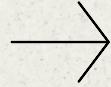# FIFA Unleashed: Eight Seasons of FIFA Player Trends (2015 – 23)

Group 7
Niharika Ahirekar
Rutvi Bheda
Mohammed Raeesul Irfan
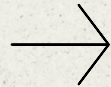
# TABLE OF CONTENTS

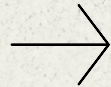O1    Phase I            →            About our dataset, how we formed and loaded it.

O2    Phase II           →            Interesting queries over our dataset, Indexing our dataset for faster query execution, Functional Dependencies, Normalization.

O3    Phase III          →            Document-oriented model Vs Relational model, Operations performed on the dataset:
Data Cleaning, Data Integration, Itemset Mining

# FIFA Unleashed: Eight Seasons of Player Trends

- ❏ Based on the videogame FIFA
- ❏ Allows comparisons for the same players across the last 8 version of the videogame (FIFA'15 – FIFA'23)
- ❏ Basically the data on the player cards

# Player Stats:

- ❏ Player stats are the real life data of the players analysed by EA
- ❏ The 6 main stats mentioned in these cards are:
    Pace,
    Shooting,
    Passing,
    Dribbling,
    Defending,
  and   Physical

However, these stats are an average of the complete player stats



PLAYER BIO  PLAYER DETAILS  **ATTRIBUTE DETAILS**  PLAYER TRAITS

HAALAND

Player Chemistry ▬▬▬▬▬▬▬ 3   Chemistry Style   **HAWK**

88
ST

🇳🇴

HAALAND

| 89 PAC | 80 DRI |
| 91 SHO | 49 DEF |
| 65 PAS | 87 PHY |

**GOLD RARE**

| PACE | SHOOTING | PASSING | DRIBBLING | DEFENDING | PHYSICAL |
|---|---|---|---|---|---|
| **89** | **91** | **65** | **80** | **49** | **87** |

| PACE | | SHOOTING | | PASSING | | DRIBBLING | | DEFENDING | | PHYSICAL | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Acceleration | 82 +4 | Att. Position | 89 +4 | Vision | 74 | Agility | 76 | Interceptions | 43 | Jumping | 74 +8 |
| Sprint Speed | 94 +4 | Finishing | 94 +4 | Crossing | 47 | Balance | 72 | Heading Acc. | 79 | Stamina | 81 |
| | | Shot Power | 94 +5 | FK Acc. | 62 | Reactions | 88 | Def. Aware | 44 | Strength | 93 +4 |
| | | Long Shots | 87 +8 | Short Pass | 74 | Ball Control | 82 | Stand Tackle | 53 | Aggression | 85 +8 |
| | | Volleys | 88 | Long Pass | 53 | Dribbling | 78 | Slide Tackle | 29 | | |
| | | Penalties | 76 +4 | Curve | 77 | Composure | 85 | | | | |

L2  R2 Change Player  ◄ R ► Views  ○ Back

| FIFA 15 | FIFA 16 | FIFA 17 | FIFA 18 | FIFA 19 | FIFA 20 |
|---------|---------|---------|---------|---------|---------|

**RONALDO**

| FIFA 15 (92 LW) | FIFA 16 (93 LW) | FIFA 17 (94 LW) | FIFA 18 (94 LW) | FIFA 19 (94 ST) | FIFA 20 (93 ST) |
|---|---|---|---|---|---|
| 93 PAC  91 DRI | 92 PAC  90 DRI | 92 PAC  91 DRI | 90 PAC  90 DRI | 90 PAC  90 DRI | 90 PAC  89 DRI |
| 93 SHO  32 DEF | 93 SHO  33 DEF | 92 SHO  33 DEF | 93 SHO  33 DEF | 93 SHO  35 DEF | 93 SHO  35 DEF |
| 81 PAS  79 PHY | 80 PAS  78 PHY | 81 PAS  80 PHY | 82 PAS  80 PHY | 81 PAS  79 PHY | 82 PAS  78 PHY |

H / L   5★ SM   4★ WF (FIFA 18)
5★ SM   4★ WF (FIFA 19)

**MESSI**

| FIFA 15 (93 CF) | FIFA 16 (94 RW) | FIFA 17 (93 RW) | FIFA 18 (93 RW) | FIFA 19 (94 CF) | FIFA 20 (94 RW) |
|---|---|---|---|---|---|
| 93 PAC  96 DRI | 92 PAC  95 DRI | 89 PAC  96 DRI | 89 PAC  95 DRI | 88 PAC  96 DRI | 87 PAC  96 DRI |
| 89 SHO  27 DEF | 88 SHO  24 DEF | 90 SHO  26 DEF | 90 SHO  26 DEF | 91 SHO  32 DEF | 92 SHO  39 DEF |
| 86 PAS  62 PHY | 86 PAS  62 PHY | 86 PAS  61 PHY | 86 PAS  61 PHY | 88 PAS  61 PHY | 92 PAS  66 PHY |

M / M   4★ SM   4★ WF (FIFA 18)
4★ SM   4★ WF (FIFA 19)

Player stats change over the years:

# Kaggle Data Set

The data set on Kaggle had many tables.
Due to the accumulation of stats of players, coaches and the teams over the years, the data was HUGE!

We decided to narrow our focus onto Male players and coaches.

We chose 2 tables to extract our data from:
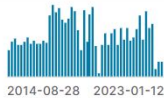- ❏ male_players_23.csv
- ❏ male_coaches_23.csv

female_coaches_23.csv
female_players (legacy)_2
female_players_23.csv
female_teams_23.csv
male_coaches_23.csv
male_players (legacy)_23.
male_players_23.csv
male_teams_23.csv

# Working on the CSV files



❏ 110 columns on unfiltered raw data
❏ Contained player personal details as well as the performance statistics all in one csv file

# DIVIDING THE CONTENT INTO TABLES:

## Players

24 columns
Contains only personal details
such as name, jersey number,
club, nationality_id, dob, etc
★     playerid : gives the
      unique id of the players
      over the years

## Metrics

72 columns
Contains attributes such as
positions (like center back,
right wing, etc), goalkeeping,
pace, etc
★     metrics_id: contains the
      unique id of the stats
      over the years

## Player Metrics

2 columns
Contains playerid
and metrics_id of the
corresponding
players

# MORE TABLES:

## Coach

4 columns
Contains personal details of the coach
★ Coach_id: unique id given to each coach

## Nationality

2 columns
Contains nationality_id and nationality
★ Nationality_id: unique id given to each country

# Document Oriented Modeling

❏ Document-oriented modeling is a methodology employed in document-oriented databases, a category of NoSQL databases.

❏ This approach is centered around the concept of documents, which are independent data entities akin to JSON or XML objects.

❏ Unlike traditional relational databases that require a predefined schema and store data in rows across multiple interconnected tables, document-oriented databases are either schema-less or have very flexible schemas. This means that each document can have a unique structure, allowing for a wide variety of data types and structures within the same database.

❏ Some examples are - MongoDB, RavenDB , Amazon DynamoDB, Couchbase, etc.

# Example of data in MongoDB



**fifa.coach**

Documents  Aggregations  Schema  Indexes  Validation

Filter  🕐 ▾  Type a query: { field: 'value' } or **Generate query** ◀

⊕ ADD DATA ▾   ☁ EXPORT DATA ▾

```
_id: ObjectId('654eb48d02bcefd6570ab792')
coach_id: 24
short_name: "D. Unsworth"
long_name: "David Unsworth"
dob: "1973-10-16"
nationality_name: "England"
year: "23"
```

**fifa.metrics**

Documents  Aggregations  Schema  Indexes  Validation

Filter  🕐 ▾  Type a query: { field: 'value' } or **Generate query** ◀

⊕ ADD DATA ▾   ☁ EXPORT DATA ▾

```
st: "87+3"
rs: "87+3"
lw: "90"
lf: "91"
cf: "91"
rf: "91"
rw: "90"
lam: "92-1"
cam: "92-1"
ram: "92-1"
lm: "89+2"
lcm: "85+3"
cm: "85+3"
rcm: "85+3"
rm: "89+2"
lwb: "64+3"
ldm: "63+3"
```

# Five Interesting Queries

1. Identify the particular sports club exhibiting a high count of players surpassing a weight threshold of 75 kilograms.

2. Determine the specific national origin attributed to the player with the most substantial compensation package whose overall performance registers above the threshold of 63.

3. Enlist all coaching personnel hailing from Germany, specifically those whose surname initiates with the consonantal sequence 'Kl' and whose date of birth falls after the year 1955.

4. Identify top 2 athletes within the Arsenal roster who exhibits the most elevated measure of their shooting skill.
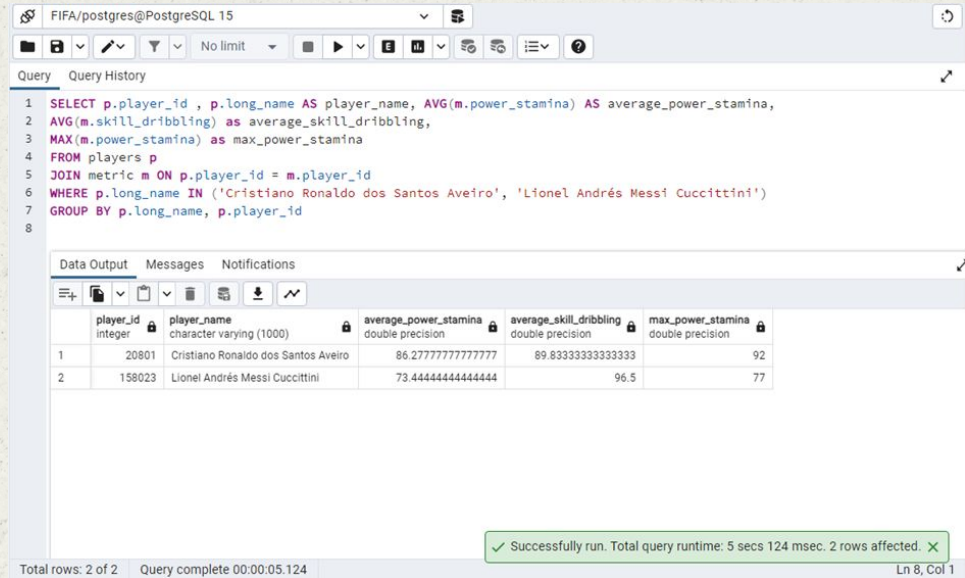


```
1  SELECT m.player_id, p.long_name, p.club_name, MAX(m.shooting) as max_shooting
2  FROM players p
3  JOIN metric m ON p.player_id = m.player_id
4  WHERE p.club_name = 'Arsenal'
5  GROUP BY p.club_name, m.player_id, p.long_name
6  LIMIT 2;
7
8
9
```

| player_id integer | long_name character varying (1000) | club_name character varying (100) | max_shooting double precision |
|---|---|---|---|
| 8473 | Tomáš Rosický | Arsenal | 70 |
| 45119 | Mikel Arteta Amatriain | Arsenal | 73 |

Total rows: 2 of 2    Query complete 00:00:17.187

5. Comparing the statistics for Lionel Messi and Cristiano Ronaldo's performance in the last 8 years.

# Indexing

❏ Indexing in databases is a technique that improves the speed of data retrieval operations.

❏ The primary advantage of indexing is significantly faster data retrieval compared to scanning the entire table, especially in large databases.

❏ This speed-up is crucial for performance, particularly in search and join operations.

| Queries | Time without Indexing (seconds) | Time after Indexing (seconds) |
|---------|---------------------------------|-------------------------------|
| **Query 1** | 8.459 | 6.1 |
| **Query 2** | 6.947 | 0.099 |
| **Query 3** | 0.107 | 0.075 |
| **Query 4** | 17.187 | 0.714 |
| **Query 5** | 5.124 | 3.8 |

# Functional Dependencies

❏ A functional dependency is a relationship between two sets of attributes in a relational database.
❏ It represents the fact that the values of one set of attributes uniquely determine the values of another set. In other words, if we have a functional dependency A -> B, it means that for every unique combination of values in A, there is a unique corresponding combination of values in B.
❏ Types: Trivial, Non-trivial, Transitive.
❏ In the players table, several functional dependencies can be identified based on the relationships between different attributes. The primary identity dependency is established by the player ID, which uniquely determines all other attributes associated with a player, including personal details, club affiliations, and skill-related information.
❏ In the nationality table, the primary key is nationality_id, and it uniquely identifies each nationality. Therefore, the functional dependency can be expressed as: nationality_id → nationality_name.
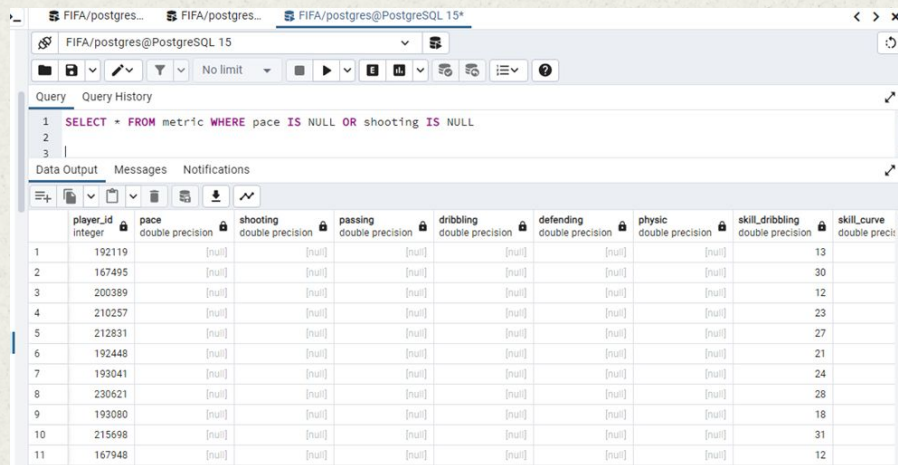
# Normalization

❏ Normalization involves breaking down a large table into smaller, more organized groups of related information to reduce redundancy and improve data integrity.
❏ It ensures that data is neatly arranged and that identical information isn't repeated in multiple places.
❏ In the players table, the foreign key nationality_id referencing another table, aids to normalization. This design choice contributes to reducing redundancy in the database, as detailed information about nationalities is stored in separate tables, and player records only contain references to these entities.
❏ The use of foreign keys also ensures data integrity, as it enforces referential consistency between the players table and the referenced table nationality.
❏ The coach table has a unique identifier, coach_id and individual columns for coach details. It follows the basics of the First Normal Form (1NF) by keeping each piece of information separate.

# Data Cleaning

❏   Data cleaning in SQL involves manipulating and modifying the data in a database to handle issues such as missing or inconsistent values. It aims to enhance the reliability and accuracy of datasets.

❏   We had some null values in some of the column, hence we would be eliminating those.

❏   Some ways of cleaning data:

   ❏   Handling missing values
   ❏   Removing duplicates
   ❏   Data Validation
   ❏   Data Type Conversion
   ❏   Data Integration
   ❏   Error Correction

# Data Cleaning



Before



After

# Data Integrity

❏ Maintaining data integrity in SQL, especially when joining tables, involves enforcing constraints, relationships, and actions that ensure the accuracy and consistency of the data.

❏ Additionally, integrating tables can significantly enhance the accuracy of the data. By breaking down data into tables and defining clear relationships, reducing the risk of data redundancy and inconsistencies.

❏ Hence we had integrated some tables in our database for ease of access and getting accuracy.

# Relational Model vs Document Oriented Model

| Based on | Relational Model | Document Oriented Model |
|---|---|---|
| **Schema** | Consists of a fixed schema | It is schema less consisting of documents of different data structures. |
| **Data Structure** | Tables with rows and columns | The documents are of JSON,XML type |
| **Query Language** | SQL | Differs system by system |
| **Normalization** | Highly normalized data | Often denormalized, data might be duplicated. |
| **Data Integrity** | Good data integrity consisting of ACID properties. | Mostly focused on performance and scalability. |

# Itemset Mining

❏ In itemset mining, the concept of lattices plays a pivotal role, particularly in how itemsets are organized and systematically explored.

❏ A lattice in this context represents the complete set of all possible itemsets, structured in a way that mirrors their subset-superset relationships.

❏ The lattice starts with the smallest possible itemsets (individual items) and expands to include larger and more complex combinations of items.

❏ Algorithms like Apriori leverage this lattice structure to efficiently mine for frequent itemsets.

❏ In our project, we implemented itemset mining and got a total of 5 lattices.

❏ Lattice 1 - 5132 combinations of players who have been in the same club for a long time.

❏ Lattice 2 - 741 combination of players.

❏ Lattice 3 - 39 combination of players.

❏ Lattice 4 - 7 combination of players - meaning 7 combinations of 4 players.

❏ And lastly, lattice 5 had no combination of players.

# Conclusion

❏ Through this project, we have discovered and understood various concepts of databases and querying.

❏ We started with getting the dataset, processing it , exploring SQL and querying it.

❏ Followed by generating interesting queries and understanding various clauses and features in SQL.

❏ Various concepts like indexing, functional dependencies, normalization were implemented.

❏ We introduced document oriented system and explored the loading of data to it.

❏ The last phase consisted of data cleaning, data integrity and itemset mining.

❏ Overall, this project helped us to dive deep into various database concepts and also taught us the ways to handle huge datasets.

# THANK YOU