

Fire! Firing Inductive Rules from Economic Geography for Fire Risk Detection

David Vaz¹, Vítor Santos Costa², and Michel Ferreira³

¹ LIACC- DCC/FCUP, University of Porto, Portugal

² CRACS- DCC/FCUP, University of Porto, Portugal

³ IT- DCC/FCUP, University of Porto, Portugal

Abstract. Wildfires can importantly affect the ecology and economy of large regions of the world. Effective prevention techniques are fundamental to mitigate their consequences. The design of such preemptive methods requires a deep understanding of the factors that increase the risk of fire, particularly when we can intervene on these factors. This is the case for the maintenance of ecological balances in the landscape that minimize the occurrence of wildfires. We use an inductive logic programming approach over detailed spatial datasets: one describing the landscape mosaic and characterizing it in terms of its use; and another describing polygonal areas where wildfires took place over several years. Our inductive process operates over a logic term representation of vectorial geographic data and uses spatial predicates to explore the search space, leveraging the framework of Spatial-Yap, its multi-dimensional indexing and tabling extensions. We show that the coupling of a logic-based spatial database with an inductive logic programming engine provides an elegant and powerful approach to spatial data mining.

1 Introduction

Wildfires are an unavoidable event in Nature and play an important role in wildland ecosystems. Naturally caused wildfires are, however, a small percentage of all the wildland fires. Preventing and mitigating the consequences of wildfires that result from the increased pressure of human activity in wildland areas has been the goal of fire control programs for more than a century. Prevention techniques range from measures aiming to reduce human infractions, to the altering of stored fuels, through controlled burns, in wildlands to affect future fire risk and behavior. In addition to the straightforward impact of fuels and weather conditions in the occurrence of fires, the role of topography is also relevant. Here we understand topography in a broader sense, as a discipline concerned with local detail of space, including not only relief but also vegetation and human-made features. In areas where human intervention has importantly reshaped this topography, as happens in many European regions where native forests have been replaced by fast-growing trees and pasture areas, the impact of this human-designed organization of landscape can potentially affect the occurrence and behavior of wildfires. In this paper we focus on this landscape organization

factor, which has been given little attention by fire control programs. While some aspects of the human-made organization of landscape can obviously affect fire behavior, such as the existence of major roads cutting through forest areas, that act as barriers to the propagation of fire, there are also potentially less obvious correlations between this local organization of the landscape and the occurrence of fires, which can profit from machine learning techniques. Understanding such correlations can then guide the human intervention in the landscape towards more efficient prevention and mitigation of the consequences of wildfires.

In this paper we propose the use of inductive logic programming (ILP) to design logic theories that correlate the local organization of the landscape with the occurrence of fires, using the framework of Spatial-Yap [1] to have term-based representation of vectorial geographic data. In addition, Spatial-Yap provides a logic-based approach to a geographic information system and is able to render the inductive engine with spatial relationship predicates that are used to formulate theories *on-the-fly* based on spatial reasoning. We use detailed spatial databases that contain vectorial representations of the landscape organization in the north of Portugal. We overlay these databases with another spatial database that keeps historical records of wildfires taking place over several years in the same region.

The remainder of this paper is organized as follows. The next section presents related work on the use of ILP with spatial datasets. Section 3 describes our fire dataset and the methodology we used. In Section 4 we report preliminary results, together with a discussion of these results. Finally, Section 5 ends the paper.

2 Related Work

Over the past years, the use of spatial data has increased in many areas of computer science. Relational Database Management Systems (RDBMS) were among the first systems to tackle this kind of data, both through extensions to support spatial data, and by providing functions to manipulate the data.

The Open Geospatial Consortium (OGC) proposed a standard to extend SQL-92 in “OpenGis Simple Features Specification for SQL” (OGC99) [2]. The purpose of this specification is to define a standard SQL schema that supports storage, retrieval, query and update of simple geospatial feature collections. Examples of RDBMS systems that conform to this standard are Oracle Spatial and PostgreSQL, through the PostGIS module.

The development of such sophisticated geographical databases has led to interest in *spatial data mining*, defined to be the branch of data-mining where the spatial neighbors of an object may have an influence on the object [3]. A typical task would be to find clusters of correlated objects [4], but a large number of different applications are possible.

Arguably, spatial learning can be considered as an instance of multi-relational learning, with a very specific type of domain knowledge, and should be an important application for ILP. Malerba [3] and his group have exploited this approach with very interesting results. In their approach, multi-relational data-mining techniques are applied by working at a higher conceptual level of the geographic

information [5]. Their approach follows a two step algorithm. First, system such as INGENS [6] extract relevant concepts and features from a spatial database, by applying and expanding on standard GIS tools. Second, this relational representation of spatial data can be mined by ILP techniques: ATRE [7] implements a sequence coverage algorithm that learns a classifier, and SPADA [8] is an association-rule learner that can find strong spatial association rules.

The INGENS work raises a number of interesting questions. One important problem, discussed by Malerba [5], is the computational cost of performing feature extraction: although spatial facts are rarely updated, attribute expansion can be expensive in terms of time and space, with often time being spent computing unnecessary attributes. One would expect this problem to grow as databases grow in size and complexity.

One possible approach to this problem is to couple a database to a deductive system: MYDDAS [9] couples YapTab [10] and MySQL extended with geometry types to form Spatial-Yap [1] (unfortunately MySQL has never evolved to conform with OGC99). In this paper we take the next step and actually couple tightly Prolog inference with the geographical data itself. In order to perform inference with logic programming, we need to address well the three key components of geographical data-mining:

1. Spatial terms for representing and storing spatial objects.
2. Spatial predicates, to manipulate (e.g., intersect two spatial terms) and to query spatial terms toward finding interesting properties such as area or distance between two spatial terms.
3. Effective indexing of spatial terms, not only because of the usual mammoth size of such terms, but also because of the number of different terms in the database and the complexity of spatial predicates.

We address the first problem by simply using Prolog terms to represent geometry types in a representation similar to the Well Known Text of OGC99.

We further define a set of spatial predicates that provide an interface to the GEOS API, that conforms to the OGC99 standard and is also used by PostGIS.

This machinery provides the foundation for a logic programming geographical information system. The next step was based on the observation that spatial data does not benefit from most of the traditional indexing techniques (namely the ones used in the logic programming), as most of them are based on single dimension indexing structures. The RDBMS community addressed this problem by proposing novel data-structures, namely R-Trees which have become standard [11].

We extended Prolog indexing through *User Defined Indexing* (UDI) [12], a new extension to Prolog indexing where the programmer is able to define the indexing mechanism based on *what* the terms in the arguments of a predicate are meant to represent. This allows users to provide an indexing function that selects a subset of the clauses of a predicate, given a set of constrained variables or bound Prolog terms. Our tests showed that UDI allows to operate effectively with spatial data, pruning the dataset prior to the application of expensive spatial operators.

3 The Fire DataSet

Portugal being the smallest of the five southern Europe countries, is the most affected by fire in terms of occurrences and relative burnt area. From 1980 to 2004, 30% of the country was burnt (equivalent to 1 fire per 20ha). The closest cases (Italy and Spain) present values of fire occurrence, density and burnt area inferior by 1/3 and 1/5 respectively.

Given the increasing trend of burnt area and with the increasing changes in temperature and precipitation, it is of the outmost importance to narrow the problematic areas in order to effectively promote fire control and landscape management. Recent efforts have provided detailed information of burnt areas between 1990 and 2007.

The starting point for our work was the COS'90 database a detailed land cover map, produced by the National Center for Geographic Information, by visual interpretation of aerial photographs from 1990 followed by polygon vectorization, with 3 digit nomenclature for each polygon. The nomenclature describes the principal and secondary type of use, e.g., PE2 would express a polygon with a mixed forest based on Pine Tree and Eucalyptus covering up to 75% of the area. The order of the first letters informs that Pine Tree is dominant, the digit informs that we have between 30% to 50% of coverage.

We further remark that polygons vary widely in size. Moreover, polygons do overlap with each other, and we have cases of polygons that are contained in other polygons. We do not exploit such overlaps in this work.

Given the fine grained level of this dataset and the absence of detail in burnt areas polygons, e.g., a given burnt area polygon may represent several fires occurring with a spawn of several months within a year, we have abstracted the burnt area by tagging the COS'90 polygons with a `burnt` label for each year, making this layer our base layer. We have only used burnt information from 1991 to 1999, an acceptable spawn given the base date of COS'90.

Additional data information can be obtained by considering a second layer with socioeconomic information. In Portugal this data is available through statistics taken over parishes.

We focus on the Viana do Castelo district (county) - see Figure 1(a). This district is one of the most heavily forested in Portugal, and has suffered from a wide diversity of fires. Previous work indicates that different regions have very different patterns: Viana do Castelo is typical of the North of Portugal and is one of the most affected sub-regions of the country. The district has 290 parishes and 15091 COS'90 polygons.

3.1 Methodology

In this work we are interested in exploring spatial predicates on-the-fly during the ILP search process. We thus rely on spatial indexing to obtain more efficient execution of queries involving spatial predicates. Even with spatial indexing, geographical queries are typically very expensive. We further use *tabling* to reduce

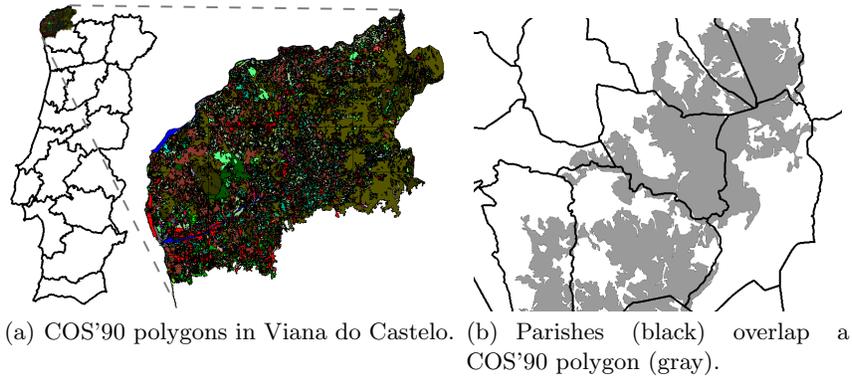


Fig. 1. Polygons and parishes in Viana do Castelo

recomputation to the minimum. As an example of this type of optimization, consider the *neighbor* relation:

```
:- table neighbor/2.
neighbor(ID1, ID2) :-
    cos90(ID1, R1, P1), R2 && R1, cos90(ID2, R2, P2), ogc_touches(P1, P2).
```

Logically, it would be sufficient to express the `cos90` polygons and then use `ogc_touches`. Assuming that R_1 is bound, the `&&` constraint selects polygons such that R_2 overlaps R_1 , R_1 and R_2 are the Minimum Bounding Rectangle (MBR) for the polygons P_1 and P_2 respectively. This is a necessary but not sufficient condition for the actual polygons to touch. The relation `ogc_touches` holds true if and only if the two polygons touch.

Notice that `&&` does not introduce any new logical information. On the other hand, `&&` is implemented very effectively by using UDI with R-Trees. In contrast, `ogc_touches` is extremely expensive: it needs to compare two complex polygons. Our approach reduces very significantly the number of calls to `ogc_touches` and makes the whole computation feasible.

However, it is clear that the `neighbor` operation will be used quite often, and may have to be recomputed every time we run a rule. We use tabling to avoid this problem. More precisely, we use tabling with *local* scheduling so that *all* solutions to the query are computed the first time the query is run.

A second interesting problem arises from the need to join the two base layers: how do parishes match COS'90? Both cover the same area, but they have different granularity and they have different information associated with it. They even overlap each other as seen in the example in Figure 1(b). In general, we do not expect to find an ideal solution to this problem: but in order to use this data we have used a weighted average based on the area of the intersection of both layers.

Spatial distance between two spatial objects corresponds to the minimum distance between any two points of each spatial object. Hence to find the mini-

minimum distance to a water class polygon, for example, we would need to calculate the distance to each water class polygon. In this case, indexing is not straightforward, but is still worthwhile. The R-Tree indexing structure abstracts spatial objects to its Minimum Bounding Rectangle, and is stored in a form that allows us to discard spatial objects far from the search rectangle. Using the indexing structure we can speedup minimum distance calculations by expanding gradually the search rectangle, starting from the MBR of the base polygon, until a matching polygon is found.

4 Results and Discussion

Our task is to predict whether a polygon will catch fire. We recall that fire is a random event: thus, we would not expect to be able to predict exactly which polygons will take fire. On the other hand, it is worthwhile to find rules that are highly indicative of vulnerability to fire. Due to space considerations, we can only present the main results here.

We use ILP system Aleph running under the Prolog system YAP-6 to search for fire risk areas. As discussed above we performed this study on the years from 1991 to 1999. In each year, positive examples (COS'90 polygons with fire occurrence) range from 180 to 1834 occurrences. We used as negative examples the remainder of 15091 polygons in the dataset. The dataset is therefore highly skewed.

We follow two different types of approaches: first, we use cross-validation over the different years; second, we try to predict the *next* year. In the latter case, we can learn with multiple years: we use up to three consecutive years. To evaluate runs on the same year we used stratified 10-fold cross validation. Results can be seen in Figure 2.

Figure 2(a) shows how our system performance in every year. Given that the dataset is very skewed, we use precision and recall as a measure of performance. Recall performance on the test set tends to range around 50%, and precision around 10%. We find these values quite acceptable, given the nature of the problem and the skew of the dataset (only about 2% of the examples are positives). Notice that the results vary significantly according to each year. In general, precision tends to be best for the years with most fires. In contrast, recall tends to be worse for these years: this is because we learn less rules in these years. The results for 1998 are quite interesting. This year about 1800 polygons burned (10% of COS'90), and the following single rule is highly predictive:

```
burnt(A,E) :-
    burnt_before(A,E).
```

Figure 2(b), 2(c), and 2(d) show next year validation performance with one year, two year, and three year training. Because years are widely different, testing the rules on the next year tends to have poorer performance than using the same year. On the other hand, as we use more years to train the system, recall and precision improve and approach same year training. Moreover, performance

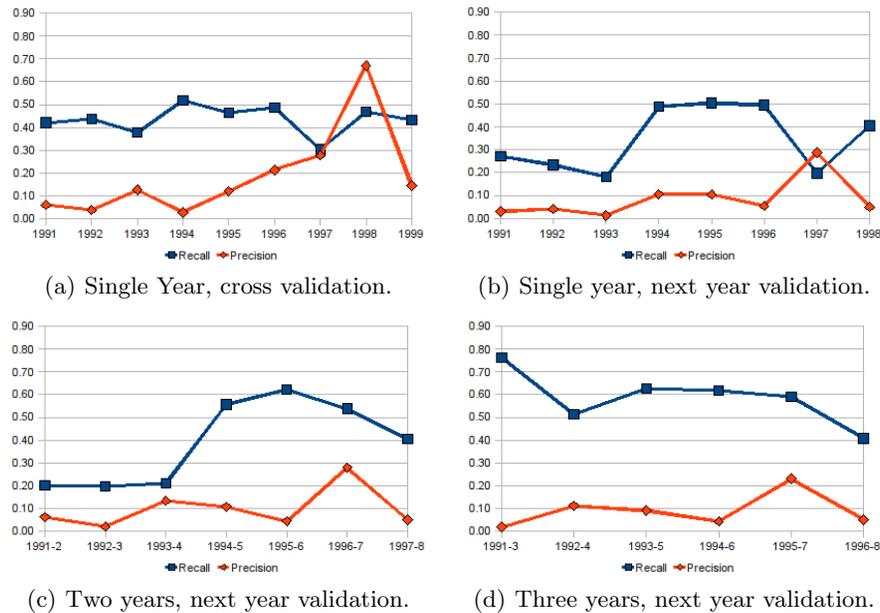


Fig. 2. Results

becomes more stable and less sensitive to variations in a year (on the other hand, we cannot take advantage of special years such as 1998). In general, with 3 year training we get a recall over 60%, with a precision of about 10%.

Two examples that give a flavor of the rules learned by our system:

```

burnt(A,_) :-
    class(A,'II'), water_pol(A,_,B), B >= 6736.85994035496.
burnt(A,E) :-
    parish(A,B), sheep(B,_,_,C), C >= 34,
    neighbor(A,D), burnt_last_year(D,E), class(D,'II').

```

The first rule refers to a polygon classified as “improductive” i.e. fallow land. The rule states that such land is quite likely to burn if more than 6Km away from a water source. The second rule applies to a polygon that is in a rural area with a high increase in sheep population, and close to fallow land that often burns. Rules most often refer to previous fire activity, to types of vegetative cover such as fallow lands, brush, pine and oak, to herding and to distance to water. Also, a high percentage of rules refer to neighboring polygons.

5 Conclusions

In this paper we have presented an ILP approach to spatial data mining, addressing the pressing problem of wildfire prevention through the understanding

of the impact of landscape organization. Our work leverages the machinery we developed in previous research, namely in the construction of an OGC-compliant logic-based geographic information system. A fundamental contribution of this work results from the coupling of an ILP system with a logic-based geographic information system. This coupling avoids the off-line materialization step of spatial features using external geographic information systems, allowing the search process to dynamically explore spatial relationship predicates in the formulation of clauses. The use of multi-dimensional indexing and tabling prove to be also crucial for the computational feasibility of our approach, providing an additional contribution for the use of ILP in the context of spatial data mining with real-world datasets.

References

1. Vaz, D., Ferreira, M., Lopes, R.: Spatial-yap: A logic-based geographic information system. In: ICLP '07: Proceedings of the 23rd International Conference on Logic Programming, Berlin, Heidelberg, Springer-Verlag (2007) 195–208
2. Open GIS Consortium, I.: OpenGIS Simple Features Specifications For SQL (1999) Available from <http://www.opengis.org/docs/99-049.pdf>.
3. Ceci, M., Appice, A., Loglisci, C., Caruso, C., Fumarola, F., Malerba, D.: Novelty detection from evolving complex data streams with time windows. In: ISMIS '09: Proceedings of the 18th International Symposium on Foundations of Intelligent Systems, Berlin, Heidelberg, Springer-Verlag (2009) 563–572
4. Ng, R.T., Han, J.: Efficient and effective clustering methods for spatial data mining. In Bocca, J.B., Jarke, M., Zaniolo, C., eds.: VLDB'94, Proceedings of 20th International Conference on Very Large Data Bases, September 12-15, 1994, Santiago de Chile, Chile, Morgan Kaufmann (1994) 144–155
5. Malerba, D., Lanza, A., Appice, A.: 10. In: Geographic Knowledge Discovery and Data Mining. 2nd Edition. CRC Press - Taylor and Francis (2009) 258–291
6. Malerba, D., Esposito, F., Lanza, A., Lisi, F.A., Appice, A.: Empowering a gis with inductive learning capabilities: the case of ingens. *Computers, Environment and Urban Systems* **27**(3) (2003) 265–281
7. Malerba, D.: Learning recursive theories in the normal ilp setting. *Fundam. Inf.* **57**(1) (2003) 39–77
8. Lisi, F.A., Malerba, D.: Inducing multi-level association rules from multiple relations. *Mach. Learn.* **55**(2) (2004) 175–210
9. Soares, T., Ferreira, M., Rocha, R.: The MYDDAS Programmer's Manual. Technical Report DCC-2005-10, Department of Computer Science, University of Porto (2005)
10. Rocha, R., Silva, F., Santos Costa, V.: YapTab: A Tabling Engine Designed to Support Parallelism. In: Conference on Tabulation in Parsing and Deduction. (2000) 77–87
11. Guttman, A.: R-trees: A dynamic index structure for spatial searching. In Yorlmark, B., ed.: SIGMOD'84, Proceedings of Annual Meeting, Boston, Massachusetts, June 18-21, 1984, ACM Press (1984) 47–57
12. Vaz, D., Santos Costa, V., Ferreira, M.: User defined indexing. In: ICLP '09: Proceedings of the 25th International Conference on Logic Programming, Berlin, Heidelberg, Springer-Verlag (2009) 372–386