

HỌC VIỆN CÔNG NGHỆ BƯU CHÍNH VIỄN THÔNG
BỘ MÔN LẬP TRÌNH PYTHON



NHÓM 11

BÁO CÁO BÀI TẬP LỚN PYTHON

NGUYỄN VIỆT TUẤN ANH – B22DCCN038

Hà Nội – 2024

Mục lục

I. Công nghệ và thuật toán sử dụng.....	3
1. Công nghệ và các thư viện sử dụng.....	3
2. Các thuật toán và thư viện sử dụng.	4
II. Thực hiện	6
1. Thu thập dữ liệu thống kê của tất cả các cầu thủ có số phút thi đấu nhiều hơn 90 phút tại giải bóng đá ngoại hạng Anh mùa 2023-2024.....	6
2. Tìm dữ liệu theo các yêu cầu sau :	10
a. tìm top ba cầu thủ có điểm cao nhất và thấp nhất ở mỗi chỉ số.	10
b. Tìm trung vị của mỗi chỉ số. Tìm trung bình và độ lệch chuẩn của mỗi chỉ số cho các cầu thủ trong toàn giải và của mỗi đội. Ghi kết quả ra file results2.csv.	12
c. Vẽ histogram phân bố của mỗi chỉ số của các cầu thủ trong toàn giải và mỗi đội.....	14
d. Tìm đội bóng có chỉ số điểm số cao nhất ở mỗi chỉ số. Theo bạn đội nào có phong độ tốt nhất giải ngoại Hạng Anh mùa 2023-2024.....	16
3. Sử dụng thuật toán K-means để phân loại các cầu thủ thành các nhóm có chỉ số giống nhau. Sử dụng thuật toán PCA, giảm số chiều dữ liệu xuống 2 chiều, vẽ hình phân cụm các điểm dữ liệu trên mặt 2D. Viết chương trình python vẽ biểu đồ rada (radar chart) so sánh cầu thủ.....	17
III. Tài liệu tham khảo.....	21

I. Công nghệ và thuật toán sử dụng.

1. Công nghệ và các thư viện sử dụng

Selenium là một bộ công cụ tự động hóa trình duyệt web, cho phép người dùng tự động hóa các thao tác trên trình duyệt, như nhấp chuột, điền biểu mẫu và thu thập dữ liệu từ các trang web.

Chức năng:

- Tự động hóa kiểm tra giao diện người dùng.
- Lấy dữ liệu từ các trang web mà không có API.
- Hỗ trợ đa trình duyệt, bao gồm Chrome, Firefox, và Safari.
- Chạy JavaScript và xử lý phần tử động: Selenium có thể thực thi mã JavaScript trên trang web và xử lý các phần tử động được tải qua AJAX.

Beautifulsoup: là một thư viện Python dùng để phân tích cú pháp HTML và XML, giúp dễ dàng trích xuất dữ liệu từ các trang web. Nó rất hữu ích khi bạn cần lấy dữ liệu từ trang web mà không cần phải điều khiển trình duyệt (như Selenium). BeautifulSoup đặc biệt phù hợp với các trang tĩnh hoặc khi đã có mã HTML của trang cần xử lý.

Chức năng:

- Tìm kiếm và lấy dữ liệu từ các tài liệu HTML/XML.
- Tích hợp tốt với các thư viện phân tích cú pháp khác như `html.parser` (mặc định), `lxml`, hoặc `html5lib` để phân tích cú pháp nhanh hơn và chính xác hơn.
- Cho phép tìm kiếm các thẻ HTML bằng nhiều cách như tên thẻ, id, class, hoặc các thuộc tính tùy ý khác.
- Cung cấp các phương thức như `find()`, `find_all()`, và `select()` để tìm và lấy dữ liệu một cách hiệu quả.

Pandas: là một thư viện phân tích và xử lý dữ liệu mạnh mẽ trong Python, rất phổ biến cho các tác vụ liên quan đến dữ liệu trong khoa học dữ liệu và học máy. Với Pandas, bạn có thể dễ dàng thao tác với dữ liệu dạng bảng, tính toán, lọc, và chuyển đổi dữ liệu thành các định dạng khác.

Đặc điểm nổi bật của Pandas:

- Cấu trúc dữ liệu linh hoạt:
 - **DataFrame**: Cấu trúc dữ liệu hai chiều (giống như bảng) có thể chứa nhiều kiểu dữ liệu khác nhau.

- **Series:** Cấu trúc dữ liệu một chiều, thường được dùng để biểu diễn một cột hoặc một hàng của DataFrame.
- Dễ dàng xử lý dữ liệu: Pandas cung cấp các hàm tiện lợi để lọc, gom nhóm, tính toán thống kê, xử lý dữ liệu thiếu, và chuyển đổi dữ liệu.
- Đọc và ghi dữ liệu:
 - Pandas hỗ trợ nhiều định dạng dữ liệu, bao gồm CSV, Excel, SQL, JSON, HTML, và các định dạng khác.
 - Cho phép lưu và tải dữ liệu từ các nguồn khác nhau, phù hợp cho các tác vụ trích xuất dữ liệu.
- Chuyển đổi dữ liệu nhanh chóng: Dễ dàng thay đổi định dạng dữ liệu, chuyển từ bảng thành mảng, từ mảng thành bảng, gộp hoặc chia nhỏ dữ liệu.

Matplotlib: là thư viện trực quan hóa dữ liệu mạnh mẽ trong Python, thường được sử dụng để tạo các loại biểu đồ khác nhau như biểu đồ đường, cột, tròn, tán xạ, và nhiều loại khác. Thư viện này rất hữu ích cho việc trình bày và phân tích dữ liệu, giúp bạn dễ dàng nhận ra các xu hướng và mẫu trong dữ liệu.

Đặc điểm nổi bật của Pandas:

- Đa dạng kiểu biểu đồ: Matplotlib cung cấp nhiều loại biểu đồ như biểu đồ đường, biểu đồ cột, biểu đồ tròn, biểu đồ tán xạ, và các biểu đồ dạng 3D.
- Tùy chỉnh linh hoạt: có thể tùy chỉnh màu sắc, kiểu đường, kiểu điểm, trục tọa độ, nhãn, tiêu đề, chú thích và nhiều yếu tố khác.
- Hoạt động tốt với các thư viện khác: Matplotlib có thể dễ dàng tích hợp với các thư viện như Pandas và NumPy để trực quan hóa dữ liệu từ các DataFrame hoặc mảng số.

2. Các thuật toán và thư viện sử dụng.

Elbow: Phương pháp Elbow (khủy tay) là một kỹ thuật thường được sử dụng trong KMeans để xác định số lượng cụm (k) tối ưu. Ý tưởng là chạy thuật toán KMeans với các giá trị khác nhau của k và sau đó tính toán tổng lỗi bình phương nội cụm (Within-Cluster Sum of Squares - WCSS) cho mỗi giá trị k. Tổng lỗi bình phương nội cụm thể hiện khoảng cách từ các điểm trong cụm đến trung tâm của cụm đó.

Điểm *khủy tay* là điểm mà ở đó tốc độ suy giảm của *hàm biến dạng* sẽ thay đổi nhiều nhất.

Phương pháp Elbow là một phương pháp thường được sử dụng để lựa chọn số lượng cụm phân chia hợp lý dựa trên biểu đồ.

Tuy nhiên có một số trường hợp chúng ta sẽ không dễ dàng phát hiện vị trí của Elbow, đặc biệt là đối với những bộ dữ liệu mà quy luật phân cụm không thực sự dễ dàng được phát hiện.

PCA: hay còn gọi là phân tích thành phần chính, là một phương pháp giảm chiều dữ liệu rất phổ biến trong học máy và thống kê. PCA giúp bạn giảm số chiều của dữ liệu nhiều chiều xuống ít chiều hơn trong khi vẫn giữ lại phần lớn thông tin quan trọng, từ đó đơn giản hóa dữ liệu và giúp việc phân tích dễ dàng hơn.

Thuật toán và cách hoạt động:

Tính vector kỳ vọng của toàn bộ dữ liệu:

$$\mathbf{x} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

Trừ mỗi điểm dữ liệu với vector kỳ vọng của toàn bộ dữ liệu.

$$\hat{\mathbf{x}}_n = \mathbf{x}_n - \mathbf{x}$$

Tính ma trận hiệp phương sai.

$$\mathbf{S} = \frac{1}{N} \hat{\mathbf{X}} \hat{\mathbf{X}}^T$$

Tính các giá trị riêng và vector riêng có norm bằng 1 của ma trận hiệp phương sai, sắp xếp theo thứ tự giảm dần của giá trị riêng.

Chọn K vector riêng ứng với K giá trị riêng lớn nhất để xây dựng ma trận \mathbf{U}_K có các cột tạo thành một hệ trực giao. K vectors này được gọi là các thành phần chính, tạo thành một không gian con gần với phân bố của dữ liệu ban đầu.

Chiều dữ liệu ban đầu theo không gian con tìm được.

Dữ liệu mới là tọa độ của các điểm dữ liệu trên không gian mới.

$$\mathbf{Z} = \mathbf{U}_K^T \hat{\mathbf{X}}$$

KMeans: là thuật toán phân cụm dữ liệu được đặc trưng bởi một *tâm* (*centroid*). *tâm* là điểm đại diện nhất cho một cụm và có giá trị bằng

trung bình của toàn bộ các quan sát nằm trong cụm. Chúng ta sẽ dựa vào khoảng cách từ mỗi quan sát tới các *tâm* để xác định nhãn cho chúng trùng thuộc về *tâm* gần nhất.

Cụ thể các bước của thuật toán k-Means được tóm tắt như sau:

Khởi tạo ngẫu nhiên k tâm cụm $k_1, k_2, k_3 \dots$

Lặp lại quá trình cập nhật tâm cụm cho tới khi dừng:

- Xác định nhãn cho từng điểm dữ liệu c_i dựa vào khoảng cách tới từng tâm cụm:

$$c_i = \arg \min_j \|\mathbf{x}_i - \mu_j\|_2^2$$

- Tính toán lại tâm cho từng cụm theo trung bình của toàn bộ các điểm dữ liệu trong một cụm:

$$\mu_j := \frac{\sum_{i=1}^n \mathbf{1}(c_i = j) \mathbf{x}_i}{\sum_{i=1}^n \mathbf{1}(c_i = j)}$$

II. Thực hiện

1. Thu thập dữ liệu thống kê của tất cả các cầu thủ có số phút thi đấu nhiều hơn 90 phút tại giải bóng đá ngoại hạng Anh mùa 2023-2024.

```
def openDriver(url):  
    driver.get(url)  
    time.sleep(5)  
    html = driver.page_source  
    soup = BeautifulSoup(html, 'html.parser')  
    return soup
```

Hàm OpenDriver nhận vào một url được khai báo ở trên như: urlRoot, goalKeepingUrl, shootingUrl, passingUrl, passTypeUrl, goalAndShotUrl, defensiveUrl, possessionUrl, playingTimeUrl, miscellaneousUrl. Sẽ mở các url bằng Selenium, tải trang, và trả về mã HTML của trang dưới dạng một đối tượng BeautifulSoup.

```
def findInStats(texts, s):  
    player = []
```

```

td = texts.find_all('td', attrs = {"data-stat":s})
for x in td:
    if len(str(x)) != 0:
        player.append(x.text.strip())
return player

```

Hàm findInStats nhận vào một hai tham số texts và s. Texts là dữ liệu HTML bảng số liệu thống kê tương ứng với từng url ở trên. S là một chuỗi kí tự thuộc tính cột cần tìm trong bảng. findInStats sẽ tìm toàn bộ dữ liệu của cột có thuộc tính là S và sẽ trả về một danh sách dữ liệu của thuộc tính đó ở dạng chuỗi ký tự.

```

res['player'] = findInStats(texts, 'player')
res['nationality'] = findInStats(texts, 'nationality')
res['team'] = findInStats(texts, 'team')
res['position'] = findInStats(texts, 'position')
res['age'] = findInStats(texts, 'age')

```

Res, res1, res2,..., res9 là các Dictionary dùng để lưu dữ liệu lấy từ từng url. Với khóa là tên thuộc tính hay tên cột và giá trị là danh sách các giá trị của thuộc tính đó.

```

df_stats = pd.DataFrame(res)
df_goalKeeper = pd.DataFrame(res1)
df_shooting = pd.DataFrame(res2)
df_passing = pd.DataFrame(res3)
df_typePass = pd.DataFrame(res4)
df_goalAndShot = pd.DataFrame(res5)
df_defensive = pd.DataFrame(res6)
df_possession = pd.DataFrame(res7)
df_playingTime = pd.DataFrame(res8)
df_miscellaneous = pd.DataFrame(res9)

```

Sử dụng Pandas để tạo mới các dataframe từ các dictionary ở trên lần lượt có chín dataframe: df_stats, df_goalKeeper, df_shooting, df_passing, df_typePass, df_goalAndShot, df_defensive, df_possession, df_playingTime, df_miscellaneous.

```
df_result = df_stats.merge(df_goalKeeper, on = ['player','team'], how='outer')
df_result =df_result.merge(df_shooting, on=['player','team'], how = 'outer')
df_result =df_result.merge(df_passing, on=['player','team'], how = 'outer')
df_result = df_result.merge( df_typePass, on=['player','team'], how = 'outer')
df_result = df_result.merge( df_goalAndShot, on=['player','team'], how = 'outer')
df_result = df_result.merge( df_defensive, on=['player','team'], how = 'outer')
df_result = df_result.merge( df_possession, on=['player','team'], how = 'outer')
df_result = df_result.merge(df_playingTime, on=['player','team'], how = 'outer')
df_result = df_result.merge( df_miscellaneous, on=['player','team'], how = 'outer')
```

Hợp nhất các dataframe ở trên thành một dataframe là df_result. Dòng đầu tiên sẽ hợp nhất df_stat và df_goalKeeper dựa trên hai cột là cột player và cột team. Sử dụng how với giá trị outer có nghĩa là một giá trị của cặp player và team trong hai dataframe sẽ được giữ lại, ngay cả khi một số giá trị chỉ có mặt ở một trong hai dataframe.

Tiếp tục hợp nhất các dataframe tiếp theo vào df_result. Nếu cầu thủ ko có giá trị trong một dataframe thì các cột tương ứng sẽ có giá trị NaN. Df_result sẽ là dataframe tổng hợp chứa tất cả thông tin từ các dataframe ban đầu. Mỗi hàng trong df_result đại diện cho một cầu thủ và mỗi cột là một số liệu khác nhau thuộc về cầu thủ này.

Các đoạn code tiếp theo sẽ thực hiện việc chuẩn bị và xử lý dữ liệu trước khi lưu vào file csv.

```
df_result['minutes'] = df_result['minutes'].str.replace(',','')
df_result['minutes'] = pd.to_numeric(df_result['minutes'], errors='coerce')
df_result = df_result[df_result['minutes'] >= 90]
```

Chức năng của đoạn code trên là lọc ra các cầu thủ có số phút thi đấu lớn hơn hoặc bằng chín mươi phút. Dòng đầu tiên giúp chuẩn hóa dữ liệu của cột minutes về dạng mỗi chuỗi ký tự và không có dấu phẩy. Dòng tiếp theo làm

chức năng chuyển đổi cột minutes sang kiểu số (numeric) với errors có giá trị là coerce. Nếu có giá trị nào không thể chuyển được nó sẽ thành NaN. Cuối cùng lọc ra các cầu thủ có số phút thi đấu lớn hơn hoặc bằng chín mươi phút.

```
df_result.fillna('n/a', inplace=True)
```

```
df_result.replace(", 'n/a', inplace=True)
```

Hàm fillna sẽ điền giá trị 'n/a' và các ô có giá trị NaN. Sau đó hàm replace sẽ thay thế các ô trống mà không phải NaN thành n/a.

```
df_result = df_result.sort_values(by=['player', 'age'], ascending=[True, False])
```

```
df_result.insert(0, 'STT', range(1, len(df_result) + 1))
```

```
df_result.to_csv('C:/result.csv', index=False)
```

Đây là bước sắp xếp và chuẩn hóa cuối cùng trước khi lưu thành csv. Hàm `sort_values` sẽ có hai tham số là `by`: cột dữ liệu muốn sắp xếp và `ascending`: True là sắp xếp tăng dần và False là sắp xếp giảm dần.

Sau khi hợp nhất các dataframe thì số thứ tự khi chuyển qua file csv ko đúng trật tự. Hàm insert sẽ có chức năng thêm một cột là STT vào dataframe bắt đầu từ vị trí 0 và có giá trị tăng dần bắt đầu từ 1.

Cuối cùng lưu df_result vào file csv thêm đường dẫn trong hàm to_csv và thêm tham số index = false để bỏ cột chỉ mục khi lưu file.

<

hình 1: file result.csv

2. Tìm dữ liệu theo các yêu cầu sau :

a. tìm top ba cầu thủ có điểm cao nhất và thấp nhất ở mỗi chỉ số.

```
def load_data(filepath):  
    data = pd.read_csv(filepath)  
    return data
```

Lấy dữ liệu từ hàm load_data với tham số filepath là đường dẫn đến nơi lưu file results.csv. Trong đó dùng hàm pd.read_csv(filepath) có chức năng đọc dữ liệu từ một file csv.

```
def find_top_3(data, columns):  
    for col in columns:  
        if pd.api.types.is_numeric_dtype(data[col]):  
            print(f"Top 3 highest players for {col}:")  
            print(data.nlargest(3, col)[['player', col]])  
            print()  
            print(f"Top 3 lowest players for {col}:")  
            print(data.nsmallest(3, col)[['player', col]])  
            print("\n" + "-" * 50 + "\n")
```

Hàm trên nhận hai tham số là data: tất cả dữ liệu từ file csv và columns: tất cả các cột trừ bốn cột name, team, nation, position không có dữ liệu thuộc kiểu numeric bằng lệnh `columns = data.columns[3:]`. Tạo ra hai dictionary để lưu lần lượt 3 cầu thủ có điểm cao nhất và thấp nhất ở mỗi chỉ số theo kiểu key = dictionary và value = dataframe với ba hàng dữ liệu.

Duyệt lần lượt qua các cột và kiểm tra xem cột đó có kiểu dữ liệu là số không thông qua hàm của Pandas: `pd.api.types.is_numeric_dtype()`. Nếu thỏa mãn top ba cầu thủ có điểm cao nhất và thấp nhất được in ra màn hình qua hai hàm sau:

`data.nlargest(3, col)[['player', col]]` : là hàm của Pandas dùng để tìm ra n hàng có giá trị lớn nhất theo cột columns và trả về dữ liệu dạng một dataframe. Trong trường hợp này, n=3 nghĩa là lấy 3 hàng có giá trị lớn nhất. `[['player',`

`col]]` có chức năng lọc lại chỉ lấy hai cột để giữ lại thông tin về tên cầu thủ và chỉ số đang xét làm cho kết quả gọn gàng và dễ đọc hơn.

Tương tự với hàm bên trên thì hàm `data.nsmallest(3, col)[['player', col]]` sẽ lấy top ba cầu thủ được chấm điểm thấp nhất ở chỉ đang xét.

```
-----
Top 3 highest players for sca_per90:
      player  sca_per90
302  Manor Solomon    8.59
13   Alex Iwobi      7.77
264  Kevin De Bruyne  7.36

Top 3 lowest players for sca_per90:
      player  sca_per90
2   Aaron Ramsdale    0.0
35  Angelo Ogbonna    0.0
120 David Ozoh        0.0

-----

Top 3 highest players for sca_passes_live:
      player  sca_passes_live
313  Martin Ødegaard    178
70   Bruno Fernandes    143
74   Bukayo Saka        137

Top 3 lowest players for sca_passes_live:
      player  sca_passes_live
2   Aaron Ramsdale    0
35  Angelo Ogbonna    0
120 David Ozoh        0

-----

Top 3 highest players for sca_passes_dead:
      player  sca_passes_dead
385  Pascal Groß        54
212  James Ward-Prowse  51
28   Andreas Pereira    48
```

hình 2.1: top three highest and lowest

```
-----

Top 3 highest players for gca:
      player  gca
39   Anthony Gordon    29
101  Cole Palmer       26
379  Ollie Watkins     26

Top 3 lowest players for gca:
      player  gca
1   Aaron Hickey    0
2   Aaron Ramsdale  0
3   Aaron Ramsey    0

-----

Top 3 highest players for gca_per90:
      player  gca_per90
381  Oscar Bobb      1.79
191  Ivan Perišić    1.75
85   Carney Chukwuemeka 1.59

Top 3 lowest players for gca_per90:
      player  gca_per90
1   Aaron Hickey    0.0
2   Aaron Ramsdale  0.0
3   Aaron Ramsey    0.0

-----

Top 3 highest players for gca_passes_live:
      player  gca_passes_live
101  Cole Palmer      20
313  Martin Ødegaard  18
410  Rodri            17
```

hình 2.2: top three highest and lowest

b. Tìm trung vị của mỗi chỉ số. Tìm trung bình và độ lệch chuẩn của mỗi chỉ số cho các cầu thủ trong toàn giải và của mỗi đội. Ghi kết quả ra file results2.csv.

Hàm `calculate_statistics(data)` gồm hai phần: tìm trung vị, trung bình, độ lệch chuẩn của tất cả các cầu thủ và tìm trung vị, trung bình, độ lệch chuẩn cho từng đội bóng.

```
tmp_df = data.copy()

club_list = data['team'].unique().tolist()

numeric_df = data.apply(pd.to_numeric, errors='coerce')

summary_stats = pd.DataFrame()

overall_stats = { }
```

`club_list` lưu một danh sách tên các câu lạc bộ. `numeric_df` là một dataframe chứa các dữ liệu đã được chuyển về kiểu số nếu ko phải kiểu số nó sẽ được chuyển thành NaN qua hàm `apply` của Pandas. Khởi tạo một dataframe rỗng `summary_stats` để lưu dữ liệu cuối cùng trước khi chuyển sang file results2.csv.

Tạo một từ điển rỗng `overall_stats` để lưu trữ các thông kê như trung vị (Median), trung bình (Mean), và độ lệch chuẩn (Standard Deviation) của tất cả các cầu thủ trong toàn bộ giải đấu cho từng chỉ số.

```
for col in numeric_df.columns:

    if not numeric_df[col].isnull().all():

        overall_stats[f'{col}_Median'] = numeric_df[col].median()

        overall_stats[f'{col}_Mean'] = numeric_df[col].mean()

        overall_stats[f'{col}_Std'] = numeric_df[col].std()
```

Duyệt lần lượt qua các cột chỉ số và tính trung vị, trung bình và độ lệch chuẩn bằng hàm `median()`, `mean()`, `std()` của thư viện pandas rồi lưu vào `overall_stats`.

```
overall_stats_df = pd.DataFrame([overall_stats], index=["all"])
summary_stats = pd.concat([summary_stats, overall_stats_df])
```

Tạo mới một dataframe chứa một phần tử với dữ liệu là `overall_stats` và chỉ định tên cho chỉ mục của dataframe này là 'all' để biểu diễn là đây là tổng

kê cho tất cả các đội. Dùng hàm `concat` của thư viện Pandas để nối hai dataframe vào nhau.

```
col_to_convert = numeric_df.columns.difference(['team', 'STT'])

numeric_df[col_to_convert]=numeric_df[col_to_convert].apply(pd.to_numeric, errors='coerce')
```

Do cần tính toán trung vị , trung bình và độ lệch chuẩn cho từng đội bóng nên hai dòng trên sẽ chịu trách nhiệm tạo ra một dataframe kiểu numeric chứa tất cả dữ liệu và thêm cột team để xác định các đội. `col_to_convert` sẽ chứa các cột cần chuyển sang kiểu số ngoại trừ hai cột ‘team’ và ‘stt’. Numeric_df sẽ chuyển tất cả các cột ngoại trừ hai cột trên sang kiểu số bằng hàm `apply`.

Duyệt lần lượt trong danh sách `club_list` chứa tên các đội bóng và lọc lấy ra dữ liệu của từng đội bóng `numeric_df[tmp_df['team'] == club]` kiểu dataframe . Tạo mới một biến kiểu từ điển `club_stats` để lưu từng giá trị trung vị, trung bình và độ lệch chuẩn của từng cột dữ liệu của mỗi đội bóng.

```
for col in club_data.columns:

    if not club_data[col].isnull().all():

        club_stats[f'{col}_Median'] = club_data[col].median()

        club_stats[f'{col}_Mean'] = club_data[col].mean()

        club_stats[f'{col}_Std'] = club_data[col].std()

club_stats_df = pd.DataFrame([club_stats], index=[club])

summary_stats = pd.concat([summary_stats, club_stats_df])

summary_stats.to_csv('C: /results2.csv', index_label="Team")
```

Duyệt qua từng cột để lấy các giá trị yêu cầu và lưu vào `club_stats`. Tiếp tục chuyển `club_stats` thành một dataframe với một trường dữ liệu duy nhất và ghép dataframe này vào `summary_stats`. Chuyển `summary_stats` thành file csv:

File

Home

Insert

Page Layout

Formulas

Data

Review

View

Help

Tell me what you want to do

</

hình 2.3: file results2.csv

c. Vẽ histogram phân bố của mỗi chỉ số của các cầu thủ trong toàn giải và mỗi đội.

Hàm `plot_histograms` nhận vào hai tham số là một dataframe và một danh sách các cột cần vẽ histogram. Columns ở đây đã được loại bỏ cột name, team, nation, position.

```
def plot_histograms(data, columns):
```

```
    for col in columns:
```

```
        if pd.api.types.is_numeric_dtype(data[col]):
```

```
            plt.figure()
```

```
            plt.hist(data[col].dropna(), bins=20, color='lightblue')
```

```
            plt.title(f'Histogram of {col}')
```

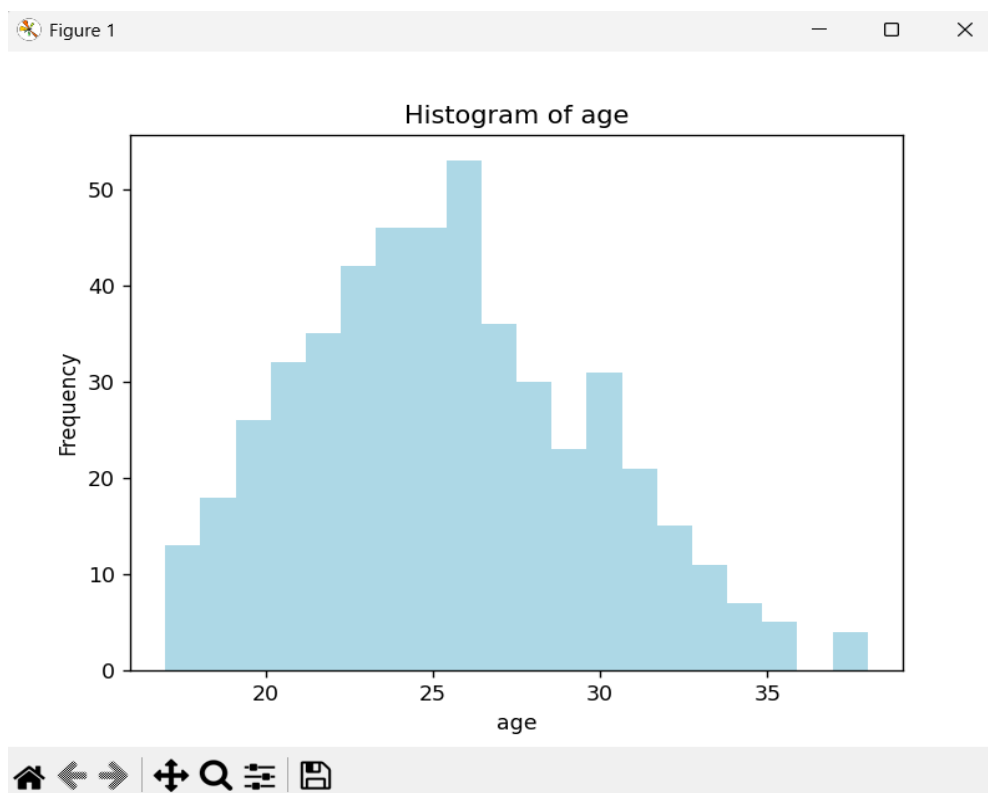
```
            plt.xlabel(col)
```

```
            plt.ylabel('Frequency')
```

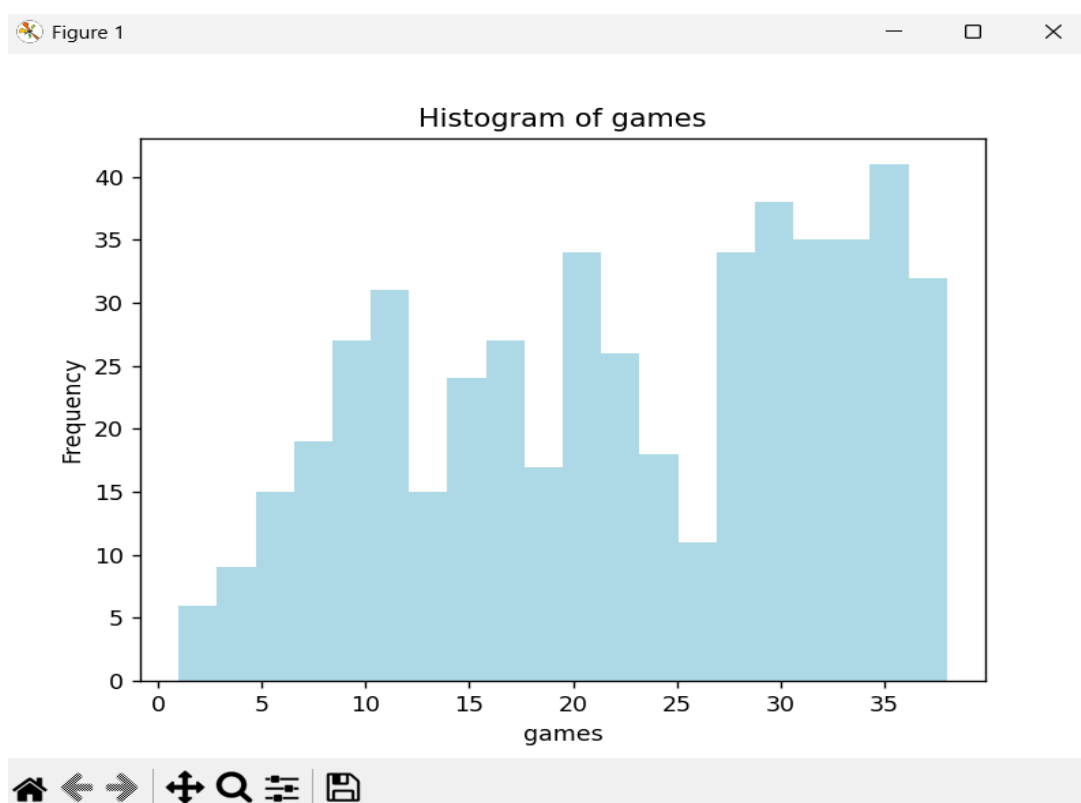
```
            plt.show()
```

Trong hàm `plot_histogram` duyệt qua tất cả các cột bỏ qua các cột không phải kiểu số và tạo mới một biểu đồ mới cho mỗi cột bằng `plt.figure()` của thư viện Matplotlib. Cài đặt các thông số cho biểu đồ với 20 khoảng dữ liệu, màu xanh nhạt và đồng thời loại bỏ các giá trị rỗng. Đặt nhãn cho trục x và trục y lần lượt là Frequency và tên cột của cột đang xét.

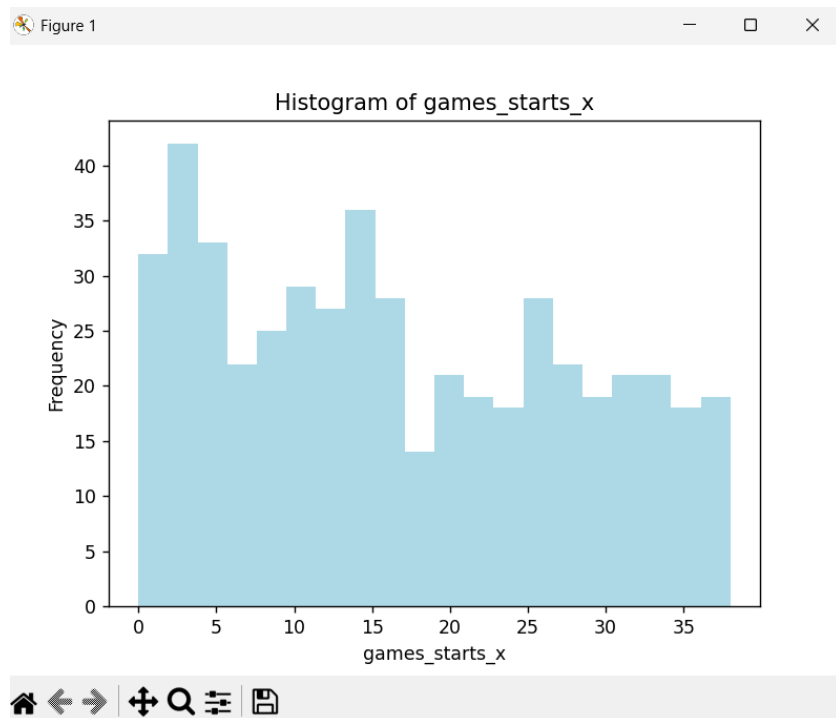
Một số biểu đồ trong tất cả dữ liệu:



hình 2.4: histogram of age



hình 2.5: histogram of games



hình 2.6: histogram of games starts

d. Tìm đội bóng có chỉ số điểm số cao nhất ở mỗi chỉ số. Theo bạn đội nào có phong độ tốt nhất giải ngoại Hạng Anh mùa 2023-2024.

```
def calculate_team_performance(data):
```

```
    numeric_data = data.select_dtypes(include='number')
```

```
    numeric_data['Team'] = data['team']
```

```
    team_performance = numeric_data.groupby('Team').mean()
```

```
    overall_mean = team_performance.mean(axis=1)
```

```
    best_team = overall_mean.idxmax()
```

```
    return best_team
```

Hàm `calculate_team_performance` nhận vào một tham số là một dataframe chứa toàn bộ dữ liệu. `Numeric_data` là một dataframe chỉ chứa các cột dữ liệu số từ tham số đầu vào và được thêm cột 'team'. Điều này cho phép ghép nhóm dữ liệu theo đội. `numeric_data.groupby('Team').mean()` sẽ tạo ra một dataframe mới với chỉ số là tên các đội và giá trị trung bình của các chỉ số tương ứng.

```
    overall_mean = team_performance.mean(axis=1)
```



```
best_team = overall_mean.idxmax()
```

Ở bước này, hàm tính trung bình của các trung bình của từng đội (theo hàng) để có được một giá trị trung bình tổng thể cho mỗi đội. Kết quả là một Series `overall_mean`, trong đó mỗi giá trị tương ứng với một đội. Đội có phong độ tốt nhất sẽ được xác định bằng hàm `idxmax()` hay đội có giá trị trung bình lớn nhất trong `overall_mean`.

Đội có phong độ cao nhất màu giải 2023-2024 giải bóng đá ngoại hạng anh là câu lạc bộ Manchester City.

```
44 Arijanet Muric          0

-----

Top 3 highest players for aerals_won_pct:
      player  aerals_won_pct
44  Arijanet Muric          100.0
50  Bart Verbruggen         100.0
81  Caoimh  n Kelleher       100.0

Top 3 lowest players for aerals_won_pct:
      player  aerals_won_pct
26  Anass Zaroury            0.0
31  Andros Townsend          0.0
60  Benson Manuel            0.0

-----

Best performing team: Manchester City
PS C:\Users\84976\Desktop\BTL_PTIT\PYTHON\bt1> |
```

h  nh 2.7: best performing team

3. Sử dụng thuật toán K-means để phân loại các cầu thủ thành các nhóm có chỉ số giống nhau. Sử dụng thuật toán PCA, giảm số chiều dữ liệu xuống 2 chiều, vẽ hình phân cụm các điểm dữ liệu trên mặt 2D. Viết chương trình python vẽ biểu đồ rada (radar chart) so sánh cầu thủ.

```
def load_data(filepath):

    data = pd.read_csv(filepath)

    return data

data = load_data('C:/Users/results.csv')
```

Hàm `load_data` nhận vào một đường dẫn đến tệp csv và tải dữ liệu vào một dataframe ở đây là `data`. Sau đó chuẩn hóa dữ liệu trước khi phân cụm:

```
numeric_data = data.select_dtypes(include='number')  
  
imputer = SimpleImputer(strategy='mean')  
  
numeric_data = imputer.fit_transform(numeric_data)
```

Tạo ra một dataframe `numeric_data` để lưu dữ liệu sau khi đã lọc ra các cột có kiểu dữ liệu số và đối tượng `imputer` từ thư viện `scikit-learn` để thay thế các giá trị NaN bằng giá trị trung bình của từng cột tránh gây lỗi khi phân cụm. `Fit_transform()` là hàm dùng để áp dụng các phương pháp thay thế này vào `numeric_data`.

```
scaler = StandardScaler()  
  
scaled_data = scaler.fit_transform(numeric_data)
```

Khởi tạo đối tượng `StandardScaler` từ thư viện `scikit-learn` sẽ chuẩn hóa dữ liệu đưa tất cả các biến cần so sánh về một thang đo và lưu vào `scaled_data`.

```
wcss = []  
  
for k in range(1, 11):  
  
    kmeans = KMeans(n_clusters=k, random_state=42)  
  
    kmeans.fit(scaled_data)  
  
    wcss.append(kmeans.inertia_)
```

Đoạn mã trên sử dụng phương pháp Elbow để xác định số cụm tối ưu cho thuật toán K-means. WCSS (within cluster sum of squares) là tổng khoảng cách bình phương của mỗi điểm đến tâm cụm tương ứng. Danh sách `wcss` lưu giá trị WCSS cho từng số lượng cụm `k` từ 1 đến 10:

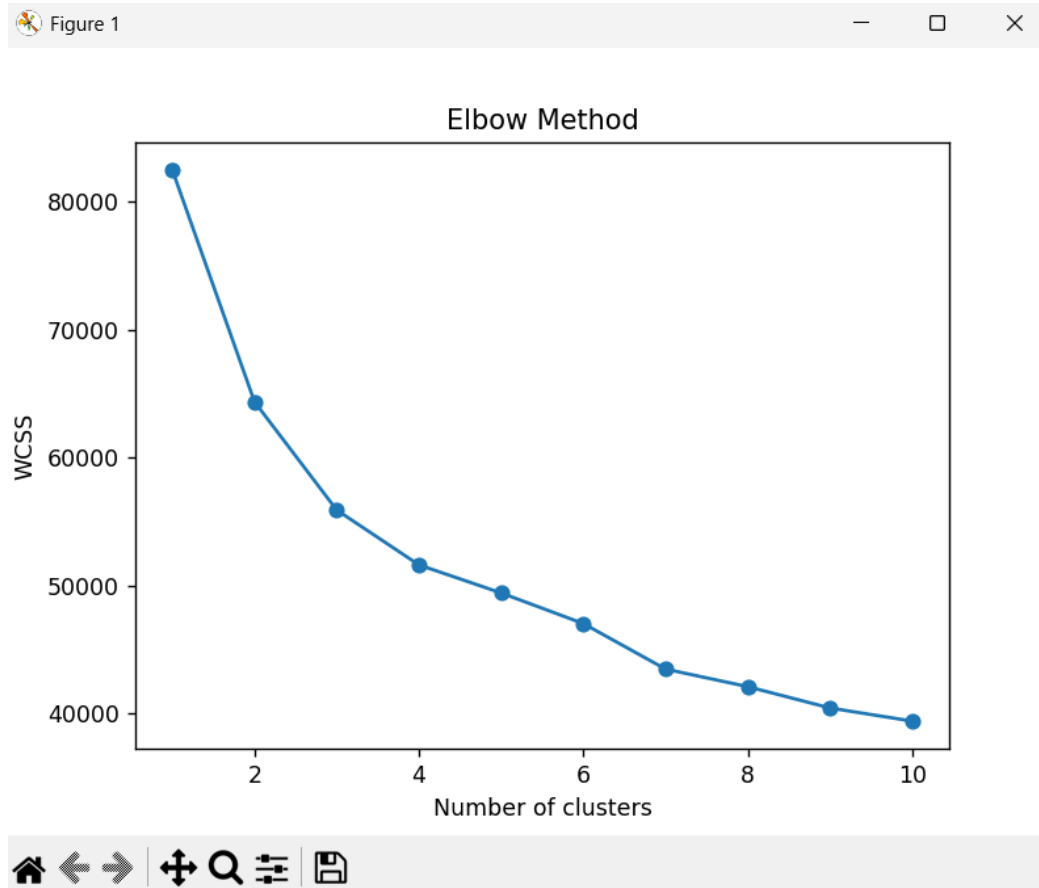
`KMeans(n_clusters=k, random_state=42)` tạo một mô hình `KMeans` với `k` cụm và thiết lập `random_state=42` để kết quả phân cụm là nhất quán. `kmeans.fit(scaled_data)`: Huấn luyện mô hình `KMeans` trên dữ liệu đã được chuẩn hóa (`scaled_data`). Trong quá trình này, thuật toán K-means sẽ tìm kiếm `k` tâm cụm và phân chia các điểm dữ liệu vào các cụm tương ứng.

`kmeans.inertia_`: Đây là giá trị WCSS cho số lượng cụm `k` hiện tại. `inertia_` trả về tổng các khoảng cách bình phương từ mỗi điểm dữ liệu đến tâm cụm của nó. Lưu giá trị này vào danh sách `wcss` để vẽ biểu đồ Elbow:

```
plt.plot(range(1, 11), wcss, marker='o')
```

```
plt.title('Elbow Method')
plt.xlabel('Number of clusters')
plt.ylabel('WCSS')
plt.show()
```

Hàm plot nhận vào ba tham số lần lượt: trục x biểu diễn số lượng cụm từ một đến 10, trục y là danh sách giá trị WCSS đã tính ở trên, marker là kí hiệu để biểu thị từng điểm dữ liệu trên biểu đồ.



hình 3.1: Elbow Method

Dựa trên biểu đồ Elbow thì điểm gấp khúc rõ rệt nhất trong biểu đồ này là lúc số lượng cụm là 3 vì từ số cụm 3 trở đi thì mức giảm của WCSS không đáng kể, nghĩa là việc thêm cụm sẽ không làm giảm nhiều độ phân tán trong mỗi cụm. Số cụm tối ưu ở đây sẽ là 3.

```
optimal_k = 3
kmeans = KMeans(n_clusters=optimal_k, random_state=42)
pca = PCA(n_components=2)
```

```
pca_data = pca.fit_transform(scaled_data)
```

Áp dụng thuật toán K-means để phân cụm dữ liệu và sử dụng PCA để giảm chiều dữ liệu xuống 2D nhằm dễ dàng trực quan hóa. `optimal_k` là số lượng cụm tối ưu đã được chọn nhờ phương pháp Elbow.

`KMeans(n_clusters=optimal_k, random_state=42)` tạo mô hình K-means với ba cụm và thiết lập `random_state=42` để đảm bảo tính nhất quán trong các lần chạy. `fit_predict(scaled_data)`: Huấn luyện mô hình K-means trên dữ liệu đã được chuẩn hóa (`scaled_data`) và trả về nhãn cụm cho từng điểm dữ liệu. Lưu nhãn cụm vào cột 'Cluster' trong dataframe `data` để có thể phân tích và trực quan hóa.

Khởi tạo một đối tượng PCA với `n_component = 2` nghĩa là giảm chiều dữ liệu xuống 2 thành phần. `pca_data` sẽ lưu dữ liệu mới sau khi đã giảm chiều dữ liệu thông qua hàm `fit_transform` áp dụng lên dữ liệu đã chuẩn hóa.

➔ Sau đoạn mã trên thì dữ liệu đã được phân vào 3 cụm, và chiều của dữ liệu đã giảm xuống còn 2 chiều để dễ dàng vẽ biểu đồ và quan sát sự phân cụm.

```
plt.scatter(pca_data[:, 0], pca_data[:, 1], c=data['Cluster'],
            cmap='viridis')

plt.title('K-means Clustering with PCA')

plt.xlabel('PCA 1')

plt.ylabel('PCA 2')

plt.show()
```

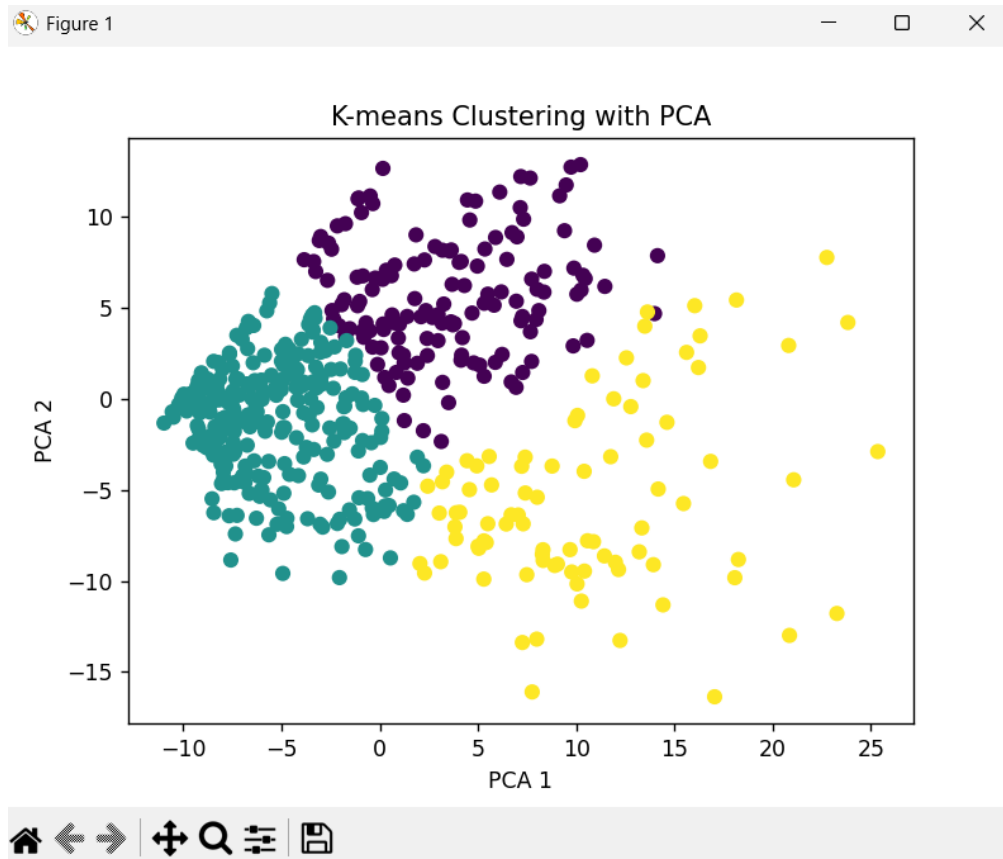
Hàm `scatter` từ thư viện Matplotlib cho ra một biểu đồ phân tán với bốn tham số lần lượt là:

`pca_data[:, 0]`: Trục x của biểu đồ là thành phần chính thứ nhất (PCA 1).

`pca_data[:, 1]`: Trục y của biểu đồ là thành phần chính thứ hai (PCA 2).

`c=data['Cluster']`: Màu sắc của các điểm dữ liệu sẽ được đặt dựa trên nhãn cụm của chúng (cột Cluster trong `data`).

`cmap='viridis'`: Sử dụng bảng màu "viridis" để mã hóa các cụm khác nhau bằng màu sắc.



Hình 3.2: K-means Clustering with PCA

III. Tài liệu tham khảo.

1. <https://scikitlearn.org/1.5/modules/generated/sklearn.cluster.KMeans>.
2. <https://codelearn.io/sharing/xu-ly-du-lieu-voi-pandas-trong-python>
3. https://www.tutorialspoint.com/beautiful_soup/index.htm
4. <https://viblo.asia/p/cao-du-lieu-doanh-nghiep-voi-beautiful-soup-mot-cach-cuc-ky-don-gian-V3m5WjgblO7>
5. <https://pandas.pydata.org/docs/>