
BỘ THÔNG TIN VÀ TRUYỀN THÔNG
HỌC VIỆN CÔNG NGHỆ BƯU CHÍNH VIỄN THÔNG



BÀI BÁO CÁO

MÔN: TOÁN RỜI RẠC 2

Giảng viên: Kim Ngọc Bách
Tên sinh viên: Bùi Quang Thắng
Mã sinh viên: B23DCCN751
Lớp: D23CQCE06 - B
Môn: Lập trình Python

Hà Nội, Tháng 5/2025

Lời mở đầu

Trong bối cảnh công nghệ thông tin và khoa học dữ liệu ngày càng phát triển, Python đã khẳng định vị thế là một trong những ngôn ngữ lập trình mạnh mẽ và phổ biến, đặc biệt trong việc xử lý dữ liệu, phân tích thống kê và học máy. Bài tập lớn này được thực hiện nhằm đáp ứng các yêu cầu của Assignment 1, môn Lập trình Python, với mục tiêu áp dụng kiến thức lý thuyết vào việc thu thập, phân tích và trực quan hóa dữ liệu thống kê của các cầu thủ tham gia mùa giải English Premier League 2024-2025.

Dự án tập trung vào việc thu thập dữ liệu từ các nguồn uy tín, phân tích các chỉ số hiệu suất, phân loại cầu thủ bằng thuật toán học máy, và đề xuất phương pháp định giá cầu thủ dựa trên giá trị chuyển nhượng. Thông qua quá trình thực hiện, chúng tôi không chỉ củng cố kỹ năng lập trình mà còn rèn luyện khả năng tư duy phân tích, giải quyết vấn đề thực tiễn và trình bày kết quả một cách khoa học.

Chúng tôi xin gửi lời cảm ơn chân thành đến giảng viên đã hướng dẫn tận tình, cùng các nguồn tài liệu tham khảo từ <https://foref.com> và <https://www.footballtransfers.com> đã hỗ trợ chúng tôi hoàn thành bài tập. Mọi ý kiến đóng góp sẽ là động lực quý báu để chúng tôi cải thiện và phát triển trong tương lai.

Trân trọng,

Bài 1:

Mã Nguồn Đầy Đủ

Dưới đây là mã nguồn đầy đủ, được định dạng bằng gói listings để dễ đọc:

```
1 # =====
2 # 1. SETUP ON COLAB / LOCAL MACHINE
3 # =====
4 !pip install selenium
5 !apt-get update -y
6 !apt install -y chromium-chromedriver
7 !cp /usr/lib/chromium-browser/chromedriver /usr/bin
8
9 # =====
10 # 2. IMPORT LIBRARIES
11 # =====
12 import time
13 import pandas as pd
14 import re
15 from bs4 import BeautifulSoup
16 from selenium import webdriver
17 from selenium.webdriver.chrome.options import Options
18
19 # =====
20 # 3. CONFIGURE SELENIUM CHROME
21 # =====
22 options = Options()
23 options.add_argument('--headless')
24 options.add_argument('--no-sandbox')
25 options.add_argument('--disable-dev-shm-usage')
26 driver = webdriver.Chrome(options=options)
27
28 # =====
29 # 4. DEFINE TABLES TO SCRAPE
30 # =====
31 tables = [
32     {
33         "url":
34             "https://fbref.com/en/comps/9/stats/Premier-League-Stats",
35         "id": "stats_standard",
36         "cols": ["Nation", "Squad", "Position", "Current age",
37                 "Matches Played", "Minutes",
38                 "Goals", "Assists", "Yellow Cards", "Red Cards",
39                 "xG: Expected Goals", "xAG: Exp. Assisted Goals",
40                 "Progressive Carries", "Progressive Passes",
41                 "Progressive Passes Rec",
42                 "Goals/90", "Assists/90", "xG/90", "xAG/90"]
43     },
44     {
45         "url":
```

```

43         "https://fbref.com/en/comps/9/keepers/Premier-League-Stats",
44         "id": "stats_keeper",
45         "cols": ["Goals Against/90", "Save Percentage", "Clean
46             Sheet Percentage", "Save% (Penalty Kicks)"]
47     },
48     {
49         "url":
50             "https://fbref.com/en/comps/9/shooting/Premier-League-Stats",
51         "id": "stats_shooting",
52         "cols": ["Shots on Target %", "Shots on target/90",
53             "Goals/Shot", "Average Shot Distance"]
54     },
55     {
56         "url":
57             "https://fbref.com/en/comps/9/passing/Premier-League-Stats",
58         "id": "stats_passing",
59         "cols": ["Passes Completed", "Pass Completion %", "Total
60             Passing Distance",
61             "Passes Completed (Short)", "Pass Completion %
62             (Medium)", "Pass Completion % (Long)",
63             "Key Passes", "Passes into Final Third", "Passes
64             into Penalty Area",
65             "Crosses into Penalty Area", "Progressive
66             Passes"]
67     },
68     {
69         "url":
70             "https://fbref.com/en/comps/9/gca/Premier-League-Stats",
71         "id": "stats_gca",
72         "cols": ["Shot-Creating Actions", "Shot-Creating
73             Actions/90",
74             "Goal-Creating Actions", "Goal-Creating
75             Actions/90"]
76     },
77     {
78         "url":
79             "https://fbref.com/en/comps/9/defense/Premier-League-Stats",
80         "id": "stats_defense",
81         "cols": ["Tackles", "Tackles Won", "Dribbles Challenged",
82             "Challenges Lost",
83             "Blocks", "Shots Blocked", "Passes Blocked",
84             "Interceptions"]
85     },
86     {
87         "url":
88             "https://fbref.com/en/comps/9/possession/Premier-League-Stats",
89         "id": "stats_possession",
90         "cols": ["Touches", "Touches (Def Pen)", "Touches (Def
91             3rd)", "Touches (Mid 3rd)",
92             "Touches (Att 3rd)", "Touches (Att Pen)",
93             "Take-Ons Attempted",

```

```

76         "Successful Take-On %", "Tackled During Take-On
77         Percentage",
78         "Carries", "Progressive Carrying Distance",
79         "Progressive Carries",
80         "Carries into Final Third", "Carries into
81         Penalty Area",
82         "Miscontrols", "Dispossessed", "Passes
83         Received", "Progressive Passes Rec"]
84     },
85     {
86         "url":
87             "https://fbref.com/en/comps/9/misc/Premier-League-Stats",
88         "id": "stats_misc",
89         "cols": ["Yellow Cards", "Red Cards", "Fouls Committed",
90                 "Fouls Drawn",
91                 "Offsides", "Crosses", "Ball Recoveries",
92                 "Aerials Won",
93                 "Aerials Lost", "% of Aerials Won"]
94     }
95 ]
96
97 # =====
98 # 5. FUNCTION TO SCRAPE A TABLE
99 # =====
100 def scrape_table(url, table_id, required_cols):
101     print(f"    Loading table {table_id}    ")
102     driver.get(url)
103     time.sleep(3)
104     soup = BeautifulSoup(driver.page_source, 'html.parser')
105     tbl = soup.find('table', {'id': table_id})
106     if tbl is None:
107         print(f"    Table {table_id} not found!")
108         return None
109
110     # Header
111     header_cells =
112         tbl.find('thead').find_all('tr')[-1].find_all('th')
113     headers = []
114     for th in header_cells:
115         labels = th.get('aria-label') or th.text.strip()
116         headers.append(labels)
117     headers[0] = 'Player'
118
119     # Data
120     rows = []
121     for tr in tbl.find('tbody').find_all('tr'):
122         if tr.get('class') == ['thead']:
123             continue
124         cols = tr.find_all(['th', 'td'])
125         if len(cols) != len(headers):
126             continue

```

```

119     vals = []
120     for i, cell in enumerate(cols):
121         txt = cell.get_text(strip=True)
122         if i == 0:
123             txt = re.sub(r'^[\d\.\s]+', '', re.sub(r'[()],' ,
124             txt)).strip()
125             vals.append(txt)
126         rows.append(vals)
127
128     df = pd.DataFrame(rows, columns=headers).set_index('Player')
129     df = df.add_prefix(f"{table_id}_")
130
131     # Select required columns
132     want = [f"{table_id}_{c}" for c in required_cols]
133     have = [c for c in want if c in df.columns]
134     missing = set(want) - set(have)
135     if missing:
136         print(f"      Table {table_id} missing columns:
137             {sorted(missing)}")
138         return df[have]
139
140     # =====
141     # 6. SCRAPE ALL TABLES
142     # =====
143     dfs = {}
144     for t in tables:
145         df = scrape_table(t["url"], t["id"], t["cols"])
146         dfs[t["id"]] = df
147
148     driver.quit()
149
150     # =====
151     # 7. MERGE DATA
152     # =====
153     std = dfs['stats_standard']
154     keepers = dfs['stats_keeper']
155     others = [df for key, df in dfs.items() if key not in
156             ['stats_standard', 'stats_keeper']]
157
158     merged = std.copy()
159     for o in others:
160         merged = merged.join(o, how='left')
161
162     # GK: Merge keeper stats for goalkeepers
163     if keepers is not None:
164         is_gk = merged['stats_standard_Position'].str.contains('GK',
165             na=False)
166         for col in keepers.columns:
167             merged[col] = pd.NA
168             merged.loc[is_gk, col] =
169                 keepers.reindex(merged.index).loc[is_gk, col]

```

```

165
166 # =====
167 # 8. CLEAN PLAYER COLUMN
168 # =====
169 merged = merged.reset_index()
170 merged['Player'] = merged['Player'].apply(lambda x: x[1] if
171     isinstance(x, tuple) else x)
172 player_col = merged.pop('Player')
173 merged.insert(0, 'Player', player_col)
174
175 # =====
176 # 9. FILTER >90 MINUTES
177 # =====
178 merged['stats_standard_Minutes'] = (
179     merged['stats_standard_Minutes']
180     .str.replace(',', '', regex=False)
181     .astype(float, errors='ignore')
182 )
183 merged = merged[ merged['stats_standard_Minutes'] > 90 ]
184
185 # =====
186 # 10. REMOVE DUPLICATE PLAYERS
187 # =====
188 merged = merged.drop_duplicates(subset=['Player'])
189
190 # =====
191 # 11. REPLACE EMPTY VALUES WITH 'N/a' (BUT KEEP 0)
192 # =====
193 merged = merged.replace(['', ''], pd.NA)
194 merged = merged.fillna('N/a')
195
196 # =====
197 # 12. SORT AND SAVE TO CSV
198 # =====
199 merged = merged.sort_values('Player')
200 print(f"Total players: {len(merged)}")
201 merged.to_csv('results.csv', index=False)
202 print(" Results saved to results.csv")

```

Phân Tích Chi Tiết Mã Nguồn

Mã nguồn được chia thành 12 phần chính, mỗi phần có chức năng riêng. Dưới đây là phân tích chi tiết:

1. Cài đặt trên Colab / Máy Local

Phần này cài đặt các thư viện và công cụ cần thiết:

- `!pip install selenium`: Cài đặt thư viện selenium để tự động hóa trình duyệt.

-
- `!apt-get update -y` và `!apt install -y chromium-chromedriver`: Cập nhật danh sách gói và cài đặt ChromeDriver.
 - `!cp /usr/lib/chromium-browser/chromedriver /usr/bin`: Sao chép ChromeDriver vào thư mục hệ thống.

2. Import Library

Nhập các thư viện cần thiết:

- `time`: Để thêm độ trễ trong quá trình tải trang web.
- `pandas`: Để xử lý dữ liệu dạng bảng.
- `re`: Để sử dụng biểu thức chính quy làm sạch dữ liệu.
- `BeautifulSoup`: Để phân tích HTML.
- `selenium`: Để điều khiển trình duyệt Chrome.

3. Cấu hình Selenium Chrome

Cấu hình trình duyệt Chrome chạy ở chế độ không giao diện (headless):

- Sử dụng `Options()` để tùy chỉnh Chrome.
- Các tham số như `-headless`, `-no-sandbox`, và `-disable-dev-shm-usage` đảm bảo hoạt động ổn định trên Colab.

4. Định nghĩa các bảng cần lấy

Xác định danh sách các bảng để thu thập dữ liệu:

- Mỗi bảng được định nghĩa bằng một từ điển với `url`, `id`, và `cols`.
- Ví dụ: Bảng `stats_standard` lấy các cột như "Goals", "Assists", "Minutes", v.v.

5. Hàm scrape 1 bảng

Hàm `scrape_table` thu thập dữ liệu từ một bảng:

- Sử dụng Selenium để tải trang và BeautifulSoup để phân tích HTML.
- Trích xuất tiêu đề và dữ liệu, làm sạch tên cầu thủ bằng regex.
- Tạo DataFrame và lọc các cột cần thiết.

6. Scrape toàn bộ

Lặp qua danh sách `tables` để thu thập tất cả các bảng và lưu vào từ điển `dfs`.

7. Gộp dữ liệu

Gộp các DataFrame:

- Bắt đầu với `stats_standard`, sau đó gộp các bảng khác bằng `join`.
- Gộp dữ liệu thủ môn (`stats_keeper`) chỉ cho các cầu thủ có vị trí GK.

8. Chỉnh lại cột Player

Làm sạch cột `Player` và di chuyển nó lên đầu DataFrame.

9. Filter >90 phút

Lọc các cầu thủ có thời gian thi đấu trên 90 phút.

10. Xóa trùng Player

Loại bỏ các cầu thủ trùng lặp dựa trên cột `Player`.

11. Thay các giá trị rỗng bằng 'N/a'

Xử lý các giá trị rỗng hoặc không hợp lệ, thay bằng 'N/a'.

12. Sắp xếp và lưu CSV

Sắp xếp theo tên cầu thủ và lưu kết quả vào `results.csv`.

Bài 2:

Mã Nguồn Đầy Đủ

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import os
5 import re
6 import random
7
8 # Read data from CSV file
9 data = pd.read_csv('results.csv')
10
11 # Replace "N/a" with NaN
12 data.replace('N/a', np.nan, inplace=True)
13
14 # Identify columns intended to be numeric, excluding irrelevant
   columns
15 intended_numeric_cols = [
16     col for col in data.columns
```

```

17     if col != 'stats_standard_Current age'
18         and col not in ['Player', 'stats_standard_Nation',
19             'stats_standard_Squad', 'stats_standard_Position']
20 ]
21 # Convert intended numeric columns to numeric type, coercing
    errors to NaN
22 for col in intended_numeric_cols:
23     data[col] = pd.to_numeric(data[col], errors='coerce')
24
25 # Filter to keep only fully numeric columns
26 numeric_cols = [col for col in intended_numeric_cols if
    data[col].dtype in ['float64', 'int64']]
27
28 # Fill NaN with 0 for numeric columns
29 data[numeric_cols] = data[numeric_cols].fillna(0)
30
31 # Categorize attack and defense metrics
32 attack_keywords = ['Goals', 'Assists', 'Shots', 'xG', 'xAG']
33 defense_keywords = ['Tackles', 'Interceptions', 'Blocks',
    'Clearances']
34 attack_cols = [col for col in numeric_cols if any(keyword in col
    for keyword in attack_keywords)]
35 defense_cols = [col for col in numeric_cols if any(keyword in col
    for keyword in defense_keywords)]
36
37 # Randomly select 3 metrics from each category (or all if fewer
    than 3)
38 selected_attack_cols = random.sample(attack_cols, min(3,
    len(attack_cols)))
39 selected_defense_cols = random.sample(defense_cols, min(3,
    len(defense_cols)))
40
41 # 1. Identify top 3 highest and lowest for each numeric metric
42 with open('top_3.txt', 'w', encoding='utf-8') as f:
43     for col in numeric_cols:
44         f.write(f"Metric: {col}\n")
45         f.write("Top 3 highest:\n")
46         top_3 = data[['Player',
47             col]].dropna().sort_values(by=col,
48             ascending=False).head(3)
49         for _, row in top_3.iterrows():
50             f.write(f"- {row['Player']}: {row[col]:.2f}\n")
51         f.write("Top 3 lowest:\n")
52         bottom_3 = data[['Player',
53             col]].dropna().sort_values(by=col).head(3)
54         for _, row in bottom_3.iterrows():
55             f.write(f"- {row['Player']}: {row[col]:.2f}\n")
56         f.write("\n")

```

```

55 # 2. Calculate median, mean, and standard deviation for all
    players
56 all_stats = pd.DataFrame({
57     '': ['all'],
58     'Team': ['all']
59 })
60 for col in numeric_cols:
61     all_stats[f'Median of {col}'] = [data[col].median()]
62     all_stats[f'Mean of {col}'] = [data[col].mean()]
63     all_stats[f'Std of {col}'] = [data[col].std()]
64
65 # Calculate statistics per team
66 team_stats =
67     data.groupby('stats_standard_Squad')[numeric_cols].agg(['median',
68         'mean', 'std']).reset_index()
69
70 # Flatten MultiIndex for columns
71 flattened_columns = [('Team', '')]
72 for col in numeric_cols:
73     flattened_columns.extend([(col, 'median'), (col, 'mean'),
74         (col, 'std')])
75 team_stats.columns = pd.MultiIndex.from_tuples(flattened_columns)
76
77 # Rename columns
78 team_stats.columns = [f'Median of {col[0]}' if col[1] == 'median'
79     else
80         f'Mean of {col[0]}' if col[1] == 'mean' else
81         f'Std of {col[0]}' if col[1] == 'std' else
82         'Team' for col in team_stats.columns]
83
84 # Add index column
85 team_stats.insert(0, '', range(1, len(team_stats) + 1))
86
87 # Ensure matching columns and combine
88 common_columns =
89     all_stats.columns.intersection(team_stats.columns)
90 team_stats = team_stats.reindex(columns=common_columns,
91     fill_value=np.nan)
92 result_stats = pd.concat([all_stats, team_stats],
93     ignore_index=True)
94 result_stats.to_csv('results2.csv', index=False)
95
96 # 3. Identify team with the highest score for each metric
97 team_totals =
98     data.groupby('stats_standard_Squad')[numeric_cols].sum()
99 top_teams = team_totals.idxmax()
100
101 print("Teams with the highest score for each metric:")
102 for col, team in top_teams.items():
103     print(f"- {col}: {team}")

```

```

97 # Evaluate the best performing team based on key metrics
98 key_stats = [
99     'stats_standard_Goals', 'stats_standard_Assists',
100     'stats_standard_xG: Expected Goals', 'stats_standard_xAG:
        Exp. Assisted Goals',
101     'stats_keeper_Save Percentage'
102 ]
103 key_stats = [stat for stat in key_stats if stat in numeric_cols]
104 team_means =
105     data.groupby('stats_standard_Squad')[key_stats].mean()
106 team_scores = team_means.sum(axis=1)
107 best_team = team_scores.idxmax()
108 best_score = team_scores.max()
109 print(f"\nBest performing team: {best_team}")
110 print(
111     f"Reason: {best_team} leads with a total average score of
        {best_score:.2f} across key metrics (goals, assists, xG,
        xAG, Save%). "
112     f"They demonstrate consistent offensive and defensive
        performance.")
113
114 # 4. Plot histograms
115 # Create histograms directory if it doesn't exist
116 if not os.path.exists('histograms'):
117     os.makedirs('histograms')
118
119 # Function to sanitize filenames
120 def sanitize_filename(name):
121     return re.sub(r'[^w\-\_\.\ ]', '_', name)
122
123 # Plot histogram for all players
124 plt.figure(figsize=(15, 10))
125 plt.suptitle("Distribution of Selected Metrics for All Players",
126             fontsize=16)
127 for i, col in enumerate(selected_attack_cols, 1):
128     plt.subplot(2, 3, i)
129     plt.hist(data[col].dropna(), bins=20, color='blue', alpha=0.7)
130     plt.title(f'Attack: {col}')
131     plt.xlabel(col)
132     plt.ylabel('Count')
133
134 for i, col in enumerate(selected_defense_cols, 1):
135     plt.subplot(2, 3, i + 3)
136     plt.hist(data[col].dropna(), bins=20, color='green',
137             alpha=0.7)
138     plt.title(f'Defense: {col}')
139     plt.xlabel(col)
140     plt.ylabel('Count')
141
142 plt.tight_layout(rect=[0, 0, 1, 0.95])

```

```

141 plt.savefig('histograms/hist_all_players.png')
142 plt.show()
143
144 # Plot histograms for each team
145 teams = data['stats_standard_Squad'].unique()
146 for team in teams:
147     team_data = data[data['stats_standard_Squad'] == team]
148     safe_team = sanitize_filename(team)
149
150     plt.figure(figsize=(15, 10))
151     plt.suptitle(f"Distribution of Selected Metrics for {team}",
152                fontsize=16)
153     for i, col in enumerate(selected_attack_cols, 1):
154         plt.subplot(2, 3, i)
155         plt.hist(team_data[col].dropna(), bins=20, color='blue',
156                alpha=0.7)
157         plt.title(f'Attack: {col}')
158         plt.xlabel(col)
159         plt.ylabel('Count')
160
161     for i, col in enumerate(selected_defense_cols, 1):
162         plt.subplot(2, 3, i + 3)
163         plt.hist(team_data[col].dropna(), bins=20, color='green',
164                alpha=0.7)
165         plt.title(f'Defense: {col}')
166         plt.xlabel(col)
167         plt.ylabel('Count')
168
169     plt.tight_layout(rect=[0, 0, 1, 0.95])
170     plt.savefig(f'histograms/hist_{safe_team}.png')
171     plt.show()

```

Phân Tích Chi Tiết Mã Nguồn

Mã nguồn này được chia thành các phần chính, mỗi phần thực hiện một nhiệm vụ cụ thể. Dưới đây là phân tích chi tiết:

1. Nhập thư viện và đọc dữ liệu

Nhập các thư viện cần thiết và đọc dữ liệu từ tệp CSV:

- `pandas`, `numpy`: Để xử lý dữ liệu.
- `matplotlib.pyplot`: Để vẽ biểu đồ.
- `os`, `re`, `random`: Để quản lý tệp, xử lý chuỗi, và chọn ngẫu nhiên.
- Đọc dữ liệu từ `results.csv` vào `DataFrame data`.

2. Làm sạch dữ liệu

Thực hiện các bước làm sạch dữ liệu:

- Thay thế "N/a" bằng `np.nan`.
- Xác định các cột nên là số (loại trừ các cột như `Player`, `Nation`, v.v.).
- Chuyển đổi các cột số sang kiểu số, ép lỗi thành `NaN`.
- Lọc các cột hoàn toàn là số và điền `NaN` bằng 0.

3. Phân loại chỉ số tấn công và phòng thủ

Phân loại các cột số thành chỉ số tấn công và phòng thủ:

- Chỉ số tấn công: Các cột chứa từ khóa như `Goals`, `Assists`, `xG`.
- Chỉ số phòng thủ: Các cột chứa từ khóa như `Tackles`, `Interceptions`.
- Chọn ngẫu nhiên tối đa 3 cột từ mỗi loại để phân tích.

4. Xác định top 3 cao nhất và thấp nhất

Tìm top 3 cầu thủ có giá trị cao nhất và thấp nhất cho mỗi chỉ số số:

- Sắp xếp dữ liệu theo từng cột số.
- Ghi kết quả vào tệp `top_3.txt`.

5. Tính thống kê cho tất cả cầu thủ và theo đội

Tính trung vị, trung bình, và độ lệch chuẩn:

- Tính cho tất cả cầu thủ và lưu vào `all_stats`.
- Tính theo từng đội (`stats_standard_Squad`) và lưu vào `team_stats`.
- Làm phẳng `MultiIndex`, đổi tên cột, và gộp dữ liệu vào `results2.csv`.

6. Xác định đội có điểm số cao nhất cho mỗi chỉ số

Tính tổng các chỉ số theo đội và xác định đội có giá trị cao nhất cho mỗi chỉ số.

`stats_standard_Matches Played`: Southampton
`stats_standard_Minutes`: Leicester City
`stats_standard_Goals`: Liverpool
`stats_standard_Assists`: Liverpool
`stats_standard_Yellow Cards`: Chelsea
`stats_standard_Red Cards`: Arsenal
`stats_standard_xG (Expected Goals)`: Liverpool
`stats_standard_xAG (Exp. Assisted Goals)`: Liverpool
`stats_standard_Progressive Carries`: Manchester City

stats_standard_Progressive Passes: Liverpool
stats_standard_Progressive Passes Rec: Liverpool
stats_standard_Goals/90: Manchester City
stats_standard_Assists/90: Chelsea
stats_standard_xG/90: Aston Villa
stats_standard_xAG/90: Chelsea
stats_shooting-Shots on Target %: Tottenham
stats_shooting-Shots on target/90: Everton
stats_shooting_Goals/Shot: Manchester City
stats_shooting_Average Shot Distance: Southampton
stats_passing_Passes Completed: Manchester City
stats_passing_Pass Completion %: Southampton
stats_passing_Total Passing Distance: Manchester City
stats_passing_Passes Completed (Short): Manchester City
stats_passing_Pass Completion % (Medium): Southampton
stats_passing_Pass Completion % (Long): Manchester Utd
stats_passing_Key Passes: Liverpool
stats_passing_Passes into Final Third: Manchester City
stats_passing_Passes into Penalty Area: Arsenal
stats_passing_Crosses into Penalty Area: Bournemouth
stats_passing_Progressive Passes: Liverpool
stats_gca_Shot-Creating Actions: Liverpool
stats_gca_Shot-Creating Actions/90: Brighton
stats_gca_Goal-Creating Actions: Liverpool
stats_gca_Goal-Creating Actions/90: Brighton
stats_defense_Tackles: Manchester Utd
stats_defense_Tackles Won: Manchester Utd
stats_defense_Dribbles Challenged: Manchester Utd
stats_defense_Challenges Lost: Brighton
stats_defense_Blocks: Leicester City
stats_defense-Shots Blocked: Brentford
stats_defense_Passes Blocked: Crystal Palace
stats_defense_Interceptions: Manchester Utd
stats_possession_Touches: Manchester City
stats_possession_Touches (Def Pen): Brentford
stats_possession_Touches (Def 3rd): Manchester Utd
stats_possession_Touches (Mid 3rd): Manchester City
stats_possession_Touches (Att 3rd): Manchester City
stats_possession_Touches (Att Pen): Manchester City
stats_possession_Take-Ons Attempted: Tottenham
stats_possession_Successful Take-On %: Aston Villa
stats_possession_Tackled During Take-On %: Manchester Utd
stats_possession_Carries: Manchester City
stats_possession_Progressive Carrying Distance: Manchester City
stats_possession_Progressive Carries: Manchester City
stats_possession_Carries into Final Third: Manchester City
stats_possession_Carries into Penalty Area: Manchester City
stats_possession_Miscontrols: Bournemouth

stats_possession_Dispossessed: Newcastle Utd
stats_possession_Passes Received: Manchester City
stats_possession_Progressive Passes Rec: Liverpool
stats_misc_Yellow Cards: Chelsea
stats_misc_Red Cards: Arsenal
stats_misc_Fouls Committed: Bournemouth
stats_misc_Fouls Drawn: Tottenham
stats_misc_Offsides: Nott'ham Forest
stats_misc_Crosses: Fulham
stats_misc_Ball Recoveries: Bournemouth
stats_misc_Aerials Won: Brentford
stats_misc_Aerials Lost: Crystal Palace
stats_misc_% of Aerials Won: Southampton
stats_keeper_Goals Against/90: Leicester City
stats_keeper_Save Percentage: Bournemouth
stats_keeper_Clean Sheet Percentage: Brentford
stats_keeper_Save% (Penalty Kicks): Everton

7. Đánh giá đội thể hiện tốt nhất

Dựa trên các chỉ số chính (Goals, Assists, xG, xAG, Save Percentage):

- Tính trung bình các chỉ số theo đội.
- Xác định đội có tổng điểm trung bình cao nhất.

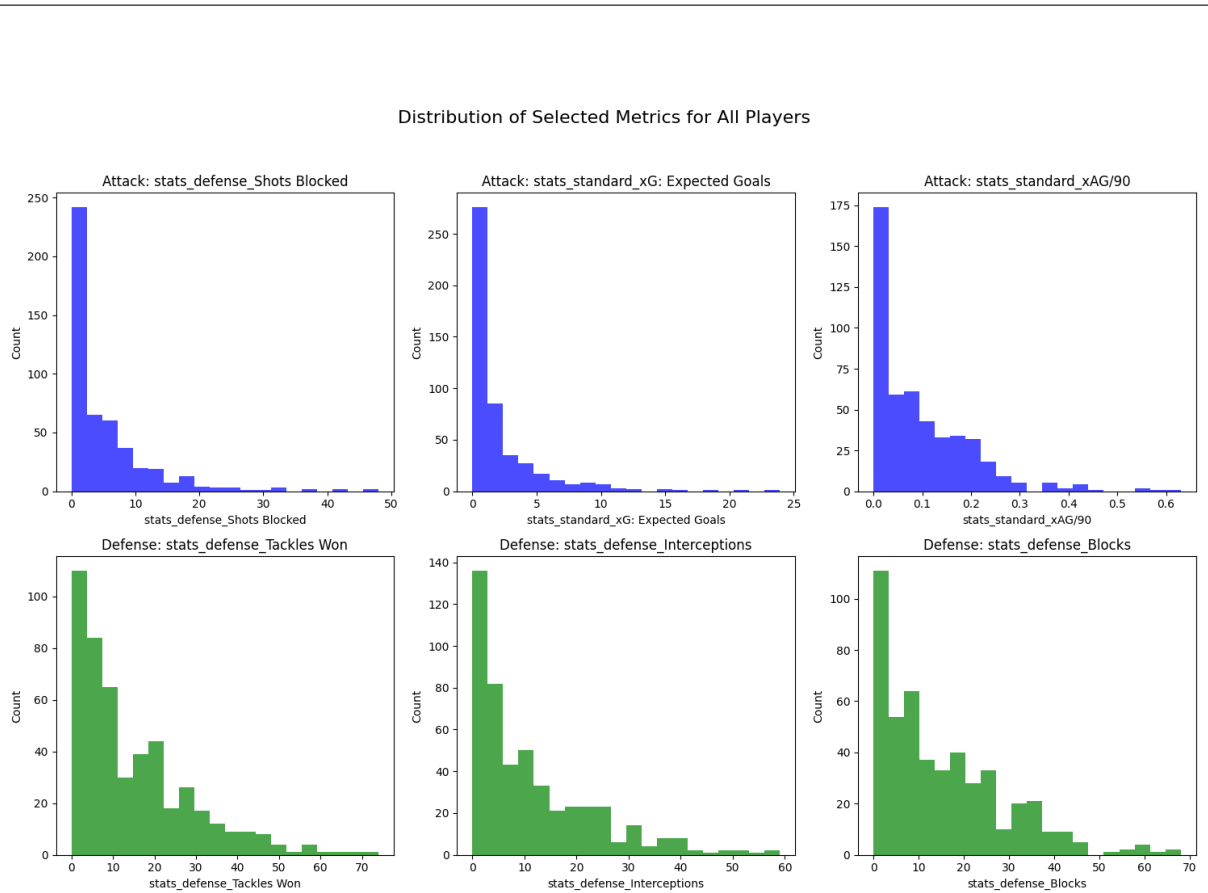
Best Performing Team: Liverpool

Reason: Liverpool leads with a total average score of **19.47** across key metrics including *goals*, *assists*, *expected goals (xG)*, *expected assisted goals (xAG)*, and *save percentage (Save%)*. This demonstrates Liverpool's consistent excellence in both offensive and defensive performance.

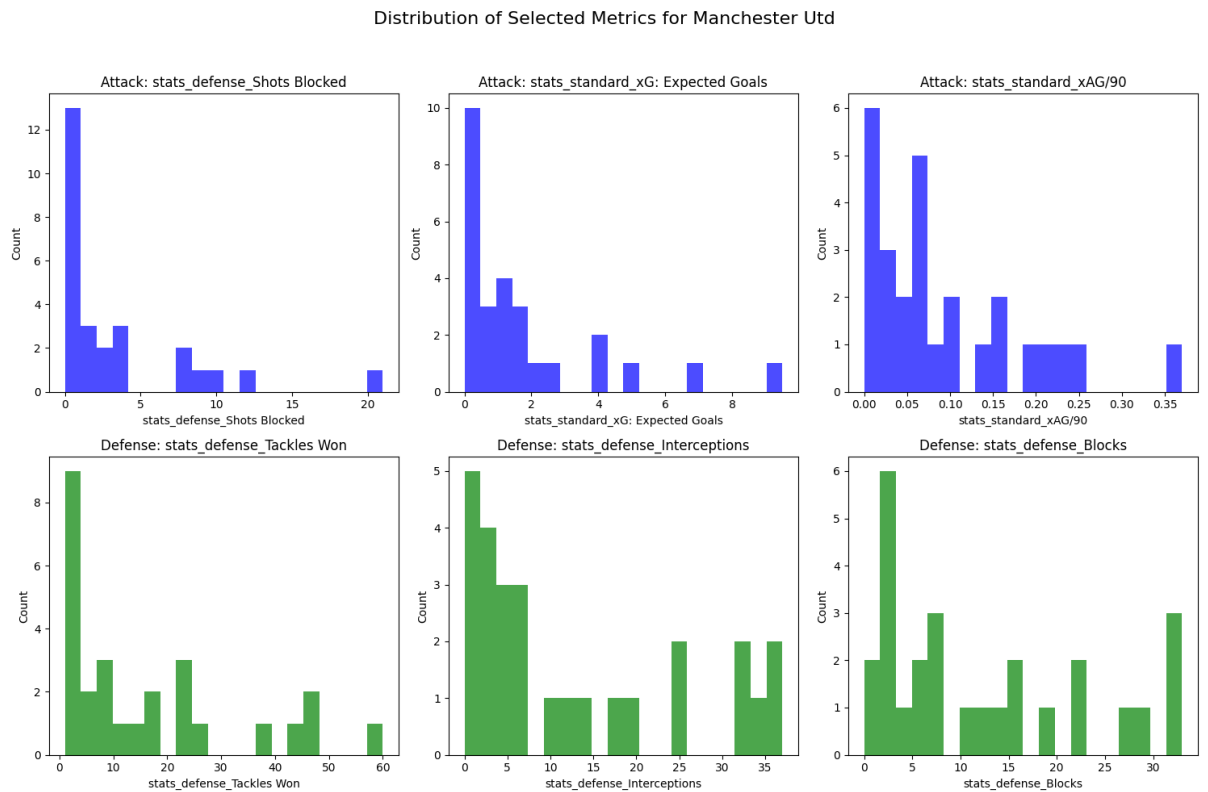
8. Vẽ biểu đồ histogram

Vẽ biểu đồ histogram cho các chỉ số tấn công và phòng thủ:

- Tạo thư mục `histograms` nếu chưa tồn tại.
- Hàm `sanitize_filename` làm sạch tên tệp.
- Vẽ histogram cho tất cả cầu thủ và lưu vào `histograms/hist_all_players.png`.
- Vẽ histogram cho từng đội và lưu vào các tệp riêng (ví dụ: `histograms/hist_[team].png`).



Hình 1: Histogram cho tất cả cầu thủ



Hình 2: Histogram cho cầu thủ của Manchester Utd

Bài 3:

Mã nguồn đầy đủ

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.cluster import KMeans
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.decomposition import PCA
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8
9 # Load data from 'results.csv' file
10 data = pd.read_csv('results.csv')
11
12 # Select features for clustering
13 features = [
14     'stats_passing_Passes Completed', # Passing control
15     'stats_standard_Goals',           # Number of goals
16     'stats_defense_Tackles',          # Tackling ability
17     'stats_standard_Assists',         # Number of assists
18     'stats_shooting-Shots on target/90' # Shots on target per 90
19     minutes
20 ]
21
22 # Extract features and handle missing values
23 data_features = data[features].fillna(0)
24
25 # Standardize the data
26 scaler = StandardScaler()
27 data_scaled = scaler.fit_transform(data_features)
28
29 # Elbow method to determine optimal number of clusters
30 inertia = []
31 k_range = range(1, 11)
32 for k in k_range:
33     kmeans = KMeans(n_clusters=k, random_state=42)
34     kmeans.fit(data_scaled)
35     inertia.append(kmeans.inertia_)
36
37 # Plot the Elbow curve
38 plt.figure(figsize=(10, 5))
39 plt.subplot(1, 2, 1)
40 plt.plot(k_range, inertia, marker='o')
41 plt.title('Elbow Method')
42 plt.xlabel('Number of clusters (k)')
43 plt.ylabel('Inertia')
44 plt.grid(True)
45 plt.savefig('elbow_plot.png')
```

```

46 # Assume optimal number of clusters is 4 (to be confirmed from
    the plot)
47 optimal_k = 4
48 kmeans = KMeans(n_clusters=optimal_k, random_state=42)
49 clusters = kmeans.fit_predict(data_scaled)
50
51 # Reduce dimensions using PCA to 2D
52 pca = PCA(n_components=2)
53 data_pca = pca.fit_transform(data_scaled)
54
55 # Create DataFrame for PCA data
56 df_pca = pd.DataFrame(data_pca, columns=['PC1', 'PC2'])
57 df_pca['Cluster'] = clusters
58
59 # Plot 2D clustering visualization
60 plt.subplot(1, 2, 2)
61 sns.scatterplot(x='PC1', y='PC2', hue='Cluster', data=df_pca,
    palette='viridis', s=100)
62 plt.title('2D Clustering Visualization')
63 plt.xlabel('Principal Component 1 (PC1)')
64 plt.ylabel('Principal Component 2 (PC2)')
65 plt.savefig('cluster_plot.png')
66
67 # Save the combined plots
68 plt.tight_layout()
69 plt.savefig('combined_plots.png')

```

Phân tích chi tiết từng bước

Bước 1: Nhập các thư viện cần thiết

Đoạn mã bắt đầu bằng việc nhập các thư viện:

- pandas và numpy: Dùng để xử lý và tính toán dữ liệu.
- sklearn.cluster.KMeans: Thư viện cung cấp thuật toán K-Means để phân cụm.
- sklearn.preprocessing.StandardScaler: Chuẩn hóa dữ liệu để đảm bảo các đặc trưng có cùng thang đo.
- sklearn.decomposition.PCA: Giảm chiều dữ liệu để trực quan hóa.
- matplotlib.pyplot và seaborn: Dùng để vẽ biểu đồ.

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.cluster import KMeans
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.decomposition import PCA
6 import matplotlib.pyplot as plt
7 import seaborn as sns

```

Bước 2: Đọc dữ liệu từ tệp

Dữ liệu được đọc từ tệp `results.csv` bằng `pd.read_csv()`. Tệp này được giả định chứa dữ liệu thống kê hiệu suất của các cầu thủ bóng đá.

```
1 # Load data from 'results.csv' file
2 data = pd.read_csv('results.csv')
```

Bước 3: Chọn các đặc trưng để phân cụm

Mã chọn 5 đặc trưng từ tập dữ liệu để phân cụm:

- `stats_passing_Passes Completed`: Số đường chuyền hoàn thành (kiểm soát bóng).
- `stats_standard_Goals`: Số bàn thắng.
- `stats_defense_Tackles`: Số lần tắc bóng.
- `stats_standard_Assists`: Số đường kiến tạo.
- `stats_shooting-Shots on target/90`: Số cú sút trúng đích mỗi 90 phút.

Các đặc trưng này được lưu vào danh sách `features`, sau đó dữ liệu tương ứng được trích xuất và các giá trị thiếu được thay bằng 0 bằng `fillna(0)`.

```
1 # Select features for clustering
2 features = [
3     'stats_passing_Passes Completed', # Passing control
4     'stats_standard_Goals',          # Number of goals
5     'stats_defense_Tackles',         # Tackling ability
6     'stats_standard_Assists',        # Number of assists
7     'stats_shooting-Shots on target/90' # Shots on target per 90
8     minutes
9 ]
10 # Extract features and handle missing values
11 data_features = data[features].fillna(0)
```

Bước 4: Chuẩn hóa dữ liệu

Dữ liệu được chuẩn hóa bằng `StandardScaler` để đưa các đặc trưng về cùng thang đo (trung bình = 0, độ lệch chuẩn = 1). Điều này rất quan trọng vì K-Means dựa trên khoảng cách Euclidean, và các đặc trưng có thang đo khác nhau có thể làm sai lệch kết quả.

```
1 # Standardize the data
2 scaler = StandardScaler()
3 data_scaled = scaler.fit_transform(data_features)
```

Bước 5: Xác định số cụm tối ưu bằng phương pháp Elbow

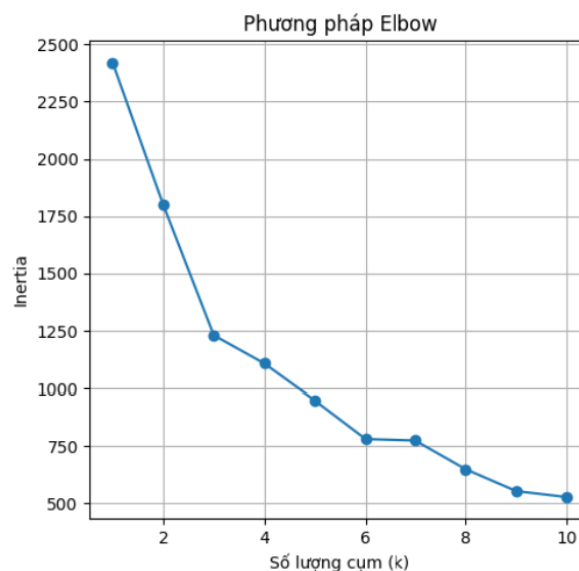
Phương pháp Elbow được sử dụng để xác định số lượng cụm (k) tối ưu. Mã lập qua các giá trị k từ 1 đến 10, huấn luyện mô hình K-Means cho từng k , và tính giá trị *inertia* (tổng bình phương khoảng cách từ mỗi điểm đến tâm cụm gần nhất). Giá trị *inertia* được lưu vào danh sách để vẽ biểu đồ.

```
1 # Elbow method to determine optimal number of clusters
2 inertia = []
3 k_range = range(1, 11)
4 for k in k_range:
5     kmeans = KMeans(n_clusters=k, random_state=42)
6     kmeans.fit(data_scaled)
7     inertia.append(kmeans.inertia_)
```

Bước 6: Vẽ biểu đồ Elbow

Biểu đồ Elbow được vẽ để trực quan hóa *inertia* theo k . Điểm "khủy tay" (nơi *inertia* giảm chậm lại) sẽ chỉ ra số cụm tối ưu. Biểu đồ được lưu vào tệp `elbow_plot.png`.

```
1 # Plot the Elbow curve
2 plt.figure(figsize=(10, 5))
3 plt.subplot(1, 2, 1)
4 plt.plot(k_range, inertia, marker='o')
5 plt.title('Elbow Method')
6 plt.xlabel('Number of clusters (k)')
7 plt.ylabel('Inertia')
8 plt.grid(True)
9 plt.savefig('elbow_plot.png')
```



Hình 3: Biểu đồ Elbow

Bước 7: Áp dụng K-Means với số cụm tối ưu

Giả sử số cụm tối ưu là 4 (dựa trên biểu đồ Elbow), mã khởi tạo và huấn luyện mô hình K-Means với $k=4$, sau đó dự đoán nhãn cụm cho từng điểm dữ liệu.

```
1 # Assume optimal number of clusters is 4 (to be confirmed from
   the plot)
2 optimal_k = 4
3 kmeans = KMeans(n_clusters=optimal_k, random_state=42)
4 clusters = kmeans.fit_predict(data_scaled)
```

Bước 8: Giảm chiều dữ liệu bằng PCA

Để trực quan hóa dữ liệu trong không gian 2 chiều, mã sử dụng PCA để giảm chiều dữ liệu từ 5 chiều (5 đặc trưng) xuống 2 chiều (PC1 và PC2).

```
1 # Reduce dimensions using PCA to 2D
2 pca = PCA(n_components=2)
3 data_pca = pca.fit_transform(data_scaled)
```

Bước 9: Tạo DataFrame cho dữ liệu PCA

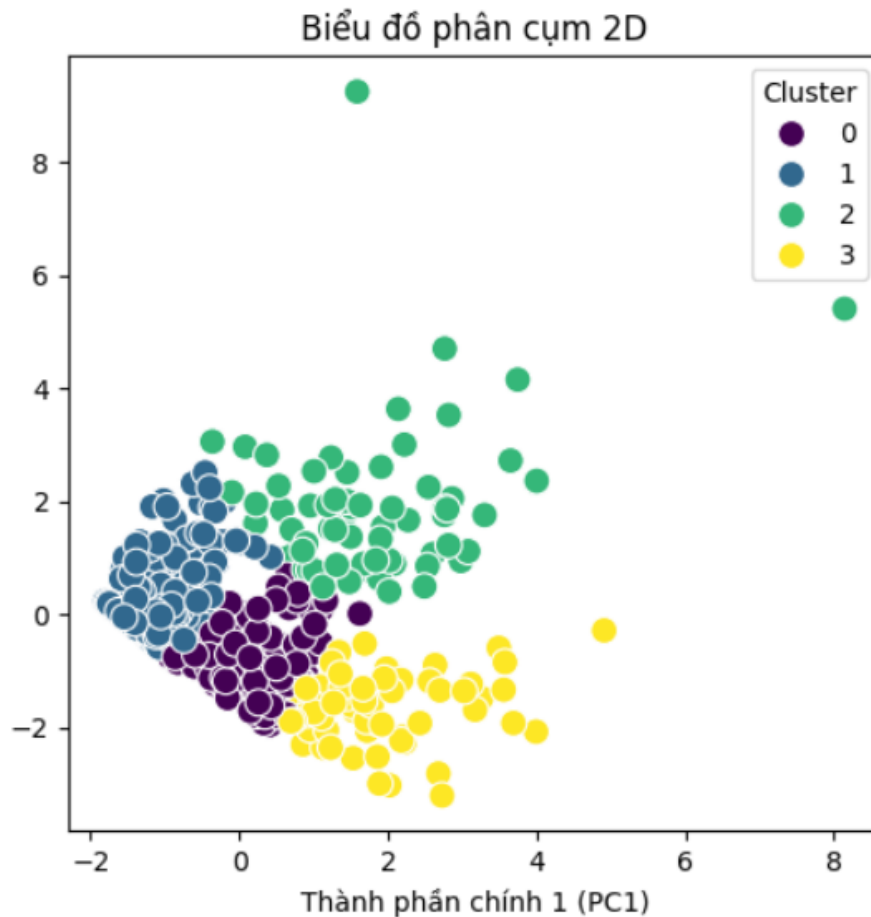
Dữ liệu sau khi giảm chiều được lưu vào một DataFrame với các cột PC1, PC2, và thêm cột **Cluster** chứa nhãn cụm.

```
1 # Create DataFrame for PCA data
2 df_pca = pd.DataFrame(data_pca, columns=['PC1', 'PC2'])
3 df_pca['Cluster'] = clusters
```

Bước 10: Vẽ biểu đồ phân cụm 2D

Biểu đồ phân tán được vẽ bằng `seaborn.scatterplot`, trong đó trục x là PC1, trục y là PC2, và màu sắc biểu thị các cụm khác nhau. Biểu đồ được lưu vào `cluster_plot.png`.

```
1 # Plot 2D clustering visualization
2 plt.subplot(1, 2, 2)
3 sns.scatterplot(x='PC1', y='PC2', hue='Cluster', data=df_pca,
   palette='viridis', s=100)
4 plt.title('2D Clustering Visualization')
5 plt.xlabel('Principal Component 1 (PC1)')
6 plt.ylabel('Principal Component 2 (PC2)')
7 plt.savefig('cluster_plot.png')
```



Hình 4: Biểu đồ phân cụm 2D

Bài 4:

Mã nguồn lấy data cho chuyển nhượng cầu thủ

```

1 # Function to scrape player data from a given URL
2 def scrape_page(driver, url):
3     driver.get(url) # Navigate to the URL
4     time.sleep(3) # Wait 3 seconds for page to load
5     soup = BeautifulSoup(driver.page_source, 'html.parser') #
6         Parse HTML content
7     # Find the table containing player data
8     table = soup.find('table', class_='table table-hover
9         no-cursor table-striped leaguetable mvp-table
10        similar-players-table mb-0')
11     if not table: # Return empty list if table is not found
12         return []
13
14     data = [] # List to store player data
15     rows = table.find('tbody').find_all('tr') # Get all rows in
16         table body

```

```

14     for row in rows: # Process each row
15         try:
16             # Extract skill and potential values
17             skill_div = row.find('div',
18                                   class_='table-skill__skill')
19             pot_div = row.find('div', class_='table-skill__pot')
20             skill = float(skill_div.text.strip()) if skill_div
21                 else None
22             pot = float(pot_div.text.strip()) if pot_div else None
23             skill_pot = f"{skill}/{pot}" if skill and pot else
24                 None # Combine skill/potential
25
26             # Extract player name
27             player_span = row.find('span', class_='d-none')
28             player_name = player_span.text.strip() if player_span
29                 else None
30
31             # Extract team name
32             team_span = row.find('span',
33                                   class_='td-team__teamname')
34             team = team_span.text.strip() if team_span else None
35
36             # Extract estimated transfer value (ETV)
37             etv_span = row.find('span', class_='player-tag')
38             etv = etv_span.text.strip() if etv_span else None
39
40             # Append data as a dictionary
41             data.append({
42                 'player_name': player_name,
43                 'team': team,
44                 'price': etv,
45                 'skill/pot': skill_pot
46             })
47         except: # Skip row if any error occurs
48             continue
49
50     return data # Return list of player data
51
52 [title={Scrape data from all pages}]
53 # Define base URL and number of pages to scrape
54 base_url =
55     "https://www.footballtransfers.com/en/players/uk-premier-league"
56 total_pages = 22
57 all_data = [] # List to store data from all pages
58
59 driver = setup_driver() # Initialize Chrome driver
60
61 try:
62     # Iterate through all pages
63     for page in range(1, total_pages + 1):

```

```
59     if page == 1:
60         url = base_url # Use base URL for first page
61     else:
62         url = f"{base_url}/{page}" # Append page number for
        subsequent pages
63     page_data = scrape_page(driver, url) # Scrape data from
        current page
64     all_data.extend(page_data) # Add page data to main list
65 finally:
66     driver.quit() # Ensure driver is closed
67     # Convert data to DataFrame and save to CSV
68 df_final = pd.DataFrame(all_data) # Create DataFrame from
        collected data
69 df_final.to_csv('football_transfers_players.csv', index=False) #
        Save to CSV without index
```

Mã nguồn chính

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.ensemble import RandomForestRegressor
4 from sklearn.model_selection import train_test_split
5 from sklearn.preprocessing import StandardScaler
6 from sklearn.metrics import mean_squared_error, r2_score
7 import matplotlib.pyplot as plt
8 import seaborn as sns
9
10 # Read data from CSV files
11 results_df = pd.read_csv('results.csv')
12 transfer_df = pd.read_csv('football_transfers_players.csv')
13
14 # Filter players with more than 900 minutes played
15 results_df = results_df[results_df['stats_standard_Minutes'] >
16     900]
17
18 # Standardize player names for merging
19 results_df['Player'] =
20     results_df['Player'].str.lower().str.strip()
21 transfer_df['player_name'] =
22     transfer_df['player_name'].str.lower().str.strip()
23
24 # Merge the two DataFrames based on player names, keeping only
25     matches
26 merged_df = pd.merge(results_df, transfer_df, left_on='Player',
27     right_on='player_name', how='inner')
28
29 # Drop rows with missing transfer prices
30 merged_df = merged_df.dropna(subset=['price'])
31
32 # Save the merged DataFrame to a CSV file
33 merged_df.to_csv('merged_results_transfer.csv', index=False,
34     encoding='utf-8-sig')
35
36 # Function to convert transfer price strings to floats
37 def convert_price(price_str):
38     if pd.isna(price_str) or not isinstance(price_str, str):
39         return np.nan
40     price_str = price_str.replace(' ', '').strip()
41     if 'M' in price_str:
42         return float(price_str.replace('M', ''))
43     elif 'K' in price_str:
44         return float(price_str.replace('K', '')) / 1000
45     return np.nan
46
47 # Apply the conversion function to the 'price' column
48 merged_df['price'] = merged_df['price'].apply(convert_price)
49
```

```

44 # Select feature columns (stats_ columns excluding certain ones)
45 feature_cols = [col for col in results_df.columns if
46                 col.startswith('stats_') and col not in [
47                     'stats_standard_Nation', 'stats_standard_Squad',
48                     'stats_standard_Position'
49 ]]
50
51 # Handle the age column: extract the numeric age
52 def extract_age(age_str):
53     if pd.isna(age_str) or not isinstance(age_str, str):
54         return np.nan
55     try:
56         return float(age_str.split('-')[0])
57     except:
58         return np.nan
59
60 if 'stats_standard_Current age' in feature_cols:
61     merged_df['stats_standard_Current age'] =
62         merged_df['stats_standard_Current age'].apply(extract_age)
63
64 # Extract features
65 features = merged_df[feature_cols]
66
67 # Replace 'N/a' with NaN and convert to numeric
68 features = features.replace('N/a', np.nan)
69 features = features.apply(pd.to_numeric, errors='coerce')
70 features = features.fillna(0)
71
72 # Target variable is the converted price
73 target = merged_df['price']
74
75 # Scale the features
76 scaler = StandardScaler()
77 features_scaled = scaler.fit_transform(features)
78
79 # Split data into training and testing sets
80 X_train, X_test, y_train, y_test =
81     train_test_split(features_scaled, target, test_size=0.2,
82                       random_state=42)
83
84 # Train the Random Forest model
85 rf_model = RandomForestRegressor(n_estimators=100,
86                                 random_state=42)
87 rf_model.fit(X_train, y_train)
88
89 # Make predictions on the test set
90 y_pred = rf_model.predict(X_test)
91
92 # Calculate loss (Mean Squared Error) and accuracy (R-squared)
93 mse = mean_squared_error(y_test, y_pred)
94 r2 = r2_score(y_test, y_pred)

```

```

89
90 # Print results
91 print(f"Mean Squared Error (MSE): {mse:.2f}")
92 print(f"R-squared Score (R2): {r2:.2f}")
93
94 # Plotting results
95 plt.figure(figsize=(12, 5))
96
97 # Plot 1: Actual vs Predicted values
98 plt.subplot(1, 2, 1)
99 plt.scatter(y_test, y_pred, color='blue', alpha=0.5)
100 plt.plot([y_test.min(), y_test.max()], [y_test.min(),
    y_test.max()], 'r--', lw=2)
101 plt.xlabel('Actual Price ( M )')
102 plt.ylabel('Predicted Price ( M )')
103 plt.title('Actual vs Predicted Prices')
104 plt.grid(True)
105
106 # Plot 2: Error distribution
107 plt.subplot(1, 2, 2)
108 errors = y_pred - y_test
109 sns.histplot(errors, kde=True, color='purple')
110 plt.xlabel('Prediction Error ( M )')
111 plt.ylabel('Frequency')
112 plt.title('Prediction Error Distribution')
113 plt.grid(True)
114
115 plt.tight_layout()
116 plt.savefig('prediction_results.png')
117 plt.close()
118
119 # Save predictions to a CSV file
120 predictions_df = pd.DataFrame({
121     'Player': merged_df.loc[y_test.index, 'Player'],
122     'Actual_Price': y_test,
123     'Predicted_Price': y_pred
124 })
125 predictions_df.to_csv('price_predictions.csv', index=False,
    encoding='utf-8-sig')

```

Phân tích chi tiết từng bước

Bước 1: Nhập các thư viện cần thiết

Đoạn mã bắt đầu bằng việc nhập các thư viện cần thiết để xử lý dữ liệu, huấn luyện mô hình, và trực quan hóa kết quả:

- `pandas` và `numpy`: Thư viện xử lý và tính toán dữ liệu.
- `sklearn.ensemble.RandomForestRegressor`: Mô hình Random Forest cho bài toán hồi quy.

-
- `sklearn.model_selection.train_test_split`: Chia dữ liệu thành tập huấn luyện và kiểm tra.
 - `sklearn.preprocessing.StandardScaler`: Chuẩn hóa dữ liệu số.
 - `sklearn.metrics`: Các hàm tính toán độ lỗi (MSE) và độ chính xác (R-squared).
 - `matplotlib.pyplot` và `seaborn`: Trực quan hóa dữ liệu qua biểu đồ.

Bước 2: Đọc dữ liệu từ tệp CSV

Dữ liệu được đọc từ hai tệp CSV: `results.csv` chứa thông tin thống kê cầu thủ và `football_transfers_players.csv` chứa giá trị chuyển nhượng.

Bước 3: Lọc cầu thủ có thời gian thi đấu trên 900 phút

Lọc dữ liệu để chỉ giữ lại các cầu thủ có thời gian thi đấu (`stats_standard_Minutes`) lớn hơn 900 phút, đảm bảo dữ liệu đủ ý nghĩa để phân tích.

Bước 4: Chuẩn hóa tên cầu thủ

Tên cầu thủ trong cả hai bảng được chuyển thành chữ thường và loại bỏ khoảng trắng thừa, giúp việc hợp nhất dữ liệu chính xác hơn.

Bước 5: Hợp nhất hai DataFrame

Hai bảng dữ liệu được hợp nhất dựa trên tên cầu thủ bằng phép nối trong (`inner join`), chỉ giữ lại các cầu thủ xuất hiện ở cả hai tập dữ liệu.

Bước 6: Loại bỏ các hàng thiếu giá trị chuyển nhượng

Các hàng không có giá trị trong cột `price` bị loại bỏ để đảm bảo mô hình hoạt động với dữ liệu đầy đủ.

Bước 7: Lưu DataFrame hợp nhất

Dữ liệu sau khi hợp nhất được lưu vào tệp `merged_results_transfer.csv` với mã hóa `utf-8-sig` để hỗ trợ ký tự tiếng Việt.

Bước 8: Chuyển đổi giá trị chuyển nhượng từ chuỗi sang số

Hàm `convert_price` xử lý chuỗi giá trị như `'€50M'` hoặc `'€500K'`, loại bỏ ký tự `'€'`, chuyển `'M'` thành triệu và `'K'` thành nghìn, trả về số thực.

Bước 9: Chọn các cột đặc trưng

Các cột bắt đầu bằng `stats_` được chọn làm đặc trưng, ngoại trừ các cột không phải số như quốc gia, đội bóng, và vị trí.

Bước 10: Xử lý cột tuổi

Hàm `extract_age` tách tuổi từ chuỗi như `'26-352'`, chỉ lấy phần số năm (26) để sử dụng làm đặc trưng.

Bước 11: Trích xuất đặc trưng

Các cột đặc trưng được trích xuất từ DataFrame hợp nhất để chuẩn bị cho việc huấn luyện mô hình.

Bước 12: Thay thế 'N/a' và chuyển sang kiểu số

Giá trị 'N/a' được thay bằng NaN, các cột được chuyển thành kiểu số, và các giá trị NaN còn lại được điền bằng 0.

Bước 13: Xác định biến mục tiêu

Biến mục tiêu là cột `price` đã được chuyển đổi sang dạng số thực.

Bước 14: Chuẩn hóa đặc trưng

Sử dụng `StandardScaler` để chuẩn hóa các đặc trưng, giúp mô hình học tốt hơn trên dữ liệu có thang đo khác nhau.

Bước 15: Chia dữ liệu thành tập huấn luyện và kiểm tra

Dữ liệu được chia thành tập huấn luyện (80%) và tập kiểm tra (20%) với `random_state=42` để đảm bảo kết quả tái lập.

Bước 16: Huấn luyện mô hình Random Forest

Mô hình Random Forest với 100 cây quyết định được khởi tạo và huấn luyện trên tập huấn luyện.

Bước 17: Dự đoán trên tập kiểm tra

Mô hình dự đoán giá trị chuyển nhượng trên tập kiểm tra, kết quả được lưu vào `y_pred`.

Bước 18: Tính toán hàm loss và độ chính xác

Tính toán Mean Squared Error (MSE) để đo lường độ lỗi và R-squared (R2) để đánh giá độ chính xác của mô hình.

Bước 19: Vẽ biểu đồ kết quả

Hai biểu đồ được tạo ra:

- Biểu đồ phân tán so sánh giá trị thực tế và dự đoán, với đường tham chiếu màu đỏ.
- Biểu đồ phân phối sai số dự đoán, sử dụng histogram và đường cong mật độ.

Kết quả được lưu vào tệp `prediction_results.png`.

Bước 20: Lưu kết quả dự đoán

Dữ liệu dự đoán (tên cầu thủ, giá trị thực tế, giá trị dự đoán) được lưu vào tệp `price_predictions.csv` với mã hóa `utf-8-sig`.

