

Induction of Optimal Semantic Semi-distances for Clausal Knowledge Bases

Claudia d'Amato, Nicola Fanizzi, Floriana Esposito

LACAM – Dipartimento di Informatica
Università degli Studi di Bari
Campus Universitario, Via Orabona 4 – 70125 Bari, Italy
{claudia.damato|fanizzi|esposito}@di.uniba.it

Abstract. Several activities related to semantically annotated resources can be enabled by a notion of similarity, spanning from clustering to retrieval, matchmaking and other forms of inductive reasoning. We propose the definition of a family of semi-distances over the set of objects in a knowledge base which can be used in these activities. In the line of works on distance-induction on clausal spaces, the family is parameterized on a committee of concepts expressed with clauses. Hence, we also present a method based on the idea of simulated annealing to be used to optimize the choice of the best concept committee.

1 Introduction

Assessing semantic similarity between objects can support a wide variety of instance-based tasks spanning from *case-based reasoning* and *retrieval* to *inductive generalization* and *clustering*.

As pointed out in related surveys [11], initially, most of the proposed similarity measures for concept descriptions focus on the similarity of atomic concepts within simple concept hierarchies or are strongly based on the structure of the terms for specific FOL fragments [4]. Alternative approaches are based on related notions of *feature* similarity or *information content*. All these approaches have been specifically aimed at assessing similarity between concepts (see also [7]). In the perspective of exploiting similarity measures in inductive (instance-based) tasks like those mentioned above, the need for a definition of a semantic similarity measure for *instances* arises [1, 2, 10].

Recently, semantic dissimilarity measures for specific FOL fragments have been proposed which turned out to be practically effective for the targeted inductive tasks. Although these measures ultimately rely on the semantics of primitive concepts as elicited from the knowledge base, still they are partly based on structural criteria (a notion of normal form) which determine also their main weakness: they are hardly portable to deal with other FOL fragments.

Therefore, we have devised a new family of dissimilarity measures for semantically annotated resources, which can overcome the aforementioned limitations. Our measures are mainly based on Minkowski's measures for Euclidean spaces

defined by means of the *hypothesis-driven* distance induction method [12]. Another source of inspiration was provided by the *indiscernibility* relationships [3] investigated *rough sets* theory [9].

Namely, the proposed measures are based on the degree of discernibility of the input objects with respect to a committee of features, which are represented by concept descriptions. As such, these new measures are not absolute, since they depend on both the choice (and cardinality) of the features committee and the knowledge base they are applied to. Rather, they rely on statistics on objects that are likely to be maintained by the knowledge base management system, which can determine a potential speed-up in the measure computation during knowledge-intensive tasks. Differently from the original idea [12], we give a definition of the notion of projections which is based on model-theory in LP.

Furthermore, we also propose ways to extend the presented measures to the case of assessing concept similarity by considering concepts as represented by their extension, i.e. the set of their instances. Specifically, we recur to notions borrowed from clustering [5] such as the *medoid*, the most centrally located instance in a concept extension w.r.t. a given metric.

Experimentally¹, it may be shown that the measures induced by large committees (e.g. including all primitive and defined concepts) can be sufficiently accurate when employed for classification tasks even though the employed committee of features were not the optimal one or if the concepts therein were partially redundant. Nevertheless, this has led us to investigate on a method to optimize the committee of features that serve as dimensions for the computation of the measure. To this purpose, the employment of genetic programming and randomized search procedures was considered. Finally we opted for an optimization search procedure based on *simulated annealing* [6], a randomized approach that can overcome the problem of the search being caught in local minima.

The remainder of the paper is organized as follows. The definition of the family of measures is proposed in Sect. 2, where we prove them to be semi-distances and extend their applicability to the case of concept similarity. In Sect. 3, we illustrate and discuss the method for optimizing the choice of concepts for the committee of features which induces the measures. Possible developments are finally examined in Sect. 4.

2 A Family of Semi-distances for Instances

In the following, we assume that objects (instances), concepts and relationships among them may be defined in terms of a function-free (yet not constant-free) clausal language such as DATALOG, endowed with the standard semantics (see [8] for reference).

We will consider a *knowledge base* $\mathcal{K} = \langle \mathcal{P}, \mathcal{D} \rangle$, where \mathcal{P} is a logic program representing the *schema*, with concepts (entities) and relationships defined

¹ Such experiments, regarding a nearest neighbor search task, are not further commented here for the sake of brevity.

through definite clauses, and the *database* \mathcal{D} is a set of ground facts concerning the world state. In this context, without loss of generality, we will consider concepts as described by unary atoms. *Primitive* concepts are defined in \mathcal{D} extensionally by means of ground facts only, whereas *defined* concepts will be defined in \mathcal{P} by means of clauses. The set of the objects occurring in \mathcal{K} is denoted with $\text{const}(\mathcal{D})$.

As regards the necessary inference services, our measures will require performing *instance-checking*, which amounts to determining whether an object belongs (is an instance) of a concept in a certain interpretation.

2.1 Basic Measure Definition

It can be observed that instances lack a syntactic structure that may be exploited for a comparison. However, on a semantic level, similar objects should *behave* similarly with respect to the same concepts, i.e. similar assertions (facts) should be shared. Conversely, dissimilar instances should likely instantiate disjoint concepts.

Therefore, we introduce novel dissimilarity measures for objects, whose rationale is the comparison of their semantics w.r.t. a fixed number of dimensions represented by concept descriptions (predicate definitions). Namely, instances are compared on the grounds of their behavior w.r.t. a reduced (yet not necessarily disjoint) committee of features, represented by a collection of concept descriptions, say $F = \{F_1, F_2, \dots, F_m\}$, which stands as a group of discriminating *features* expressed in the language taken into account. In this case, we will consider unary predicates which have a definition in the knowledge base.

Following [12], a family of totally semantic distance measures for objects can be defined for clausal representations. In its simplest formulation, inspired by Minkowski's metrics, these functions can be defined as follows:

Definition 2.1 (family of measures). *Let \mathcal{K} be a knowledge base. Given a set of concept descriptions $F = \{F_1, F_2, \dots, F_m\}$, a family $\{d_p^F\}_{p \in \mathbb{N}}$ of functions $d_p^F : \text{const}(\mathcal{D}) \times \text{const}(\mathcal{D}) \mapsto [0, 1]$ is defined as follows:*

$$\forall a, b \in \text{const}(\mathcal{D}) \quad d_p^F(a, b) := \frac{1}{m} \left[\sum_{i=1}^m |\pi_i(a) - \pi_i(b)|^p \right]^{1/p}$$

where $\forall i \in \{1, \dots, m\}$ the i -th projection function π_i is defined by:

$$\forall a \in \text{const}(\mathcal{D}) \quad \pi_i(a) = \begin{cases} 1 & \mathcal{K} \vdash F_i(a) \\ 0 & \text{otherwise} \end{cases}$$

The superscript F will be omitted when the set of features is fixed.

2.2 Discussion

We can prove that these functions have the standard properties for semi-distances:

Proposition 2.1 (semi-distance). For a fixed feature set and $p \in \mathbb{N}$, function d_p is a semi-distance.

Proof. In order to prove the thesis, given any three objects $a, b, c \in \text{const}(\mathcal{D})$ it must hold that:

1. $d_p(a, b) \geq 0$ (positivity)
2. $d_p(a, b) = d_p(b, a)$ (symmetry)
3. $d_p(a, c) \leq d_p(a, b) + d_p(b, c)$ (triangular inequality)

Now, we observe that:

1. trivial, by definition
2. trivial, for the commutativity of the operators involved
3. it follows from the properties of the power function:

$$\begin{aligned}
d_p(a, c) &= \frac{1}{m} \left[\sum_{i=1}^m |\pi_i(a) - \pi_i(c)|^p \right]^{1/p} \\
&= \frac{1}{m} \left[\sum_{i=1}^m |\pi_i(a) - \pi_i(b) + \pi_i(b) - \pi_i(c)|^p \right]^{1/p} \\
&\leq \frac{1}{m} \left[\sum_{i=1}^m |\pi_i(a) - \pi_i(b)|^p + \sum_{i=1}^m |\pi_i(b) - \pi_i(c)|^p \right]^{1/p} \\
&= \frac{1}{m} \left[\sum_{i=1}^m |\pi_i(a) - \pi_i(b)|^p + \sum_{i=1}^m |\pi_i(b) - \pi_i(c)|^p \right]^{1/p} \\
&\leq \frac{1}{m} \left[\sum_{i=1}^m |\pi_i(a) - \pi_i(b)|^p \right]^{1/p} + \frac{1}{m} \left[\sum_{i=1}^m |\pi_i(b) - \pi_i(c)|^p \right]^{1/p} \\
&= d_p(a, b) + d_p(b, c)
\end{aligned}$$

As such, these are only a semi-distances. Namely, it cannot be proved that $d_p(a, b) = 0$ iff $a = b$. This is the case of *indiscernible* instances with respect to the given set of hypotheses \mathbf{F} [3].

Here, we make the assumption that the feature-set \mathbf{F} may represent a sufficient number of (possibly redundant) features that are able to discriminate really different objects. As hinted in [12], redundancy may help appreciate the relative differences in similarity.

Compared to other proposed distance (or dissimilarity) measures, the presented functions are not based on structural (syntactical) criteria; namely, they require only deciding whether an object can be an instance of the concepts in the committee.

Note that the computation of projection functions can be performed in advance (with the support of suitable DBMSs) thus determining a speed-up in the actual computation of the distance measure. This is very important for the integration of these measures in instance-based methods which massively use distances, such as in case-based reasoning and clustering.

2.3 Extensions

The definition above might be further refined and extended by recurring to model theory. Namely, the set of Herbrand models of the knowledge base $\mathcal{M}_{\mathcal{K}} \subseteq 2^{|\mathcal{B}_{\mathcal{K}}|}$ may be considered, where $\mathcal{B}_{\mathcal{K}}$ stands for its Herbrand base.

Now, given two instances a and b to be compared w.r.t. a certain feature F_i , $i = 1, \dots, m$, we might check whether they can be distinguished in the world represented by a Herbrand interpretation $\mathcal{I} \in \mathcal{M}_{\mathcal{K}}$: $\mathcal{I} \models F_i(a)$ and $\mathcal{I} \models F_i(b)$. Hence, a distance measure should count the cases of disagreement, varying the Herbrand models of the knowledge base: The resulting definition for a dissimilarity measure is the following:

$$\forall a, b \in \text{const}(\mathcal{D}) \quad d_p^{\mathcal{F}}(a, b) := \frac{1}{m \cdot |\mathcal{M}_{\mathcal{K}}|} \left[\sum_{\mathcal{I} \in \mathcal{M}_{\mathcal{K}}} \sum_{i=1}^m |\pi_i^{\mathcal{I}}(a) - \pi_i^{\mathcal{I}}(b)|^p \right]^{1/p}$$

where the projections are computed for a specific world state as encoded by a Herbrand interpretation \mathcal{I} :

$$\forall a \in \text{const}(\mathcal{D}) \quad \pi_i^{\mathcal{I}}(a) = \begin{cases} 1 & F_i(a) \in \mathcal{I} \\ 0 & \text{otherwise} \end{cases}$$

Following the rationale of the average link criterion used in clustering [5], the measures can be extended to the case of concepts, by recurring to the notion of medoids. The *medoid* of a group of objects is the object that has the highest similarity w.r.t. the others. Formally, given a group $G = \{a_1, a_2, \dots, a_n\}$, the medoid is defined:

$$m = \text{medoid}(G) = \underset{a \in G}{\text{argmin}} \sum_{j=1}^n d_p^{\mathcal{F}}(a, a_j)$$

Now, given two concepts C_1, C_2 , we can consider the two corresponding groups of objects obtained by retrieval $R_i = \{a \in \text{const}(\mathcal{D}) \mid \mathcal{K} \models C_i(a)\}$, and their resp. medoids $m_i = \text{medoid}(R_i)$ for $i = 1, 2$ w.r.t. a given measure $d_p^{\mathcal{F}}$ (for some $p > 0$ and committee \mathcal{F}). Then we can define the function for concepts as follows:

$$d_p^{\mathcal{F}}(C_1, C_2) := d_p^{\mathcal{F}}(m_1, m_2)$$

Alternatively, a metric can be defined based on the single-link and complete-link principles [5]:

$$d_p^{\mathcal{F}}(C_1, C_2) = \frac{\min\{d_p^{\mathcal{F}}(a, b) \mid \mathcal{K} \models C_1(a) \wedge C_2(b)\}}{\max\{d_p^{\mathcal{F}}(a, b) \mid \mathcal{K} \models C_1(a) \wedge C_2(b)\}}$$

3 Optimization

Although the measures could be implemented according to the definitions, their effectiveness and also the efficiency of their computation strongly depends on the

choice of the feature committee (*feature selection*). Indeed, various optimizations of the measures can be foreseen as concerns their parametric definition.

Among the possible committees, those that are able to better discriminate the objects in the ABox ought to be preferred:

Definition 3.1 (good feature set). *Let $F = \{F_1, F_2, \dots, F_m\}$ be a set of concept descriptions. We call F a good feature set for the knowledge base $\mathcal{K} = \langle \mathcal{T}, \mathcal{A} \rangle$ iff $\forall a, b \in \text{const}(\mathcal{D}) \exists i \in \{1, \dots, m\} : \pi_i(a) \neq \pi_i(b)$.*

Note that, when the function defined in the previous section adopts a good feature set, it has the properties of a metric on the related instance-space.

Since the function strongly depends on the choice of concepts included in the committee of features F , two immediate heuristics can be derived:

1. controlling the number of concepts of the committee (which has an impact also on efficiency), including especially those that are endowed with a real discriminating power;
2. finding optimal sets of discriminating features of a given cardinality, by allowing also their composition employing the specific refinement operators.

Both these heuristics can be enforced by means of suitable ILP techniques especially when knowledge bases with large sets of instances are available. Namely, part of the entire data can be drawn in order to induce optimal F sets, in advance with respect to the application of the measure for other specific purposes as those mentioned above. The adoption of genetic programming has been considered for constructing optimal sets of features. Yet these algorithms are known to suffer from being possibly caught in local minima. An alternative may consist in employing a different probabilistic search procedure which aims at a global optimization. Thus a method based on simulated annealing [6] has been devised, whose algorithm is reported in Fig. 1.

Essentially the algorithm searches the space of all possible feature committees starting from an initial guess (determined by $\text{MAKEINITIALFS}(\mathcal{K})$) based on the concepts (both primitive and defined) currently referenced in the knowledge base. The loop controlling the search is repeated for a number of times that depends on the temperature which gradually decays to 0, when the current committee can be returned. The current feature set is iteratively refined calling a suitable procedure $\text{RANDOMSUCCESSOR}()$. Then the fitness of the new feature set is compared to that of the current one determining the increment of energy ΔE . If this is positive then the candidate committee replaces the current one. Otherwise it will be replaced with a probability that depends on ΔE .

As regards the heuristic $\text{FITNESSVALUE}(F)$, it can be computed as the average *discernibility factor* [9, 3] of the objects w.r.t. the feature set. For example, given a set of objects $IS = \{a_1, \dots, a_n\} \subseteq \text{const}(\mathcal{D})$ the fitness function may be defined:

$$\text{FITNESSVALUE}(F) = k \cdot \sum_{1 \leq i < j \leq n} \sum_{h=1}^m | \pi_h(a_i) - \pi_h(a_j) |$$

```

FeatureSet OPTIMIZEFEATURESET( $\mathcal{K}$ ,  $\Delta T$ )
input  $\mathcal{K}$ : Knowledge base
         $\Delta T$ : function controlling the decrease of temperature
output FeatureSet
local  currentFS: current Feature Set
        nextFS: next Feature Set
        Temperature: controlling the probability of downward steps
begin
currentFS  $\leftarrow$  MAKEINITIALFS( $\mathcal{K}$ )
for  $t \leftarrow 1$  to  $\infty$  do
    Temperature  $\leftarrow$  Temperature  $- \Delta T(t)$ 
    if (Temperature = 0)
        return currentFS
    nextFS  $\leftarrow$  RANDOMSUCCESSOR(currentFS, $\mathcal{K}$ )
     $\Delta E \leftarrow$  FITNESSVALUE(nextFS)  $-$  FITNESSVALUE(currentFS)
    if ( $\Delta E > 0$ )
        currentFS  $\leftarrow$  nextFS
    else // replace FS with given probability
        REPLACE(currentFS, nextFS,  $e^{\Delta E}$ )
end

```

Fig. 1. Feature Set optimization based on a Simulated Annealing procedure.

where k is a normalization factor which may be set to: $(1/m)(n \cdot (n - 1)/4 - n)$, depending on the number of couples of different instances that really determine the fitness measure.

As concerns finding candidates to replace the current committee (RANDOMSUCCESSOR()), the function was implemented by recurring to simple transformations of a feature set:

- adding (resp. removing) a concept C : $\text{nextFS} \leftarrow \text{currentFS} \cup \{C\}$
(resp. $\text{nextFS} \leftarrow \text{currentFS} \setminus \{C\}$)
- randomly choosing one of the current concepts from currentFS , say C , and replacing it with one of its refinements $C' \in \text{REF}(C)$

Refining concept descriptions is language-dependent. For the adopted clausal logic, various refinement operators have been proposed in the literature [8]. *Complete* operators are to be preferred to ensure exploring the whole search-space.

4 Conclusions and Ongoing Work

In the line of past works on distance-induction, we have proposed the definition of a family of semi-distances over the instances in a clausal knowledge base. The measures are parameterized on a committee of concepts. Therefore, we have also presented a randomized search method to find optimal committees.

Possible subsumption relationships between clauses in the committee may be explicitly exploited in the measure for making the relative distances more accurate. The extension to the case of concept distance may also be improved.

The measures may have a wide range of application in distance-based methods to knowledge bases. They have been integrated in an instance-based learning system implementing a nearest-neighbor learning algorithm: an experimentation on performing semantic-based retrieval proved the effectiveness of the new measures.

The next step will concern exploiting the measures in a conceptual clustering algorithm where clusters will be formed by grouping instances on the grounds of their similarity assessed through the measure, triggering the induction of new emerging concepts.

References

- [1] G. Bisson. Learning in fol with a similarity measure. In W. Swartout, editor, *Proceedings of the 10th National Conference on Artificial Intelligence*, pages 82–87. MIT Press, 1992.
- [2] W. Emde and D. Wettschereck. Relational instance-based learning. In L. Saitta, editor, *Proceedings of the 13th International Conference on Machine Learning, ICML96*, pages 122–130. Morgan Kaufmann, 1996.
- [3] S. Hirano and S. Tsumoto. A knowledge-oriented clustering technique based on rough sets. In *Proceedings of the 25th Annual International Computer Software and Applications Conference, COMPSAC01*, pages 632–635. IEEE Computer Society, 2001.
- [4] A. Hutchinson. Metrics on terms and clauses. In M. van Someren and G. Widmer, editors, *Proceedings of the 9th European Conference on Machine Learning, ECML97*, volume 1224 of *LNAI*, pages 138–145. Springer, 1997.
- [5] A.K. Jain, M.N. Murty, and P.J. Flynn. Data clustering: A review. *ACM Computing Surveys*, 31(3):264–323, 1999.
- [6] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi. Optimization by simulated annealing. *Science*, 220(4598):671–680, 1983.
- [7] S.-H. Nienhuys-Cheng. Distances and limits on herbrand interpretations. In D. Page, editor, *Proceedings of the 8th International Workshop on Inductive Logic Programming, ILP98*, volume 1446 of *LNAI*, pages 250–260. Springer, 1998.
- [8] S.-H. Nienhuys-Cheng and R. de Wolf. *Foundations of Inductive Logic Programming*, volume 1228 of *LNAI*. Springer, 1997.
- [9] Z. Pawlak. *Rough Sets: Theoretical Aspects of Reasoning About Data*. Kluwer Academic Publishers, 1991.
- [10] I. Ramon and M. Bruynooghe. A framework for defining distances between first-order logic objects. Technical Report CW 263, Department of Computer Science, Katholieke Universiteit Leuven, 1998.
- [11] A. Rodriguez. *Assessing semantic similarity between spatial entity classes*. PhD thesis, University of Maine, 1997.
- [12] M. Sebag. Distance induction in first order logic. In S. Džeroski and N. Lavrač, editors, *Proceedings of the 7th International Workshop on Inductive Logic Programming, ILP97*, volume 1297 of *LNAI*, pages 264–272. Springer, 1997.