

# Mining over Football Match Data: Seeking Associations among Explicit and Implicit Events

Jan Nemrava<sup>1,2</sup>, Vojtěch Svátek<sup>1</sup>, Milan Šimůnek<sup>1</sup>, Paul Buitelaar<sup>2</sup>

<sup>1</sup> Department of Information and Knowledge Engineering,  
University of Economics, Prague, W. Churchill Sq. 4, 130 67 Praha 3, Czech Republic  
e-mail: {nemrava,svatek,simunek}@vse.cz

<sup>2</sup> DFKI (German Research Center for Artificial Intelligence) GmbH,  
Stuhlsatzenhausweg 3, D-66123 Saarbrücken, Germany  
e-mail: paulb@dfki.de

**Abstract.** Football match data potentially cover a wide scope of resource modalities, such as video/audio streams, free-text news reports, structured online reports and structured data. One of the partial tasks in the *K-Space* project is to analyse multimedia resources together with resources that are complementary to them: texts and structured data. Eventually, the events (such as goals, fouls or substitutions) contained or automatically detected in different resources should be mapped to each other, so as to ease automated retrieval from ‘difficult’ sources via the ‘easier’ ones. In this paper we concentrate on the possibility of using association mining over each resource separately, however, with the aim of detecting relationships that could help an expert formulate criteria for future mapping across the resources.

## 1 Introduction

Football match data of various types, including the actual video, relevant news articles or tables with statistics, are of high interest for a large community of sports fans, and, consequently, a large quantity of such data is generated and consumed. This makes the football domain well suitable for experiments with analysing similar information in different modalities, such as video, text and structured data. This analysis can result e.g. in faster retrieval of a desired match, sequence of video or the like.

One of the partial tasks in the *K-Space* project<sup>3</sup> is to analyse multimedia resources together with resources that are complementary to them: texts and structured data. An ongoing effort described in [2] among other deals with matching the events detected in *video*, which are represented at a relatively low level (such as audio peaks, close shots etc.) but with fine-grained timing (in seconds), with events described in *textual* and *structured* (XML) data, which are more ‘semantic’ (e.g. goals, yellow cards) but with coarse-grained timing (in minutes). This paper however focuses at a side-issue of this temporal matching: detecting typical relationships among different attributes that describe the events in a single resource—video, structured data, and texts, respectively.

<sup>3</sup> <http://www.kspace-noe.org>

## 2 Datasets Available

The first resource involved is the *SmartWeb* dataset<sup>4</sup> compiled at DFKI, Germany and consisting of:

- *Structured* match protocols in XML format.
- A corpus of *textual* match reports (free-text or semi-structured documents in German and English) that are derived from freely available web sources. The bilingual documents are not translations, but are aligned on the level of a particular match (i.e. they are about the same match).
- An *ontology* on football that is integrated with foundational (DOLCE), general (SUMO) and task-specific (discourse, navigation) ontologies.
- A *knowledge base* of events and entities in the World Cup domain that have been automatically extracted from the German documents.

In the structured reports, there are on average 12 events per match; these consist in *goals*, *substitutions*, and *yellow/red cards*. They are assumed to contain official, approved data, and thus can act as ‘gold standard’. An excerpt from a structured report (listing the cards given in one match) is in Fig. 1.

```
<Cards>
  <YellowCard>
    <Player>BOROWSKI Tim</Player>
    <Team>GER</Team>
    <Minute>40</Minute>
  </YellowCard>
  <YellowCard>
    <Player>METZELDER Christoph</Player>
    <Team>GER</Team>
    <Minute>56</Minute>
  </YellowCard>
  <YellowCard>
    <Player>CAMORANESI Mauro</Player>
    <Team>ITA</Team>
    <Minute>90</Minute>
  </YellowCard>
</Cards>
```

Fig. 1. Fragment of structured protocol

The second major resource consists in the results of *football video analysis* in terms of identified candidate events with time points and confidence values. Both resources have an intersection in terms of matches covered, among other, those of World Cup 2006. There are typically between 200-300 events per match. They have been detected using Machine Learning techniques (SVM) based on the values of six low-level detectors focussing on, respectively: *crowd detection*, *speech-band audio activity*, *on-screen graphics tracking* (dealing in particular with score-board absence detection), *motion activity measure* (possibly detecting e.g. a celebrating player), *close-up shot detection*, and *field line orientation* (possibly detecting events that take place near the end of the field). The analysis, carried out at Dublin City University, has been described in [4].

<sup>4</sup> <http://smartweb.dfki.de/>

### 3 First Mining Experiments on Timeless Data

A straightforward transformation of structured XML data into relational form suitable for tabular data mining was relatively easy: we created one record per event, resulting, for example, for WC2006 qualification matches, in about 15 000 records (of which about 2 500 goals, 3 500 cards and 9 000 substitutions). We then could apply the *LISp-Miner* tool [3], which can find various types of associations in such data, possibly with complex combinations of allowed values (subsets, intervals, cuts etc.) within one element (literal) of an association.

Using *LISp-Miner*, obvious relationships have primarily been found, such as “a player who was substituted cannot obtain red card (and vice versa)” or “goalkeepers are much less frequently substituted”. There were also a few less obvious ones, which might deserve expert interpretation, such as “players who scored a goal are more often substituted”.

### 4 Towards Using LISp-Miner for Temporal Mining

All events in data, both the *explicit* ones (in structured and textual data) and the *implicit* ones (in video analysis data) are of *temporal* nature. It is thus desirable to capture the temporal aspects of these events in a way that would make them exploitable by the data mining tool.

We are currently preparing a transformation of *structured data* such that it would also include in each record (corresponding to individual event) the information about the *minute-offset* from the previous event of a certain type (e.g. “with respect to the given substitution event, the last red card event was 8 minutes ago, and it was a card to the team of the substituted player”). When mining on such data, we could take advantage of the capability of *LISp-Miner* to suggest value intervals at runtime, possibly arriving at hypotheses such as “teams more often substitute within 5-10 minutes after getting a red card”.

The same approach could be taken with *textual data*, in principle. In this way we could obtain richer information including match events that are not part of formal protocols such as rushes, corners, shots, passes or (keeper’s) saves. In our ongoing research, we managed to successfully apply *ontology-driven information extraction* on this data, resulting in event lists such as in Fig. 2 (coverage of one minute of match). However, the data obtained via textual information extraction is often noisy and there are discrepancies in temporal stamps across different resources, so further pre-processing is needed. In addition, having possibly several tens of event types instead of 4–5 as in structured data would lead to a significant increase of the number of attributes, and thus require a more sophisticated method of building analytic questions to be tested, so as to avoid high computational overhead.

The situation is even more difficult with *video analysis* data, where we have a huge number of ‘events’ but these lack clear semantics. We plan to pre-process the data via assigning to each event (i.e. data record) specific attributes for each of the aforementioned detectors indicating how much time (in seconds) passed since the detector reached a certain pre-defined value threshold.

```

city-gazetteer
player "Kehl"
player_action freekick - "freistoß""Kopfballverlängerung"
amount
playeraction shot - "schuss"
player_action balldeflection - "abfälschen""Valente"
player_action balldeflection - "abfälschen""Valente"
playeraction scoregoal - "tor"
player "Valente"
player_action shot - "schuss""Ball"
playeraction penaltykick - "elfmeter"

```

Fig. 2. Sample of data extracted from an online textual report

## 5 Possible Applications of Discovered Associations

The analysis of ‘clean’ structured data, and perhaps also of textual data, could be seen as method for discovering new relationships to be inserted into the *domain ontology* or *knowledge base*, respectively. Similar relationships have been used in the same (football) domain to correct errors and to solve inconsistencies during the analysis and merging of match reports [1], which is also our interest.

In contrast, mining over video analysis results should rather give us a better insight into these data (otherwise rather obscure for a non-expert in signal analysis) and possibly even give feedback to the developers of the underlying video analysis tools

*The authors cordially thank James Lanagan from the Dublin City University for preparing the video analysis event data, and Jan Rauch from University of Economics, Prague, for providing his expertise with the LISp-Miner tool.*

*The research leading to this paper was supported by the European Commission under contract FP6-027026, “Knowledge Space of Semantic Inference for Automatic Annotation and Retrieval of Multimedia Content” - K-Space, by the IGA VSE grant no.11/06, and by the grant no.201/05/0325 of the Czech Science Foundation, “New methods and tools for knowledge discovery in databases”.*

## References

1. Kuper, J., Saggion, H., Cunningham, H., Declerck, T., de Jong, F., Reidsma, D., Wilks, Y., Wittenburg, P.: Intelligent Multimedia Indexing and Retrieval through Multi-source Information Extraction and Merging. In: Proceedings of IJCAI, 409-414, 2003.
2. Nemrava, J., Svátek, V., Declerck, T., Buitelaar, P., Zeiner, H., Alcantara, M.: Report on algorithms for mining complementary sources. Deliverable D5.4, IST FP6-027026 K-Space.
3. Rauch, J., Šimůnek, M.: An Alternative Approach to Mining Association Rules. In: Lin, T. Y., Ohsuga, S., Liau, C. J., Tsumoto, S. (eds.), Data Mining: Foundations, Methods, and Applications, Springer-Verlag, 2005, pp. 211–232
4. Sadlier, D., O’Connor, N.: Event Detection in Field Sports Video using Audio-Visual Features and a Support Vector Machine. *IEEE Transactions on Circuits and Systems for Video Technology*, October 2005.