

Bài 1:

1. Khai báo thư viện

```
1 from selenium import webdriver
2 from selenium.webdriver.common.by import By
3 from selenium.webdriver.chrome.service import Service
4 from selenium.webdriver.chrome.options import Options
5 import pandas as pd
6 import time
```

- **selenium**: Thư viện Selenium giúp tự động hóa việc điều khiển trình duyệt web. Selenium sẽ giúp chúng ta duyệt trang web, tìm các phần tử trên trang và lấy dữ liệu từ chúng.
- **webdriver**: Đây là module chính trong Selenium, được sử dụng để điều khiển trình duyệt (như Chrome, Firefox).
- **By**: Cung cấp các phương thức để xác định các phần tử trên trang web (ví dụ: qua ID, class, XPath).
- **Service**: Dùng để chỉ định đường dẫn đến ChromeDriver (trình điều khiển của trình duyệt Chrome).
- **Options**: Cho phép cấu hình các tùy chọn cho trình duyệt, ví dụ như chế độ không giao diện người dùng (headless).
- **pandas**: Là thư viện xử lý và phân tích dữ liệu trong Python, rất hữu ích để lưu trữ và xử lý bảng dữ liệu (DataFrame).
- **time**: Thư viện giúp chúng ta tạm dừng chương trình trong một khoảng thời gian (ví dụ: để đợi trang web tải xong).

2. Cấu hình Selenium

```
7 chrome_options = Options()
8 chrome_options.add_argument("--headless")
9 chrome_options.add_argument("--no-sandbox")
10 chrome_options.add_argument("--disable-dev-shm-usage")
11 service = Service(executable_path=r"C:\Users\VIETTELSTORE\Downloads\chromedriver-win64\chromedriver-win64\chromedriver.exe")
12 driver=webdriver.Chrome(options=chrome_options, service=Service)
```

- **chrome_options**: Tạo một đối tượng cấu hình trình duyệt Chrome. Các tùy chọn được thiết lập ở đây:
 - **--headless**: Chạy trình duyệt mà không mở giao diện người dùng (giúp tiết kiệm tài nguyên và tăng tốc độ khi lấy dữ liệu).
 - **--no-sandbox** và **--disable-dev-shm-usage**: Các tùy chọn này giúp tránh lỗi khi chạy trong môi trường máy chủ hoặc các máy tính có hạn chế về tài nguyên.

- `Service`: Chỉ định đường dẫn đến `chromedriver.exe`, công cụ dùng để điều khiển Chrome. Đây là phần bắt buộc khi sử dụng Selenium với Chrome.
- `webdriver.Chrome`: Khởi tạo đối tượng `driver` để điều khiển trình duyệt Chrome theo cấu hình đã thiết lập.

3. URL và Danh mục thống kê

```
13 url='https://fbref.com/en/comps/9/{}/Premier-League-Stats'
14 categories = ["stats", "keepers", "shooting","passing", "gca", "defense", "possession", "misc"]
15 path=["standard", 'keeper', 'shooting','passing', 'gca', 'defense', 'possession', 'misc']
16 dataframes={}
17 n=[0,1,2,3,4,5,6,7]
```

- `url`: Đây là URL của trang FBref nơi chứa dữ liệu về cầu thủ Premier League. `{}` sẽ được thay thế bằng từng mục trong danh sách `categories`.
- `categories`: Các danh mục thống kê mà bạn sẽ thu thập dữ liệu từ trang web, bao gồm các chỉ số như: thống kê chung, thủ môn, sút bóng, chuyền bóng, v.v.
- `path`: Là các phần của URL đại diện cho từng loại dữ liệu cụ thể. Mỗi mục trong `categories` tương ứng với một phần tử trong `path`.
- `dataframes`: Tạo một từ điển để lưu trữ dữ liệu từ mỗi danh mục thống kê.
- `n`: Một danh sách chỉ số để lặp qua tất cả các danh mục thống kê (từ 0 đến 7).

4. Cấu hình ánh xạ các cột dữ liệu

```
5. mapping_config = {
6.     "stats": {
7.         "Player": "player",
8.         "Nation": "nationality",
9.         "Team": "team",
10.        "Position": "position",
11.        "Age": "age",
12.        "Matches Played": "games",
13.        "Starts": "games_starts",
14.        "Minutes": "minutes",
15.        "Goals": "goals",
16.        "Assists": "assists",
17.        "Yellow Cards": "cards_yellow",
18.        "Red Cards": "cards_red",
19.        "xG": "xg",
20.        "xAG": "xg_assist",
21.        "PrgC": "progressive_carries",
22.        "PrgP": "progressive_passes",
23.        "PrgR": "progressive_passes_received",
24.        "Gls90": "goals_per90",
25.        "Ast90": "assists_per90",
26.        "xG90": "xg_per90",
```

```
27.     "xAG90": "xg_assist_per90"
28. },
29. "keepers": {
30.     "Player": "player",
31.     "Nation": "nationality",
32.     "Team": "team",
33.     "Position": "position",
34.     "Age": "age",
35.     "Matches Played": "gk_games",
36.     "Starts": "gk_games_starts",
37.     "Minutes": "gk_minutes",
38.     "GA90": "gk_goals_against_per90",
39.     "Save%": "gk_save_pct",
40.     "CS%": "gk_clean_sheets_pct",
41.     "PK Save%": "gk_pens_save_pct"
42. },
43. "shooting": {
44.     "Player": "player",
45.     "Nation": "nationality",
46.     "Team": "team",
47.     "Position": "position",
48.     "Age": "age",
49.     "SoT%": "shots_on_target_pct",
50.     "SoT/90": "shots_on_target_per90",
51.     "G/Sh": "goals_per_shot",
52.     "Dist": "average_shot_distance"
53. },
54. "passing": {
55.     "Player": "player",
56.     "Nation": "nationality",
57.     "Team": "team",
58.     "Position": "position",
59.     "Age": "age",
60.     "Cmp": "passes_completed",
61.     "Cmp%": "passes_pct",
62.     "TotDist": "passes_total_distance",
63.     "Cmp% (Short)": "passes_pct_short",
64.     "Cmp% (Medium)": "passes_pct_medium",
65.     "Cmp% (Long)": "passes_pct_long",
66.     "KP": "assisted_shots",
67.     "Passes 1/3": "passes_into_final_third",
68.     "PPA": "passes_into_penalty_area",
69.     "CrsPA": "crosses_into_penalty_area"
70. },
71. "gca": {
72.     "Player": "player",
73.     "Nation": "nationality",
74.     "Team": "team",
```

```

75.     "Position": "position",
76.     "Age": "age",
77.     "SCA": "sca",
78.     "SCA90": "sca_per90",
79.     "GCA": "gca",
80.     "GCA90": "gca_per90"
81. },
82. "defense": {
83.     "Player": "player",
84.     "Nation": "nationality",
85.     "Team": "team",
86.     "Position": "position",
87.     "Age": "age",
88.     "Tkl": "tackles",
89.     "TklW": "tackles_won",
90.     "Att": "challenges",
91.     "Challenges Lost": "challenges_lost",
92.     "Blocks": "blocks",
93.     "Sh": "blocked_shots",
94.     "Pass": "blocked_passes",
95.     "Int": "interceptions"
96. },
97. "possession": {
98.     "Player": "player",
99.     "Nation": "nationality",
100.     "Team": "team",
101.     "Position": "position",
102.     "Age": "age",
103.     "Touches": "touches",
104.     "Def Pen": "touches_def_pen_area",
105.     "Def 3rd": "touches_def_3rd",
106.     "Mid 3rd": "touches_mid_3rd",
107.     "Att 3rd": "touches_att_3rd",
108.     "Att Pen": "touches_att_pen_area",
109.     "Att (TO)": "take_ons",
110.     "Succ%": "take_ons_won_pct",
111.     "Tkld%": "take_ons_tackled_pct",
112.     "Carries": "carries",
113.     "ProDist": "carries_progressive_distance",
114.     "ProgC": "progressive_carries",
115.     "Carries 1/3": "carries_into_final_third",
116.     "CPA": "carries_into_penalty_area",
117.     "Mis": "miscontrols",
118.     "Dis": "dispossessed",
119.     "Rec": "passes_received"
120. },
121. "misc": {
122.     "Player": "player",

```

```

123.         "Nation": "nationality",
124.         "Team": "team",
125.         "Position": "position",
126.         "Age": "age",
127.         "Fls": "fouls",
128.         "Fld": "fouled",
129.         "Off": "offsides",
130.         "Crs": "crosses",
131.         "Recov": "ball_recoveries",
132.         "Won": "aerials_won",
133.         "Aerials Lost": "aerials_lost",
134.         "Won%": "aerials_won_pct"
135.     }
136. }

```

- `mapping_config`: Đây là cấu hình ánh xạ giữa các tên cột trong dữ liệu lấy từ FBref và các tên cột chuẩn trong DataFrame của bạn. Ví dụ, từ "Goals" trong dữ liệu gốc sẽ được ánh xạ thành "goals" trong DataFrame.

5. Hàm `add`

```

150 def add(all_data, col1, col2, col3):
151     col_index = all_data.columns.get_loc(col1)
152     all_data[col3] = all_data[col1].combine_first(all_data[col2])
153     all_data.drop(columns=[col1, col2], inplace=True)
154     cols = list(all_data.columns)
155     cols.remove(col3)
156     cols.insert(col_index, col3)
157     return all_data[cols]

```

- Hàm này dùng để hợp nhất hai cột có tên là `col1` và `col2` thành một cột mới tên là `col3`.
- Nếu một trong các giá trị ở `col1` là NaN, giá trị từ `col2` sẽ thay thế.
- Sau đó, hàm sẽ xóa đi hai cột ban đầu và giữ cột mới tại đúng vị trí.

6. Hàm `restructure_row`

```

158 def restructure_row(row_data, mapping):
159     new_data = {}
160     for field, header in mapping.items():
161         new_data[field] = row_data.get(header, "N/A")
162     return new_data

```

- Hàm này sẽ nhận vào một dòng dữ liệu và ánh xạ các giá trị từ `row_data` theo cấu hình trong `mapping`. Nếu một giá trị không có, nó sẽ gán "N/A".

- Kết quả trả về là một từ điển với các trường đã được chuẩn hóa.

7. Thu thập và xử lý dữ liệu

```
for i in n:
    player = []
    driver.get(url.format(categories[i]))
    time.sleep(3)
    header_elements = driver.find_elements(By.XPATH, '//*[@id="stats_{}"]//thead//tr[2]//th'.format(path[i]))
    headers = [h.get_attribute('data-stat') for h in header_elements]
    row_elements = driver.find_elements(By.XPATH, '//*[@id="stats_{}"]//tbody//tr'.format(path[i]))
    headers.pop(0)
    for row in row_elements:
        cell_elements = row.find_elements(By.TAG_NAME, 'td')
        row_data = {header: cell.text.strip() for header, cell in zip(headers, cell_elements)}
        restructured = restructure_row(row_data, mapping_config[categories[i]])
        player.append(restructured)
    df = pd.DataFrame(player)
    dataframes[categories[i]] = df
    del df
all_data=dataframes[categories[0]]
for i in range(1,8):
    df = dataframes[categories[i]]
    df = df.groupby("Player").first().reset_index() # <- dòng thêm vào
    all_data = pd.merge(all_data, df, on=["Player","Nation","Team","Position","Age"], how='outer')
    del df
```

- Mã này lặp qua từng danh mục trong `categories`, lấy dữ liệu từ trang web bằng cách sử dụng `driver.get` để tải trang.
- Sau khi trang web tải xong, sử dụng `driver.find_elements` để lấy các phần tử cần thiết từ bảng (ví dụ: các tiêu đề cột và các dòng dữ liệu).
- Các dữ liệu thu thập được sẽ được xử lý và chuẩn hóa thông qua hàm `restructure_row`, rồi lưu vào `DataFrame`.

8. Kết hợp dữ liệu từ các danh mục thống kê

```
180 for i in range(1,8):
181     df = dataframes[categories[i]]
182     df = df.groupby("Player").first().reset_index() # <- dòng thêm vào
183     all_data = pd.merge(all_data, df, on=["Player","Nation","Team","Position","Age"], how='outer')
184     del df
```

- Tạo `all_data` từ dữ liệu của danh mục đầu tiên (`stats`).
- Sau đó, lặp qua các danh mục còn lại, nhóm dữ liệu theo "Player" (cầu thủ) và kết hợp với `all_data` bằng phương thức `merge`.

9. Xử lý các cột và lọc dữ liệu

```
186 all_data= add(all_data, 'Minutes_x', 'Minutes_y', 'Minutes')
187 all_data= add(all_data, 'Matches Played_x', 'Matches Played_y', 'Matches Played')
188 all_data= add(all_data, 'Starts_x', 'Starts_y', 'Starts')
189 all_data['Minutes'] = pd.to_numeric(all_data['Minutes']).astype(str).str.replace(",",""), errors='coerce')
190 all_data=all_data[all_data['Minutes'] > 90]
```

- Dữ liệu từ các cột "Minutes" và "Matches Played" được hợp nhất bằng hàm `add`.
- Cột "Minutes" được chuyển sang kiểu số thực (`float`), và chỉ giữ lại các cầu thủ có trên 90 phút thi đấu.

10. Tiền xử lý dữ liệu và lưu vào file CSV

```
191 all_data['First Name'] = all_data['Player'].str.split().str[0]
192 all_data['Last Name'] = all_data['Player'].str.split().str[-1]
193 all_data = all_data.sort_values(by="First Name")
194 all_data = all_data.drop(columns=['First Name', 'Last Name'])
195 all_data = all_data.fillna("N/a")
196 all_data = all_data.replace("", "N/a")
197 all_data.to_csv('results.csv', index=False)
```

- Tách tên cầu thủ thành "First Name" và "Last Name", sau đó sắp xếp dữ liệu theo "First Name".
- Thay thế các giá trị thiếu hoặc rỗng bằng "N/a".
- Lưu dữ liệu vào file `results.csv`.

11. Thoát khỏi trình duyệt

```
driver.quit()
```

- Sau khi hoàn tất thu thập và xử lý dữ liệu, đóng trình duyệt để giải phóng tài nguyên.

Bài 2: Phân tích thống kê dữ liệu các chỉ số cầu thủ

1. Đọc dữ liệu và xử lý chỉ số cần phân tích

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import os
4
5 # ===== 1. Đọc dữ liệu và xử lý chỉ số cần phân tích =====
6 df = pd.read_csv("results.csv")
7 columns = ['Age', 'Matches Played', 'Starts', 'Minutes', 'Goals', 'Assists', 'Yellow Cards', 'Red Cards',
8            'xG', 'xAG', 'PrgC', 'PrgP', 'PrgR', 'Gls90', 'Ast90', 'xG90', 'xAG90', 'GA90', 'Save%',
9            'CS%', 'PK Save%', 'SoT%', 'SoT/90', 'G/Sh', 'Dist', 'Cmp', 'Cmp%', 'TotDist', 'Cmp% (Short)',
10           'Cmp% (Medium)', 'Cmp% (Long)', 'KP', 'Passes 1/3', 'PPA', 'CrsPA', 'SCA', 'SCA90', 'GCA', 'GCA90',
11           'Tkl', 'TklW', 'Att', 'Challenges Lost', 'Blocks', 'Sh', 'Pass', 'Int', 'Touches', 'Def Pen',
12           'Def 3rd', 'Mid 3rd', 'Att 3rd', 'Att Pen', 'Att (TO)', 'Succ%', 'Tkld%', 'Carries', 'ProDist',
13           'ProgC', 'Carries 1/3', 'CPA', 'Mis', 'Dis', 'Rec', 'Fls', 'Fld', 'Off', 'Crs', 'Recov', 'Won',
14           'Aerials Lost', 'Won%']
15 teams = df['Team'].unique()
16
```

- Đoạn mã này sử dụng thư viện `pandas` để đọc dữ liệu từ file CSV (`results.csv`) vào một `DataFrame`.
- Dữ liệu này bao gồm thông tin về các cầu thủ bóng đá, với nhiều chỉ số thống kê như "Goals", "Assists", "Minutes", "xG", và các chỉ số khác.
- Mảng `columns` chứa danh sách các chỉ số mà bạn muốn phân tích.

- Biến `teams` chứa danh sách các đội bóng tham gia trong dữ liệu, giúp phân nhóm các cầu thủ theo đội sau này.

2. Chuyển đổi giá trị "Age" từ dạng sang số năm

```
# ===== 2. Chuyển đổi giá trị "Age" sang số năm =====
def convert_age_to_years(age_str):
    try:
        # Tách tuổi và số ngày
        age, days = age_str.split('-')
        age = int(age)
        days = int(days)
        # Tính tuổi tính theo năm, cộng thêm phần năm từ số ngày
        return age + (days / 365.25)
    except Exception as e:
        print(f"Lỗi khi chuyển đổi tuổi {age_str}: {e}")
        return 0 # Trả về 0 nếu không thể chuyển đổi

# Áp dụng hàm chuyển đổi cho cột 'Age'
df['Age'] = df['Age'].apply(convert_age_to_years)
```

- Cột `Age` chứa giá trị tuổi của các cầu thủ dưới dạng `xx-yyy`, trong đó `xx` là số tuổi và `yyy` là số ngày trong năm tính từ ngày sinh.
- Đoạn mã này định nghĩa một hàm `convert_age_to_years` để tách phần tuổi và phần ngày, sau đó tính toán tuổi chính xác theo năm.
- Việc chia cho 365.25 giúp tính đúng số ngày trong một năm tính cả năm nhuận.
- Hàm này được áp dụng cho toàn bộ cột `Age` của `DataFrame`, chuyển đổi tất cả các giá trị thành kiểu dữ liệu số (float), giúp dễ dàng xử lý và phân tích dữ liệu.

3. Chuyển đổi các cột còn lại về kiểu số

```
33 # ===== 3. Chuyển đổi các cột còn lại về kiểu số =====
34 for col in columns:
35     df[col] = pd.to_numeric(df[col].astype(str).replace(",","").replace("N/a", "0"), errors='coerce')
```

- Sau khi xử lý cột `Age`, đoạn mã này tiếp tục chuyển đổi các cột còn lại trong `columns` về kiểu dữ liệu số.
- Đối với mỗi cột, hàm `apply` sẽ loại bỏ dấu phẩy (,), chuyển đổi giá trị thành số và thay thế các giá trị "N/a" bằng 0.
- `errors='coerce'` đảm bảo rằng những giá trị không hợp lệ sẽ bị thay thế bằng `NaN`, tránh lỗi trong quá trình tính toán sau này.

4. Top 3 cao nhất và thấp nhất mỗi chỉ số


```

36 # ===== 4. Top 3 cao nhất và thấp nhất mỗi chỉ số =====
37 with open("top_3.txt", "w", encoding="utf-8") as f:
38     for col in columns:
39         if df[col].dropna().empty:
40             continue
41         top_3 = df[['Player', col]].sort_values(by=col, ascending=False).head(3)
42         bottom_3 = df[['Player', col]].sort_values(by=col, ascending=True).head(3)
43
44         f.write(f"=== Top 3 for {col} ===\n")
45         for _, row in top_3.iterrows():
46             f.write(f"{row['Player']}: {row[col]}\n")
47         f.write(f"\n=== Bottom 3 for {col} ===\n")
48         for _, row in bottom_3.iterrows():
49             f.write(f"{row['Player']}: {row[col]}\n")
50         f.write("\n" + "="*50 + "\n\n")
51
52 print("✅ Đã tạo file top_3.txt")
53

```

- Đoạn mã này tìm ra các cầu thủ có chỉ số cao nhất và thấp nhất trong mỗi cột trong danh sách `columns`.
- Đối với mỗi cột, các cầu thủ được sắp xếp theo giá trị của cột đó và lấy 3 cầu thủ có giá trị cao nhất và thấp nhất.
- Kết quả được ghi vào file `top_3.txt` với tên cầu thủ và giá trị tương ứng.
- Đây là cách để so sánh các cầu thủ trong các chỉ số quan trọng, giúp bạn tìm ra những cầu thủ nổi bật hoặc thiếu sót trong từng chỉ số.

5. Tính median, mean, std theo đội bóng

```

54 # ===== 5. Median/Mean/Std theo đội bóng =====
55 results = []
56
57 # Kiểm tra giá trị hợp lệ trong mỗi đội
58 for team in teams:
59     row = {'Team': team}
60     for col in columns:
61         team_df = df[df['Team'] == team]
62
63         # Kiểm tra xem dữ liệu trong cột có hợp lệ
64         valid_values = team_df[col].dropna()
65         if not valid_values.empty:
66             # Tính các chỉ số nếu có giá trị hợp lệ
67             row[f'Median of {col}'] = valid_values.median()
68             row[f'Mean of {col}'] = valid_values.mean()
69             row[f'Std of {col}'] = valid_values.std()
70         else:
71             # Nếu không có giá trị hợp lệ, gán giá trị mặc định
72             row[f'Median of {col}'] = 0.0
73             row[f'Mean of {col}'] = 0.0
74             row[f'Std of {col}'] = 0.0
75
76     results.append(row)
77
78 # Lưu kết quả vào file CSV
79 pd.DataFrame(results).round(3).fillna(0.0).to_csv("results2.csv", index=False)
80 print("✅ Đã tạo file results2.csv")

```

- Đoạn mã này tính toán các chỉ số thống kê cơ bản (median, mean, và standard deviation) cho từng chỉ số trong `columns`, nhưng phân nhóm theo đội bóng.
- Kết quả được lưu vào một file `team_stats.csv`, giúp bạn phân tích các đội bóng dựa trên các chỉ số của các cầu thủ trong đội đó.
- Các phép tính thống kê này sẽ giúp bạn so sánh các đội bóng về các chỉ số hiệu suất và ra quyết định trong việc phân tích.

6. Tìm đội có chỉ số trung bình cao nhất cho từng tiêu chí

```

82 # ===== 6. Tìm đội có chỉ số cao nhất mỗi tiêu chí =====
83 best_teams = {}
84 for col in columns:
85     team_means = df.groupby('Team')[col].mean()
86     best_team = team_means.idxmax()
87     best_teams[col] = best_team
88
89 with open("best_teams.txt", "w") as f:
90     for stat, team in best_teams.items():
91         f.write(f'Top team for {stat}: {team}\n')
92
93 print("✅ Đã tạo file best_teams.txt")

```

- Đoạn mã này xác định đội bóng nào có **giá trị trung bình cao nhất** cho từng chỉ số trong danh sách `columns`.
- Dữ liệu được nhóm theo `Team`, sau đó tính trung bình (`mean`) cho từng chỉ số.
- Đội có giá trị trung bình cao nhất (`idxmax()`) được ghi nhận là đội “tốt nhất” trong chỉ số đó.
- Kết quả được lưu vào file `best_teams.txt`, giúp bạn dễ dàng xác định đội nổi bật nhất theo từng tiêu chí.

7. Vẽ biểu đồ histogram cho 3 chỉ số tấn công và phòng thủ toàn giải và từng đội

```

95 # ===== 7. Vẽ biểu đồ histogram từng chỉ số toàn giải và từng đội =====
96 df = pd.read_csv("results.csv")
97 columns = ['PrgC', 'PrgR', 'PrgP', 'Save%', 'PK Save%', 'CS%']
98 teams = df['Team'].unique()
99 os.makedirs("histograms", exist_ok=True)
100
101 for col in columns:
102     # Toàn giải
103     plt.figure(figsize=(10,6))
104     plt.hist(df[col].dropna(), bins=20, alpha=0.7, color='blue')
105     plt.title(f'Distribution of {col} (All Players)')
106     plt.xlabel(col)
107     plt.ylabel('Frequency')
108     plt.grid(True)
109     plt.savefig(f'histograms/{col}_all.png')
110     plt.close()
111
112     # Từng đội
113     for team in teams:
114         team_df = df[df['Team'] == team]
115         plt.figure(figsize=(10,6))
116         plt.hist(team_df[col].dropna(), bins=20, alpha=0.7, color='green')
117         plt.title(f'Distribution of {col} ({team})')
118         plt.xlabel(col)
119         plt.ylabel('Frequency')
120         plt.grid(True)
121         plt.savefig(f'histograms/{col}_{team}.png')
122         plt.close()
123
124 print("✅ Đã tạo biểu đồ histogram trong thư mục histograms/")

```

- Biểu đồ histogram giúp trực quan hóa phân bố giá trị của từng chỉ số, từ đó hiểu được độ trải đều hay chênh lệch.
- Đầu tiên, tạo thư mục `histograms/` để lưu các hình ảnh biểu đồ.
- Với mỗi chỉ số trong `columns`, vẽ:
 - Một biểu đồ histogram cho toàn bộ giải đấu (tất cả cầu thủ).
 - Một biểu đồ riêng cho từng đội bóng, giúp so sánh phân bố nội bộ.
- Điều này hỗ trợ phát hiện các đội có xu hướng vượt trội hoặc tụt hậu ở một chỉ số cụ thể.

Bài 3: Phân cụm cầu thủ bằng K-means và giảm chiều PCA

1. Đọc dữ liệu và xử lý giá trị số

```
1  import pandas as pd
2  import numpy as np
3  from sklearn.preprocessing import StandardScaler
4  from sklearn.cluster import KMeans
5  from sklearn.decomposition import PCA
6  import matplotlib.pyplot as plt
7  import seaborn as sns
8
9  # Đọc dữ liệu
10 df = pd.read_csv('results.csv')
11
12 # Chuyển 'N/a' thành NaN và ép kiểu numeric cho các cột số
13 df.replace('N/a', np.nan, inplace=True)
14
```

- Đọc dữ liệu từ file `results.csv`
- Thay thế các giá trị "N/a" bằng `NaN` để xử lý dữ liệu thiếu

2. Xác định và chuyển đổi các cột số

```
15 # Xác định các cột số
16 numeric_cols = df.select_dtypes(include=[np.number]).columns.tolist()
17
18 # Nếu chưa có numeric cols – ép kiểu thử các cột còn lại
19 if not numeric_cols:
20     for col in df.columns:
21         try:
22             df[col] = pd.to_numeric(df[col])
23         except:
24             continue
25     numeric_cols = df.select_dtypes(include=[np.number]).columns.tolist()
26
```

- Tự động xác định các cột dạng số (`float/int`)
- Nếu tất cả cột đang ở dạng chuỗi (`object`), cố gắng chuyển đổi sang kiểu số

3. Xử lý thiếu dữ liệu và chuẩn hóa

```

27 # Loại bỏ các hàng bị thiếu dữ liệu
28 df.dropna(subset=numeric_cols, inplace=True)
29
30 # Chuẩn hóa dữ liệu
31 scaler = StandardScaler()
32 x_scaled = scaler.fit_transform(df[numeric_cols])
33

```

- Bỏ các cầu thủ thiếu dữ liệu ở bất kỳ cột chỉ số nào
- Chuẩn hóa dữ liệu về cùng thang đo (mean=0, std=1) để phân cụm hiệu quả

4. Elbow Method tìm số lượng cụm tối ưu (k)

```

34 # Tìm số cluster tối ưu bằng Elbow Method
35 inertia = []
36 K_range = range(1, 11)
37 for k in K_range:
38     kmeans = KMeans(n_clusters=k, n_init=10, random_state=42)
39     kmeans.fit(x_scaled)
40     inertia.append(kmeans.inertia_)
41
42 # Vẽ biểu đồ Elbow
43 plt.figure(figsize=(8, 5))
44 plt.plot(K_range, inertia, 'bo-')
45 plt.xlabel('Number of Clusters (k)')
46 plt.ylabel('Inertia')
47 plt.title('Elbow Method for Optimal k')
48 plt.grid()
49 plt.show()
50

```

- Chạy thuật toán K-means với các giá trị $k = 1$ đến 10
- Ghi lại độ biến thiên (*inertia*) của từng mô hình
- Vẽ biểu đồ Elbow để chọn số lượng cụm hợp lý (nơi "gập khúc")

5. Tiến hành phân cụm với k cụ thể (ví dụ: k=3)

```

51 # Chọn số cluster – ví dụ k=3 từ Elbow Method
52 k = 3
53 kmeans = KMeans(n_clusters=k, n_init=10, random_state=42)
54 clusters = kmeans.fit_predict(x_scaled)
55
56 # Thêm cột cluster vào dataframe
57 df['cluster'] = clusters
58

```

- Dùng K-means để phân chia dữ liệu thành $k=3$ nhóm

- Gán kết quả vào cột `Cluster` trong dataframe để tiện theo dõi nhóm của từng cầu thủ

6. Giảm chiều dữ liệu với PCA xuống 2 chiều

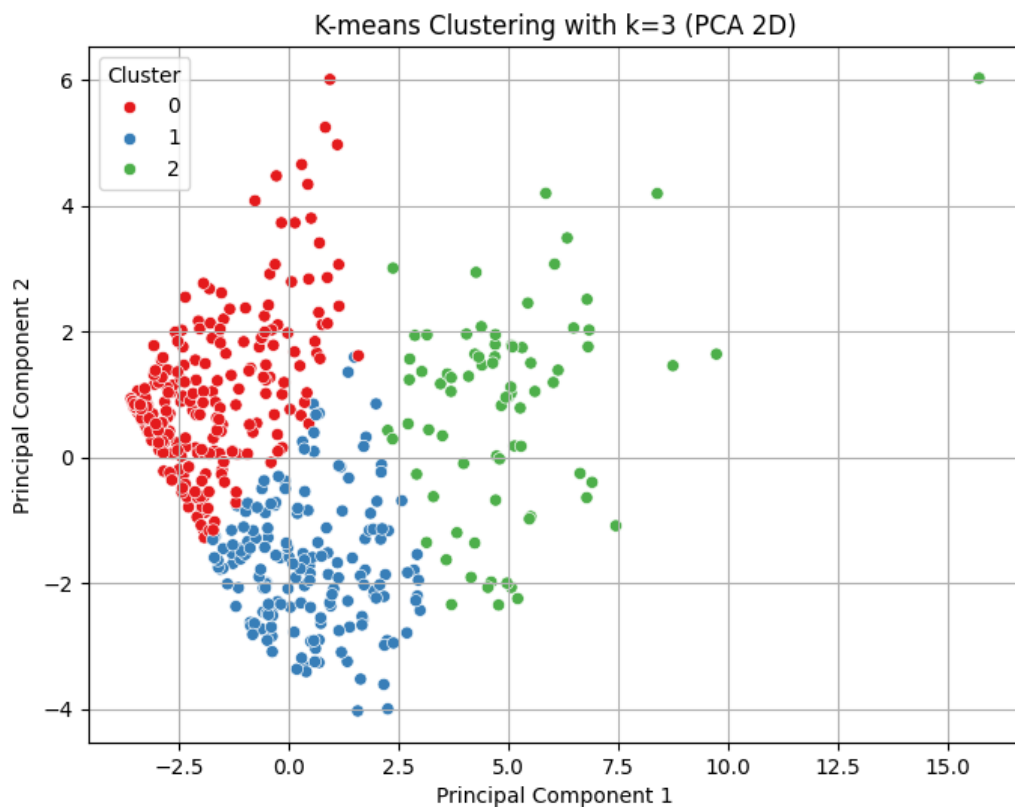
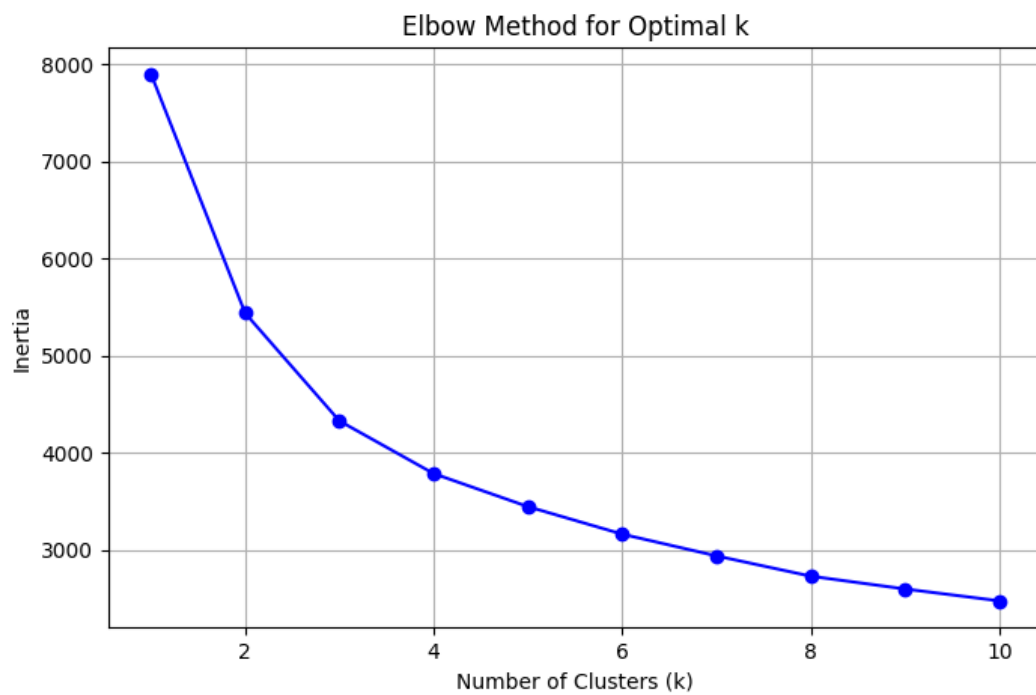
```
58
59 # Giảm chiều dữ liệu xuống 2D với PCA
60 pca = PCA(n_components=2)
61 X_pca = pca.fit_transform(X_scaled)
62
63 # Thêm PCA vào dataframe
64 df['PCA1'] = X_pca[:, 0]
65 df['PCA2'] = X_pca[:, 1]
66
```

- Dùng PCA để chuyển dữ liệu đa chiều thành 2 chiều (`PCA1`, `PCA2`)
- Giúp trực quan hóa dữ liệu phân cụm dễ hiểu hơn

7. Vẽ biểu đồ phân cụm cầu thủ theo PCA

```
67 # Vẽ scatter plot các nhóm
68 plt.figure(figsize=(8, 6))
69 sns.scatterplot(data=df, x='PCA1', y='PCA2', hue='Cluster', palette='Set1')
70 plt.title(f'K-means Clustering with k={k} (PCA 2D)')
71 plt.xlabel('Principal Component 1')
72 plt.ylabel('Principal Component 2')
73 plt.legend(title='Cluster')
74 plt.grid(True)
75 plt.show()
76 |
```

- Vẽ biểu đồ phân tán (`scatter plot`) thể hiện các cụm cầu thủ trên mặt phẳng 2D
- Mỗi màu đại diện cho 1 cụm → dễ so sánh cách các cầu thủ được nhóm lại theo phong cách, vai trò, hoặc năng lực



Bài 4 – Phần 1: Thu thập giá trị chuyển nhượng cầu thủ

1. Import thư viện cần thiết

```
1 import pandas as pd
2 from selenium import webdriver
3 from selenium.webdriver.common.by import By
4 from selenium.webdriver.chrome.service import Service
5 from selenium.webdriver.chrome.options import Options
6 import time
```

- **pandas**: xử lý dữ liệu
- **selenium**: tự động hóa việc truy cập và lấy dữ liệu từ website
- **time.sleep**: tạm dừng để website kịp tải nội dung

2. Đọc dữ liệu đầu vào và lọc cầu thủ hợp lệ

```
8 # --- Phần 1: Thu thập giá trị chuyển nhượng cầu thủ ---
9
10 # Đọc dữ liệu từ results.csv
11 df = pd.read_csv("results.csv")
12 df.columns = df.columns.str.lower() # Đảm bảo cột được viết thường
13
14 # Lọc cầu thủ có thời gian thi đấu > 900 phút
15 df_filtered = df[df["minutes"] > 900]
16 players_needed = set(df_filtered["player"].str.strip())
17
```

- Đọc dữ liệu từ file `results.csv` đã có từ các phần trước
- Lọc các cầu thủ có số phút thi đấu > 900 phút
- Tạo danh sách tên cầu thủ cần thu thập giá trị chuyển nhượng

3. Cấu hình Selenium với Chrome headless

```
18 # Cấu hình Selenium
19 chrome_options = Options()
20 chrome_options.add_argument("--headless")
21 chrome_options.add_argument("--no-sandbox")
22 chrome_options.add_argument("--disable-dev-shm-usage")
23 service = Service(executable_path=r"C:\Users\VIETTELSTORE\Downloads\chromedriver-win64\chromedriver-win64\chromedriver.exe")
24 driver = webdriver.Chrome(options=chrome_options, service=service)
25
```

- Dùng Chrome ở chế độ **headless** để không cần mở trình duyệt
- Chỉ định đường dẫn đến **ChromeDriver** để Selenium điều khiển trình duyệt

4. Truy cập và thu thập dữ liệu từ nhiều trang


```

26 # URL của các trang cần truy cập
27 base_url = "https://www.footballtransfers.com/us/values/players/most-valuable-soccer-players/playing-in-uk-premier-league/"
28
29 # Danh sách để lưu kết quả
30 transfer_data = []
31
32 # Lấy dữ liệu từ các trang 1 đến 22
33 for page in range(1, 23):
34     driver.get(f"{base_url}{page}")
35     time.sleep(3)
36
37 # Lấy tất cả tên cầu thủ và giá trị chuyển nhượng
38 players = driver.find_elements(By.XPATH, "//td[contains(@class, 'td-player')]/a")
39 values = driver.find_elements(By.XPATH, "//td[contains(@class, 'text-center')]/span[contains(@class, 'player-tag')]")
40
41 for player, value in zip(players, values):
42     player_name = player.text.strip()
43     transfer_value = value.text.strip()
44
45     # Kiểm tra nếu cầu thủ có trong danh sách đã lọc
46     if player_name.lower() in [p.lower() for p in players_needed]:
47         transfer_data.append({
48             "Player": player_name,
49             "Value": transfer_value
50         })
51

```

- Duyệt qua 22 trang danh sách cầu thủ Premier League trên FootballTransfers
- Mỗi trang:
 - Lấy danh sách tên cầu thủ (td-player)
 - Lấy giá trị chuyển nhượng (player-tag)
- Chỉ giữ lại cầu thủ có tên trong danh sách lọc `players_needed`
- Thêm dữ liệu vào danh sách `transfer_data`

5. Lưu dữ liệu vào file CSV và đóng trình duyệt

```

52 # Lưu kết quả vào file CSV
53 transfer_df = pd.DataFrame(transfer_data)
54 transfer_df.to_csv("transfer_values.csv", index=False)
55
56 # Đóng trình duyệt
57 driver.quit()
58
59 print("Lấy dữ liệu thành công và lưu vào file 'transfer_values.csv'.")

```

Câu 4 – Phần 2 - Xây dựng mô hình định giá cầu thủ:

1. Cài đặt thư viện và môi trường

```

1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  import warnings
6  import os
7  import logging
8  from sklearn.model_selection import train_test_split, RandomizedSearchCV
9  from sklearn.preprocessing import StandardScaler
10 from sklearn.compose import ColumnTransformer
11 from sklearn.pipeline import Pipeline
12 from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
13 from scipy.stats import spearmanr, randint, uniform
14 from math import sqrt
15 import xgboost as xgb
16

```

- Gọi các thư viện cần thiết để xử lý dữ liệu (pandas, numpy), vẽ biểu đồ (matplotlib, seaborn), và kiểm soát môi trường thực thi (warnings, logging, os).
- Import các thành phần học máy từ `scikit-learn` để:
 - Tiền xử lý dữ liệu.
 - Xây dựng pipeline học máy.
 - Tối ưu mô hình.
 - Đánh giá độ chính xác dự đoán.
 - Sử dụng mô hình hồi quy XGBoost.

```

17  # Cài đặt môi trường và hiển thị
18  warnings.filterwarnings("ignore")
19  os.environ["TF_CPP_MIN_LOG_LEVEL"] = "3"
20  logging.getLogger('sklearn').setLevel(logging.CRITICAL)
21  plt.style.use('ggplot')
22

```

Giảm bớt cảnh báo không cần thiết, tối ưu hiển thị và chọn phong cách vẽ biểu đồ.

2. Load dữ liệu

```

23  # --- Load dữ liệu ---
24  def load_data():
25      features = pd.read_csv("results.csv")
26      values = pd.read_csv("transfer_values.csv")
27
28      # Chuẩn hóa tên cầu thủ
29      for df in [features, values]:
30          df["Player"] = df["Player"].str.strip().str.lower()
31
32      # Merge theo "Player"
33      merged = pd.merge(features, values, on="Player", how="inner", suffixes=('', '_drop'))
34      merged = merged.filter(regex='^(?!.*_drop)') # Bỏ các cột bị trùng
35
36      return merged
37

```

Hàm này đọc dữ liệu từ hai file CSV, chuẩn hóa tên cầu thủ để tránh lỗi khi merge. Kết quả là một DataFrame đã ghép nối đầy đủ thông tin.

3. Tiền xử lý dữ liệu

```
39 def preprocess_data(df):
40     # Chuyển đổi giá trị từ dạng "€10m" -> 10_000_000
41     def convert_value(value):
42         try:
43             value = str(value).lower().replace("€", "").replace(",", "").strip()
44             if 'm' in value:
45                 return float(value.replace('m', '')) * 1e6
46             if 'k' in value:
47                 return float(value.replace('k', '')) * 1e3
48             return float(value)
49         except:
50             return np.nan
51
52     df["Value"] = df["Value"].apply(convert_value)
53     df = df[df["Value"] > 1000] # Loại bỏ cầu thủ không có giá trị thực
54
55     # Chuyển các cột bắt buộc sang số nếu chưa đúng kiểu
56     numeric_cols = df.select_dtypes(include=np.number).columns.tolist()
57     for col in ['Goals', 'Assists', 'Age']:
58         if col in df.columns and col not in numeric_cols:
59             df[col] = pd.to_numeric(df[col], errors='coerce')
60             numeric_cols.append(col)
61
62     # Thêm đặc trưng mới
63     df['Goals_per_Age'] = df['Goals'] / df['Age']
64     df['Assists_per_Age'] = df['Assists'] / df['Age']
65     df['Log_Transfer_Value'] = np.log1p(df['Value'])
66
67     # Xử lý missing value
68     df[['Goals_per_Age', 'Assists_per_Age']] = df[['Goals_per_Age', 'Assists_per_Age']].fillna(0)
69     df[numeric_cols] = df[numeric_cols].fillna(df[numeric_cols].median())
70
71     return df
72
```

Hàm xử lý dữ liệu như chuyển "€10m" thành 10000000, lọc giá trị bất hợp lệ, chuẩn hóa số liệu, thêm đặc trưng chia theo độ tuổi và logarit hóa giá trị chuyển nhượng để giúp mô hình học tốt hơn.

4. Xây dựng pipeline mô hình

```
73 # --- Xây dựng pipeline model ---
74 def build_model(numeric_features):
75     preprocessor = ColumnTransformer(transformers=[
76         ('num', StandardScaler(), numeric_features)
77     ])
78
79     model = Pipeline([
80         ('preprocessor', preprocessor),
81         ('regressor', xgb.XGBRegressor(objective='reg:squarederror', random_state=42, n_jobs=-1, verbosity=0))
82     ])
83
84     return model
85
```

Pipeline gồm hai bước: (1) Chuẩn hóa dữ liệu số bằng StandardScaler, (2) Huấn luyện mô hình hồi quy XGBoost để dự đoán log giá trị chuyển nhượng.

5. Đánh giá mô hình

```

86 # --- Đánh giá mô hình ---
87 def evaluate_model(model, X, y, log_target=True):
88     pred = model.predict(X)
89     if log_target:
90         y_true = np.expm1(y)
91         pred = np.expm1(pred)
92     else:
93         y_true = y
94
95     rmse = sqrt(mean_squared_error(y_true, pred))
96     mae = mean_absolute_error(y_true, pred)
97     r2 = r2_score(y_true, pred)
98     spear_corr, _ = spearmanr(y_true, pred)
99
100     print(f" R2 Score: {r2:.4f}")
101     print(f" RMSE: {rmse / 1e6:.4f} triệu €")
102     print(f" MAE: {mae / 1e6:.4f} triệu €")
103     print(f" Spearman Corr: {spear_corr:.4f}")
104

```

Hàm đánh giá mô hình bằng các chỉ số:

- R^2 : mức độ giải thích biến mục tiêu.
- RMSE & MAE: sai số dự đoán.
- Hệ số tương quan Spearman: độ liên hệ giữa giá trị thật và dự đoán.

6. Huấn luyện và tối ưu mô hình

```

106 if __name__ == "__main__":
107     df = load_data()
108     df = preprocess_data(df)
109
110     # Tách biến đầu vào và đầu ra
111     X = df.drop(columns=['Value', 'Log_Transfer_Value', 'Player'])
112     y = df['Log_Transfer_Value']
113     numeric_features = X.select_dtypes(include=np.number).columns.tolist()
114
115     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
116
117     model = build_model(numeric_features)
118
119     # Random search để tối ưu tham số
120     param_dist = {
121         'regressor__n_estimators': randint(100, 400),
122         'regressor__learning_rate': uniform(0.05, 0.2),
123         'regressor__max_depth': randint(4, 8),
124         'regressor__subsample': uniform(0.7, 0.3),
125         'regressor__colsample_bytree': uniform(0.7, 0.3)
126     }
127
128     search = RandomizedSearchCV(model, param_dist, n_iter=15, cv=3,
129                                scoring='neg_root_mean_squared_error',
130                                random_state=42, n_jobs=-1)
131     search.fit(X_train, y_train)
132     best_model = search.best_estimator_
133
134     print("\nĐánh giá trên tập huấn luyện:")
135     evaluate_model(best_model, X_train, y_train)
136
137     print("\nĐánh giá trên tập kiểm tra:")
138     evaluate_model(best_model, X_test, y_test)
139

```

- Tách tập dữ liệu thành huấn luyện và kiểm tra (80/20).
- Áp dụng **RandomizedSearchCV** để chọn tham số tốt nhất cho mô hình XGBoost.
- In kết quả đánh giá cả trên tập huấn luyện và kiểm tra.