# AAMNA NAZ AWAN

# BIG DATA ANALYTICS

# FINAL EXAM

**DATASET:** Football data from transfer markt

## ABOUT DATA SET:

Clean, structured and automatically updated football data from Transfermarkt, including

60,000+ games from many seasons on all major competitions

400+ clubs from those competitions

30,000+ players from those clubs

400,000+ player market valuations historical records

1,200,000+ player appearance records from all games

## FILE NAMES:

1. appearances.csv
2. club_games.csv
3. clubs.csv
4. competitions.csv
5. game_events.csv
6. game_lineups.csv
7. games.csv
8. player_valuations.csv
9. players.csv
10. transfers.csv

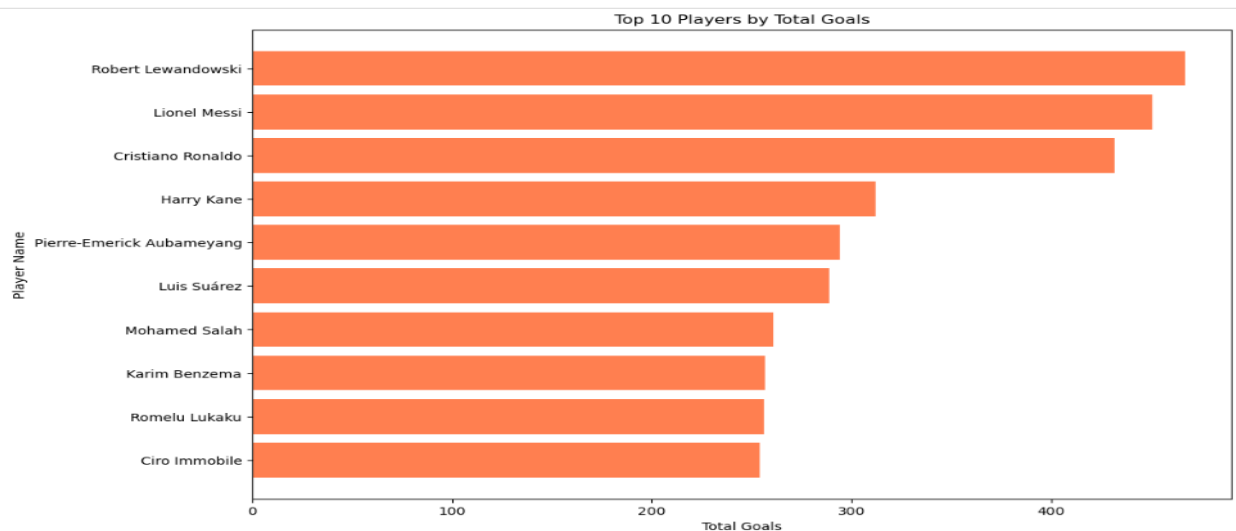QUERY NO:1    **Top Players by Total Goals**

```
1    import pandas as pd
2    df = pd.read_csv('abfss://aamna@youexcel.dfs.core.windows.net/appearances.csv',
3    storage_options = {'account_key':'reQHz4k3MhhRRj1l79ewjF/vDUYcDaQmZ3E07rUQENBdNWIiZ+wiu6VfvrWu+czw
4    display(df)
5    |
```

✓  9 sec - Command executed in 8 sec 990 ms by aamnanaz403 on 2:42:34 AM, 8/22/24

```
1
2    import matplotlib.pyplot as plt
3
4    # Aggregate total goals by player
5    top_scorers = df.groupby('player_id')['goals'].sum().reset_index()
6    top_scorers = top_scorers.sort_values(by='goals', ascending=False).head(10)
7
8    # Load player names for better readability
9    players = pd.read_csv('abfss://aamna@youexcel.dfs.core.windows.net/players.csv',
10   storage_options = {'account_key':'reQHz4k3MhhRRj1l79ewjF/vDUYcDaQmZ3E07rUQENBdNWIiZ+wiu6VfvrWu+czw
11   display(df)
12
13   top_scorers = top_scorers.merge(players[['player_id', 'name']], on='player_id')
14
15   # Visualization
16   plt.figure(figsize=(12, 8))
17   plt.barh(top_scorers['name'], top_scorers['goals'], color='coral')
18   plt.title('Top 10 Players by Total Goals')
19   plt.xlabel('Total Goals')
20   plt.ylabel('Player Name')
21   plt.gca().invert_yaxis()
22   plt.show()
23
```

✓  3 sec - Command executed in 3 sec 169 ms by aamnanaz403 on 2:47:21 AM, 8/22/24
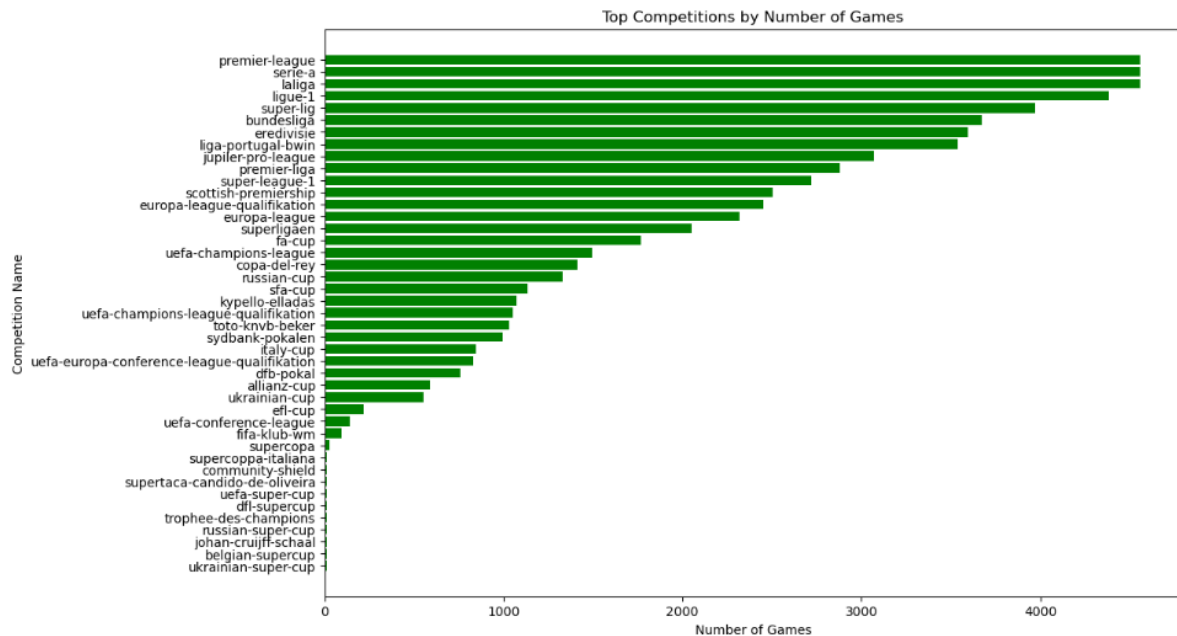
OUTPUT:-

QUERY NO:2 **Top Competitions by Number of Games**

```
1    import seaborn as sns
2
3    # Load datasets
4    games = pd.read_csv('abfss://aamna@youexcel.dfs.core.windows.net/games.csv',
5        storage_options = {'account_key':'reQHz4k3MhhRRj1l79ewjF/vDUYcDaQmZ3E07rUQENBdNWIiZ+wiu6VfvrWu+czw
6    competitions = pd.read_csv('abfss://aamna@youexcel.dfs.core.windows.net/competitions.csv',
7        storage_options = {'account_key':'reQHz4k3MhhRRj1l79ewjF/vDUYcDaQmZ3E07rUQENBdNWIiZ+wiu6VfvrWu+czw
8
9    # Count games per competition
10   games_per_competition = games['competition_id'].value_counts().reset_index()
11   games_per_competition.columns = ['competition_id', 'number_of_games']
12
13   # Merge with competition names
14   games_per_competition = games_per_competition.merge(competitions[['competition_id', 'name']], on='
15
16   # Visualization
17   plt.figure(figsize=(12, 8))
18   plt.barh(games_per_competition['name'], games_per_competition['number_of_games'], color='green')
19   plt.title('Top Competitions by Number of Games')
20   plt.xlabel('Number of Games')
21   plt.ylabel('Competition Name')
22   plt.gca().invert_yaxis()
23   plt.show()
24
```

✓ 2 sec - Command executed in 2 sec 1 ms by aamnanaz403 on 2:52:56 AM, 8/22/24
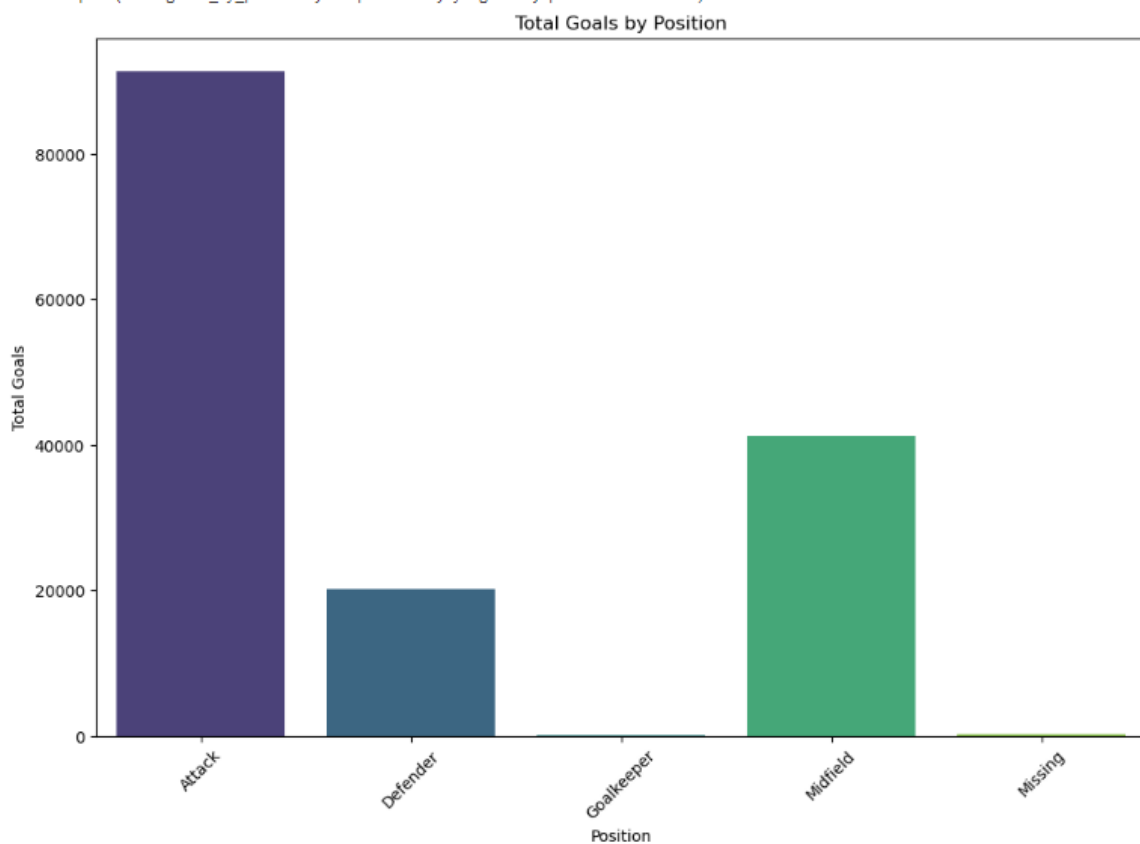
OUTPUT:-

QUERY : 3 **Goal Distribution by Position**

```
1
2     # Merge to get player positions
3     player_goals = df.merge(players[['player_id', 'position']], on='player_id')
4
5     # Aggregate goals by position
6     goals_by_position = player_goals.groupby('position')['goals'].sum().reset_index()
7
8     # Visualization
9     plt.figure(figsize=(12, 8))
10    sns.barplot(data=goals_by_position, x='position', y='goals', palette='viridis')
11    plt.title('Total Goals by Position')
12    plt.xlabel('Position')
13    plt.ylabel('Total Goals')
14    plt.xticks(rotation=45)
15    plt.show()
16
```

✓ 2 sec - Command executed in 2 sec 1 ms by aamnanaz403 on 2:54:49 AM, 8/22/24
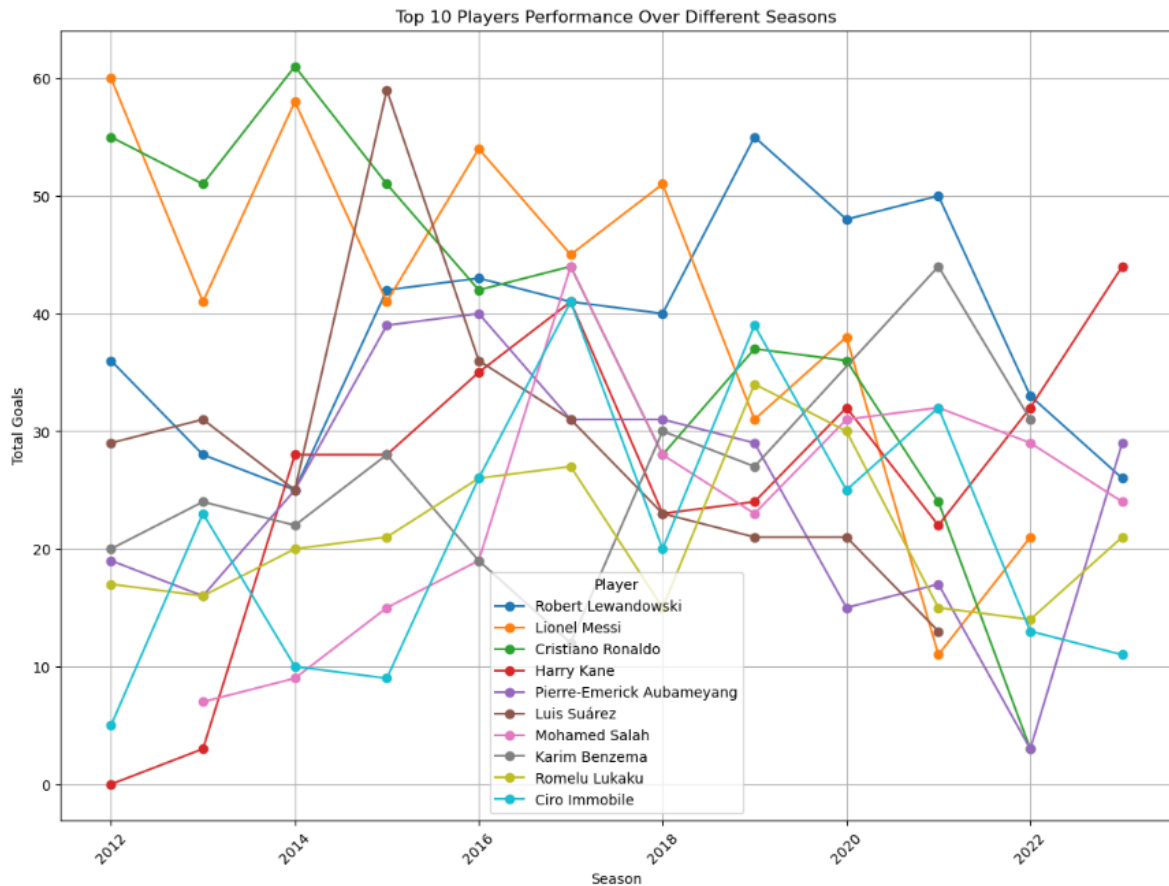
OUTPUT:-

QUERY : 4 **Player Performance Over Different Seasons**

```
1
2    # Merge datasets to get season information
3    player_performance = df.merge(games[['game_id', 'season']], on='game_id')
4    player_performance = player_performance.groupby(['player_id', 'season'])['goals'].sum().reset_index()
5
6    # Get top 10 players by total goals
7    top_players = player_performance.groupby('player_id')['goals'].sum().nlargest(10).index
8    top_players_performance = player_performance[player_performance['player_id'].isin(top_players)]
9
10   # Merge with player names
11   top_players_performance = top_players_performance.merge(players[['player_id', 'name']], on='player_id')
12
13   # Visualization
14   plt.figure(figsize=(14, 10))
15   for player in top_players:
16       player_data = top_players_performance[top_players_performance['player_id'] == player]
17       plt.plot(player_data['season'], player_data['goals'], marker='o', label=player_data['name'].iloc[0])
18
19   plt.title('Top 10 Players Performance Over Different Seasons')
20   plt.xlabel('Season')
21   plt.ylabel('Total Goals')
22   plt.legend(title='Player')
23   plt.grid(True)
24   plt.xticks(rotation=45)
25   plt.show()
26
```

✓ 2 sec - Command executed in 2 sec 71 ms by aamnanaz403 on 2:55:39 AM, 8/22/24
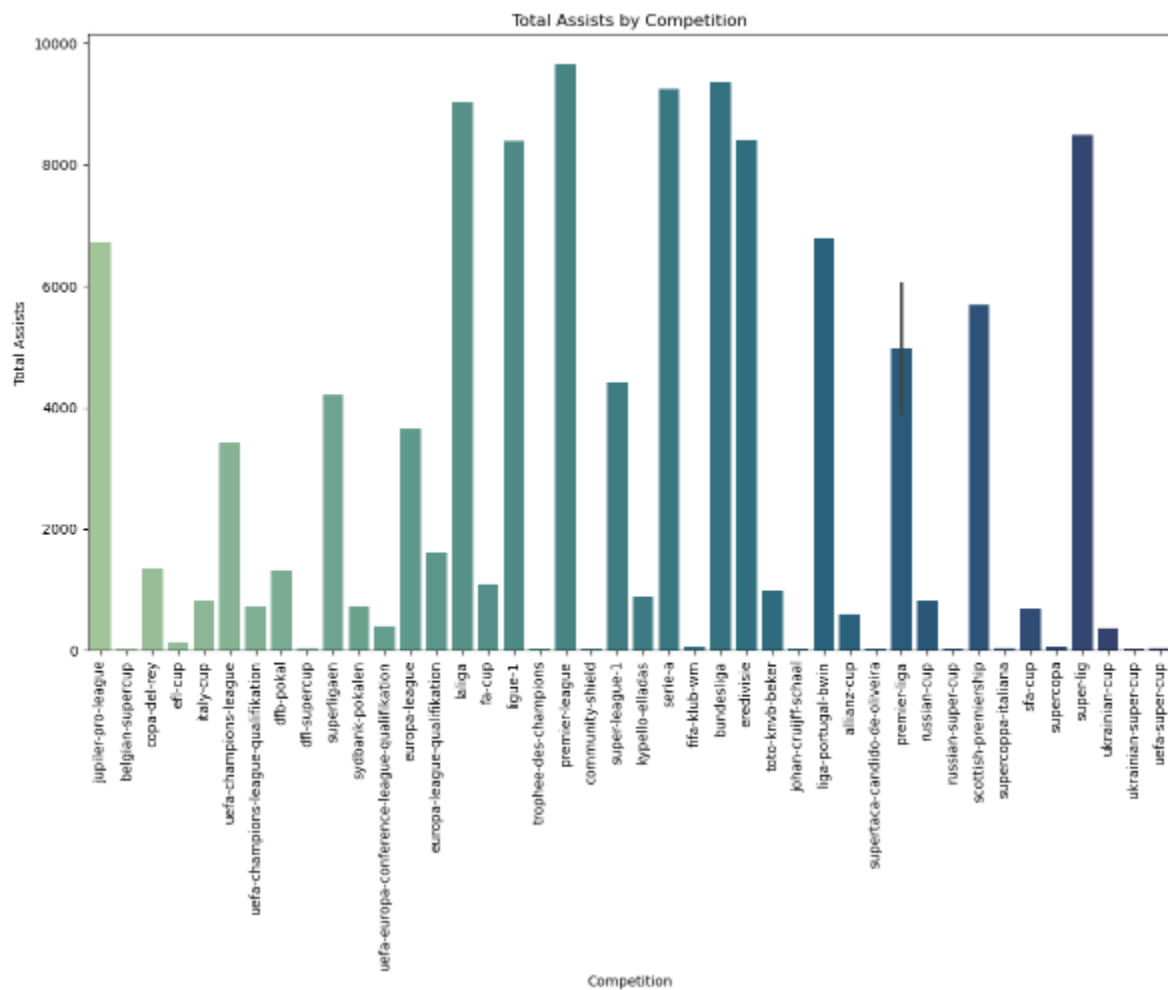
OUTPUT:-

QUERY: 5 **Comparison of Total Assists by Competition**

```
1
2    # Merge to get competition names
3    assist_per_competition = df[['competition_id', 'assists']].groupby('competition_id').sum().reset_index()
4    assist_per_competition = assist_per_competition.merge(competitions[['competition_id', 'name']], on='competition_id')
5
6    # Visualization
7    plt.figure(figsize=(14, 8))
8    sns.barplot(data=assist_per_competition, x='name', y='assists', palette='crest')
9    plt.title('Total Assists by Competition')
10   plt.xlabel('Competition')
11   plt.ylabel('Total Assists')
12   plt.xticks(rotation=90)
13   plt.show()
14
```

✓ 2 sec - Command executed in 1 sec 997 ms by aamnanaz403 on 2:56:17 AM, 8/22/24
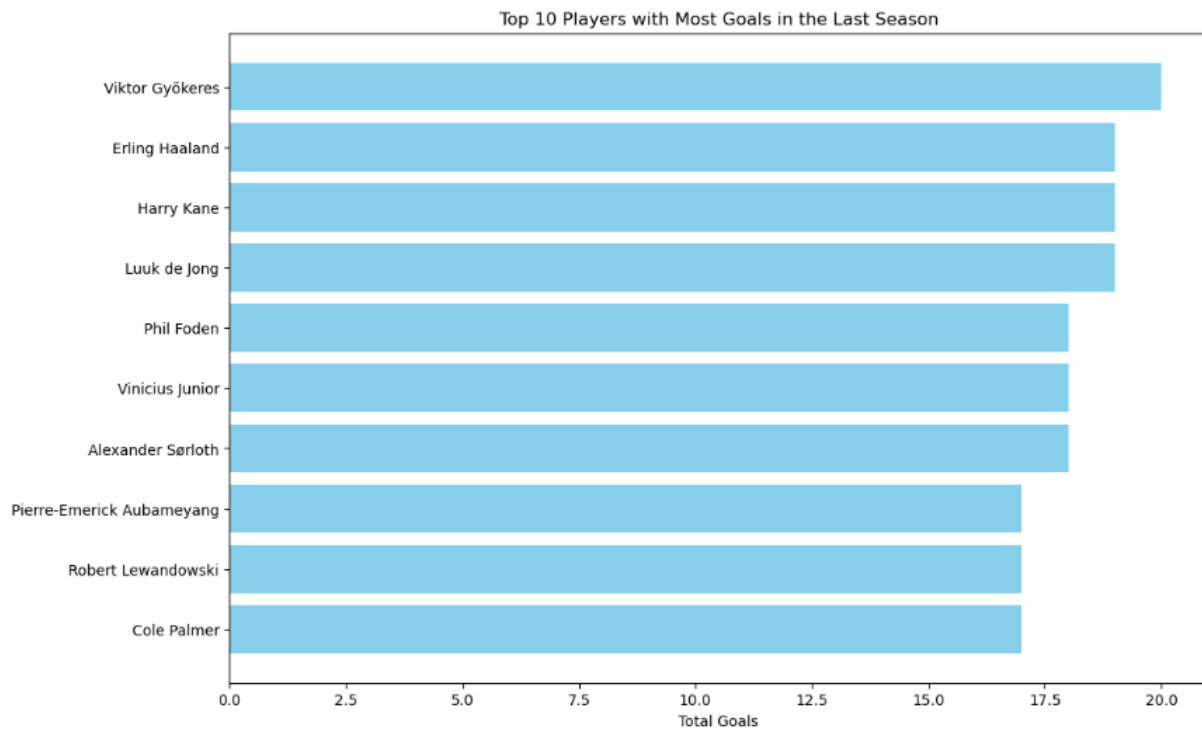
OUTPUT:-

QUERY: 6 **Top 10 Players with Most Goals in the Last Season**

```
1
2
3    # Filter for the last season
4    last_season = df['date'].max().split('-')[0]   # assuming the latest date corresponds to the last season
5    appearances_last_season = df[df['date'].str.startswith(last_season)]
6
7    # Aggregate goals by player
8    player_goals = appearances_last_season.groupby('player_id')['goals'].sum().reset_index()
9    player_goals = player_goals.sort_values(by='goals', ascending=False).head(10)
10
11   # Merge with players data for names
12   player_goals = player_goals.merge(players[['player_id', 'name']], on='player_id')
13
14   # Plot
15   plt.figure(figsize=(12, 8))
16   plt.barh(player_goals['name'], player_goals['goals'], color='skyblue')
17   plt.xlabel('Total Goals')
18   plt.title('Top 10 Players with Most Goals in the Last Season')
19   plt.gca().invert_yaxis()
20   plt.show()
21
```
✓  1 sec - Command executed in 1 sec 950 ms by aamnanaz403 on 2:57:19 AM, 8/22/24

OUTPUT:-

QUERY: 7 **Comparison of Average Player Market Value by Nationality**

```
1
2    # Aggregate average market value by nationality
3    market_value_by_nationality = players.groupby('country_of_citizenship')['market_value_in_eur'].mean().reset_index()
4
5    # Sort and plot
6    market_value_by_nationality = market_value_by_nationality.sort_values(by='market_value_in_eur', ascending=False).head(10)
7
8    plt.figure(figsize=(14, 8))
9    sns.barplot(data=market_value_by_nationality, x='country_of_citizenship', y='market_value_in_eur', palette='coolwarm')
10   plt.title('Comparison of Average Player Market Value by Nationality')
11   plt.xlabel('Nationality')
12   plt.ylabel('Average Market Value (EUR)')
13   plt.xticks(rotation=45)
14   plt.show()
15
```
✓  <1 sec - Command executed in 627 ms by aamnanaz403 on 3:01:00 AM, 8/22/24

OUTPUT:-



Comparison of Average Player Market Value by Nationality