

European Soccer Database

BY: Anik Shikarpuri

<https://www.kaggle.com/datasets/hugomathien/soccer>

The database I selected for this assignment is on European soccer. It has information about players, teams, matches, player attributes, and team attributes. This information is really interesting because we can dissect each player individually. And it is a great tool for teams looking for new players.

The dataset has very specific details about player attributes such as speed, dominant foot, shooting power, freekick rating, penalties scored, tackles won, player overall rating, and much more. There are eight different entries for this database. There is also data on scores for games between the European top teams from 2008 - 2020. The reason this information is relevant is because we can take what we know from the team's performance and find new players on the market that could improve their team as a whole. For example, in soccer you want the left wing position to usually be left footed so they can cross the ball in more easily. So using the database we can find a left winger who is left footed with a high player rating. Which will make it much easier for teams to find good replacements for their current players. The data was retrieved from FIFA (Federation Internationale de Football Association). Which is the organization that runs all these games and teams. So the data given is accurate because it was taken from the source itself.

When I was loading the data there were a lot of challenges I faced. First of which was finding which data I actually need. There was so much data given that it was hard to find what I needed. As you can see in the image below the database was divided into 8 different tables and each of those tables were then divided into several parts. I loaded the data by just looking going to the 'browse data' tab and that allowed me to get a look at what I want to use.

Name	Type	Schema
Tables (8)		
Country	CREATE	
League	CREATE	
Match	CREATE	
Player	CREATE	
Player_Attributes	CREATE	
Team	CREATE	
Team_Attributes	CREATE	
sql_sequence	CREATE	
Indices (0)		
Views (0)		
Triggers (0)		

Database Structure Browse Data Edit Pragmas Execute SQL												
Tables: Player_Attribu												
	jumping	stamina	strength	long_shots	aggression	interceptions	positioning	vision	penalties	marking	sta	
	Filter	Filter	Filter	Filter	Filter	Filter	Filter	Filter	Filter	Filter	Filter	Filter
2	58	54	76	35	71	70	45	54	48	65		
3	58	54	76	35	63	41	45	54	48	65		
4	58	54	76	34	62	40	44	53	47	62		
5	58	54	76	34	62	40	44	53	47	62		
6	85	79	56	62	68	67	60	66	59	76		
7	85	79	56	60	68	67	60	66	59	76		
8	84	79	56	59	67	66	58	65	59	76		
9	84	79	56	58	67	66	58	65	59	76		
10	84	79	56	58	67	66	58	65	59	76		
11	84	79	56	58	67	66	58	65	59	76		
12	84	79	56	58	67	66	58	65	59	76		
13	84	79	56	58	67	66	58	65	59	76		
14	84	79	50	56	66	65	57	64	58	73		
15	84	79	50	56	66	65	57	64	58	73		
16	84	79	50	56	66	65	57	64	58	73		
17	84	79	50	56	66	62	57	64	58	73		
18	84	80	50	56	66	62	57	64	58	73		
19	84	80	50	56	66	62	57	64	58	73		
20	84	79	49	55	66	62	57	61	58	73		
21	84	79	49	55	66	62	57	61	58	73		
22	84	79	49	55	66	62	57	61	58	73		
23	84	77	48	56	63	62	56	58	58	73		
24	84	76	48	55	63	62	56	58	58	73		
25	77	74	48	55	63	62	56	58	58	72		

The picture above shows the initial table for the Player Attributes table. Which was one of the tables I used for my queries.

I created three queries to answer the question of which teams really need new players, and what players should they look for. I divided this into three parts.

1. Which teams have lost by the biggest goal differentials at home games?

New Database

Open Database

Write Changes

Revert Changes

Open Project

Save Project

Attach Databases

Database Structure

Browse Data

Edit Pragmas

Execute SQL

1

2

3

4

5

6

7

8

9

10

11

12

13

14

```

select R.away_team_goal, R.home_team_goal,
S.team_long_name
FROM Match AS R LEFT OUTER JOIN Team AS S
ON (R.home_team_api_id = S.team_api_id)
where home_team_goal < away_team_goal
order by away_team_goal DESC

```

away_team_goal	home_team_goal	team_long_name
9	0	ES Troyes AC
8	1	FC St.Pauli
8	0	UD Almeria
8	0	Córdoba CF
8	2	RC Deportivo de La Coruña
8	0	RC Deportivo de La Coruña
7	1	Beerschot AC
7	1	KV Oostende

Execution finished without errors.

Result: 7466 rows returned in 248ms

At line 1:

```

select R.away_team_goal, R.home_team_goal,
S.team_long_name
FROM Match AS R LEFT OUTER JOIN Team AS S
ON (R.home_team_api_id = S.team_api_id)
where home_team_goal < away_team_goal
order by away_team_goal DESC

```

Mode: Text

NULL

Type of data currently in
0 byte(s)

Identity Select an id

Name

2. Which attackers would be best suited for these teams?

The screenshot shows a database application interface with a menu bar (New Database, Open Database, Write Changes, Revert Changes, Open Project, Save Project, Attach) and a toolbar. The 'Execute SQL' tab is active. The SQL editor contains the following query:

```
1 select R.player_name, S.attacking_work_rate, S.finishing, S.heading_accuracy, S.dribbling, S.acceleration, S.sprint_speed, S.shot_power, S.overall_rating
2 FROM Player AS R LEFT OUTER JOIN Player_Attributes AS S
3 ON (R.player_fifa_api_id = S.player_fifa_api_id)
4
5
6
7 WHERE
8 attacking_work_rate = 'high' and
9 overall_rating > 85
10
11
12
13
14
```

The results are displayed in a table with 10 columns: player_name, attacking_work_rate, finishing, heading_accuracy, dribbling, acceleration, sprint_speed, shot_power, overall_rating. The results show 15 rows, all for 'Andres Iniesta' with an overall_rating of 89.

	player_name	attacking_work_rate	finishing	heading_accuracy	dribbling	acceleration	sprint_speed	shot_power	overall_rating
8	Alexis Sanchez	high	76	66	90	92	89	79	86
9	Andres Iniesta	high	73	54	90	76	75	65	88
10	Andres Iniesta	high	73	54	90	76	75	65	88
11	Andres Iniesta	high	73	54	92	76	75	65	89
12	Andres Iniesta	high	73	54	92	76	75	65	89
13	Andres Iniesta	high	73	54	92	76	75	65	89
14	Andres Iniesta	high	73	54	91	76	75	65	89
15	Andres Iniesta	high	73	54	91	76	75	65	89

Execution finished without errors.
Result: 618 rows returned in 68ms
At line 2:
select R.player_name,
S.attacking_work_rate, S.finishing, S.heading_accuracy, S.dribbling, S.acceleration, S.sprint_speed, S.shot_power, S.overall_rating
FROM Player AS R LEFT OUTER JOIN Player_Attributes AS S
ON (R.player_fifa_api_id = S.player_fifa_api_id)

WHERE
attacking_work_rate = 'high' and
overall_rating > 85

3. Which defensive players are best suited for these teams?

The screenshot shows the same database application interface. The SQL editor contains the following query:

```
1 select
2 R.overall_rating, R.defensive_work_rate, R.standing_tackle, R.sliding_tackle, R.interceptions,
3 S.player_name, S.weight
4 FROM Player AS S LEFT OUTER JOIN Player_Attributes AS R
5 ON (R.player_fifa_api_id = S.player_fifa_api_id)
6
7 WHERE
8 defensive_work_rate = 'high'
9 AND overall_rating > 85
10
11
12
13
14
```

The results are displayed in a table with 7 columns: overall_rating, defensive_work_rate, standing_tackle, sliding_tackle, interceptions, player_name, weight. The results show 18 rows, all for 'Andres Iniesta' with an overall_rating of 89.

	overall_rating	defensive_work_rate	standing_tackle	sliding_tackle	interceptions	player_name	weight
11	86	high	42	36	42	Alexis Sanchez	137
12	87	high	74	64	76	Andres Iniesta	150
13	87	high	74	64	76	Andres Iniesta	150
14	87	high	74	64	76	Andres Iniesta	150
15	87	high	74	64	76	Andres Iniesta	150
16	87	high	72	64	76	Andres Iniesta	150
17	87	high	70	66	76	Andres Iniesta	150
18	86	high	63	61	66	Angel Di Maria	165

Execution finished without errors.
Result: 271 rows returned in 96ms
At line 1:
select
R.overall_rating, R.defensive_work_rate, R.standing_tackle, R.sliding_tackle, R.interceptions,
S.player_name, S.weight
FROM Player AS S LEFT OUTER JOIN Player_Attributes AS R
ON (R.player_fifa_api_id = S.player_fifa_api_id)

WHERE
defensive_work_rate = 'high'
AND overall_rating > 85

As you can see I loaded the data by taking different parts from different tables and creating a correlation between them. In the three tables I created, I connected them through their similar columns which allowed me to create a connection. One of the challenges I had was finding the correlation and how to connect the two tables together. And I also had an issue with my table

having repeating columns. To summarize the three tables, in table 1 I made a query that finds the teams with the biggest home game losses. I did this by finding the team api ID from the 'player' and 'player attributes' dataset. That way I can find the connection and then give the names of the team with the biggest home loss. For the second table I found a player with player ratings higher than 85. And gave all the attributes of an attacking player. So the teams can look for good attacking players with all the relevant attributes that could help them find the right player. I then did the same thing but for defensive players in table three. So this will allow teams to see which of them are performing the worst. And this will lead them to wanting new players which table 2 and 3 help them find. I was able to pull the same data from different tables. But I think no relation is better because it has lower data storage and has a higher scalability.

Part 2: NoSQL

So now we will take two of the questions I used for SQL, and try to find the same data in MongoDB. For MongoDB I will try to find solid defensive and offensive players that can make a significant impact to the team. In other words, which attackers would be best suited for these teams? And which defenders would be best suited for these teams? Like I stated previously, the NoSQL data tool I used was MongoDB. The difference is MongoDB uses python, therefore using a different language than SQL. It also uses a dynamic schema and is best suited for hierarchical data storage. However, it is not good for complex queries. I had a lot of issues when trying to load the data. The problem I had was trying to convert a .zip to a .csv. The data I used was downloaded as a .zip file, which needed to be converted to .csv. I tried different website converters but nothing was working. So, I opened the data in SQLite and then I exported the project as a .csv file so I could then open it in MongoDB. This solved the problem and allowed me to begin writing the queries in MongoDB. Below are my two queries:

1) Which attacking player will be best suited for a team?

My Query:

```
Filter: { _rating: { $gt: "85" }, preferred_foot: "left", attacking_work_rate: "high", shot_power: { $gt: "85" } }
Project: { field: 0 }
Sort: { field: -1 } or [{"field": -1}]
Collection: { locale: 'simple' }
```

ADD DATA EXPORT COLLECTION

```
{ "_id": "ObjectID('639253f07d170bafa73f9b91')",
  "id": "13139",
  "player_fifa_api_id": "183898",
  "player_api_id": "48589",
  "date": "2014-05-02 00:00:00",
  "overall_rating": "86",
  "potential": "86",
  "preferred_foot": "left",
  "attacking_work_rate": "high",
  "defensive_work_rate": "medium",
  "crossing": "90",
  "finishing": "75",
  "heading_accuracy": "53",
  "short_passing": "91",
  "volleys": "77",
  "dribbling": "87",
  "curve": "83",
  "free_kick_accuracy": "72",
  "long_passing": "81",
  "ball_control": "86",
  "acceleration": "88",
  "sprint_speed": "88",
  "agility": "90",
  "reactions": "79",
  "balance": "79" }
```

```
1 overall_rating: {
2   $gt: "85"
3 },
4 preferred_foot: 'left',
5 attacking_work_rate: 'high',
6 shot_power: {
7   $gt: "85"
8 }
9
10
```

2) Which defensive players are best for the team?

My Query:

```
1 overall_rating: {
2   $gt: "85"
3 },
4 defensive_work_rate: 'high',
5 standing_tackle: {
6   $gt: "85"
7 },
8 },
9 sliding_tackle: {
10   $gt: "85"
11 }
12
```

FIFA.FIFA DB

Documents Aggregations Schema Explain Plan Indexes Validation

Filter: { \$gt:"85"},defensive_work_rate:'high',standing_tackle:{ \$gt:"85"},sliding_tackle:{ \$gt:"85"} Reset Find

Project: { field: 0 }

Sort: { field: -1 } or [{"field": -1}] MaxTime

Collection: { locale: 'simple' } Skip Lin

ADD DATA EXPORT COLLECTION

```
heading_accuracy: "74"
short_passing: "84"
volleys: "78"
dribbling: "83"
curve: "79"
free_kick_accuracy: "68"
long_passing: "80"
ball_control: "83"
acceleration: "79"
sprint_speed: "79"
agility: "75"
reactions: "87"
balance: "77"

_id: "ObjectID('639253f07d170bafa73f9b94b')",
id: "131421",
player_fifa_api_id: "181872",
player_api_id: "49939",
date: "2014-05-02 00:00:00",
overall_rating: "86",
potential: "88",
preferred_foot: "right",
attacking_work_rate: "high",
defensive_work_rate: "high"
```

So here we see that I was able to get the same data from MongoDB as SQLite. However, there are some issues. First of all it did find information I wanted, but it also gave all the information that falls under the data table with it. So it did give me the information I asked for but it also gave me data I didn't want. The other problem I had was that I couldn't use the 'match' table because when I tried adding the data it said it was "750,000" documents. Which meant it took way too long for it to import the dataset. So I chose not to use it. Lastly, the programming in itself is much more complicated for MongoDB than SQLite. After using both I would say SQLite would be preferable when trying to create complex queries. And the NoSQL data tool is more applicable when working with hierarchical data and if you want it to run well on the cloud. When working with large scale data the NoSQL tools are much better than SQL. This is because NoSQL was designed to be used across large distributed systems due to the dynamic

schema. The other benefit being large amounts of data can be stored in NoSQL databases. So to conclude, NoSQL is much more suitable for large data sets but, when working with more complicated or challenging queries SQL is more suitable. After working on this project I would say if you are looking at data for personal use, SQL is much more user friendly and easier to find the data you are looking for.