

# Approximate Measures of Semantic Dissimilarity under Uncertainty

Nicola Fanizzi, Claudia d'Amato and Floriana Esposito

*Lacam Lab – Department of Computer Science  
University of Bari – Italy*

**Uncertainty Reasoning for the Semantic Web**

Workshop  $\diamond$  ISWC, 12 November 2007

# Contents

- 1 Introduction & Motivation
- 2 Semantic Distance Measures
  - Overall Idea
  - The Projection Function
  - Measure Definition
  - Measure Optimization: Feature Selection
- 3 Semantic Dissimilarity under Uncertainty
  - Overall Idea
  - Probability Masses: Computation
  - The Discernibility Function
  - Measure Definition
- 4 Concept Dissimilarity under Uncertainty
- 5 Conclusions & Future Works

## Introduction & Motivations

- In the context of Reasoning in the SW, it is growing the interest in alternative *inductive procedures* (i.e. case-based reasoning, retrieval, conceptual clustering, ontology matching...)
  - Many of them are based on the notion of *similarity*
- Most of the measures able to assess similarity in DL representation focus on *similarity between atomic concepts*
  - Inductive learning methods often need for a notion of **similarity among individuals**
- **A new family of dissimilarity measures for semantically annotated resources has been devised**

# Knowledge Base Representation

**Assumption:** resources, concepts and relationships are defined in terms of a representation that can be mapped to some DL language (with the standard model-theoretic semantics)

$$\mathcal{K} = \langle \mathcal{T}, \mathcal{A} \rangle$$

- *T-box*  $\mathcal{T}$  is a set of definitions  $C \equiv D$
- *A-box*  $\mathcal{A}$  contains extensional assertions on concepts and roles e.g.  $C(a)$  and  $R(a, b)$
- The set of the individuals (resources) occurring in  $\mathcal{A}$  will be denoted  $\text{Ind}(\mathcal{A})$

*Instance checking* and *retrieval* inference services will be used

## Semantic Distance Measure: Main Idea

- **IDEA:** *on a semantic level, similar individuals should behave similarly w.r.t. the same concepts*
- Following HDD [**Sebag 1997**]: individuals can be compared on the grounds of their behavior w.r.t. a given set of hypotheses  $F = \{F_1, F_2, \dots, F_m\}$ , that is a collection of (primitive or defined) concept descriptions
  - $F$  stands as a group of *discriminating features* expressed in the considered language
- As such, the new measure *totally depends on semantic* aspects of the individuals in the KB

# The Projection Function

## Projection Function

Given a concept  $F_i \in F$ , the related *projection function*  $\pi_i : \text{Ind}(\mathcal{A}) \mapsto \{0, 1/2, 1\}$  is defined,  $\forall a \in \text{Ind}(\mathcal{A})$

$$\pi_i(a) := \begin{cases} 1 & \mathcal{K} \models F_i(a) \\ 0 & \mathcal{K} \models \neg F_i(a) \\ 1/2 & \textit{otherwise} \end{cases}$$

- Case:  $\pi_i(a) = 1/2 \Rightarrow$  the reasoner cannot give the truth value for a certain membership query
  - This is due to the *OWA* normally made in this context
- Hence, as in the classic probabilistic models, *uncertainty* is coped with by considering a *uniform distribution* over the possible cases.

## Semantic Distance Measure: Definition

[Fanizzi et al. @ DL 2007] Let  $\mathcal{K} = \langle \mathcal{T}, \mathcal{A} \rangle$  be a KB. Given sets of concept descriptions  $F = \{F_1, F_2, \dots, F_k\}$ , a *family of semi-distance functions*  $d_p^F : \text{Ind}(\mathcal{A}) \times \text{Ind}(\mathcal{A}) \mapsto \mathbb{R}$ , inspired to Minkowski's distance, is defined as follows:

$$\forall a, b \in \text{Ind}(\mathcal{A}) \quad d_p^F(a, b) := \frac{1}{k} \sqrt[p]{\sum_{i=1}^k \delta_i(a, b)^p}$$

where  $p > 0$  and  $\delta_i : \text{Ind}(\mathcal{A}) \times \text{Ind}(\mathcal{A}) \mapsto [0, 1]$  is the **discernibility function**:  $\forall (a, b) \in \text{Ind}(\mathcal{A}) \times \text{Ind}(\mathcal{A})$

$$\delta_i(a, b) = |\pi_i(a) - \pi_i(b)|$$

that compare two individuals  $(a, b)$  w.r.t. a feature concept  $F_i \in F$

## Distance Measure: Example

$$\mathcal{T} = \{ \text{Female} \equiv \neg\text{Male}, \text{Parent} \equiv \forall\text{child.}\text{Being} \sqcap \exists\text{child.}\text{Being}, \\ \text{Father} \equiv \text{Male} \sqcap \text{Parent}, \\ \text{FatherWithoutSons} \equiv \text{Father} \sqcap \forall\text{child.}\text{Female} \}$$

$$\mathcal{A} = \{ \text{Being}(\text{ZEUS}), \text{Being}(\text{APOLLO}), \text{Being}(\text{HERCULES}), \text{Being}(\text{HERA}), \\ \text{Male}(\text{ZEUS}), \text{Male}(\text{APOLLO}), \text{Male}(\text{HERCULES}), \\ \text{Parent}(\text{ZEUS}), \text{Parent}(\text{APOLLO}), \neg\text{Father}(\text{HERA}), \\ \text{God}(\text{ZEUS}), \text{God}(\text{APOLLO}), \text{God}(\text{HERA}), \neg\text{God}(\text{HERCULES}), \\ \text{hasChild}(\text{ZEUS}, \text{APOLLO}), \text{hasChild}(\text{HERA}, \text{APOLLO}), \\ \text{hasChild}(\text{ZEUS}, \text{HERCULES}), \}$$

Suppose  $F = \{F_1, F_2, F_3, F_4\} = \{\text{Male}, \text{God}, \text{Parent}, \text{FatherWithoutSons}\}$ .

Let us compute the distances (with  $p = 1$ ):

$$d_1^F(\text{HERCULES}, \text{ZEUS}) =$$

$$(|1 - 1| + |0 - 1| + |1/2 - 1| + |1/2 - 0|) / 4 = 1/2$$

$$d_1^F(\text{HERA}, \text{HERCULES}) =$$

$$(|0 - 1| + |1 - 0| + |1 - 1/2| + |0 - 1/2|) / 4 = 3/4$$



## Distance Measure: Discussion

- The measure is a semi-distance (i.e. it does not guarantee that if  $d_p^F(a, b) = 0 \Rightarrow a = b$ )
- *More similar* the considered *individuals are*, more similar the projection function values are  $\Rightarrow d_p^F \simeq 0$
- *More different* the considered *individuals are*, more different the projection values are  $\Rightarrow$  the value of  $d_p^F$  will increase
- The measure complexity mainly depends from the complexity of the *Instance Checking* operator for the chosen DL
  - $Compl(d_p^F) = |F| \cdot 2 \cdot Compl(IChk)$
- **Optimal discriminating feature set could be learned**

# Measure Optimization: Feature Selection

- **Assumption:**  $F$  represents a sufficient number of (possibly redundant) features able to really discriminate individuals.
  - The choice of the features – *feature selection* – may be crucial
- **Proposal of optimization algorithms** that are able to find/build optimal discriminating concept committees  
**[Fanizzi et al. @ DL 2007 and @ ICSC 2007]**
  - **Idea:** Optimization of a *fitness function* that is based on the *discernibility factor of the committee*, namely
  - Given  $\text{Ind}(\mathcal{A})$  (or just a hold-out sample)  $HS \subseteq \text{Ind}(\mathcal{A})$  find the subset  $F$  that maximize the following function:

$$\text{DISCERNIBILITY}(F, HS) := \sum_{(a,b) \in HS^2} \sum_{i=1}^k \delta_i(a, b)$$

- The results obtained with KSs drawn from ontology libraries show that (a selection) of the (primitive and defined) concepts is often sufficient to induce satisfactory dissimilarity measures

## Dissimilarity under Uncertainty: Motivation

- The defined measure deals with uncertainty in a uniform way
  - the degree of discernibility of two individuals is null when they have the same behavior w.r.t. the same feature, even in the presence of total uncertainty of class-membership for both
  - When uncertainty regards only one projection, then they are considered partially (possibly) similar
  - **GOAL:** makes this uncertainty more explicit
- **New Proposal:** The dissimilarity between two individuals is assessed as a *combination* of **degree of evidence** that they differ w.r.t. a feature set
  - The measure is again based on the *degree of belief of discernibility* of individuals w.r.t. the features
    - the notion of *probability masses* of the basic events (class-membership) is exploited

## Computing the Probability Masses

Given the feature set  $F = \{F_1, F_2, \dots, F_k\}$ , the *probability mass* of the basic events "class-membership",  $\forall a \in \text{Ind}(\mathcal{A})$  and  $i \in \{1, 2, \dots, k\}$  *in case of uncertainty* is given by:

$$m_i(\mathcal{K} \models F_i(a)) \approx |\text{retrieval}(F_i, \mathcal{K})| / |\text{Ind}(\mathcal{A})|$$

$$m_i(\mathcal{K} \models \neg F_i(a)) \approx |\text{retrieval}(\neg F_i, \mathcal{K})| / |\text{Ind}(\mathcal{A})|$$

$$m_i(\mathcal{K} \models F_i(a) \vee \mathcal{K} \models \neg F_i(a)) \approx 1 - m_i(\mathcal{K} \models F_i(a)) - m_i(\mathcal{K} \models \neg F_i(a))$$

**Rationale:** the larger the (estimated) extension the more likely is for individuals to belong to the concept.

**In case of a certain answer received from the reasoner, the probability mass amounts to 0 or 1.**

# The Discernibility Function

The discernibility function (w.r.t. a concept) measures the amount of evidence that two input individuals are separated by that concept

## Discernibility Function

Given  $F_i \in F$ , the *discernibility function*  $\delta_i : \text{Ind}(\mathcal{A}) \times \text{Ind}(\mathcal{A}) \mapsto [0, 1]$  is defined,  $\forall (a, b) \in \text{Ind}(\mathcal{A}) \times \text{Ind}(\mathcal{A})$ , as follows:

$$\delta_i(a, b) := \begin{cases} 0 & \text{if } \mathcal{K} \models F_i(a) \wedge \mathcal{K} \models F_i(b) \\ 1 & \text{if } \mathcal{K} \models F_i(a) \wedge \mathcal{K} \models \neg F_i(b) \text{ or viceversa} \\ m_i(\mathcal{K} \models \neg F_i(b)) & \text{else if } \mathcal{K} \models F_i(a) \\ m_i(\mathcal{K} \models F_i(b)) & \text{else if } \mathcal{K} \models \neg F_i(a) \\ \delta_i(b, a) & \text{else if } \mathcal{K} \models F_i(b) \vee \mathcal{K} \models \neg F_i(b) \\ 2 \cdot m_i(\mathcal{K} \models F_i(a)) \cdot m_i(\mathcal{K} \models \neg F_i(b)) & \text{otherwise} \end{cases}$$

## Discernibility Function: Interpretation

- The extreme values  $\{0, 1\}$  are returned when the answers from the instance-checking service are certain for both individuals.
- If  $a$  is an instance of  $F_i$  (resp., its complement)  $\Rightarrow$  the discernibility depends on the belief of class-membership to the complement concept of  $b$ .
- If there is uncertainty for  $a$  but not for  $b$ , the function is computed swapping the roles of the two individuals.
- In case of uncertainty for both individuals, the discernibility is computed as the chance that they may belong one to  $F_i$  and one to its complement

## Dissimilarity Measure under Uncertainty: Definition

Following the *mixing combination rule*, the degree of belief can be combined for assessing a dissimilarity measure between individuals:

### Dissimilarity Measure under Uncertainty

Given an ABox  $\mathcal{A}$ , a dissimilarity measure  $d_{avg}^F : \text{Ind}(\mathcal{A}) \times \text{Ind}(\mathcal{A}) \mapsto [0, 1]$ ,  $\forall (a, b) \in \text{Ind}(\mathcal{A}) \times \text{Ind}(\mathcal{A})$ , is defined as follows:

$$d_{avg}^F(a, b) := \sum_{i=1}^k w_i \delta_i(a, b)$$

Where the weights can be defined uniform as:

- $w_i = 1/k$  or
- $w_i = \frac{u_i}{u}$  where  $u_i = \frac{1}{|\text{Ind}(\mathcal{A}) \setminus \text{retrieval}(F_k, \mathcal{K})|}$  and  $u = \sum_{i=1}^k u_i$

## Concept Dissimilarity under Uncertainty

The measures can be extended to the case of concepts, by recurring to the notion of *medoids*.

- The *medoid* of a group of individuals  $G = \{a_1, a_2, \dots, a_n\}$  is the individual that has the highest similarity w.r.t. the others i.e.  $\text{medoid}(G) = \operatorname{argmin}_{a \in G} \sum_{j=1}^n d(a, a_j)$

### Concept Dissimilarity

Given  $C_1$  and  $C_2$  concepts, let  $R_i = \{a \in \text{Ind}(\mathcal{A}) \mid \mathcal{K} \models C_i(a)\}$  be groups of individuals for  $i = 1, 2$  and  $m_i = \text{medoid}(R_i)$  their resp. medoids w.r.t. a given measure  $d_p^F$ . Then the function for concepts can be defined as:

$$d_p^F(C_1, C_2) := d_p^F(m_1, m_2)$$

Similarly, the distance of an individual  $a$  to a concept  $C$  can be defined:

$$d_p^F(a, C) := d_p^F(a, m)$$



# Conclusions

- The definition of dissimilarity measures over the spaces of individuals in a KB have been proposed
  - The measures are totally semantic (not language-dependent)
  - The measures are parameterized on a committees of concepts
- Optimal committees can be found maximizing a discernibility function, by the use of randomized search methods
- Dissimilarity measures able to cope with cases of uncertainty have been defined
  - based on a simple evidence combination method

## Future Works

Embedding the presented measures in distance-based methods to apply to KBs for:

- setting up logic approaches to ontology matching
- supporting a process of *(semi-)automatic classification* of new data (also as a first step towards ontology evolution)
- ranking the answers provided by a matchmaking algorithm on the ground of the similarity between the query concept and the retrieved individuals

# The End

That's all!