

SemMF: A Framework for Calculating Semantic Similarity of Objects Represented as RDF Graphs

Radoslaw Oldakowski¹, Christian Bizer²

Freie Universität Berlin, Institut für Produktion, Wirtschaftsinformatik und OR

Garystr. 21, D-14195 Berlin, Germany

¹cax@zedat.fu-berlin.de, ²chris@bizer.de

Abstract

In this paper we introduce an easy-to-use, flexible, open source framework for computing semantic similarity between objects represented as arbitrary RDF graphs, based on similarities of specified object properties. We outline the general architecture of the framework and describe three kinds of concept matching techniques implemented: the string matcher, the numeric matcher, and the taxonomic matcher exploiting a concept hierarchy. In a use case example we utilize the framework to match job postings with applicants' profiles.

1 Introduction

Semantic similarity measures play an important role in information retrieval by providing means to improve recall and precision. They are used in various application domains ranging from product comparison to job recruitment.

Whereas traditional similarity measures treat objects as sets of concepts drawn from a flat domain and calculate their similarity based on set intersection (shared concepts) more sophisticated approaches enhance these object models by adding a taxonomy describing the relationships among concepts leading to improved measures of similarity [3].

In this paper we introduce an easy-to-use, flexible framework for computing semantic similarity between objects represented as arbitrary RDF graphs, based on similarities of specified object properties. The user is given the choice to apply either hierarchical or non-hierarchical concept matching techniques for different kinds of object properties.

2 Framework Architecture

The matching engine takes as input a query object and a collection of resource objects to be matched against the query object (see Figure 1). Both are represented in RDF and may utilize different schema vocabularies. If they use concepts from a common taxonomy, an RDFS or OWL representation of this taxonomy has to be provided.

For accessing and querying of resource and query graphs as well as the underlying taxonomies the engine utilizes the Jena2 Semantic Web Framework [2].

2.1 Matching Description

In most cases not every object property is relevant for the similarity computation. For example, an object representing a certain product may contain manufacturer's phone number which may be irrelevant for comparing product's characteristics with customer's preferences. Thus, SemMF gives users the opportunity to explicitly specify each relevant property in the query RDF graph and to map it to the semantically corresponding property (i.e. holding the same kind of information, e.g. price information) in a resource RDF graph. Each mapping is assigned a concept matcher implementing a certain matching technique (see Section 2.2).

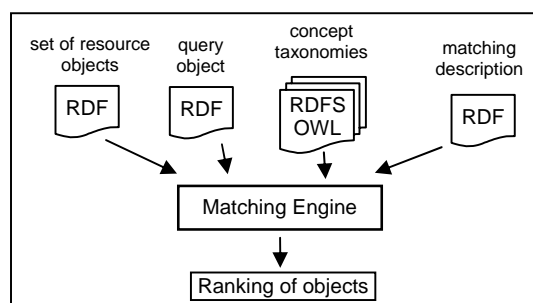


Figure 1. Matching Process

In the matching description the importance of each property can be specified by assigning it a certain *weight*. Moreover, properties to be matched can be grouped into thematic clusters, for example all properties describing the technical specification of a product. The property clustering enables to sort the matching result by cluster similarities. The matching description is represented in RDF using SemMF vocabulary provided with the framework.

2.2 Matchers

Owing to its flexibility by being able to work with different kinds of property matchers implementing a simple interface the framework can be utilized in various use case scenarios. In this paper we present three example implementations of different property matching techniques operating on either concepts from a taxonomy or RDF literals. Each matcher takes as input an RDF node from the query graph and a se-

mentally corresponding node from the resource graph. The value returned is their similarity score. There are a number of parameters for customizing each matcher.

Taxonomic Matcher

The taxonomic matcher computes the similarity between two concepts c_1 and c_2 based on the distance $d_c(c_1, c_2)$ between them, which reflects their respective position in the concept hierarchy [4]. The matcher is able to handle multiple inheritance of concepts at the leaf level of a taxonomy.

The concept similarity is defined as: $sim_c(c_1, c_2) = 1 - d_c(c_1, c_2)$. Every concept in a taxonomy is assigned a milestone value. Since the distance between two given concepts in a hierarchy represents the path over the closest common parent (ccp), the distance is calculated as:

$$d_c(c_1, c_2) = d_c(c_1, ccp) + d_c(c_2, ccp)$$

where $d_c(c, ccp) = milestone(ccp) - milestone(c)$

The milestone values of concepts in a taxonomy are calculated with (as set the matching description) either:

- a *linear milestone calculator* using the formula:

$$milestone(n) = 1 - [l(n) / l(N)]$$

where $l(n)$ is the depth of the node n in the hierarchy and $l(N)$ represents the deepest hierarchy level.

- or an *exponential milestone calculator* using the formula from [4]:

$$milestone(n) = \frac{1/2}{k^{l(n)}}$$

where k is a factor larger than 1 indicating the rate at which the milestone values decrease along the hierarchy. It can be assigned different values depending on the hierarchy depth.

This formula implies two assumptions: the semantic differences between upper level concepts are bigger than those between lower level concepts (in other words: two general concepts are less similar than two specialized ones) and that the distance between ‘brothers’ is greater than the distance between ‘parent’ and ‘child’.

Numeric Matcher

The numeric matcher is used to determine the similarity of two numeric values. A good application example for this matching technique is the comparison of product price some person is willing to pay (p_q) with the actual product price (p_r). For all $p_r > p_q$ the similarity shall decrease with increasing p_r . However, for all $p_r > upper\ bound$ specified in the matching description the similarity shall equal 0.

String Matcher

This matcher calculates the similarity of two given RDF nodes based on their string serialization. Thus, in the case of RDF literals not only their string value but also their language and datatype are compared.

2.3 Matching Process

Inside each thematic cluster the engine calculates the similarity between each query property and the corresponding resource property. These similarities are multiplied by the indicated weights and summed up yielding the cluster simi-

larity. All cluster similarities, in turn, multiplied by the specified cluster weights yield the object similarity.

However, if for a given query property value there is more than one semantically corresponding resource property value (e.g. a product may be available in different colors) the engine chooses the one with the highest similarity.

2.4 Matching Result

The output of the matching engine is a ranking of objects by their similarity values. The engine also provides a detailed description of the matching process (i.e. similarity values for all clusters and for each single property within a cluster together with the information about matchers used, property values, weights, etc.) which can be used to generate explanations for the calculated object similarity.

3 Use Case

The framework presented in this paper has been utilized as a foundation in a prototypical implementation of a semantic job portal [1]. It was used for matching job postings with applicants’ profiles (and vice versa) represented in RDF. All three kinds of matching techniques were applied: the string matcher for literal values (e.g. job duration [full time | part time]), numeric matcher for salary information, and the taxonomic matcher for concepts from taxonomies like skills, industry, and occupational classifications.

4 Conclusion

Computation of semantic similarity of objects is an often occurring problem. In this paper we present an easy-to-use flexible framework for calculating semantic similarity between objects represented as arbitrary RDF graphs, based on similarities of specified properties. The framework can be used in semantic web applications working with RDF graphs containing concepts drawn from both a flat and a hierarchical domain making it applicable in a wide range of different use cases.

SemMF is available under GNU Lesser General Public License (LGPL) from:

<http://www.wiwiss.fu-berlin.de/suhl/radek/semmf/>

References

- [1] Ch. Bizer et al. The Impact of Semantic Web Technologies on Job Recruitment Processes. *7. Internationale Tagung Wirtschaftsinformatik (WI2005)*, Bamberg, February 2005.
- [2] J. Carroll et al. Jena: Implementing the Semantic Web Recommendations. *The 14th International World Wide Web Conference (WWW2004)*, New York, May 2004.
- [3] P. Ganesan et al. Exploiting hierarchical domain structure to compute similarity. *ACM Transactions on Information Systems*, 21(1):63-94, 2003
- [4] J. Zhong et al. Conceptual Graph Matching for Semantic Search. *The 2002 International Conference on Computational Science (ICCS2002)*, Amsterdam, April 2002.

Live Demonstration of SemMF

We demonstrate the application of the Semantic Matching Framework (SemMF) in the online job recruitment domain. The framework presented in this paper has been utilized as a foundation in a prototypical implementation of the semantic job portal [1]. SemMF is used for matching job postings with applicants' profiles (and vice versa) represented in RDF using concepts from a flat domain as well as concepts from different taxonomies. As matching result a ranked list of best matching candidates for a given open position (or alternatively a list of best matching open positions for a given applicant's profile) is generated.

The *properties* to be matched describing a job posting and an applicant's profile are grouped into 3 thematic clusters:

- job details (weight value 0.2)
 - *job duration* [part time | full time]
 - *travel desired* [true | false]
 - *salary* [integer value]
- job specification (weight value 0.2)
 - *industry* [concepts from a taxonomy based on the *Classification of Industry Sectors (WZ2003)* – a German classification standard for economic activities]
 - *occupation* [concepts from a taxonomy based on the *Occupation Codes (Berufskennziffer* –

BKZ) – a German version of the SOC (*Standard Occupational Classification System*) which classifies workers into 5597 occupational categories according to occupational definitions]

- skill specification (weight value 0.6)
 - *skill* [concepts from a skill classification] (weight value 0.8)
 - *skill level* [concepts from a rating taxonomy] (weight value 0.2)

By presenting this use case of SemMF in the recruitment domain we:

- show an example of the matching description (see Section 2.1) created for this particular use case, serialized in RDF using SemMF vocabulary provided with the framework
- show possible parameter settings for the kinds of matchers chosen and explain their influence on the calculated concept similarity
- explain how weights at property and cluster level are applied to indicate the importance of certain properties
- generate ranked list of best matching candidates, and justify the usefulness of property clustering
- demonstrate how the matching information returned by the engine is used to generate explanations for the similarities computed (see Figure 2).

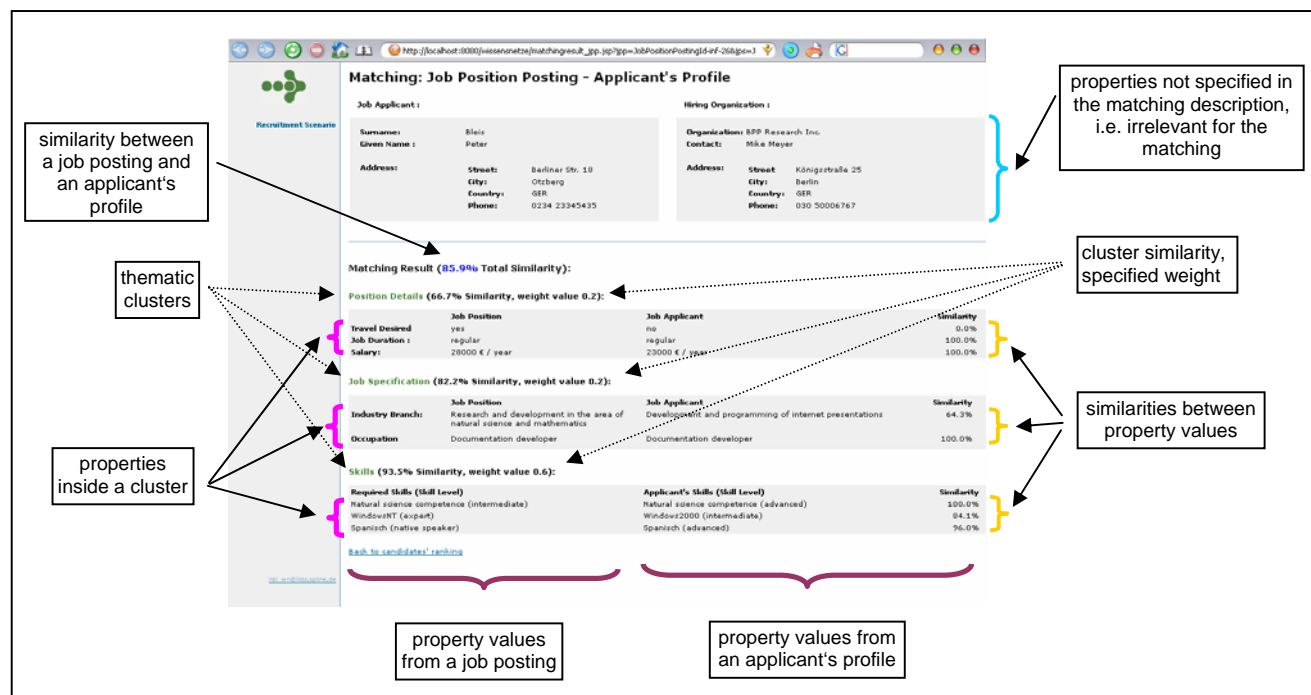


Figure 2. An Explanation of the Calculated Similarity Between a Job Posting and an Applicant's Profile.