

CM3010: Databases and Advanced Data Techniques

Midterm Assignment

Arjun Muralidharan

24th December 2021

Contents

1	Find and critique a dataset	3
1.1	Finding the dataset	3
1.2	Critiquing the dataset	3
1.3	Interest in the data set	3
1.4	Questions to answer	4
1.5	Model your data	4
1.6	Create the database	5
1.7	Create a simple web application	6

1 Find and critique a dataset

1.1 Finding the dataset

I will use a dataset of sports matches, specifically test cricket matches from 2004 until 2021, obtained from [1].

The data describes ball-by-ball statistics and can be considered **pre-existing external data**. It is assembled with automated **data curation** to assemble the data from various sources. Some general points on the structure of the data:

- Each match's data is organised in a separate JSON file.
- Each file includes some meta data about the creation and the event that this match was part of.
- All ball-by-ball data is stored in **nested objects** within the JSON structure.
- This leads to lots of repetition of information as this follows a **tree-like** structure rather than a relational one.

Hence the data needs to be normalized and mapped to a relational model in order to use it in a relational database.

1.2 Critiquing the dataset

One advantage of the dataset is that it is very **up to date** as new matches are added every day. This makes the data set highly usable and accurate, at least starting from 2004 onwards. The dataset is also **lightweight to use** given that each match has it's own file; new data does not change or affect the size and integrity of previous data.

However, the data is distributed and hence **harder to consolidate** into a single repository, and it contains lots of repeated data as player names occur again for each individual delivery (a ball played) where the player was involved. Matches that belong to the same event or test match series are in separate files and the event name is repeated in all the files related to an event.

Event names are further not unique, e.g. "New Zealand tour of India" may have happened many times in the past 17 years, so the year of the event needs to be considered when defining what a unique event is.

1.3 Interest in the data set

I chose this data set as I personally enjoy watching Cricket as a sport, and believe the data set to be interesting for the following reasons:

- Cricket is a statistics-heady sport, with lots of data recorded officially on a ball-by-ball basis

- Generally, sports data is rarely available in an openly available format and often hidden behind paywalls or graphical interfaces, gatekept by sports companies and sports news sites. This dataset made openly available a lot of information about this fascinating sport.
- The data is structured in an interesting manner that in its current form is not suited for relational representation and many questions cannot be answered without transforming the data into a relational form.

1.4 Questions to answer

By transferring a part of the data into relational form, I aim to answer the following questions which cannot be answered with the current source structure.

1. How many test matches have the test-playing nations played since 2004?
2. Who are the 10 batsmen with the highest number of runs scored since 2004?
3. Which 10 players are the most successful bowlers in terms of wickets taken and runs conceded since 2004?
4. What is the most common type of dismissal since 2004?

1.5 Model your data

I have built a tentative ER diagram shown in [Figure 1](#). Some notes on normalization:

1. The core entities are *players* and *matches*. A match in cricket is played through *deliveries*, or individual balls. Hence, players are involved in deliveries, and called up to line-ups to play in matches. These relationships are modelled with linking entities that would eventually result in link tables.
2. This ER diagram is not a complete representation, as further, more granular normalization might be warranted based on the use case. For example, the entity **delivery role** might need an entity that stores the possible roles, from which an attribute is drawn (e.g. bowler, batter, fielder, etc.).

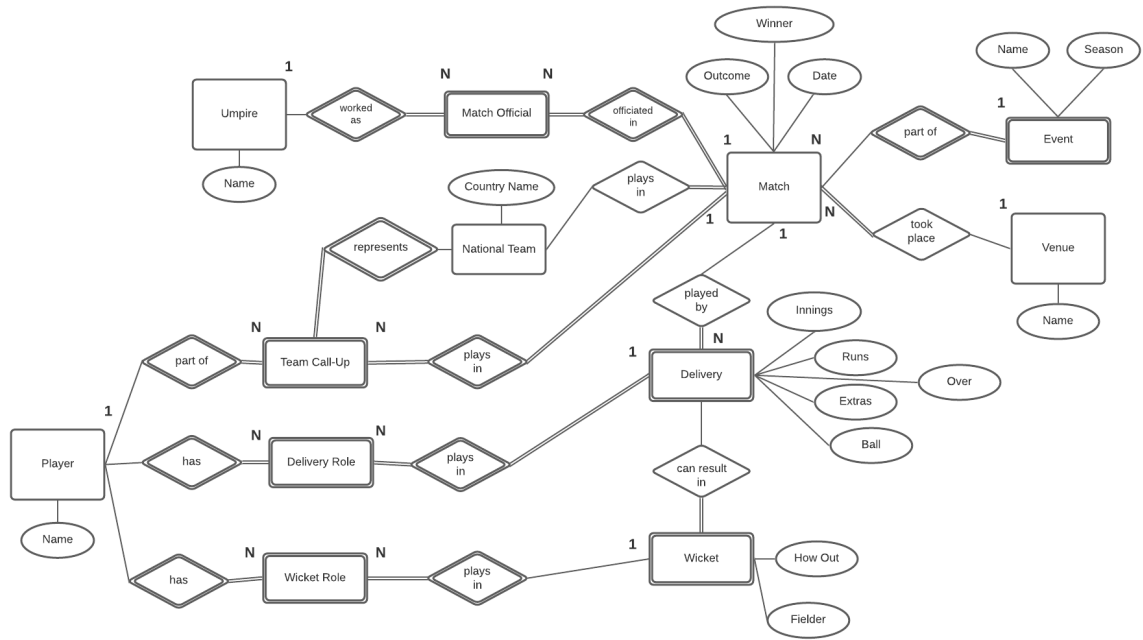


Figure 1. *ER diagram for Cricsheet data*

1.6 Create the database

The database implementation does not completely reflect the ER diagram exactly, because the scope far exceeds the current needs of answering the questions. Therefore, I took a slightly simplified approach and did not implement every link table shown above. I recognize this shortcoming but want to demonstrate that implementing additional link tables would be trivial.

Listing 1 Listing: Database Implementation

```
1 CREATE TABLE `player` (  
2   `name` VARCHAR(100),  
3   PRIMARY KEY (`name`)  
4 );  
5  
6 CREATE TABLE `team` (  
7   `country_name` VARCHAR(50),  
8   PRIMARY KEY (`country_name`)  
9 );  
10  
11 CREATE TABLE `match` (  
12   `match_id` INT NOT NULL AUTO_INCREMENT,  
13   `team_A` VARCHAR(50) NOT NULL,  
14   `team_B` VARCHAR(50) NOT NULL,  
15   `date` DATE NOT NULL,  
16   `winner` VARCHAR(50),  
17   `outcome` VARCHAR(100) NOT NULL,  
18   PRIMARY KEY (`match_id`),  
19   FOREIGN KEY (`team_A`) REFERENCES `team`(`country_name`) ON DELETE CASCADE,  
20   FOREIGN KEY (`team_B`) REFERENCES `team`(`country_name`) ON DELETE CASCADE,  
21   FOREIGN KEY (`winner`) REFERENCES `team`(`country_name`) ON DELETE CASCADE,  
22   CONSTRAINT unique_match UNIQUE (`team_A`, `team_B`, `date`)  
23 );  
24  
25 CREATE TABLE `delivery` (  
26   `match_id` INT NOT NULL,  
27   `innings` INT NOT NULL,  
28   `over` INT NOT NULL,  
29   `ball` INT NOT NULL,  
30   `runs` INT NOT NULL,  
31   `extras` INT NOT NULL,  
32   `batter` VARCHAR(50) NOT NULL,  
33   `bowler` VARCHAR(50) NOT NULL,  
34   `nonstriker` VARCHAR(50) NOT NULL,  
35   PRIMARY KEY (`match_id`, `innings`, `over`, `ball`),  
36   FOREIGN KEY (`match_id`) REFERENCES `match`(`match_id`),  
37   FOREIGN KEY (`batter`) REFERENCES `player`(`name`),  
38   FOREIGN KEY (`bowler`) REFERENCES `player`(`name`),  
39   FOREIGN KEY (`nonstriker`) REFERENCES `player`(`name`)  
40 );  
41  
42 CREATE TABLE `wicket`(  
43   `playerout` VARCHAR(50) NOT NULL,  
44   `fielder` VARCHAR(50),  
45   `kind` VARCHAR(50),  
46   `match_id` INT NOT NULL,  
47   `innings` INT,  
48   `over` INT,  
49   `ball` INT,  
50   PRIMARY KEY (`match_id`, `innings`, `over`, `ball`),  
51   FOREIGN KEY (`match_id`, `innings`, `over`, `ball`) REFERENCES `delivery`(`match_id`, `innings`, `over`, `ball`),  
52   FOREIGN KEY (`playerout`) REFERENCES `player`(`name`)  
53 );
```

Reflecting on this implementation, for a more scalable version I would implement more linking tables. Also, I have not implemented tables unneeded for the questions posed, such as Umpires, Events or Venues. If I were building a more full-fledged cricket scoring app, I would implement this analogously.

1.7 Create a simple web application

The web application is hosted in the lab environment. Some notes in running it:

1. It can be started with the `npm start` command.

2. The database is already populated, but can be re-populated by running `node db/cric.js`. This instantiates a class that processes the JSON files in the `matches` folder and adds them to the database. Because this operation can be long and heavy on the database, I decided to decouple this from running the application.
3. The application requires about a minute (in the lab environment) to query the database, as I have not implemented any caching of query results. Because some queries need to sift through every ball every played since 2004, this takes a moment. The app presents a warning and asks you to refresh the page after waiting a minute.

The app displays all the associated SQL queries on-demand. This functionality may not work in the lab environment as the browser preview can be slow and tedious to interact with; hence I recommend installing the app locally if desired.

Databases & Advanced Data Techniques

CricApp

This small application displays basic queries against a database of Cricket match data obtained from [CricSheet](#). The following sections outline some interesting questions we aimed to ask this dataset.

Most matches played by a team

What are the test-playing nations in our dataset, and how many matches has each of them played?

- It's clear that the traditional nations such as England and Australia have played far more matches than other teams
- New countries who have become test playing nations in recent countries rank at the bottom (Afghanistan, Ireland)
- There has been a match where an "ICC World XI" played. This match was played against [Australia in 2005](#).

Team	Matches
England	211
Australia	174
South Africa	149
Sri Lanka	146
West Indies	135
New Zealand	132
Pakistan	126
Bangladesh	88
India	88
Zimbabwe	37
Afghanistan	6
Ireland	3
ICC World XI	1

View SQL

```

SELECT
  Team,
  sum(matches) as Matches
FROM
(
  (
    select
      team_A as team,
      count(*) as matches
    from
      `match`
    GROUP BY

```

Figure 2. Screenshot of the working app

References

- [1] S. Rushe, “Cricsheet Matches,” Dec 2021. [Online]. Available: <https://cricsheet.org/matches/>