423 total_goals = df['TotalGoals'].sum()  
424 total_shots = df['TotalShots'].sum()  
425 goal_conversion_rate = total_goals / total_shots  
    if total_shots > 0 else 0
```

```
426
427
428 print("Match Statistics Summary")
429 print("-" * 40)
430 print(f"Average total shots per match : {average_shots:.2f}")
431 print(f"Average shots on target per match : {average_shots_on_target:.2f}")
432 print(f"Average goals per match : {average_goals:.2f}")
433 print(f"Goal conversion rate (Goals/Shots): {goal_conversion_rate:.2%}")
434


---


435 #%% md
436 ## Average distribution of fouls, yellow cards,
437 and red cards (Home vs Away)


---


438 # Module: [Advanced Analysis2]
439 # Contributor: [Jun Sun]
440 import matplotlib.pyplot as plt
441 import numpy as np
442
443 # Compute averages
444 categories = ['Fouls', 'Yellow Cards', 'Red Cards']
445 home_values = [
446     df['HomeFouls'].mean(),
447     df['HomeYellowCards'].mean(),
448     df['HomeRedCards'].mean()
449 ]
450 away_values = [
451     df['AwayFouls'].mean(),
452     df['AwayYellowCards'].mean(),
453     df['AwayRedCards'].mean()
454 ]
455
456 x = np.arange(len(categories))
457 bar_width = 0.35
458
459 fig, ax = plt.subplots(figsize=(10, 6))
460 bars1 = ax.bar(x - bar_width / 2, home_values,
```

```
460 bar_width, label='Home Team', color='skyblue')
461 bars2 = ax.bar(x + bar_width / 2, away_values,
   bar_width, label='Away Team', color='salmon')
462
463
464 ax.bar_label(bars1, fmt='%.2f', padding=3)
465 ax.bar_label(bars2, fmt='%.2f', padding=3)
466
467 ax.set_ylabel('Average per Match')
468 ax.set_title('Average Fouls and Cards per Match (
   Home vs Away)')
469 ax.set_xticks(x)
470 ax.set_xticklabels(categories)
471 ax.set_ylim(0, max(max(home_values), max(
   away_values)) + 2)
472 ax.grid(axis='y', linestyle='--', alpha=0.7)
473 ax.legend()
474
475 plt.tight_layout()
476 plt.show()
477 #%% md
478 # Correlation between corner kicks and win rate
479 #%%
480 # Module: [Advanced Analysis3]
481 # Contributor: [Jun Sun]
482 import pandas as pd
483
484
485 df = pd.read_csv("cleaned_football_data.csv")
486
487 df['HomeWin'] = df['FullTimeResult'].apply(lambda
   x: 1 if x == 'H' else 0)
488 df['AwayWin'] = df['FullTimeResult'].apply(lambda
   x: 1 if x == 'A' else 0)
489
490 home_correlation = df[['HomeCorners', 'HomeWin']].
   corr().iloc[0, 1]
491 away_correlation = df[['AwayCorners', 'AwayWin']].
   corr().iloc[0, 1]
492
493 print("Correlation between Home Corners and Home
```

```
493 Win:", home_correlation)
494 print("Correlation between Away Corners and Away
      Win:", away_correlation)
495 #%% md
496 # Analysis of the relationship between goal
      difference and match outcome
497 #%%
498 # Module: [Advanced Analysis4]
499 # Contributor: [Jun Sun]
500 import matplotlib.pyplot as plt
501 import pandas as pd
502
503 required_columns = ['FullTimeHomeGoals', '
      FullTimeAwayGoals', 'FullTimeResult']
504 if not all(col in df.columns for col in
      required_columns):
505     missing = [col for col in required_columns if
      col not in df.columns]
506     raise ValueError(f"Missing required columns: {'
      missing}'")
507
508 df_clean = df.dropna(subset=required_columns).copy()
509
510 df_clean['GoalDifference'] = df_clean['
      FullTimeHomeGoals'] - df_clean['FullTimeAwayGoals'
      ]
511
512 result_mapping = {'H': 1, 'D': 0, 'A': -1}
513 df_clean['MatchOutcomeScore'] = df_clean['
      FullTimeResult'].map(result_mapping)
514
515
516 df_clean = df_clean.dropna(subset=['
      MatchOutcomeScore'])
517 gd_corr = df_clean[['GoalDifference', '
      MatchOutcomeScore']].corr().iloc[0, 1]
518
519
520 print(f"  Correlation between goal difference and
      match result: {gd_corr:.3f}")
```

```
521 if abs(gd_corr) >= 0.7:
522     print("→ Strong correlation: Goal difference
523         is a good predictor of result.")
524 elif abs(gd_corr) >= 0.4:
525     print("→ Moderate correlation: Goal
526         difference has some predictive power.")
527 else:
528     print("→ Weak correlation: Goal difference
529         isn't strongly predictive of result.")
530
531 plt.figure(figsize=(10, 6))
532 plt.scatter(
533     df_clean['GoalDifference'], df_clean['
534     MatchOutcomeScore'],
535     alpha=0.6, edgecolors='k', color='
536     mediumseagreen'
537 )
538 plt.axhline(0, color='gray', linestyle='--', lw=1)
539 plt.axvline(0, color='gray', linestyle='--', lw=1)
540 plt.title('Goal Difference vs Match Outcome Score'
541 , fontsize=14)
542 plt.xlabel('Goal Difference (Home - Away)', fontsize=12)
543 plt.ylabel('Match Outcome Score (H=1, D=0, A=-1)'
544 , fontsize=12)
545 plt.grid(True, linestyle='--', alpha=0.5)
546 plt.tight_layout()
547 plt.show()
548
549 #%% md
550 # Liverpool vs Mancity
551 #%%
552 # Module: [Advanced Analysis5]
553 # Contributor: [Ruoyao Yan]
554 import matplotlib.pyplot as plt
555 import pandas as pd
556
557 df = pd.read_csv(
558     'cleaned_football_data.csv')
```

```
554
555 LvsM = []
556 for index, row in df.iterrows():
557     home = row['HomeTeam']
558     away = row['AwayTeam']
559     if (home == 'Liverpool' and away == 'Man City'
560         ) or (home == 'Man City' and away == 'Liverpool'):
561         LvsM.append(row)
562
563 liverpool_vs_mancity = pd.DataFrame(LvsM)
564
565 total_liverpool_wins = len(liverpool_vs_mancity[((liverpool_vs_mancity['HomeTeam'] == 'Liverpool') & (liverpool_vs_mancity['FullTimeResult'] == 'H')) | ((liverpool_vs_mancity['AwayTeam'] == 'Liverpool') & (liverpool_vs_mancity['FullTimeResult'] == 'A'))])
566
567 total_mancity_wins = len(liverpool_vs_mancity[((liverpool_vs_mancity['HomeTeam'] == 'Man City') & (liverpool_vs_mancity['FullTimeResult'] == 'H')) | ((liverpool_vs_mancity['AwayTeam'] == 'Man City') & (liverpool_vs_mancity['FullTimeResult'] == 'A'))])
568
569 total_draws = len(liverpool_vs_mancity[liverpool_vs_mancity['FullTimeResult'] == 'D'])
570
571 labels = ['Liverpool Wins', 'Man city Wins', 'Draws']
572 counts = [total_liverpool_wins, total_mancity_wins, total_draws]
573
574 plt.figure(figsize=(6, 4))
575 plt.pie(
576     counts,
577     labels=labels,
578     autopct='%1.2f%%',
579     startangle=90,
580     colors=colors,
```

```
581 )
582
583 plt.title('Liverpool vs Man City Match Results',
      fontsize=14)
584 plt.tight_layout()
585 plt.show()
586 #%% md
587 # Prediction
588 #%% md
589 Predict the number of goals scored by Liverpool
as the home team
590 #%%
591 # Module: [Prediction1]
592 # Contributor: [Ruoyao Yan]
593 import pandas as pd
594 from sklearn.model_selection import
train_test_split
595 from sklearn.linear_model import LinearRegression
596 from sklearn.metrics import mean_squared_error,
r2_score
597 import numpy as np
598 from collections import Counter
599 import matplotlib.pyplot as plt
600
601 liverpool_home = df[df['HomeTeam'] == 'Liverpool'
].copy()
602
603 features = ['HomeShotsOnTarget', '
AwayShotsOnTarget']
604 X = liverpool_home[features]
605 y = liverpool_home['FullTimeHomeGoals']
606
607 X_train, X_test, y_train, y_test =
train_test_split(X, y, test_size=0.2, random_state
=42)
608
609 model = LinearRegression()
610 model.fit(X_train, y_train)
611
612 y_pred = model.predict(X_test)
613
```

```
614 print("Predicted goals:", y_pred)
615 y_pred_rounded = np.round(y_pred).astype(int)
616 print("Predicted goals (rounded):", y_pred_rounded)
617 counter = Counter(y_pred_rounded)
618 most_common_pred, frequency = counter.most_common(1)[0]
619 print(f"Most frequently predicted goals: {most_common_pred} (appeared {frequency} times)")
620
621 print("Mean Squared Error (MSE):",
       mean_squared_error(y_test, y_pred))
622 print("R2 score:", r2_score(y_test, y_pred))
623 print("Maximum predicted value:", y_pred.max())
624
625 plt.scatter(y_test, y_pred, alpha=0.5)
626 plt.xlabel("Actual goals")
627 plt.ylabel("Predicted goals")
628 plt.title("Liverpool home goal prediction")
629 plt.plot([min(y_test), max(y_test)], [min(y_test),
                                         max(y_test)], 'r--') # Ideal line
630 plt.ylim(0, 5)
631 plt.show()
632
633 total_goals_method1 = most_common_pred * 19
634 print(f"Season total goals: {total_goals_method1:.1f}")
635


---


636 #%% md
637 Predict Liverpool's specific winning rate in home
   games


---


638 #%%
639 # Module: [Prediction2]
640 # Contributor: [Ruoyao Yan]
641 import pandas as pd
642 import numpy as np
643 from collections import Counter
644 from sklearn.model_selection import
   train_test_split
645 from sklearn.metrics import mean_squared_error,
   r2_score
```

```
646
647 liverpool_home = df[df['HomeTeam'] == 'Liverpool'
648 ].copy()
649 liverpool_home['HomeWin'] = (liverpool_home['
649 FullTimeResult'] == 'H').astype(int)
649 win_rate_by_shots = liverpool_home.groupby('
649 HomeShotsOnTarget')['HomeWin'].mean().reset_index
649 ()
650 win_rate_by_shots.columns = ['HomeShotsOnTarget',
650 'WinRate']
651
652 X = win_rate_by_shots[['HomeShotsOnTarget']]
653 y = win_rate_by_shots['WinRate']
654 from sklearn.linear_model import LinearRegression
655
656 model = LinearRegression()
657 model.fit(X, y)
658
659 X_test = pd.DataFrame(np.arange(0, 15).reshape(-1
659 , 1), columns=['HomeShotsOnTarget'])
660 y_pred = model.predict(X_test)
661 import matplotlib.pyplot as plt
662
663 plt.figure(figsize=(10, 6))
664 plt.scatter(X, y, color='blue', label='
664 Actual_win_rate')
665 plt.plot(X_test, y_pred, color='red', label='
665 Predicted_win_rate')
666 plt.xlabel('Liverpool home shots on target')
667 plt.ylabel('Predicted_win_rate')
668 plt.title('Liverpool home shots on target vs win
668 rate')
669 plt.legend()
670 plt.grid(True)
671 plt.show()
672 next_season_shots = 4
673 predicted_win_rate = model.predict(pd.DataFrame({'
673 HomeShotsOnTarget': [next_season_shots]}))
674 print(f"Projected average home shots on target for
674 Liverpool next season {next_season_shots},
674 Predicted win rate: {predicted_win_rate[0]:.1%}")
```

```
675
676 features = ['HomeShotsOnTarget', '
677     AwayShotsOnTarget']
677 X = liverpool_home[features]
678 y = liverpool_home['FullTimeHomeGoals']
679
680 X_train, X_test, y_train, y_test =
681     train_test_split(X, y, test_size=0.2, random_state
682 =42)
683
684 model = LinearRegression()
685 model.fit(X_train, y_train)
686
687 y_pred = model.predict(X_test)
688
689 print("Predicted goals:", y_pred)
690 y_pred_rounded = np.round(y_pred).astype(int)
691 print("Predicted goals (rounded):", y_pred_rounded)
692 counter = Counter(y_pred_rounded)
693 most_common_pred, frequency = counter.most_common(
694     1)[0]
695 print(f"Most frequently predicted goals: {most_common_pred} (appeared {frequency} times)")
696
697 print("Mean Squared Error (MSE):",
698     mean_squared_error(y_test, y_pred))
699 print("R2 score:", r2_score(y_test, y_pred))
700 print("Maximum predicted value:", y_pred.max())
701 #%% md
702 # Predicting goals using cards and corner kicks
703 #%% md
704 # Module: [Prediction3]
705 # Contributor: [Ruoyao Yan] and [Jun Sun]
706 import pandas as pd
707 import numpy as np
708 import matplotlib.pyplot as plt
709 import seaborn as sns
710
711 from sklearn.model_selection import
```

```
708 train_test_split
709 from sklearn.linear_model import LinearRegression
710 from sklearn.metrics import mean_squared_error,
    r2_score, mean_absolute_error
711
712
713 df = pd.read_csv('cleaned_football_data.csv')
714 df['TotalGoals'] = df['FullTimeHomeGoals'] + df['
    FullTimeAwayGoals']
715 df['TotalCorners'] = df['HomeCorners'] + df['
    AwayCorners']
716 df['TotalYellowCards'] = df['HomeYellowCards'] +
    df['AwayYellowCards']
717 df['TotalRedCards'] = df['HomeRedCards'] + df['
    AwayRedCards']
718
719 features = ['TotalCorners', 'TotalYellowCards', '
    TotalRedCards']
720 target = 'TotalGoals'
721
722 X = df[features]
723 y = df[target]
724
725 X_train, X_test, y_train, y_test =
    train_test_split(
726     X, y, test_size=0.2, random_state=42
727 )
728
729 model = LinearRegression()
730 model.fit(X_train, y_train)
731
732 y_pred = model.predict(X_test)
733
734 mse = mean_squared_error(y_test, y_pred)
735 mae = mean_absolute_error(y_test, y_pred)
736 r2 = r2_score(y_test, y_pred)
737
738 print(f"MSE: {mse:.2f}")
739 print(f"MAE: {mae:.2f}")
740 print(f"R2: {r2:.2f}")
741
```

```
742 coef_df = pd.DataFrame({  
743     'Feature': features,  
744     'Coefficient': model.coef_,  
745     'AbsoluteEffect': np.abs(model.coef_)  
746 }).sort_values(by='AbsoluteEffect', ascending=  
    False)  
747  
748 print("\nThe impact of each feature on the total  
    number of goals (sorted by absolute value):")  
749 print(coef_df)  
750  
751  
752 plt.figure(figsize=(18, 5))  
753  
754 for i, feature in enumerate(features):  
755     plt.subplot(1, 3, i+1)  
756     sns.regplot(x=feature, y='TotalGoals', data=df  
        , scatter_kws={'alpha': 0.4})  
757     plt.title(f"{feature} vs TotalGoals")  
758     plt.xlabel(feature)  
759     plt.ylabel("TotalGoals")  
760  
761 plt.tight_layout()  
762 plt.show()  
763  
764 # example  
765 new_data = pd.DataFrame([[12, 4, 1]], columns=  
    features) # 12角球、4黄牌、1红牌  
766 predicted_goal = model.predict(new_data)[0]  
767 print(f"\nPredict the number of goals in the match  
    : {predicted_goal:.1f}")
```