

Generalization-based Similarity for Conceptual Clustering

S. Ferilli, T.M.A. Basile, N. Di Mauro, M. Biba, and F. Esposito

Dipartimento di Informatica
Università di Bari
via E. Orabona, 4 - 70125 Bari - Italia
{ferilli, basile, ndm, biba, esposito}@di.uniba.it

Abstract. Knowledge extraction represents an important issue that concerns the ability to identify valid, potentially useful and understandable patterns from large data collections. Such a task becomes more difficult if the domain of application cannot be represented by means of an attribute-value representation. Thus, a more powerful representation language, such as First-Order Logic, is necessary. Due to the complexity of handling First-Order Logic formulæ, where the presence of relations causes various portions of one description to be possibly mapped in different ways onto another description, few works presenting techniques for comparing descriptions are available in the literature for this kind of representations. Nevertheless, the ability to assess similarity between first-order descriptions has many applications, ranging from description selection to flexible matching, from instance-based learning to clustering. This paper tackles the case of Conceptual Clustering, where a new approach to similarity evaluation, based on both syntactic and semantic features, is exploited to support the task of grouping together similar items according to their relational description. After presenting a framework for Horn Clauses (including criteria, a function and composition techniques for similarity assessment), classical clustering algorithms are exploited to carry out the grouping task. Experimental results on real-world datasets prove the effectiveness of the proposal.

1 Introduction

The large amount of information available nowadays makes more difficult the task of extracting useful knowledge, i.e. valid, potentially useful and understandable patterns, from data collections. Such a task becomes more difficult if the collection requires a more powerful representation language than simple attribute-value vectors. First-order logic (FOL for short) is a powerful formalism, that is able to express relations between objects and hence can overcome the limitations shown by propositional or attribute-value representations. However, the presence of relations causes various portions of one description to be possibly mapped in different ways onto another description, which poses problems of computational effort when two descriptions have to be compared to each other.

Specifically, an important subclass of FOL refers to sets of *Horn clauses*, i.e. logical formulæ of the form $l_1 \wedge \dots \wedge l_n \Rightarrow l_0$ where the l_i 's are *atoms*, usually represented in Prolog style as $l_0 :- l_1, \dots, l_n$ to be interpreted as “ l_0 (called *head* of the clause) is true, provided that l_1 and ... and l_n (called *body* of the clause) are all true”. Without loss of generality [16], we will deal with the case of linked Datalog clauses.

The availability of techniques for the comparison between FOL (sub-)descriptions could have many applications: helping a subsumption procedure to converge quickly, guiding a generalization procedure by focussing on the components that are more similar and hence more likely to correspond to each other, implementing flexible matching, supporting instance-based classification techniques or conceptual clustering. Cluster analysis concerns the organization of a collection of unlabeled patterns into groups (clusters) of homogeneous elements based on their similarity. The similarity measure exploited to evaluate the distance between elements is responsible for the effectiveness of the clustering algorithms. Hence, the comparison techniques are generally defined in terms of a metric that must be carefully constructed if the clustering is to be relevant. In supervised clustering there is an associated output class value for each element and the efficacy of the metric exploited for the comparison of elements is evaluated according to the principle that elements belonging to the same class are clustered together as much as possible.

In the following sections, a similarity framework for first-order logic clauses will be presented. Then, Section 5 will deal with related work, and Section 6 will show how the proposed formula and criteria are able to effectively guide a clustering procedure for FOL descriptions. Lastly, Section 7 will conclude the paper and outline future work directions.

2 Similarity Formula

Intuitively, the evaluation of similarity between two items i' and i'' might be based both on the presence of common features, which should concur in a positive way to the similarity evaluation, and on the features of each item that are not owned by the other, which should concur negatively to the whole similarity value assigned to them [10]. Thus, plausible similarity parameters are:

n , the number of features owned by i' but not by i'' (*residual* of i' wrt i'');

l , the number of features owned both by i' and by i'' ;

m , the number of features owned by i'' but not by i' (*residual* of i'' wrt i').

A novel similarity function that expresses the degree of similarity between i' and i'' based on the above parameters, developed to overcome some limitations of other functions in the literature (e.g., Tverski's, Dice's and Jaccard's), is:

$$sf(i', i'') = sf(n, l, m) = 0.5 \frac{l+1}{l+n+2} + 0.5 \frac{l+1}{l+m+2} \quad (1)$$

It takes values in $]0, 1[$, to be interpreted as the degree of similarity between the two items. A complete overlapping of the two items tends to the limit of 1

as long as the number of common features grows. The full-similarity value 1 is never reached, and is reserved to the exact identification of items, i.e. $i' = i''$ (in the following, we assume $i' \neq i''$). Conversely, in case of no overlapping the function will tend to 0 as long as the number of non-shared features grows. This is consistent with the intuition that there is no limit to the number of different features owned by the two descriptions, which contribute to make them ever different. Since each of the two terms refers specifically to one of the two clauses under comparison, a weight could be introduced to give different importance to either of the two.

3 Similarity Criteria

The main contribution of this paper is in the exploitation of the formula in various combinations that can assign a similarity degree to the different clause constituents. In FOL formulæ, terms represent specific objects; unary predicates represent term properties and n -ary predicates express relationships. Hence, two levels of similarity between first-order descriptions can be defined: the *object* level, concerning similarities between terms in the descriptions, and the *structure* one, referring to how the nets of relationships in the descriptions overlap.

Example 1. Let us consider, as a running example throughout the paper, the following toy clause (a real-world one would be too complex):

$$C : h(a) :- p(a, b), p(a, c), p(d, a), r(b, f), o(b, c), q(d, e), t(f, g), \\ \pi(a), \phi(a), \sigma(a), \tau(a), \sigma(b), \tau(b), \phi(b), \tau(d), \rho(d), \pi(f), \phi(f), \sigma(f).$$

3.1 Object Similarity

Consider two clauses C' and C'' . Call $A' = \{a'_1, \dots, a'_n\}$ the set of terms in C' , and $A'' = \{a''_1, \dots, a''_m\}$ the set of terms in C'' . When comparing a pair of objects $(a', a'') \in A' \times A''$, two kinds of object features can be distinguished: the properties they own as expressed by unary predicates (*characteristic features*), and the roles they play in n -ary predicates (*relational features*). More precisely, a *role* can be seen as a couple $R = (\text{predicate}, \text{position})$ (written compactly as $R = \text{predicate}/\text{arity.position}$), since different positions actually refer to different roles played by the objects. For instance, a characteristic feature could be **male**(X), while relational features in a **parent**(X,Y) predicate are the ‘parent’ role (*parent/2.1*) the ‘child’ role (*parent/2.2*).

Two corresponding similarity values can be associated to a' and a'' : a *characteristic similarity*,

$$\text{sf}_c(a', a'') = \text{sf}(n_c, l_c, m_c)$$

based on the set P' of properties related to a' and the set P'' of properties related to a'' , for the following parameters:

$n_c = |P' \setminus P''|$ number of properties owned by a' in C' but not by a'' in C'' (*characteristic residual of a' wrt a''*);

$l_c = |P' \cap P''|$ number of common properties between a' in C' and a'' in C'' ;
 $m_c = |P'' \setminus P'|$ number of properties owned by a'' in C'' but not by a' in C'
 (*characteristic residual of a'' wrt a'*).

and a *relational similarity*,

$$\text{sf}_r(a', a'') = \text{sf}(n_r, l_r, m_r)$$

based on the *multisets* R' and R'' of roles played by a' and a'' , respectively, for the following parameters:

$n_r = |R' \setminus R''|$ how many times a' plays in C' role(s) that a'' does not play in C'' (*relational residual of a' wrt a''*);

$l_r = |R' \cap R''|$ number of times that both a' in C' and a'' in C'' play the same role(s);

$m_r = |R'' \setminus R'|$ how many times a'' plays in C'' role(s) that a' does not play in C' (*relational residual of a'' wrt a'*).

Overall, we can define the *object similarity* between two terms as

$$\text{sf}_o(a', a'') = \text{sf}_c(a', a'') + \text{sf}_r(a', a'')$$

Example 2. Referring to clause C , the set of properties of a is $\{\pi, \phi, \sigma, \tau\}$, for b it is $\{\sigma, \tau\}$ and for c it is $\{\phi\}$. The multiset of roles of a is $\{p/2.1, p/2.1, p/2.2\}$, for b it is $\{p/2.2, r/2.1, o/2.1\}$ and for c it is $\{p/2.2, o/2.2\}$.

3.2 Structural Similarity

When checking for the structural similarity of two formulæ, many objects can be involved, and hence their mutual relationships represent a constraint on how each of them in the former formula can be mapped onto another in the latter. The structure of a formula is defined by the way in which n -ary *atoms* (predicates applied to a number of terms equal to their arity) are applied to the various objects to relate them. This is the most difficult part, since relations are specific to the first-order setting and are the cause of indeterminacy in mapping (parts of) a formula into (parts of) another one. In the following, we will call *compatible* two FOL (sub-)formulæ that can be mapped onto each other without yielding inconsistent term associations (i.e., a term in one formula cannot correspond to different terms in the other formula).

Given an n -ary literal, we define its *star* as the multiset of n -ary predicates corresponding to the literals linked to it by some common term (a predicate can appear in multiple instantiations among these literals). The *star similarity* between two compatible n -ary literals l' and l'' having stars S' and S'' , respectively, can be computed for the following parameters:

$n_s = |S' \setminus S''|$ how many more relations l' has in C' than l'' has in C'' (*star residual of l' wrt l''*);

$l_s = |S' \cap S''|$ number of relations that both l' in C' and l'' in C'' have in common;
 $m_s = |S'' \setminus S'|$ how many more relations l'' has in C'' than l' has in C' (*star residual* of l'' wrt l').

by taking into account also the object similarity values for all pairs of terms included in the association θ that map l' onto l'' of their arguments in corresponding positions:

$$\text{sf}_s(l', l'') = \text{sf}(n_s, l_s, m_s) + C^s(\{\text{sf}_o(t', t'')\}_{t'/t'' \in \theta})$$

where C^s is a composition function (e.g., the average).

Then, Horn clauses can be represented as a graph in which atoms are the nodes, and edges connect two nodes *iff* they share some term, as described in the following. In particular, we will deal with *linked* clauses only (i.e. clauses whose associated graph is connected). Given a clause C , we define its *associated graph* G_C , where the edges to be represented form a Directed Acyclic Graph (DAG), *stratified* in such a way that the head is the only node at level 0 and each successive level is made up by nodes not yet reached by edges that have at least one term in common with nodes in the previous level. In particular, each node in the new level is linked by an incoming edge to each node in the previous level having among its arguments at least one term in common with it.

Example 3. In the graph G_C , the head represents the 0-level of the stratification. Then directed edges may be introduced from $h(X)$ to $p(X, Y)$, $p(X, Z)$ and $p(W, X)$, which yields level 1 of the stratification. Now the next level can be built, adding directed edges from atoms in level 1 to the atoms not yet considered that share a variable with them: $r(Y, U)$ – end of an edge starting from $p(X, Y)$ –, $o(Y, Z)$ – end of edges starting from $p(X, Y)$ and $p(X, Z)$ – and $q(W, W)$ – end of an edge starting from $p(W, X)$. The third level of the graph includes the only remaining atom, $s(U, V)$ – having an incoming edge from $r(Y, U)$.

Now, all possible paths starting from the head and reaching *leaf* nodes are univoquely determined, which reduces the amount of indeterminacy in the comparison. Given two clauses C' and C'' , we define the *intersection* between two paths $p' = \langle l'_1, \dots, l'_{n'} \rangle$ in $G_{C'}$ and $p'' = \langle l''_1, \dots, l''_{n''} \rangle$ in $G_{C''}$ as the pair of longest compatible initial subsequences of p' and p'' :

$p' \cap p'' = (p_1, p_2) = (\langle l'_1, \dots, l'_k \rangle, \langle l''_1, \dots, l''_k \rangle)$ s.t.
 $\forall i = 1, \dots, k : l'_1, \dots, l'_i$ compatible with $l''_1, \dots, l''_i \wedge$
 $(k = n' \vee k = n'' \vee l'_1, \dots, l'_{k+1}$ incompatible with $l''_1, \dots, l''_{k+1})$
 and the two residuals as the incompatible trailing parts:

$$p' \setminus p'' = \langle l'_{k+1}, \dots, l'_{n'} \rangle \quad p'' \setminus p' = \langle l''_{k+1}, \dots, l''_{n''} \rangle$$

Hence, the *path similarity* between p' and p'' , $\text{sf}_s(p', p'')$, can be computed by applying (1) to the following parameters:

$n_p = |p' \setminus p''| = n' - k$ is the length of the trail incompatible sequence of p' wrt p'' (*path residual* of p' wrt p'');

$l_p = |p_1| = |p_2| = k$ is the length of the maximum compatible initial sequence of p' and p'' ;
 $m_p = |p'' \setminus p'| = n'' - k$ is the length of the trail incompatible sequence of p'' wrt p' (*path residual* of p'' wrt p').

by taking into account also the star similarity values for all pairs of literals associated by the initial compatible sequences:

$$\text{sf}_p(p', p'') = \text{sf}(n_p, l_p, m_p) + C^p(\{\text{sf}_s(l'_i, l''_i)\}_{i=1, \dots, k})$$

where C^p is a composition function (e.g., the average).

Example 4. In C , the star of $p(a, b)$ is the multiset $\{p/2, p/2, r/2, o/2\}$, while that of $p(a, c)$ is $\{p/2, p/2, o/2\}$. The paths in C (ignoring the head that, being unique, can be univoquely matched) are $\langle p(a, b), r(b, f), t(f, g) \rangle$, $\langle p(a, b), o(b, c) \rangle$, $\langle p(a, c), o(b, c) \rangle$, $\langle p(d, a), q(d, e) \rangle$.

Note that no single criterion is by itself neatly discriminant, but their cooperation succeeds in assigning sensible similarity values to the various kinds of components, and in distributing on each kind of component a proper portion of the overall similarity, so that the difference becomes ever clearer as long as they are composed one atop the previous ones.

4 Clause Similarity

Now, similarity between two (tuples of) terms reported in the head predicates of two clauses, according to their description reported in the respective bodies, can be computed based on their generalization. In particular, one would like to exploit their *least general generalization*, i.e. the most specific model for the given pair of descriptions. Unfortunately, such a generalization is not easy to find: either classical θ -subsumption is used as a generalization model, and then one can compute Plotkin's least general generalization [13], at the expenses of some undesirable side-effects concerning the need of computing its reduced equivalent (and also of some counter-intuitive aspects of the result), or, as most ILP learners do, one requires the generalization to be a subset of the clauses to be generalized. In the latter option, that we choose for the rest of the work, the θ_{OI} generalization model [5], based on the Object Identity assumption, represents a supporting framework with solid theoretical foundations to be exploited.

Given two clauses C' and C'' , call $C = \{l_1, \dots, l_k\}$ their least general generalization, and consider the substitutions θ' and θ'' such that $\forall i = 1, \dots, k : l_i \theta' = l'_i \in C'$ and $l_i \theta'' = l''_i \in C''$, respectively. Thus, a formula for assessing the overall similarity between C' and C'' , called *formulae similitudo* and denoted fs , can be computed according to the amounts of common and different literals:

$n = |C'| - |C|$ how many literals in C' are not covered by its least general generalization with respect to C'' (*clause residual* of C' wrt C'');

$l = |C| = k$ maximal number of literals that can be put in correspondence between C' and C'' according to their least general generalization;
 $m = |C''| - |C|$ how many literals in C'' are not covered by its least general generalization with respect to C' (*clause residual* of C'' wrt C').

and of common and different objects:

$n_o = |terms(C')| - |terms(C)|$ how many terms in C' are not associated by its least general generalization to terms in C'' (*object residual* of C' wrt C'');
 $l_o = |terms(C)|$ maximal number of terms that can be put in correspondence in C' and C'' as associated by their least general generalization;
 $m_o = |terms(C'')| - |terms(C)|$ how many terms in C'' are not associated by its least general generalization to terms in C' (*object residual* of C'' wrt C').

by taking into account also the star similarity values for all pairs of literals associated by the least general generalization:

$$fs(C', C'') = sf(n, l, m) \cdot sf(n_o, l_o, m_o) + C^c(\{sf_s(l'_i, l''_i)\}_{i=1, \dots, k})$$

where C^c is a composition function (e.g., the average). This function evaluates the similarity of two clauses according to the composite similarity of a maximal subset of their literals that can be put in correspondence (which includes both structural and object similarity), smoothed by adding the overall similarity in the number of overlapping and different literals and objects between the two (whose weight in the final evaluation should not overwhelm the similarity coming from the detailed comparisons, hence the multiplication).

In particular, the similarity formula itself can be exploited for computing the generalization. The path intersections are considered by decreasing similarity, adding to the partial generalization generated thus far the common literals of each pair whenever they are compatible [6]. The proposed similarity framework proves actually able to lead towards the identification of the proper sub-parts to be put in correspondence in the two descriptions under comparison, as shown indirectly by the portion of literals in the clauses to be generalized that is preserved by the generalization. More formally, the compression factor (computed as the ratio between the length of the generalization and that of the shortest clause to be generalized) should be as high as possible. Interestingly, on the document dataset (see section 6 for details) the similarity-driven generalization preserved on average more than 90% literals of the shortest clause, with a maximum of 99,48% (193 literals out of 194, against an example of 247) and just 0,006 variance. As a consequence, one would expect that the produced generalizations are least general ones or nearly so. Noteworthy, using the similarity function on the document labelling task leads to runtime savings that range from 1/3 up to 1/2, in the order of hours.

5 Related Works

Few works faced the definition of similarity or distance measures for first-order descriptions. [4] proposes a distance measure based on probability theory applied

to the formula components. Compared to that, our function does not require the assumptions and simplifying hypotheses to ease the probability handling, and no *a-priori* knowledge of the representation language is required. It does not require the user to set weights on the predicates' importance, and is not based on the presence of 'mandatory' relations, like for the *G1* subclause in [4]. *KGB* [1] uses a similarity function, parameterized by the user, to guide generalization; our approach is more straightforward, and can be easily extended to handle negative information in the clauses. In *RIBL* [3] object similarity depends on the similarity of their attributes' values and, recursively, on the similarity of the objects related to them, which poses the problem of indeterminacy. [17] presents an approach for the induction of a distance on FOL examples, that exploits the truth values of whether each clause covers the example or not as features for a distance on the space $\{0,1\}^k$ between the examples. [12] organizes terms in an importance-related hierarchy, and proposes a distance between terms based on interpretations and a level mapping function that maps every simple expression on a natural number. [14] presents a distance function between atoms based on the difference with their lgg, and uses it to compute distances between clauses. It consists of a pair where the second component allows to differentiate cases where the first component cannot.

As pointed out, we focus on the identification and exploitation of similarity measures for first-order descriptions in the clustering task. Many research efforts on data representation, elements' similarity and grouping strategies have produced several successful clustering methods (see [9] for a survey). The classical strategies can be divided in bottom-up and top-down. In the former, each element of the dataset is considered as a cluster. Successively, the algorithm tries to group the clusters that are more similar according to the similarity measure. This step is performed until the number of clusters the user requires as a final result is reached, or the minimal similarity value among clusters is greater than a given threshold. In the latter approach, known as hierarchical clustering, at the beginning all the elements of the dataset form a unique cluster. Successively, the cluster is partitioned into clusters made up of elements that are more similar according to the similarity measure. This step is performed until the number of clusters required by the user as a final result is reached. A further classification is based on whether an element can be assigned (NotExclusive or Fuzzy Clustering) or not (Exclusive or Hard Clustering) to more than one cluster. Also the strategy exploited to partition the space is a criterion used to classify the clustering techniques: in Partitive Clustering a representative point (centroid, medoid, etc.) of the cluster in the space is chosen; Hierarchical Clustering produces a nested series of partitions by merging (Hierarchical Agglomerative) or splitting (Hierarchical Divisive) clusters, Density-based Clustering considers the density of the elements around a fixed point.

Closely related to data clustering is Conceptual Clustering, a Machine Learning paradigm for unsupervised classification which aims at generating a concept description for each generated class. In conceptual clustering both the inherent structure of the data and the description language, available to the learner, drive

cluster formation. Thus, a concept (regularity) in the data could not be learned by the system if the description language is not powerful enough to describe that particular concept (regularity). This problem arises when the elements simultaneously describe several objects whose relational structures change from one element to the other. First-Order Logic representations allow to overcome these problems. However, most of the clustering algorithms and systems work on attribute-value representation (e.g., CLUSTER/2 [11], CLASSIT [8], COBWEB [7]). Other systems such as LABYRINTH [18] can deal with structured objects exploiting a representation that is not powerful enough to express the dataset in a lot of domains. There are few systems that cluster examples represented in FOL (e.g., AUTOCLASS-like [15], KBG [1]), some of which still rely on propositional distance measures (e.g., TIC [2]).

6 Experiments on Clustering

The proposed similarity framework was tested on the conceptual clustering task, where a set of items must be grouped into homogeneous classes according to the similarity between their first-order logic description. In particular, we adopted the classical K-means clustering technique. However, since first-order logic formulae do not induce an euclidean space, it was not possible to identify/build a *centroid* prototype for the various clusters according to which the next distribution in the loop would be performed. For this reason, we based the distribution on the concept of *medoid* prototypes, where a medoid is defined as the observation that actually belongs to a cluster and that has the minimum average distance from all the other members of the cluster. As to the stop criterion, it was set as the moment in which a new iteration outputs a partition already seen in previous iterations. Note that it is different than performing the same check on the set of prototypes, since different prototypes could yield the same partition, while there cannot be several different sets of prototypes for one given partition. In particular, it can happen that the last partition is the same as the last-but-one, in which case a fixed point is reached and hence a single solution has been found and has to be evaluated. Conversely, when the last partition equals a previous partition, but not the last-but-one one, a loop is identified, and one cannot focus on a single minimum to be evaluated.

Experiments on Conceptual Clustering were run on a real-world dataset¹ containing 353 descriptions of scientific papers first page layout, belonging to 4 different classes: Elsevier journals, Springer-Verlag Lecture Notes series (SVLN), Journal of Machine Learning Research (JMLR) and Machine Learning Journal (MLJ). The complexity of such a dataset is considerable, and concerns several aspects of the dataset: the journals layout styles are quite similar, so that it is not easy to grasp the difference when trying to group them in distinct classes; moreover, the 353 documents are described with a total of 67920 literals, for an average of more than 192 literals per description (some descriptions are made up of more than 400 literals); last, the description is heavily based on a *part_of*

¹ <http://lacam.di.uniba.it:8000/systems/inthelex/index.htm#datasets>

relation that increases indeterminacy. A short example of paper description (with predicate names slightly changed for the sake of brevity) is:

```

observation(d) :- num_pages(d,1), page_1(d,p1), page_w(p1,612.0), page_h(p1,792.0), last_page(p1), frame(p1,f4),
t_text(f4), w_medium_large(f4), h_very_very_small(f4), center(f4), middle(f4), frame(p1,f2), t_text(f2), w_large(f2),
h_small(f2), center(f2), upper(f2), frame(p1,f1), t_text(f1), w_large(f1), h_large(f1), center(f1), lower(f1), frame(p1,f6),
t_text(f6), w_large(f6), h_very_small(f6), center(f6), middle(f6), frame(p1,f12), t_text(f12), w_medium(f12), h_very_very_small(f12),
left(f12), middle(f12), frame(p1,f10), t_text(f10), w_large(f10), h_small(f10), center(f10), upper(f10), frame(p1,f3),
t_text(f3), w_large(f3), h_very_small(f3), center(f3), upper(f3), frame(p1,f9), t_text(f9), w_large(f9), h_medium(f9),
center(f9), middle(f9), on_top(f4,f12), to_right(f4,f12), to_right(f6,f4), on_top(f4,f6), on_top(f10,f4), to_right(f10,f4),
on_top(f2,f4), to_right(f2,f4), to_right(f1,f4), on_top(f4,f1), on_top(f3,f4), to_right(f3,f4), to_right(f9,f4), on_top(f4,f9),
on_top(f2,f12), to_right(f2,f12), on_top(f2,f6), valign_center(f2,f6), on_top(f10,f2), valign_center(f2,f10), on_top(f2,f1),
valign_center(f2,f1), on_top(f3,f2), valign_center(f2,f3), on_top(f2,f9), valign_center(f2,f9), on_top(f12,f1), to_right(f1,f12),
on_top(f6,f1), valign_center(f1,f6), on_top(f10,f1), valign_center(f1,f10), on_top(f3,f1), valign_center(f1,f3), on_top(f9,f1),
valign_center(f1,f9), on_top(f6,f12), to_right(f6,f12), on_top(f10,f6), valign_center(f6,f10), on_top(f3,f6), valign_center(f6,f3),
on_top(f9,f6), valign_center(f6,f9), on_top(f10,f12), to_right(f10,f12), on_top(f3,f12), to_right(f3,f12), on_top(f9,f12),
to_right(f9,f12), on_top(f3,f10), valign_center(f10,f3), on_top(f10,f9), valign_center(f10,f9), on_top(f3,f9), valign_center(f3,f9).

```

Since the class of each document in the dataset is known, we performed a supervised clustering; after hiding the correct class to the clustering procedure, we provided it with the ‘anonymous’ dataset, asking for a partition of 4 clusters. Then, we compared each outcoming cluster with each class, and assigned it to the best-matching class according to precision and recall. In practice, we found that for each cluster the precision-recall values were neatly high for one class, and considerably low for all the others; moreover, each cluster had a different best-matching class, so that the association and consequent evaluation became straightforward.

The clustering procedure was run first on 40 documents randomly selected from the dataset, then on 177 documents and lastly on the whole dataset, in order to evaluate its performance behaviour when tackling increasingly large data. Results are reported in Table 1: for each dataset size it reports the number of instances in each cluster and in the corresponding class, the number of matching instances between the two and the consequent precision (Prec) and recall (Rec) values, along with the overall number of correctly split documents in the dataset. Compound statistics, shown below, report the average precision and recall for each dataset size, along with the overall accuracy, plus some information about runtime and number of description comparisons to be carried out.

The overall results show that the proposed method is highly effective since it is able to autonomously recognize the original classes with precision, recall and purity (Pur) well above 80% and, for larger datasets, always above 90%. This is very encouraging, especially in the perspective of the representation-related difficulties (the lower performance on the reduced dataset can probably be explained with the lack of sufficient information for properly discriminating the clusters, and suggests further investigation). Runtime refers almost completely to the computation of the similarity between all couples of observations: computing each similarity takes on average about 2sec, which can be a reasonable time

Table 1. Experimental results

Instances	Cluster	Class	Intersection	Prec (%)	Rec (%)	Total Overlapping
40	8	Elsevier (4)	4	50	100	35
	6	SVLN (6)	5	83,33	83,33	
	8	JMLR (8)	8	100	100	
	18	MLJ (22)	18	100	81,82	
177	30	Elsevier (22)	22	73,33	100	164
	36	SVLN (38)	35	97,22	92,11	
	48	JMLR (45)	45	93,75	100	
	63	MLJ (72)	62	98,41	86,11	
353	65	Elsevier (52)	52	80	100	326
	65	SVLN (75)	64	98,46	85,33	
	105	JMLR (95)	95	90,48	100	
	118	MLJ (131)	115	97,46	87,79	
Instances	Runtime	Comparisons	Avg Runtime (sec)	Prec (%)	Rec (%)	Pur (%)
40	25'24"	780	1,95	83,33	91,33	87,5
177	9h 34' 45"	15576	2,21	90,68	94,56	92,66
353	39h 12' 07"	62128	2,27	91,60	93,28	92,35

considering the descriptions complexity and the fact that the prototype has no optimization in this preliminary version. Also the semantic perspective is quite satisfactory: an insight of the clustering outcomes shows that errors are made on very ambiguous documents (the four classes have a very similar layout style), while the induced cluster descriptions highlight interesting and characterizing layout clues. Preliminary comparisons on the 177 dataset with other classical measures report an improvement with respect to both Jaccard's, Tverski's and Dice's measures up to +5,48% for precision, up to + 8,05% for recall and up to + 2,83% for purity.

7 Conclusions

Knowledge extraction concerns the ability to identify valid, potentially useful and understandable patterns from large data collections. Such a task becomes more difficult if the domain of application requires a First-Order Logic representation language, due to the problem of indeterminacy in mapping portions of descriptions onto each other. Nevertheless, the ability to assess similarity between first-order descriptions has many applications, ranging from description selection to flexible matching, from instance-based learning to clustering.

This paper deals with Conceptual Clustering, and proposes a framework for Horn Clauses similarity assessment. Experimental results on real-world datasets prove that, endowing classical clustering algorithms with this framework, considerable effectiveness can be reached. Future work will concern fine-tuning of the similarity computation methodology, and a more extensive experimentation.

References

- [1] G. Bisson. Conceptual clustering in a first order logic representation. In *ECAI '92: Proceedings of the 10th European conference on Artificial intelligence*, pages 458–462. John Wiley & Sons, Inc., 1992.
- [2] H. Blockeel, L. De Raedt, and J. Ramon. Top-down induction of clustering trees. In J. Shavlik, editor, *Proceedings of the 15th International Conference on Machine Learning*, pages 55–63. Morgan Kaufmann, 1998.
- [3] W. Emde and D. Wettschereck. Relational instance based learning. In L. Saitta, editor, *Proc. of ICML-96*, pages 122–130, 1996.
- [4] F. Esposito, D. Malerba, and G. Semeraro. Classification in noisy environments using a distance measure between structural symbolic descriptions. *IEEE Transactions on PAMI*, 14(3):390–402, 1992.
- [5] Floriana Esposito, Nicola Fanizzi, Stefano Ferilli, and Giovanni Semeraro. A generalization model based on oi-implication for ideal theory refinement. *Fundam. Inform.*, 47(1-2):15–33, 2001.
- [6] S. Ferilli, T.M.A. Basile, N. Di Mauro, M. Biba, and F. Esposito. Similarity-guided clause generalization. In *Proc. of AI*IA-2007*, LNAI, page 12. Springer, 2007 (To appear).
- [7] D. H. Fisher. Knowledge acquisition via incremental conceptual clustering. *Machine Learning*, 2(2):139–172, 1987.
- [8] J. H. Gennari, P. Langley, and D. Fisher. Models of incremental concept formation. *Artificial Intelligence*, 40(1-3):11–61, 1989.
- [9] A. K. Jain, M. N. Murty, and P. J. Flynn. Data clustering: a review. *ACM Computing Surveys*, 31(3):264–323, 1999.
- [10] Dekang Lin. An information-theoretic definition of similarity. In *Proc. 15th International Conf. on Machine Learning*, pages 296–304. Morgan Kaufmann, San Francisco, CA, 1998.
- [11] R. S. Michalski and R. E. Stepp. Learning from observation: Conceptual clustering. In R. S. Michalski, J. G. Carbonell, and T. M. Mitchell, editors, *Machine Learning: An Artificial Intelligence Approach*, pages 331–363. Springer: Berlin, 1984.
- [12] S. Nienhuys-Cheng. Distances and limits on herbrand interpretations. In D. Page, editor, *Proc. of ILP-98*, volume 1446 of *LNAI*, pages 250–260. Springer, 1998.
- [13] G. D. Plotkin. A note on inductive generalization. *Machine Intelligence*, 5:153–163, 1970.
- [14] J. Ramon. *Clustering and instance based learning in first order logic*. PhD thesis, Dept. of Computer Science, K.U.Leuven, Belgium, 2002.
- [15] J. Ramon and L. Dehaspe. Upgrading bayesian clustering to first order logic. In *Proceedings of the 9th Belgian-Dutch Conference on Machine Learning*, pages 77–84. Department of Computer Science, K.U.Leuven, 1999.
- [16] C. Rouveirol. Extensions of inversion of resolution applied to theory completion. In *Inductive Logic Programming*, pages 64–90. Academic Press, 1992.
- [17] M. Sebag. Distance induction in first order logic. In N. Lavrač and S. Džeroski, editors, *Proc. of ILP-97*, volume 1297 of *LNAI*, pages 264–272. Springer, 1997.
- [18] K. Thompson and P. Langley. Incremental concept formation with composite objects. In *Proceedings of the sixth international workshop on Machine learning*, pages 371–374. Morgan Kaufmann Publishers Inc., 1989.