



HAL
open science

Soccer2014DS: a dataset containing player events from the 2014 World Cup

Marcos Roberto Ribeiro, Maria Camila N. Barioni, Sandra de Amo, Claudia Roncancio, Cyril Labbé

► **To cite this version:**

Marcos Roberto Ribeiro, Maria Camila N. Barioni, Sandra de Amo, Claudia Roncancio, Cyril Labbé. Soccer2014DS: a dataset containing player events from the 2014 World Cup. 32nd BRAZILIAN SYMPOSIUM ON DATABASES (SBBDD 2017), 2017, Uberlandia, Brazil. hal-01656405

HAL Id: hal-01656405

<https://hal.science/hal-01656405>

Submitted on 5 Dec 2017

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Soccer2014DS: a dataset containing player events from the 2014 World Cup

Marcos Roberto Ribeiro^{1,2}, Maria Camila N. Barioni²,
Sandra de Amo², Claudia Roncancio³, Cyril Labbé³

¹ Instituto Federal de Minas Gerais (IFMG), Bambuí, Brazil

² Universidade Federal de Uberlândia (UFU), Uberlândia, Brazil

³ Univ. Grenoble Alpes, CNRS, Grenoble INP, LIG, F-38000, Grenoble, France

marcos.ribeiro@ifmg.edu.br, {camila.barioni, deamo}@ufu.br,
{claudia.roncancio, cyril.labbe}@imag.fr

Abstract. *The player monitoring has become a common task in many sports. However, there is no public datasets of detailed soccer player events. Thus the creation of such datasets can be useful for diverse research fields such data mining, sports analytics and continuous preference queries. In this paper, we describe the construction of the dataset Soccer2014DS containing player events of the 2014 Soccer World Cup. This dataset is composed by the raw extracted data collected by a web crawler and by derived streams with new calculated attributes. We also explain how we are using this dataset in experiments related to the development of a new query language.*

1. Introduction

The monitoring of sport players during matches started in the end of XX century [Ali and Farrally 1991]. In this period, the researchers were more concerned to collect just health data from players. Since 2000, the monitoring tasks became more sophisticated due to the development of new technologies for GPS devices and softwares for video processing [Baca et al. 2009, Barris and Button 2008]. This new technologies allow to collect complex data to be used on detailed analysis of player events.

Despite the players monitoring be very common in official competitions, there are few public datasets with these information available. If we consider the soccer sport, to the best of our knowledge, there exist no public datasets with detailed player events. Thus, our main goal herein is to present the public dataset *Soccer2014DS* which contains player events of the 2014 Soccer World Cup. This dataset can be useful for diverse research fields such as data mining [Bialkowski et al. 2014, Gyarmati and Hefeeda 2015], sports analytics [Lucey et al. 2013, Perin et al. 2013] and continuous queries [Arasu et al. 2016].

This paper is organized as follows. Section 2 describes the original extracted data. Next, Section 3 presents the data streams derived from the original data. Section 4 discusses about the research opportunities and Section 5 explains the dataset limitations. Finally, Section 6 concludes the paper.

2. Data Sources

The creation of the *Soccer2014DS* dataset started with the extraction of the original data available on the Huffpost Data web site¹ [Boice et al. 2014]. This web site contains information about the 2014 Soccer World Cup provided by the company Opta Sports². These

¹<http://data.huffingtonpost.com/2014/world-cup>

²<http://www.optasports.com/>

data is used to display statistics and graphs about player events. Every match has an individual page where a user can delimit a time-line and see the details of this selection. Our first task was to study the source code of the Huffpost Data web site. Based on this study, we developed a web crawler to extract the data. The crawler starts the extraction in the page of the final match and follows the links to the remaining matches to complete the data collecting.

The extraction of the original raw data was performed in 2015. After this task, we organized the extracted data into the relations `Matches`, `Teams` and `Players` and the stream `Events`. The relations have just one instance and the duplicated data from all matches was eliminated. On the other hand, the stream `Events` has 64 instances (one instance per match) preserving all extracted event data. Appendix B presents the logical schema of the dataset (the derived data is addressed in Section 3).

The attributes of the relation `Matches` are `id` (match identifier), `date`, `time`, `venue` and `attendance`. Table 1(a) displays the attributes of the stream `Events`. The player coordinates (`x` and `y`) and final coordinates of the ball (`to_x`, `to_y` and `to_z`) are expressed as a percentage of the field dimensions. The attributes `type` and `outcome` are used to identify the move performed by players. Appendix A presents the events associated to every combination of values for these attributes. The attribute `field_pass` represents the continuous ball possession, `t` for true and `f` for false. The attribute `side` is the field side of the team, `H` for left side and `A` for right side. When the move is a pass to another player, the attribute `to` assumes the identifier of this player. For streams, we also must associate a `timestamp` for every tuple [Arasu et al. 2016], in the stream `Events` the `timestamp` is calculated using the attributes `min` and `sec` (`timestamp = min × 60 + sec`).

Table 1. Relation attributes: (a) Events (b) Players

| (a) | | (b) | |
|-------------------------------|--------------------------------|---------------------------------|-------------------------------|
| Attribute | Description | Attribute | Description |
| <code>id</code> | Event identifier | <code>id</code> | Player identifier |
| <code>period</code> | Period of math | <code>name</code> | Player name |
| <code>min, sec</code> | Minute and second of the event | <code>real_position</code> | Detailed position |
| <code>displaymin</code> | Displayed minute | <code>real_position_side</code> | Position side |
| <code>team</code> | Team identifier | <code>known_name</code> | Known name |
| <code>player_id</code> | Player identifier | <code>short_name</code> | Short name |
| <code>x, y</code> | Player coordinates | <code>last_name</code> | Last name |
| <code>type</code> | Type of event (move performed) | <code>first_name</code> | First name |
| <code>outcome</code> | Result of the event | <code>middle_name</code> | Middle name |
| <code>field_pass</code> | Continuous ball possession | <code>team_id</code> | Team identifier |
| <code>side</code> | Field side of team | <code>preferred_foot</code> | Preferred foot |
| <code>to_x, to_y, to_z</code> | Final coordinates of ball | <code>club</code> | Club |
| <code>to</code> | Player identifier of next move | <code>caps</code> | Matches played by player team |
| | | <code>goals</code> | Goals |
| | | <code>jersey_num</code> | Jersey number |
| | | <code>country</code> | Birth country |
| | | <code>birth_date</code> | Birth date |
| | | <code>position</code> | Position |

The relation `Teams` is composed by the attributes `id` (team identifier), `name` and `iso` (ISO acronym). Table 1(b) presents the attributes of the relation `Players`. The values for the attributes `real_position`, `real_position_side`, `preferred_foot` and `position` are shown in Table 2. Please see [Bakker 2015] for

more details about the data gathered by the company Opta Sports.

Table 2. Values for `Players` attributes

| Attribute | Values |
|---------------------------------|---|
| <code>real_position</code> | Attacking Midfielder, Central Defender, Central Midfielder, Defensive Midfielder, Full Back, Goalkeeper, Second Striker, Striker, Wing Back, Winger |
| <code>real_position_side</code> | Centre, Centre/Right, Left, Left/Centre, Left/Centre/Right, Left/Right, Right, Unknown |
| <code>preferred_foot</code> | Both, Left, Mostly Left, Mostly Right, Right, (empty) |
| <code>position</code> | Defender, Forward, Goalkeeper, Midfielder |

3. Derived Streams

After the data extraction described in the previous section, we created derived streams by applying cleaning and conversions over the original data. As described in the previous section, in order to know the exact player move, we must check the attributes `type` and `outcome` of the relation `Events`. In addition, the coordinates of the player and the ball, expressed by float values, could not be suitable for some applications where the user has to indicate a region of the soccer field. To deal with this situation, we decided to create new derived streams: `Moves(player_id, place, move)` for the performed moves; and `Places(player_id, place, ball, direc)` for the player positioning. The logical schema presented in Appendix B shows the relationship of the derived streams and the original relations.

The stream `Moves` contains just the moves performed by the players. So, we do not consider events without ball like cards, substitutions, etc. More precisely, the events types 17, 18, 19, 34, 43, 58, 60, 102 and the events with `(type, outcome)` equal to (5, 1), (6, 1), (53, 1) and (57, 1) are ignored. For the remaining move types, we use the mapping from `(Events.type, Events.outcome)` to `Moves.move` described in Table 3.

Table 3. Move mapping

| move | (type, outcome) |
|-------------------------------|--|
| <i>pass</i> | (1, 1), (59, 1) |
| <i>bpas</i> (bad pass) | (1, 0), (2, 1) |
| <i>lbal</i> (lost ball) | (3, 0), (7, 0), (44, 0), (50, 1), (51, 1), (57, 0), (59, 0), (61, 0) |
| <i>drib</i> (dribble) | (3, 1), (42, 1) |
| <i>foul</i> | (4, 0) |
| <i>fsuf</i> (foul suffered) | (4, 1), (55, 1) |
| <i>dled</i> (dribbled) | (45, 0) |
| <i>bout</i> (ball out) | (5, 0), (6, 0) |
| <i>brec</i> (ball recovery) | (7, 1), (44, 1), (49, 1), (56, 1), (61, 1) |
| <i>int</i> (interception) | (8, 1), (74, 1) |
| <i>gsav</i> (goalkeeper save) | (10, 1), (11, 1), (41, 1), (52, 1), (54, 1) |
| <i>clea</i> (clearance) | (12, 1) |
| <i>wsho</i> (wrong shot) | (13, 1), (14, 1), (15, 1) |
| <i>goal</i> | (16, 1) |
| <i>rec</i> (reception) | (100, 1) |
| <i>cond</i> (conduction) | (101, 1) |

The possible values for the attribute `place` are defensive area (da), defensive intermediary (di), middle-field (mf), offensive intermediary (oi) and offensive area (oa). Figure 1 displays how we compute these values according to the attribute `Events.x` represented by dashed lines. The attribute `ball` (ball possession) is a mapping from the at-

tribute `field_pass`. If `field_pass = t` then `ball = 1`, and if `field_pass = f` then `ball = 0`. In order to compute de attribute `direc` (move direction of a player), we consider the previous and the current place of each player. Next, we calculate the horizontal distance (`xdist`) and vertical distance (`ydist`) between these places. When `ydist = 0` and `xdist = 0`, the direction is *none* since the player did not move. If `ydist > xdist` then the direction is *lateral*. Otherwise, the direction is *backward* (for `xdist < 0`) or *forward* (for `xdist > 0`).

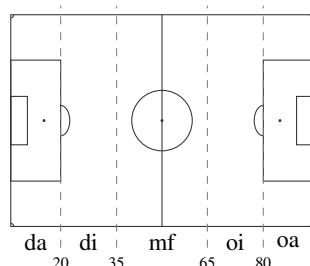


Figure 1. Soccer field division

The full dataset and the importing tool are available for download in a Github repository³. This repository also provides additional tools for benchmarking continuous queries with temporal conditional preferences. All tools were developed using the python language. The dataset is stored in the directory `data`. This directory contains the relations `Matches`, `Teams` and `Players` stored in CSV (comma separated values) format. Such relations are the union of the data extracted from all matches. The `data` directory also contains the subdirectories `raw`, `events`, `moves` and `places` to store respectively the extracted raw data, the stream `Events`, the stream `Moves` and the stream `Places`.

The subdirectory `raw` has, for every match, the JSON files `match_id.json`, `match_id-players.json` and `match_id-teams.json` where `match_id` is the match identifier. The subdirectories `events`, `moves` and `places` contain the CSV files `match_id.csv` with the stream data of the matches. Table 4 presents the information about tuples and instances for all relations and streams of the dataset. The streams `Moves` and `Places` have less tuples than stream the `Events` due to the data cleaning and the computation of the new attributes.

Table 4. Data statistics

| Relation/Stream | Tuples | Instances | Tuples/Instance |
|-----------------|---------|-----------|-----------------|
| Matches | 64 | 1 | 64 |
| Teams | 32 | 1 | 32 |
| Players | 736 | 1 | 736 |
| Events | 167,801 | 64 | 2621 |
| Moves | 130,607 | 64 | 2040 |
| Places | 137,621 | 64 | 2150 |

4. Research Opportunities

As we mentioned in the first section, our dataset is useful for a multitude of studies. This section outlines a non-exhaustive list of research fields that can benefit from use our dataset.

³<https://streampref.github.io/wcimport/>

Data Mining. In the data mining field, our dataset can be explored for the validation of new techniques aiming the detection of temporal patterns and key events. As the dataset has the player coordinates, we can also discover the special relations between this patterns and some specific field regions. Using additional information about the localization of the matches, it could be possible to find player patterns correlated to environmental variables.

Sport Analytics. Our dataset can be used to make various sport analysis over individual players or teams. Using our dataset, specialists are able to analyze the moves performed, the player positioning, pass distances and many others variables related to a specific player. In addition, the information of all team players can be combined to perform analysis over the team strategy and positioning.

Preference Continuous Queries. An interesting research topic in the field of data streams processing is to incorporate temporal conditional preferences into continuous queries [Ribeiro et al. 2017b]. This new kind of query makes use of the implicit temporal information of data streams to select sequences of elements that best fit user preferences. As the tuples of data streams has a *timestamp*, we can know the order of the tuples. So, by using temporal preferences, the user can express wishes like “if there exists a value X in the past then I prefer a value Y to a value X in the current moment”.

Unlike traditional databases, data stream applications do not store all data due to limitations of time and space. The existence of datasets for data streams scenarios are important to allow the validation of new techniques for the evaluation of continuous queries. There are works that proposed synthetic data generators [Bifet et al. 2011], but the real datasets are still useful for many specific situations.

In the work [Ribeiro et al. 2017a] we proposed a preference model for reasoning with temporal conditional preferences on data stream scenarios. The *Soccer2014DS* dataset was used in the experiments to demonstrate the effectiveness of our proposed approach. We also used the *Soccer2014DS* dataset in the work [Ribeiro et al. 2017b]. In this latter work, the dataset was used to conduct an extensive set of experiments to compare the performance and the memory usage of algorithms to process continuous queries with temporal conditional preferences.

5. Limitations

The real soccer datasets collected by Opta Sports have additional information provided by an special attribute (*qualifiers*). However, these datasets are not public. The extracted data has this attribute, but it has no meaningful values. So, we drop this attribute from our dataset.

Our dataset does not have the coordinates of all players at every second. Only events related to moves, cards, fouls and ball outs were gathered. So, positioning analysis must take this information into consideration.

We calculated just the attributes `move`, `place`, `ball` and `direc` for the derived streams `Moves` and `Places`. Although, new attributes can be computed using the available information in the stream `Events`, for example, the distance of passes and shots. In addition, we decided to keep each match into an individual stream, but it is possible to join these data into a single stream if it is a requirement of the data analysis task.

6. Conclusion

In this paper we described the creation of the public dataset *Soccer2014DS* containing player events of the 2014 Soccer World Cup. The construction of the dataset started with

the extraction of data from *Internet* using a web crawler. This extracted data was used to create the derived data by applying cleaning and conversions techniques.

The dataset and all the developed tools are available for download in a public repository. So, the dataset can be used on the development of new research works related to data mining, sports analytics and continuous queries. We already use the *Soccer2014DS* dataset on our previous works [Ribeiro et al. 2017a, Ribeiro et al. 2017b] and we are still using this dataset on new researches about continuous queries with temporal conditional preferences.

References

- Ali, A. and Farrally, M. (1991). Recording soccer players' heart rates during matches. *Journal of Sports Sciences*, 9(2):183–189.
- Arasu, A., Babcock, B., Babu, S., Cieslewicz, J., Datar, M., Ito, K., Motwani, R., Srivastava, U., and Widom, J. (2016). *STREAM: The Stanford Data Stream Management System*, pages 317–336. Springer, Berlin, Germany.
- Baca, A., Dabnichki, P., Heller, M., and Kornfeind, P. (2009). Ubiquitous computing in sports: A review and analysis. *Journal of Sports Sciences*, 27(12):1335–1346.
- Bakker, L. F. B. C. (2015). Visualizing football team strategies and player performance. Master's thesis, Eindhoven University of Technology, Eindhove, Netherlands.
- Barris, S. and Button, C. (2008). A review of vision-based motion analysis in sport. *Sports Medicine*, 38(12):1025–1043.
- Bialkowski, A., Lucey, P., Carr, P., Yue, Y., Sridharan, S., and Matthews, I. (2014). Large-scale analysis of soccer matches using spatiotemporal tracking data. In *International Conference on Data Mining (ICDM)*, Shenzhen, China.
- Bifet, A., Holmes, G., Pfahringer, B., Read, J., Kranen, P., Kremer, H., Jansen, T., and Seidl, T. (2011). MOA: A real-time analytics open source framework. In *European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML PKDD)*, pages 617–620, Athens, Greece.
- Boice, J., Fung, H., and Bycoffe, A. (2014). URL: <http://data.huffingtonpost.com/2014/world-cup> (visited on 23/04/2015).
- Gyarmati, L. and Hefeeda, M. (2015). Estimating the maximal speed of soccer players on scale. In *Machine Learning and Data Mining for Sports Analytics Workshop*, Porto, Portugal.
- Lucey, P., Oliver, D., Carr, P., Roth, J., and Matthews, I. (2013). Assessing team strategy using spatiotemporal data. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1366–1374, Chicago, IL, USA.
- Perin, C., Vuillemot, R., and Fekete, J.-D. (2013). Soccerstories: A kick-off for visual soccer analysis. *IEEE Transactions on Visualization and Computer Graphics*, 19(12):2506–2515.
- Ribeiro, M. R., Barioni, M. C. N., de Amo, S., Roncancio, C., and Labbé, C. (2017a). Reasoning with temporal preferences over data streams. In *International Florida Artificial Intelligence Research Society Conference (FLAIRS)*, Marco Island, Florida.
- Ribeiro, M. R., Barioni, M. C. N., de Amo, S., Roncancio, C., and Labbé, C. (2017b). Temporal conditional preference queries on streams. In *International Conference on Database and Expert Systems Applications (DEXA)*, Lyon, France. (to appear).

Appendix

A. Events Associated to Values of Attributes Type and Outcome

| type | outcome | Event description | type | outcome | Event description |
|------|---------|--|------|---------|--|
| 1 | 0 | Non completed pass | 44 | 0 | Lost aerial duel |
| | 1 | Completed pass | | 1 | Wined aerial duel |
| 2 | 1 | Pass to offside player | 45 | 0 | Player dribbled |
| 3 | 0 | Dribble losing ball | 49 | 1 | Ball recovery |
| | 1 | Successful dribble | | 50 | 1 |
| 4 | 0 | Foul committed | 51 | 1 | Error (causing ball dispossession) |
| | 1 | Foul suffered | | 52 | 1 |
| 5 | 0 | Ball out | 53 | 1 | Cross not claimed by goalkeeper |
| | 1 | Wined throw-in or goal kick | | 54 | 1 |
| 6 | 0 | Ball out on goal line | 55 | 1 | Offside provoked |
| | 1 | Wined corner kick | | 56 | 1 |
| 7 | 0 | Dispossessed opponent without possession | 57 | 0 | Player causes throw-in reversion |
| | 1 | Dispossessed opponent with possession | | 1 | Player wins throw-in reversion |
| 8 | 1 | Interception | 58 | 1 | Goalkeeper faced to penalty kick |
| 10 | 1 | Goalkeeper save shot | | 59 | 0 |
| 11 | 1 | Goalkeeper catches crossed ball | 60 | 1 | Goalkeeper clears ball with possession |
| 12 | 1 | Clearance (shot out defensive zone) | | 61 | 0 |
| 13 | 1 | Miss (shot out goal) | 61 | 0 | Player touches ball and without possession |
| 14 | 1 | Post (shot on goal frame) | | 1 | Player touches ball and with possession |
| 15 | 1 | Attempt saved by other player | 74 | 1 | Accidental blocking |
| 16 | 1 | Goal | | 100 | 1 |
| 17 | 1 | Card | 101 | 1 | Ball conduction |
| 18 | 1 | Player substituted | | 102 | 1 |
| 19 | 1 | Player comes on (as substitute) | | | |
| 34 | 1 | Player line up and formation | | | |
| 41 | 1 | Goalkeeper punches ball | | | |
| 42 | 1 | Skill on the ball | | | |
| 43 | 1 | Deleted event | | | |

B. Soccer2014DS Logical Schema

