

Part 1:

1. Idea:

This Python program automates the collection of detailed statistical data for English Premier League football players from FBref.com, focusing on those who have played over 90 minutes. The goal is to extract a diverse set of performance metrics and professional parameters, then compile, process, and save them into a results.csv file. This file is structured with each row representing a player, each column corresponding to a statistic, players sorted alphabetically by their first name, and missing values marked as "N/a", facilitating subsequent in-depth data analysis.

2. Code Analysis

The program is written in Python and utilizes several powerful libraries to perform tasks ranging from web browsing to data processing:

- **Selenium (webdriver):** Used to automate the Chrome browser, allowing the program to access and interact with the dynamic content of the FBref website.
- **BeautifulSoup4 (bs4):** This library parses the HTML code of web pages, helping to extract the necessary data tables.
- **Pandas (pd):** The central library for data manipulation. Pandas is used to create and manage DataFrames (data tables), and to perform operations such as merging, filtering, and cleaning data.
- **StringIO (io.StringIO):** Supports reading strings as file objects, which is necessary when converting HTML tables (as strings) into Pandas DataFrames.
- **Time (time):** The time module is used to introduce delays (time.sleep(5)) during web browsing, ensuring pages have sufficient time to load completely before data extraction, and also to avoid sending too many requests to the server 연속적으로.

Main Operational Flow of the Code:

1. **Initialize Browser:** An instance of webdriver.Chrome() is created.

```
driver=webdriver.Chrome()
```

2. **Define Links (links):** A dictionary stores target URLs on FBref and the IDs of the HTML tables containing the corresponding statistical data (e.g., Standard, Goalkeeping, Shooting, etc.).

```

10  # các liên kết đến các bảng dữ liệu
11  links = {
12      "Standard": ("https://fbref.com/en/comps/9/stats/Premier-League-Stats", "stats_standard"),
13      "Goalkeeping": ("https://fbref.com/en/comps/9/keepers/Premier-League-Stats", "stats_keeper"),
14      "Shooting": ("https://fbref.com/en/comps/9/shooting/Premier-League-Stats", "stats_shooting"),
15      "Passing": ("https://fbref.com/en/comps/9/passing/Premier-League-Stats", "stats_passing"),
16      "GnS Creation": ("https://fbref.com/en/comps/9/gca/Premier-League-Stats", "stats_gca"),
17      "Defensive act": ("https://fbref.com/en/comps/9/defense/Premier-League-Stats", "stats_defense"),
18      "Possession": ("https://fbref.com/en/comps/9/possession/Premier-League-Stats", "stats_possession"),
19      "Misc": ("https://fbref.com/en/comps/9/misc/Premier-League-Stats", "stats_misc")
20  }

```

3. scraping(url, table_id) Function:

- This function is the core of data collection from each page.
- It navigates to the url using driver.get().
- Pauses for 5 seconds for the page to load.
- Retrieves the page source (driver.page_source) and uses BeautifulSoup to parse it.
- Finds the specific HTML table based on table_id.
- Uses pd.read_html() (in conjunction with StringIO) to convert the HTML table into a Pandas DataFrame. The header row is specified by header=1.
- Performs basic cleaning steps to remove any redundant or unnecessary header rows within the table.

```

24  def (variable) driver: WebDriver
25      driver.get(url)
26      time.sleep(5)
27      soup = BeautifulSoup(driver.page_source, "html.parser")
28      table = soup.find("table", {"id": table_id})
29      df = pd.read_html(StringIO(str(table)), header=1)[0]
30      df=df[~df["Rk"].str.contains("Rk", na=False)]
31      df = df[df.apply(lambda row: not any(row.astype(str) == row.name), axis=1)]
32      return df

```

4. Data Collection and Merging (all_df):

- The program iterates through each item in the links dictionary.
- Calls the scraping function to get a DataFrame for each type of statistic.
- The first DataFrame is assigned as all_df. Subsequent DataFrames are merged (pd.merge) into all_df based on the common columns "Player" and "Squad". An outer join is used to ensure no player data is lost. Suffixes (suffixes) are added to duplicate column names to distinguish their origin.

```

33     all_df = None
34     # Duyệt qua từng liên kết và lấy dữ liệu, sau đó gộp lại
35     for name, (url, table_id) in links.items():
36         df = scraping(url, table_id)
37         if all_df is None:
38             all_df = df
39         else:
40             all_df = pd.merge(all_df, df, on=["Player", "Squad"], how="outer", suffixes=("", f"_{name}"))
41             print(f"Gộp với {name} xong, kích thước all_df: {all_df.shape}")

```

5. Post-Merge Data Processing:

- **Filter by playing time:** The "Min" column is converted to a numeric type. Only players with playing time (Min) greater than 90 minutes are retained.

```

90     if "Min" in all_df.columns:
91         all_df["Min"] = pd.to_numeric(all_df["Min"], errors="coerce")
92         all_df = all_df[all_df["Min"] > 90]
93

```

- **Handle missing values (NaN):** All NaN values in the DataFrame are replaced with the string "N/a".

```

94     # Thay thế giá trị NaN bằng "N/a"
95     all_df.fillna("N/a", inplace=True)

```

- **Create "First_name" column:** A new column containing the player's first name is created from the "Player" column to facilitate sorting.

```

97     all_df["First_name"] = all_df["Player"].apply(lambda x: x.split()[0] if isinstance(x, str) else "N/a")

```

- **Sort data:** The all_df DataFrame is sorted alphabetically by "First_name".

```

100     all_df = all_df.sort_values(by="First_name")

```

- **Select columns (columns_to_keep):** Only columns whose names are in the predefined columns_to_keep list (containing all desired statistical indicators) are retained in the final DataFrame.

```

# Các cột cần giữ lại trong DataFrame
columns_to_keep = [
    "Player", "Squad", "Nation", "Pos", "Age",
    "MP", "Starts", "Min",
    "Gls", "Ast", "CrdY", "CrdR",
    "xG", "xAG",
    "PrgC", "PrgP", "PrgR",
    "Gls.1", "Ast.1", "xG.1", "xAG.1",
    "GA90", "Save%", "CS%",
    "PKsv",
    "SoT%", "SoT/90", "G/Sh", "Dist",
    "Cmp", "Cmp%", "TotDist",
    "Cmp%.1", "Cmp%.2", "Cmp%.3",
    "KP", "1/3", "PPA", "CrsPA", "PrgP_Passing",
    "SCA", "SCA90", "GCA", "GCA90",
    "Tkl", "TklW", "Att_Defensive act", "Lost",
    "Blocks", "Sh_Defensive act", "Pass", "Int",
    "Touches", "Def Pen", "Def 3rd_Possession",
    "Mid 3rd_Possession", "Att 3rd_Possession", "Att Pen",
    "Att_Possession", "Succ%", "Tkld%",
    "Carries", "PrgDist_Possession", "PrgC_Possession",
    "1/3_Possession", "CPA", "Mis", "Dis",
    "Rec", "PrgR_Possession",
    "Fls", "Fld_Misc", "Off", "Crs", "Recov",
    "Won", "Lost_Misc", "Won%"
]

```

6. Save Results:

- The complete all_df DataFrame is saved to a results.csv file.
- The index=False parameter is used to prevent writing the DataFrame index to the file.
- encoding="utf-8-sig" ensures correct handling of special characters.
- A confirmation message indicating the file has been saved is printed to the console.

```
86 # Lưu DataFrame vào file CSV
87 all_df.to_csv("results.csv", index=False, encoding="utf-8-sig")
88 print("Dữ liệu đã được lưu vào file results.csv.")
```

Part 2:

1. Concept

This Python program is designed to perform statistical analysis on a pre-existing dataset of footballer statistics, presumably stored in a file named results.csv. The core idea is to process this data to extract meaningful insights, including identifying top and bottom performers, calculating descriptive statistics for various metrics both league-wide and فريق-wise, and visualizing the distribution of key statistics. The program aims to generate summary files (top_3.txt, results2.csv) and display histograms for selected attacking and defensive metrics.

2. Code Analysis

The program is written in Python and utilizes several common libraries for data manipulation and visualization:

- **Pandas (pd):** Used extensively for reading the input CSV file, data manipulation (handling missing values, selecting columns, grouping), and creating summary DataFrames.
- **NumPy (np):** While not directly used in many operations, Pandas relies on NumPy for its numerical operations and data structures. It's explicitly imported and used for np.nan when handling potentially missing aggregated statistics.
- **Matplotlib.pyplot (plt):** Employed for generating and displaying histograms to visualize the distribution of player statistics.
- **OS (os):** Imported but not actively used in the provided snippet (the comment indicates it's for working with directories, which might be relevant in a fuller version of the script, e.g., for creating output directories if they don't exist).

Main Operational Flow of the Code:

1. **plot_histograms(df, numeric_cols, grouped_by_squad) Function:**
 - This function is responsible for generating and displaying histograms.
 - It iterates through the provided numeric_cols.
 - For each statistic:
 - It prepares the data by dropping any missing (NaN) values.

- **League-wide Histogram:** A histogram is plotted showing the distribution of the statistic for all players in the league.
- **Team-specific Histograms:** If grouped_by_squad (a Pandas GroupBy object) is provided, the function then iterates through each team. For each team, it filters the data for that team's players and plots a separate histogram for the current statistic.
- plt.show() is called after each histogram is configured to display it.

```

6  def plot_histograms(df: pd.DataFrame, numeric_cols: list, grouped_by_squad=None): # Bỏ output_dir
7      for col_stat in numeric_cols:
8          data_for_hist = df[col_stat].dropna()
9          # 1. Histogram cho toàn giải
10         plt.figure(figsize=(10, 6))
11         plt.hist(data_for_hist, bins='auto', color='skyblue', edgecolor='black')
12         plt.title(f"Phân phối của {col_stat} - Toàn giải")
13         plt.xlabel(col_stat)
14         plt.ylabel("Số lượng cầu thủ")
15         plt.grid(axis='y', alpha=0.75)
16         plt.show() # Hiển thị biểu đồ
17
18         # 2. Histogram cho từng đội (nếu có grouped_by_squad)
19         if grouped_by_squad is not None:
20             print(f"    Đang xử lý histogram theo đội cho cột: {col_stat}")
21             for team_name, team_data_group in grouped_by_squad:
22                 team_stat_data = team_data_group[col_stat].dropna()
23
24                 if team_stat_data.empty:
25                     continue
26
27                 plt.figure(figsize=(10, 6)) # Tạo figure mới cho mỗi biểu đồ của đội
28                 plt.hist(team_stat_data, bins='auto', color='lightcoral', edgecolor='black')
29                 plt.title(f"Phân phối của {col_stat} - Đội: {team_name}")
30                 plt.xlabel(col_stat)
31                 plt.ylabel("Số lượng cầu thủ")
32                 plt.grid(axis='y', alpha=0.75)
33                 plt.show() # Hiển thị biểu đồ của đội
34             print(f"    Đã xử lý xong histogram theo đội cho cột: {col_stat}")
35
36     print("Hoàn tất việc chuẩn bị hiển thị histogram. Các cửa sổ biểu đồ sẽ lần lượt xuất hiện.")

```

2. main() Function: This is the primary execution block.

- **Read Data:** Reads data from results.csv into a Pandas DataFrame df. na_values='N/a' ensures that "N/a" strings are treated as missing values.

```
df = pd.read_csv("results.csv", na_values='N/a')
```

○ Handle Missing Values:

- Identifies all numeric columns in the DataFrame.
- Fills any NaN values in these numeric columns with the mean of their respective columns.

- Checks if essential text columns ('Player', 'Squad') are present.

```

43     numeric_cols = df.select_dtypes(include='number').columns.tolist()
44     df[numeric_cols] = df[numeric_cols].fillna(df[numeric_cols].mean())
45     print(f"Đã xử lý giá trị bị thiếu cho các cột số: {numeric_cols}.")
46     required_text_cols = ['Player', 'Squad']
47     missing_text_cols = [col for col in required_text_cols if col not in df.columns]

```

- **Write top_3.txt File:**

- Opens top_3.txt for writing.
- For each numeric column:
 - Writes a header for the statistic.
 - Determines which columns to display (Player, Squad, and the current statistic).
 - Uses nlargest(3, ...) and nsmallest(3, ...) to find the top 3 and bottom 3 players for that statistic.
 - Writes this information to the text file.

```

50  with open("top_3.txt", "w", encoding="utf-8") as f:
51    for col_stat in numeric_cols:
52      f.write(f"----- Thông kê: {col_stat} -----\\n")
53
54      display_cols_for_top_bottom = []
55      if 'Player' in df.columns:
56        display_cols_for_top_bottom.append('Player')
57      if 'Squad' in df.columns:
58        display_cols_for_top_bottom.append('Squad')
59      display_cols_for_top_bottom.append(col_stat)
60
61      if df[col_stat].notna().any():
62        top3 = df.nlargest(3, col_stat)[display_cols_for_top_bottom]
63        bot3 = df.nsmallest(3, col_stat)[display_cols_for_top_bottom]
64
65        f.write(f"Top 3 CAO NHẤT theo {col_stat}:\\n")
66        f.write(top3.to_string(index=False) + "\\n")
67        f.write(f"Top 3 THẤP NHẤT theo {col_stat}:\\n")
68        f.write(bot3.to_string(index=False) + "\\n\\n")
69      else:
70        f.write(f"Cột {col_stat} không có dữ liệu hợp lệ để tìm top/bottom 3.\\n\\n")

```

- **Group Data by Team:**

- If the 'Squad' column exists, it groups the DataFrame df by this column, creating a DataFrameGroupBy object named grouped_by_squad.

```

75      grouped_by_squad = None
76      if 'Squad' in df.columns (variable) df: DataFrame
77        grouped_by_squad = df.groupby('Squad')
78      else:
79        print("CẢNH BÁO: Không tìm thấy cột 'Squad'. Không thể nhóm theo đội.")

```

- **Calculate League-wide Statistics:**

- If numeric columns exist and are not empty, it calculates the median, mean, and standard deviation (std) for each numeric column across all players. The result is stored in all_league_stats after transposing (.T).

```

82     all_league_stats = None
83     if numeric_cols:
84         if not df[numeric_cols].empty:
85             all_league_stats = df[numeric_cols].agg(['median', 'mean', 'std']).T
86             print("Đã tính toán median, mean, std cho toàn giải.")
87         else:
88             print("Không có dữ liệu số để tính thống kê toàn giải.")
89     else:
90         print("Không có cột số nào được xác định để tính thống kê toàn giải.")
```

- **Prepare and Create results2.csv:**

- This section executes if numeric_cols, all_league_stats, and grouped_by_squad are all available.
- It prepares a list of rows (rows_for_results2) to build the new CSV.
- Defines the column structure for results2.csv: 'Squad', followed by 'Median of StatX', 'Mean of StatX', 'Std of StatX' for each numeric statistic.
- League Totals Row:** Adds a row to rows_for_results2 for the 'all' league statistics.
- Team-specific Rows:** If grouped_by_squad contains numeric data, it aggregates median, mean, and std for each team and each numeric column. Each team's aggregated stats are added as a new row. .get() is used with np.nan as a default to prevent errors if a specific statistic combination is missing.
- A new DataFrame results2_df is created from rows_for_results2 and saved to results2.csv without the index and with utf-8-sig encoding.

```

94     if numeric_cols and all_league_stats is not None and grouped_by_squad is not None:
95         print("Đang chuẩn bị dữ liệu cho 'results2.csv'...")
96         rows_for_results2 = []
97         columns_for_results2 = ['Squad']
98         for col_stat in numeric_cols:
99             columns_for_results2.extend([f"Median của {col_stat}", f"Mean của {col_stat}", f"Std của {col_stat}"])
100        all_league_values = []
101        for col_stat in numeric_cols:
102            if col_stat in all_league_stats.index:
103                all_league_values.extend([
104                    all_league_stats.loc[col_stat, 'median'],
105                    all_league_stats.loc[col_stat, 'mean'],
106                    all_league_stats.loc[col_stat, 'std']
107                ])
108            else:
109                all_league_values.extend([np.nan, np.nan, np.nan])
110        rows_for_results2.append(['all'] + all_league_values)
111    if not grouped_by_squad[numeric_cols].first().empty: # Kiểm tra nếu có dữ liệu số trong nhóm
112        team_aggregated_stats = grouped_by_squad[numeric_cols].agg(['median', 'mean', 'std'])
113        for team_name, team_specific_stats in team_aggregated_stats.iterrows():
114            team_values = []
115            for col_stat in numeric_cols:
116                team_values.extend([
117                    team_specific_stats.get(col_stat, 'median'), np.nan,
118                    team_specific_stats.get(col_stat, 'mean'), np.nan,
119                    team_specific_stats.get(col_stat, 'std'), np.nan
120                ])
121            rows_for_results2.append([team_name] + team_values)
122        else:
123            print("Cảnh báo: Không có dữ liệu số trong các nhóm đội để tổng hợp cho results2.csv.")
124        results2_df = pd.DataFrame(rows_for_results2, columns=columns_for_results2)
125        results2_df.to_csv('results2.csv', index=False, encoding='utf-8-sig') |
126        print(results2_df.head())
127    else:
128        print("Bỏ qua việc tạo 'results2.csv' do thiếu grouped_by_squad, all_league_stats hoặc numeric_cols.")

```

- **Plot Histograms (Specific Columns):**

- Defines two lists, attack_cols_defined (["Gls", "Ast", "Dist"]) and defense_cols_defined (["Tkl", "TklW", "Att_Defensive act"]).
- These are combined into target_hist_cols.
- cols_to_plot_actual is created by filtering target_hist_cols to include only those columns that actually exist in the DataFrame df and are of a numeric data type.
- If cols_to_plot_actual is not empty, the plot_histograms function is called with these specific columns.

```

131     attack_cols_defined = ["Gls", "Ast", "Dist"]
132     defense_cols_defined = ["Tk1", "Tk1W", "Att_Defensive act"]
133
134     target_hist_cols = attack_cols_defined + defense_cols_defined
135
136     cols_to_plot_actual = []
137     for col_name in target_hist_cols:
138         if col_name in df.columns:
139             if pd.api.types.is_numeric_dtype(df[col_name]):
140                 cols_to_plot_actual.append(col_name)
141
142     if cols_to_plot_actual:
143         print(f"\nSẽ vẽ histogram cho các cột: {', '.join(cols_to_plot_actual)}")
144         plot_histograms(df, cols_to_plot_actual, grouped_by_squad)
145     else:
146         print("\nKhông có cột tấn công/phòng ngự hợp lệ nào được tìm thấy để vẽ histogram.")

```

3. Execution Block (if __name__ == '__main__':)

- This standard Python construct ensures that the main() function is called only when the script is executed directly (not when imported as a module).

Part 3:

1. Concept

This Python program employs K-means clustering to classify football players into distinct groups based on their statistical performance. The objective is to identify inherent player categories—such as different roles, performance tiers, or styles of play—with pre-existing labels.

The process involves several key steps:

- Data Preprocessing:** Cleaning the statistical data, imputing any missing values, and scaling features to ensure fair comparison.
- Optimal Cluster (K) Determination:** Utilizing the "Elbow method" by plotting K-means inertia (sum of squared distances to cluster centers) for various K values. The script identifies the optimal K as the point where the most significant drop in inertia occurs, indicating a good balance for cluster formation.
- K-means Clustering:** Applying the K-means algorithm with the determined optimal K to assign players to clusters.
- Dimensionality Reduction & Visualization:** Using Principal Component Analysis (PCA) to reduce the data to two dimensions. This allows for a 2D scatter plot where players are colored by their assigned cluster, providing a visual

representation of the groupings and their separation.

2. Code Analysis

The program is written in Python and utilizes several libraries from the scikit-learn ecosystem for machine learning tasks, along with Pandas for data handling and Matplotlib for plotting:

- **Pandas (pd):** Used for reading the input CSV file (results.csv), data manipulation (dropping columns, applying cleaning functions), and storing cluster assignments.
- **NumPy (np):** Implicitly used by Pandas and scikit-learn for numerical operations.
- **Scikit-learn (sklearn):**
 - StandardScaler: For standardizing features by removing the mean and scaling to unit variance.
 - KMeans: The K-means clustering algorithm.
 - SimpleImputer: For handling missing values (NaNs) by replacing them with a specified strategy (e.g., mean).
 - PCA: Principal Component Analysis for dimensionality reduction.
- **Matplotlib.pyplot (plt):** Used for generating plots, specifically the Elbow method plot and the 2D scatter plot of the clusters.
- **OS (os):** Imported in the previous script but not explicitly used in this clustering script snippet.

Main Operational Flow of the Code:

1. **Load Data:** Player statistics are read from results.csv into a Pandas DataFrame.

```
# Read data  
df = pd.read_csv("results.csv")
```

2. **Initial Data Cleaning (Dropping Columns):** Identifier columns that are not suitable for clustering (like Player name, Squad, Nation, Position, Age) are dropped from the DataFrame. errors='ignore' ensures that if any of these columns are already missing, the code doesn't raise an error.

```
11 # Bỏ các cột định danh không dùng cho phân cụm  
12 cols_to_drop = ["Player", "Squad", "Nation", "Pos", "Age"]  
13 df_clean = df.drop(columns=cols_to_drop, errors='ignore')
```

3. **Numeric Data Cleaning and Conversion:**

- A function clean_numeric is defined to convert columns to a numeric type. It

first converts the column to string type, then removes any '%' characters (e.g., from percentage values), and finally attempts to convert to numeric. Values that cannot be converted become NaN.

- This function is applied to all columns in df_clean.

```
16  def clean_numeric(col):
17  |     return pd.to_numeric(col.astype(str).str.replace('%', '', regex=False), errors='coerce')
18
19  df_clean = df_clean.apply(clean_numeric)
```

4. **Impute Missing Values:** SimpleImputer is used to fill any NaN values that resulted from the previous cleaning step or were already present. The strategy is "mean", meaning NaNs in each column will be replaced by the mean of that column.

```
22  imputer = SimpleImputer(strategy="mean")
23  df_clean_imputed = imputer.fit_transform(df_clean)
```

5. **Data Standardization:** The features (player statistics) are standardized using StandardScaler. This process transforms the data to have a mean of 0 and a standard deviation of 1, which is important for distance-based algorithms like K-means to prevent features with larger magnitudes from dominating the results.

```
26  scaler = StandardScaler()
27  X_scaled = scaler.fit_transform(df_clean_imputed)
```

6. **Determine Optimal K (Elbow Method):**

- The script calculates the K-means inertia (sum of squared distances of samples to their closest cluster center) for a range of cluster numbers K_range (1 to 29).
- It then identifies the optimal_K by finding where the largest drop in inertia occurs between consecutive K values. This point is considered the "elbow" and suggests a good balance between the number of clusters and the compactness of each cluster.
- The Elbow curve (Inertia vs. Number of clusters) is plotted to visually aid in selecting K.

```

30     inertia = []
31     K_range = range[1, 30]
32
33     for k in K_range:
34         kmeans = KMeans(n_clusters=k, random_state=42, n_init=10)
35         kmeans.fit(X_scaled)
36         inertia.append(kmeans.inertia_)
37
38     # Tính toán số K tối ưu từ sự thay đổi lớn nhất trong inertia
39     # Tìm sự thay đổi lớn nhất giữa các inertia liên tiếp
40     inertia_diff = [inertia[i] - inertia[i+1] for i in range(len(inertia)-1)]
41     optimal_K = inertia_diff.index(max(inertia_diff)) + 2
42     print(f"Số K tối ưu là {optimal_K}, vì từ K={optimal_K-1} tới K={optimal_K} ta nhận thấy sự thay đổi lớn nhất")
43
44     # Vẽ biểu đồ Elbow
45     plt.figure(figsize=(8, 5))
46     plt.plot(K_range, inertia, 'bo-', markersize=8)
47     plt.xlabel('Số cụm (k)')
48     plt.ylabel('Inertia')
49     plt.title('Phương pháp Elbow để chọn số cụm K')
50     plt.grid(True)
51     plt.tight_layout()
52     plt.show()
53

```

7. **Apply K-means Clustering:** K-means clustering is performed on the scaled data (`X_scaled`) using the `optimal_K` found in the previous step. The cluster labels for each player are then added as a new 'Cluster' column to the `df_clean` DataFrame (note: it might be more appropriate to add it to the original `df` or a copy to retain player identifiers alongside cluster labels for easier interpretation).

```

54     # Áp dụng K-means phân cụm với số K tối ưu
55     kmeans_optimal = KMeans(n_clusters=optimal_K, random_state=42, n_init=10)
56     df_clean['Cluster'] = kmeans_optimal.fit_predict(X_scaled)
57

```

8. **Dimensionality Reduction with PCA:** PCA is used to reduce the dimensionality of the scaled data (`X_scaled`) to 2 components (`n_components=2`). This allows for visualization of the clusters in a 2D scatter plot.

```

58     # Giảm chiều dữ liệu xuống 2D với PCA
59     pca = PCA(n_components=2)
60     X_pca = pca.fit_transform(X_scaled)
61

```

9. **Visualize Clustering Results:**

- A scatter plot is created where the x-axis is the first principal component (PC1) and the y-axis is the second principal component (PC2).
- Each point represents a player, and the color of the point is determined by the cluster label assigned by K-means.
- This plot helps to visually inspect the separation and distribution of the identified player clusters.

```

63 plt.figure(figsize=(8, 5))
64 plt.scatter(X_pca[:, 0], X_pca[:, 1], c=df_clean['Cluster'], cmap='viridis', s=50)
65 plt.title('Phân cụm K-means với số K tối ưu (PCA 2D)')
66 plt.xlabel('PC 1')
67 plt.ylabel('PC 2')
68 plt.colorbar(label='Cluster')
69 plt.show()

```

Part4 :

1. Concept

This Python program automates collecting football player transfer values from footballtransfers.com and integrates them with existing statistics from a local results.csv file. Key steps include:

1. Scraping player names, teams, and transfer values.
2. Reading local player statistics.
3. Filtering players from the stats file based on a minimum playing time (900 minutes).
4. Merging scraped values with filtered stats using fuzzy name matching.
5. Saving the combined dataset (stats + transfer values) to players_900mins_transfer_values.csv.

The aim is to enrich player performance data with market valuations for UK Premier League players.

2. Code Analysis

The script uses Python libraries like Pandas, BeautifulSoup, Selenium, WebDriver Manager, Time, and Rapidfuzz for web scraping, data handling, and fuzzy string matching.

Global Configurations:

The script defines constants for URLs, file paths, and processing thresholds (e.g., NUMBER_OF_PAGES_TO_SCRAPE, MINIMUM_PLAYING_TIME_MINS,

FUZZY_MATCH_SIMILARITY_THRESHOLD_PERCENT).

Main Operational Flow of the Code:

1. `scrape_player_values(...)` Function:

- o Scraps player names, teams, and transfer values from multiple pages of footballtransfers.com using Selenium and BeautifulSoup.
- o Initializes and quits a Chrome WebDriver.
- o Includes error handling and page load delays.
- o Returns a DataFrame of scraped data.

```
24 def scrape_player_values(base_url: str, num_pages: int) -> pd.DataFrame:
25     players_data = []
26     print("Đang khởi tạo WebDriver cho việc scraping...")
27     driver = webdriver.Chrome(service=ChromeService(ChromeDriverManager().install()))
28     print(f"Bắt đầu quét {num_pages} trang từ {base_url}...")
29     for i in range(1, num_pages + 1):
30         url_to_scrape = base_url
31         if i != 1:
32             url_to_scrape += "/" + str(i)
33         print(f"Đang xử lý trang {i}: {url_to_scrape}")
34         try:
35             driver.get(url_to_scrape)
36             time.sleep(3) # Chờ trang tải - có thể cần điều chỉnh hoặc dùng WebDriverWait
37             html_source = driver.page_source
38             soup = BeautifulSoup(html_source, "html.parser")
39             player_table = soup.find("table", class_="table table-hover no-cursor table-striped leaguetable mvp-table")
40             if player_table:
41                 table_rows = player_table.find_all("tr")
42                 for row_index, row in enumerate(table_rows):
43                     if row_index == 0: continue
44                     columns = row.find_all("td")
45                     if columns and len(columns) > 4:
46                         try:
47                             player_name_tag = columns[2].find("a")
48                             player_name = player_name_tag.text.strip() if player_name_tag else "N/A"
49                             team_tag = columns[4].find("span", class_="td-team_teamname")
50                             team_name = team_tag.text.strip() if team_tag else "N/A"
51                             player_value = columns[-1].text.strip()
52                             if player_name != "N/A":
53                                 players_data.append({"Player": player_name, "Team": team_name, "Value_Scraped": player_value})
54                             except (IndexError, AttributeError) as e_parse:
55                                 print(f" Lỗi nhỏ khi phân tích hàng {row_index} ở trang {i}: {e_parse}")
56                         else:
57                             print(f" Không tìm thấy bảng dữ liệu trên trang {i}.")
58                         except Exception as e_page:
59                             print(f" Lỗi không mong muốn khi xử lý trang {i}: {e_page}")
60
```

2. `read_and_filter_player_stats(...)` Function:

- o Reads player statistics from a local CSV (stats_csv_path).
- o Handles file errors and checks for required columns ('Player', 'Min').
- o Converts 'Min' column to numeric and filters players by min_playing_time.
- o Selects specified columns_to_keep.
- o Returns the filtered DataFrame.

```

70     def read_and_filter_player_stats(stats_csv_path: str, min_playing_time: int, columns_to_keep: list) -> pd.DataFrame:
71
72         df_stats = pd.read_csv(stats_csv_path)
73         required_cols_stats = ["Player", "Min"]
74         for col in required_cols_stats:
75             if col not in df_stats.columns:
76                 print(f'Lỗi: Cột bắt buộc '{col}' không tìm thấy trong '{stats_csv_path}'.')
77                 return pd.DataFrame()
78         if df_stats['Min'].dtype == 'object': # Xử lý cột 'Min' nếu là dạng chuỗi (ví dụ: '1,234')
79             df_stats['Min'] = df_stats['Min'].astype(str).str.replace(',', '', regex=False)
80             df_stats['Min'] = pd.to_numeric(df_stats['Min'], errors='coerce')
81             df_stats.dropna(subset=['Min'], inplace=True)
82         filtered_players_df = df_stats[df_stats['Min'] > min_playing_time].copy()
83         if filtered_players_df.empty:
84             print(f'Không tìm thấy cầu thủ nào có số phút thi đấu > {min_playing_time} phút.')
85             return pd.DataFrame()
86         print(f'Tìm thấy {len(filtered_players_df)} cầu thủ có số phút thi đấu > {min_playing_time} phút.')
87         valid_columns_to_keep = [col for col in columns_to_keep if col in filtered_players_df.columns]
88         final_df = filtered_players_df[valid_columns_to_keep].reset_index(drop=True)
89         print(f'Dã chọn các cột từ file stats: {valid_columns_to_keep}')
90
91         return final_df

```

3. **combine_data_and_add_values(...)** Function:

- Merges filtered statistics with scraped transfer values using fuzzy name matching (rapiddfuzz).
- Adds 'Transfer_Value' and 'Match_Score' columns to the statistics DataFrame.
- **Drops players for whom no transfer value is found.**
- Returns the combined DataFrame.

```

93     def combine_data_and_add_values(
94         df_filtered_stats: pd.DataFrame,
95         df_scraped_values: pd.DataFrame,
96         similarity_threshold: int
97     ) -> pd.DataFrame:
98         if df_scraped_values.empty:
99             df_filtered_stats['Transfer_Value'] = None
100            df_filtered_stats['Match_Score'] = None
101            return df_filtered_stats
102        df_combined = df_filtered_stats.copy()
103        df_combined["Transfer_Value"] = None
104
105        scraped_player_names = df_scraped_values['Player'].dropna().astype(str).tolist()
106        if not scraped_player_names:
107            print("CẢNH BÁO: Không có tên cầu thủ nào trong DataFrame từ web để so khớp.")
108            return df_combined
109        print(f"Thực hiện so khớp tên cho {len(df_combined)} cầu thủ đã lọc...")
110        matched_count = 0
111        for index, row in df_combined.iterrows():
112            player_name_to_match = str(row["Player"])
113            match_result = process.extractOne(
114                player_name_to_match,
115                scraped_player_names,
116                scorer=fuzz.token_sort_ratio
117            )
118            if match_result:
119                best_match_name, score, _ = match_result
120                if score >= similarity_threshold:
121                    value_series = df_scraped_values.loc[df_scraped_values["Player"] == best_match_name, "Value_Scraped"]
122                    if not value_series.empty:
123                        df_combined.at[index, "Transfer_Value"] = value_series.iloc[0]
124                        df_combined.at[index, "Match_Score"] = score
125                        matched_count += 1
126        df_combined.dropna(subset=['Transfer_Value'], inplace=True)
127        print(f"Số cầu thủ còn lại sau khi loại bỏ những người không có giá trị chuyển nhượng: {len(df_combined)}")

```

4. main() Function:

- o Orchestrates the entire workflow: calls scraping, filtering, and combining functions.
- o Includes a fallback to create a sample results.csv if the input stats file is not found, for testing.
- o Saves the final combined DataFrame to OUTPUT_FINAL_CSV_PATH.
- o Prints progress and summary messages.

```

131 def main():
132     # 1. Thu thập dữ liệu từ web
133     df_scraped = scrape_player_values(BASE_URL_SCRAPING, NUMBER_OF_PAGES_TO_SCRAPE)
134     try:
135         with open(INPUT_STATS_CSV_PATH, 'r') as f:
136             pass # Tệp tồn tại
137     except FileNotFoundError:
138         print(f"\nCẢNH BÁO: Tệp '{INPUT_STATS_CSV_PATH}' không tìm thấy. Tạo tệp mẫu...")
139         sample_data_stats = {
140             'Player': ['Robert Lewandowski', 'Kylian Mbappe', 'Erling Haaland', 'Mohamed Salah', 'Lionel Messi Jr'],
141             'Nation': ['POL', 'FRA', 'NOR', 'EGY', 'ARG', 'BEL', 'SEN'],
142             'Pos': ['FW', 'FW', 'FW', 'FW', 'MF', 'FW'],
143             'Squad': ['Barcelona', 'PSG', 'Man City', 'Liverpool', 'Inter Miami', 'Man City', 'Al-Nassr'],
144             'Age': [35, 25, 23, 31, 36, 32, 31],
145             'Min': [3000, 2800, 2500, 3200, 850, 2700, 700], # Một số < 900
146         }
147         pd.DataFrame(sample_data_stats).to_csv(INPUT_STATS_CSV_PATH, index=False)
148         print(f'Tệp mẫu '{INPUT_STATS_CSV_PATH}' đã được tạo. Vui lòng kiểm tra và cập nhật nếu cần.')
149         df_stats_filtered = read_and_filter_player_stats(
150             INPUT_STATS_CSV_PATH,
151             MINIMUM_PLAYING_TIME_MINS,
152             COLUMNS_TO_KEEP_FROM_STATS_FILE
153         )
154     # 3. Kết hợp dữ liệu
155     df_final_output = combine_data_and_add_values(
156         df_stats_filtered,
157         df_scraped,
158         FUZZY_MATCH_SIMILARITY_THRESHOLD_PERCENT
159     )
160     print(f"\nTổng số cầu thủ trong kết quả cuối cùng: {len(df_final_output)}")
161     print("\n5 dòng đầu của dữ liệu kết hợp cuối cùng:")
162     print(df_final_output.head())
163     # 4. Lưu kết quả
164     df_final_output.to_csv(OUTPUT_FINAL_CSV_PATH, index=False, encoding='utf-8-sig')
165     print("\n--- KẾT THÚC QUY TRÌNH ---")

```

5. Execution Block (`if __name__ == "__main__":`)

- Standard Python entry point that calls `main()` when the script is run.

```

if __name__ == "__main__":
    main()

```