

Learning medical ontologies from the Web

David Sánchez, Antonio Moreno

Intelligent Technologies for Advanced Knowledge Acquisition (ITAKA) Research Group
Department of Computer Science and Mathematics
Universitat Rovira i Virgili (URV)
Avda. Països Catalans, 26. 43007 Tarragona (Spain)
{david.sanchez, antonio.moreno}@urv.net

Abstract. The development of intelligent healthcare support systems always requires a formalization of medical knowledge. Domain ontologies are especially suitable for this purpose but their construction is, in most cases, manually addressed. This results in long and tedious development processes that hamper their real applicability. This is why there is a need of ontology learning methods that aid the ontology construction process. Considering the enormous amount of medical knowledge available freely on the Web, one may consider it as a valid source for developing knowledge acquisition systems. In this paper we offer an overview of an automatic and unsupervised method for learning domain ontologies from the Web. We also introduce its application over a specific medical domain in the frame of the K4Care European project.

Keywords: Ontology learning, web mining, knowledge acquisition, medical knowledge modelling.

1 Introduction

Medical ontologies are developed to solve problems such as the demand for reusing and sharing patient data or the transmission of these data. The unambiguous communication of complex and detailed medical concepts is a crucial feature in current medical information systems. In these systems, several agents must interact in order to share their results and, thus, they must use a medical terminology with a clear and non-confusing meaning [6].

The development of these ontologies is a complex task: on the one hand, they are general enough to be required for achieving consensus between a wide community of users and, on the other hand, they are concrete enough to present an enormous diversity with hundreds of possible concepts to model.

Medical ontology engineering is typically performed manually, requiring the intervention of medical specialists (which provide the medical knowledge) and knowledge engineers (which are able to formalize that knowledge). The required consensus is typically hampered by the difficulty of translating the shared world model of a medical community to the formal and explicit knowledge representation of an ontology. This produces long and tedious development stages that delay the applicability of the resulting ontologies.

Due to all these reasons, nowadays, there is a need of methods that can perform, or at least ease, the construction of medical ontologies. In this sense, *Ontology learning* is defined as the set of methods and techniques used for building from scratch, enriching, or adapting an existing ontology in a semi-automatic fashion using distributed and heterogeneous knowledge and information sources [6]. These methods allow a reduction in the time and effort needed in the ontology development process.

This data-driven knowledge acquisition process typically uses scientific texts, electronic dictionaries or medical repositories (such as UMLS). Considering the nature of those learning corpus (reduced scope, noise-free, trusted, structured), classical ontology learning methods have been designed [6].

In the last years, the growth of the medical information available on the Web provides users with a way for fast data access and information exchange. It is an invaluable tool for researchers, information engineers and health care companies and practitioners [5] for retrieving knowledge. These characteristics have motivated researchers [17] to consider the Web as a valid repository for *Information Retrieval* and *Knowledge Acquisition*. However, the extraction of information from web resources is a difficult task, due to their lack of structure and untrustiness, in addition to the ambiguity inherent to all resources written in natural language.

Despite all these shortcomings, as the number of resources available is so vast and the amount of people generating web pages is so enormous, it has been argued that the Web information distribution approximates the real distribution as used in society [3]. From the learning point of view, this is a very interesting characteristic and our motivation for using the Web as the source for knowledge acquisition.

So, in this paper, we present an overview of a novel approach for automatic domain ontology learning from the Web. Thanks to the amount of medical information available on the Web and the structured nature of medical knowledge, our method is especially suitable for learning medical ontologies. As a result of the application of this methodology over a medical domain, we introduce a case of study framed in the scope of the K4Care European research project. At the end, the main aim of this paper is to show the usefulness of the developed automatic learning method to aid medical researchers in modelling knowledge.

The rest of the paper is organized as follows. Section 2 presents an overview of the main steps involved in the ontology construction process, introducing the learning techniques employed for knowledge acquisition. Section 3 gives a general vision of our approach for learning domain ontologies from the Web, including the acquisition of taxonomic and non-taxonomic relationships. Section 4 presents and evaluates an example of the obtained results for a medical domain in the context of the K4Care project. The final section presents the conclusions and proposes lines of future work.

2 Ontology learning overview

Ontologies are composed at least by *classes* (concepts of the domain), *relations* (different types of binary associations between concepts or data values) and *instances* (real world individuals). Formally, an ontology often boils down to an object model represented by a set of concepts or classes C , which are *taxonomically* related by the

transitive *IS-A* relation $H \in C \times C$ and *non-taxonomically* related by named object relations $R^* \in C \times C \times \text{String}$. On the basis of the object model, a set of logical axioms, A , enforce semantic constraints.

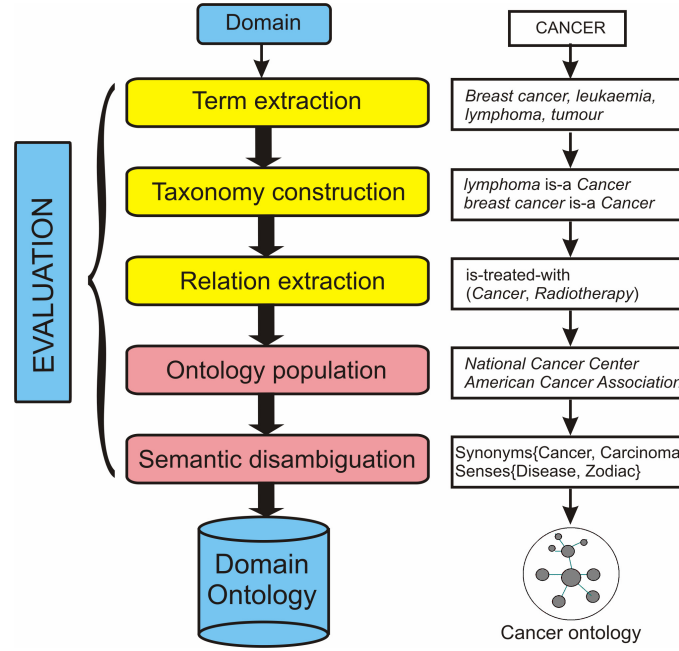


Fig. 1. General steps of the domain ontology learning process.

From the *ontology engineering* point of view [6], the main steps and ontology learning techniques employed for building ontologies are the following (see Fig.1):

- Extraction of terms that represent domain concepts, building a lexicon. In order to define an unsupervised approach, we have used statistical analysis about term co-occurrence [20]. Statistical values can be obtained in a very immediate way from hit counts of web search engines if the appropriate queries are performed [21]. In addition, thanks to the size and heterogeneity of the Web, those values are very robust, as they can approximate the true societal words usage [3].
- Construction of an initial taxonomy of concepts using *is-a* relations. Without using any kind of background knowledge, it is possible to extract hyponymy relationships using domain independent linguistic patterns [7, 8].
- Learning non-taxonomic relations. It is considered as the least tackled problem within ontology learning [10]. We have addressed the problem by extensively using verb phrases as the central point of a relation. Following the same philosophy as in the taxonomic case, we consider specific verb phrases as particular domain dependent semantic patterns that express a particular non-taxonomic relationship [18]. Lightweight analytic procedures [15] and statistics compiled from querying a web search engine [21] complete the proposed non-taxonomic learning method.

- Ontology population by the detection of instances for the discovered concepts. We have limited this stage to the discovery of named entities. We use language dependent rules (capitalization for the English language) to detect proper names.
- Optionally, we can also treat semantic ambiguity in order to improve the quality of the results. We have developed complementary mechanisms to deal with polysemy and synonymy [19].
- Evaluation of the obtained results (concepts, instances and relationships). As ontological knowledge is non-uniquely expressible, the comparative evaluation of different approaches is difficult. For that reason, ontology learning evaluation is recognized to be an open problem [6]. In our case, as the quality of the final result will depend on the performance of every step of the learning process, specific evaluation methods for each one of them have been designed. Whenever a domain standard is available (e.g. MESH for the taxonomic case), results have been carefully compared. In other cases (as for the non taxonomic relationships), an expert's opinion may be required.

3 Ontology learning methodology

The core of our novel Web based approach covers the acquisition of domain terms and the definition of taxonomic and non taxonomic relations. Its main advantage is the automatic and unsupervised operation, creating domain ontologies from scratch.

Even though we have developed individual methodologies for dealing with each of the learning steps, they have been designed to be executed in an integrated and iterative way. So, each step can be bootstrapped with the knowledge acquired up to that moment. In this manner, new concepts and relationships can be used as seeds for further analysis. Through several iterations, the system incrementally constructs the semantic network of concepts composing the domain ontology.

As shown in Fig. 2, the learning process is divided in several phases. The *Taxonomic learning* [20] starts from a user specified keyword (e.g. *cancer*) that indicates the domain for which the ontology should be constructed. The system starts by querying a web search engine to obtain a corpus of web documents to analyse. At this initial stage, only general queries using several domain independent patterns for hyponymy detection (e.g. "*cancers such as*") are constructed. Web content is parsed in order to find matchings for those patterns and extract taxonomic candidates (e.g. "*cancers such as leukaemia or breast cancer*"). Domain verbs are also retrieved at this stage (e.g. "*cancer is associated with*"). Several iterations using different patterns are performed in order to minimize language ambiguity and a final set of candidates is compiled. Each taxonomic candidate is then evaluated using web based statistical scores about term co-occurrence. New queries for web search engines are constructed in order to infer the degree of relatedness of a taxonomic candidate (e.g. "*breast cancer*") and the domain (e.g. "*cancer*"). Web search hit counts are used to compute statistical scores (more details in [20]); those candidates with the higher scores are selected as valid taxonomic specialisations for the domain. In parallel, a procedure that detects named entities using capitalization heuristics is executed. It allows filtering the retrieved candidates by including real world individuals (e.g. "*American*").

Cancer Association”) as instances –and not subclasses- of the ontology. The output of this process is a one-level taxonomy with general terms (e.g. several types of *cancer*) and a set of verbs that have appeared in the same context –sentence- as the searched domain keyword (e.g. *is associated with*, *causes*, *is treated with*, etc.).

The next stage is the *Non-taxonomic learning* [18]. This process begins with the verb list compiled in the previous step, which is used as the knowledge base for the non-taxonomic learning. Each verb can be used as a bootstrap by constructing domain related patterns (e.g. “*cancer is treated with*”) that are queried into the Web search engine. Additional web resources are retrieved and analysed to find verb-based pattern matchings (e.g. “*cancer*” “*is treated with*” “*radiotherapy*”). Again, candidates for non-taxonomic relations (e.g. “*radiotherapy*”) are ranked and selected using web-scale statistical scores (more details in [18]). Finally, the verb phrase is used to link each pair of concepts, defining a set of domain binary relations.

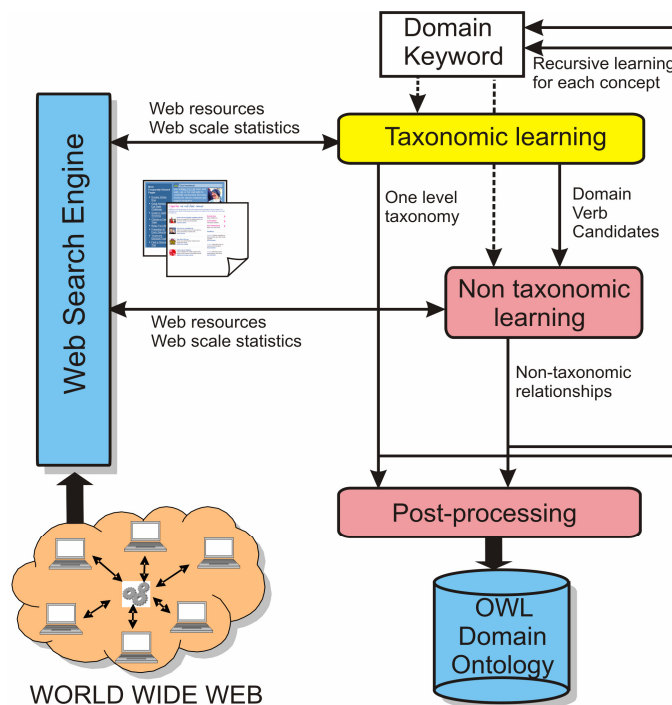


Fig. 2. Ontology learning methodology.

The two previous learning stages are *recursively* executed for each obtained concept (taxonomically –e.g. “*breast cancer*”- or non-taxonomically –e.g. “*radiotherapy*”- related). Each one becomes an individual seed for further analysis. Those new learning iterations can use the already acquired knowledge as a bootstrap to contextualize web queries and to obtain more concrete web resources. The specific number of learning iterations, the amount of resources analysed at each step and the finalization of the recursive analysis is controlled by the algorithm itself. In more detail, the system continuously monitors the learning evolution computing, at the end

of each individual learning pass (i.e. the query and processing of a specific taxonomic or non-taxonomic pattern), the percentage of selected and rejected candidates according to the statistical scores. This value measures the learning throughput of a specific concept and pattern, allowing the system to self control the learning process. On the one hand, the most productive ones are further evaluated by retrieving and analysing additional web resources. On the other hand, for the less productive ones, the process is finished and the next pattern and/or concept are taken. Thanks to the higher degree of contextualization allowed by the bootstrapped knowledge, the learning process is able to finish adequately.

At the end of this incremental process, we obtain a multi-level taxonomy in which each concept can be non-taxonomically related to other ones. An illustrative example of the kind of the structure that we are able to obtain is presented in the next section.

Finally, a *post-processing* stage is introduced in order to detect implicit relationships (such as multiple inheritances), equivalencies, avoid redundancies and to obtain a more compact structure that becomes the final domain ontology.

Note that this section only represents an overview of the learning process, as our main objective is to introduce the usefulness of the developed methodologies in modelling domain knowledge. More details are offered in [18], [19] and [20].

4 Case of study

In this section we offer an example and an evaluation of our learning method for a particular medical domain in the framework of the European project K4Care.

K4Care (<http://www.k4care.net>) aims to create, implement, and validate a knowledge-based healthcare model for the professional assistance to senior patients at home. This new Healthcare Model for home care will contribute to achieve a European standard supported by ICT technologies that improves the efficiency of the care services for all the citizens in the enlarged Europe. K4Care relies on the definition of domain ontologies, Electronic Health-Care Records and Formal Intervention Plans. In more detail, a specific Patient-Case Profile Ontology (CPO) is being constructed. It aims to structure the knowledge available about the care of patients. It combines diseases, syndromes, signs and symptoms, social issues, assessment tests, and interventions to define a knowledge model of how to deal with Home-Care Patients. More concretely [9], the information available in a patient's *Electronic Health-Care Record*, combined with the results obtained for some clinical tests (*Comprehensive Assessment*), will be processed using the CPO as the knowledge base in order to infer the patient's syndromes. Then, associated *Formal Intervention Plans* for the discovered pathologies can be used to aid the healthcare providers in specifying the patient's particular treatment (*Individual Intervention Plans*).

The CPO is being currently defined manually from scratch, from the interaction of medical experts and knowledge engineers, supposing a considerable effort. Up to this moment, the ontology models the main entities that are relevant within the project scope. This ontology is heavily focused on the taxonomic aspect of the knowledge modelling (e.g. classification of different types of diseases), and offers a very little degree of general –non-taxonomic– semantic interlinkage between concepts (e.g. the

In more detail, among the different entities modelled in the CPO, there are several diseases which are considered within the K4Care scope (senior patient typical pathologies). The most exhaustively covered one is *Dementia*, for which several specialisations have been defined.

We executed our learning methodology for that domain. As a result, we obtained a Dementia ontology covering related classes, instances and taxonomic and non-taxonomic relationships. Most of the taxonomic relationships and some of the more relevant non-taxonomic ones are presented in Fig. 3. Considering the amount of discovered ontological entities, one can imagine the degree of human effort required to compile and structure them appropriately.

In order to evaluate the quality of these results in terms of *precision*, we compared them against a widely used medical standard (MESH <http://www.nlm.nih.gov/mesh/MBrowser.html>). We have queried the MESH browser to check if a discovered concept is present or not, obtaining a precision of 74% for the taxonomic case. Non-taxonomic relationships cannot be so easily checked as they are typically not modelled in standard classifications. A manual evaluation of the 99 discovered relationships measured a precision of 71.1%. In both cases, precision is enough to consider the results as a reliable knowledge base for the domain.

Next, we compared the obtained ontology in terms of *recall* against the K4Care hand made ontological specification. Considering that mainly taxonomic relationships are modelled in the CPO, including 15 types of diseases and 7 classes of dementia, we were able to retrieve 57% of them. The non discovered ones are referred to the vaguest subclasses (e.g. *Mixed Type*, *Other Degenerative Dementia* and *Unspecified Dementia*) which are hardly distinguished from general adjectives. However, in total, we automatically discovered 25 direct types of dementia, more than 200 classes and 99 non-taxonomic relationships. Those last ontological entities are especially interesting due to the inherent difficulty of modelling non-taxonomic knowledge.

All those results were obtained in a completely automatic and unsupervised way, without requiring any kind of previous knowledge, search tuning or user's intervention. The system extensively queried a web search engine and analysed a large amount of web resources (21004). In any case, before their real application, the ontology should be checked and filtered by a medical expert.

5 Conclusions and future work

Many knowledge acquisition approaches have been developed in the past. Different methodologies have been designed according to the knowledge source [14]: texts, dictionaries, knowledge bases, semi-structured data, relation schemas, etc. Considering the nature of those classical repositories, the common characteristics of classical knowledge acquisition methods are:

- Many of them [4, 12] use as learning corpus a reduced and pre-selected set of relevant documents for the covered domain. This approach solves some problems about untrustiness, lack of structure, noise and size that arise when developing an unsupervised, domain-independent Web-based approach.

- Most of the knowledge acquisition methodologies [1, 11] use predefined knowledge to some degree, like training examples, previous ontologies or semantic repositories. This fact hampers the development of domain independent solutions, weakening the scalability and versatility of those systems in wide and heterogeneous environments like the Web.
- Most of them only cover the taxonomic aspect of the ontology learning process [10]. There have been very few attempts of non-taxonomic learning and, in many situations [13, 16], extracted relationships remain unlabelled.

On the contrary, we aim to obtain domain ontologies from scratch without any previous knowledge, adapting several classical techniques for knowledge acquisition (linguistic patterns, statistical analysis, etc.) to the special casuistry of the Web. We also cover all the main steps of the ontology learning process, configuring an integrated and intelligent learning approach.

Our approach is fully unsupervised. This is especially important due to the amount of available web resources, avoiding the need of a domain expert. The system has continuous feedback about the productivity of the learning task performed at each moment, guiding the learning to the most productive entities. In addition, the learning is automatic, allowing to easily perform executions at any time in order to retrieve updated results. This characteristic fits very well with the dynamic nature of the Web.

Domain ontologies are crucial in many knowledge intensive areas requiring interoperability such as the Semantic Web [2], e-commerce and e-medicine. From the presented results and posterior analysis, we can conclude that the use of automated ontology learning tools that are able to obtain with quite good accuracy (precision) a domain ontology in a few hours, can suppose a great help for ontology modellers. For the introduced example, the labour of specifying taxonomic entities can be reduced by more than a half. In addition, new ontological entities not yet considered (like new taxonomic terms and additional non-taxonomic relationships) are proposed.

Thanks to those advantages, ontology construction can be reduced from the fully manual ontology engineering effort -requiring an active participation of knowledge engineers- to a semi-automatic process which only requires refining a quite complete ontological structure. In this last case, ontologies can be evaluated and edited by the domain expert without advanced knowledge modelling skills.

As future research, we plan to apply our results to aid in the construction of the K4Care knowledge model. Other interesting syndromes, symptoms or diseases framed in the scope of the project can be further analysed. We would like to receive feedback for our results from expert medical partners of the K4Care project. This may give us an idea of the potential benefits and improvements that our solution may offer, such as the reduction in development time of the required knowledge structures.

Acknowledgements

The work has been supported by *Departament d'Universitats, Recerca i Societat de la Informació de la Generalitat de Catalunya i del Fons Social Europeu* of Catalonia. Authors would also like to acknowledge the support of the K4Care European research project (IST-2004-026968).

References

1. Alfonseca, E., Manandhar, S.: An unsupervised method for general named entity recognition and automated concept discovery. In Proc. of the 1st International Conference on General WordNet (2002)
2. Berners-lee, T., Hendler, J., Lassila, O.: The semantic web. Scientific American (2001)
3. Cilibrasi, R., Vitanyi, P.M.B.: Automatic meaning discovery using Google, Available at: <http://xxx.lanl.gov/abs/cs.CL/0412098> (2004)
4. Faatz, A., Steinmetz, R.: Ontology enrichment with texts from the WWW. In Proc. of Semantic Web Mining 2nd Workshop at ECML/PKDD-2002 (2002)
5. Forkner-Dunn, J.: Internet-based Patient Self-care: the Next Generation of Health Care delivery. Journal Med Internet Research 5(2) (2003)
6. Gómez-Pérez, A., Fernández-López, M., Corcho, O.: Ontological Engineering, Springer Verlag, 2nd printing. (2004)
7. Grefenstette, G.: SQLET: Short Query Linguistic Expansion Techniques: Palliating One-Word Queries by Providing Intermediate Structure to Text. In: Information Extraction: A Multidisciplinary Approach to an Emerging Information Technology, LNAI 1299(6) (1997) 97-114
8. Hearst, M.A.: Automatic acquisition of hyponyms from large text corpora. In Proceedings of the 14th International Conference on Computational Linguistics (1992) 539-545
9. Isern, D, Moreno, A., Pedone, G. and Varga, L.: An Intelligent Platform to Provide Home Care Services. In Proceedings of the Workshop From Knowledge to Global Care. 11th Conference on Artificial Intelligence in Medicine (2007).
10. Kavalec, M., Maedche, A., Skátek, V.: Discovery of Lexical Entries for Non-taxonomic Relations in Ontology Learning. SOFSEM 2004, LNCS 2932 (2004) 249-256
11. Khan, L. and Luo, F.: Ontology Construction for Information Selection. In Proc. of 14th IEEE International Conference on Tools with Artificial Intelligence (2002) 122-127
12. Lee, D., Na, J., Khoo, C.: Ontology Learning for Medical Digital Libraries. In Proc. of ICADL 2003, LNCS 2911 (2003) 302-305
13. Maedche, A., Staab, S.: Discovering Conceptual Relations from Text. In Proc. of the 14th European Conference on Artificial Intelligence, IOS Press (2000) 321-325
14. Maedche, A., Staab, S.: Ontology Learning for the Semantic Web, IEEE Intelligent Systems, S.I. on the Semantic Web, 16(2) (2001) 72-79
15. Pasca, M.: Acquisition of Categorized Named Entities for Web Search. CIKM'04 (2004) 137-145
16. Reinberger, M.L., Spyns, P.: Discovering knowledge in texts for the learning of DOGMA inspired ontologies. In Proc. of Workshop on Ontology Learning and Population, ECAI 2004 (2004) 19-24
17. Resnik, P., Smith, N.: The web as a parallel corpus, Computational Linguistics, 29(3) (2003) 349-380
18. Sánchez, D. and Moreno, A.: Discovering Non-taxonomic Relations from the Web. 7th International Conference on Intelligent Data Engineering and Automated Learning. LNCS 4224 (2006) 629-636
19. Sánchez, D., Moreno A.: Development of new techniques to improve Web Search. In Proc. of 9th International Joint Conference on Artificial Intelligence (2005) 1632-1633
20. Sánchez, D., Moreno, A.: A methodology for knowledge acquisition from the web. International Journal of Knowledge-Based and Intelligent Engineering Systems 10(6) (2006) 453-475
21. Turney, P.D.: Mining the Web for synonyms: PMI-IR versus LSA on TOEFL. In Proceedings of the Twelfth European Conference on Machine Learning (2001) 491-502