



# Data-Mining Algorithms with Semantic Knowledge

PROYECTO DE INVESTIGACIÓN PROGRAMA NACIONAL DE INVESTIGACIÓN FUNDAMENTAL, PLAN NACIONAL DE I+D+i 2008-2011 ÁREA TEMÁTICA DE GESTIÓN: Tecnologías informáticas

# Internal project report T3.1 Damask Ontology

Authored by Carlos Vicient, Universitat Rovira i Virgili

**Co-authored by** Antonio Moreno, Universitat Rovira i Virgili





# Document information

project name:	DAMASK	
Project reference:	TIN2009-11005	
type of document:	Deliverable	
file name:	ProjRepT3.1	
version:	final	
authored by:	C.Vicient	14/10/2011
co-authored by	A.Moreno	
released by:	Co-ordinator	
approved by:	Co-ordinator	A. Moreno



DAMASK

# Document history

version	date	reason of modification
0.1	14.October.2011	A preliminary release of the document.
1.0	20.Nov.2001	Final version
	200	
	200	
	200	
	200	
	200	



# **Table of Contents**

1		Introduction	3
2		Ontology	4
	2.1	Background ontologies	4
	2.2	Damask Ontology	5
3		References	6
4		Annex I – TourismOWL	7
5		Annex II – Space.owl	12
6		Annex III – PCTTO.owl	16
7		Annex IV – damaskontology.owl	20





### **1** Introduction

The first task of DAMASK, called T1 - Semantic integration of the information available in heterogeneous Web resources, includes two preliminary subtasks, tasks 1.1 and 1.2 (see Figure 1). The former task discusses all related works in information extraction from the Web distinguishing algorithms to extract structured, semi-structured and non-structured resources. Deliverable D1 was the result of the task 1.1., and its aim was to make a state-of-the-art on Information Extraction techniques applied to Web resources. An internal report, which was the output of task 1.2, lists the types of data that can be extracted using the previous algorithms. As stated in that document, DAMASK takes into account three data types: Measurement, Nominal and Semantic. Finally, Deliverable D2 was the result of task 1.3. (see Figure 1), and its aim was to present the proposed methodology for Ontology-based feature extraction. The main goal of that task was to design and implement a novel method that was able to extract relevant features from a range of textual documents going from plain textual data to semi-structured resources. The designed methodology was able to take profit from pre-processed input when it is available in order to complement its own learning algorithms[1, 2]. The key point of the work was to complement the syntactical parsing and other natural language processing techniques with the knowledge contained in an input ontology (which ideally, should model the knowledge domain in which the posterior data analysis will be focused -e.g. touristic points of interest).



Figure 1: Tasks of DAMASK



Task T2 was focused on the goal O2 of the project: design of a clustering method based on ontologies. The inputs of this task were (1) a data matrix object  $\times$  attribute (e.g. touristic destinations) and (2) a domain ontology. Based on those inputs, a method for automatically building clusters was needed. During the clustering, contextual knowledge provided by the domain ontology was used [3]. Finally, an automatic interpretation process of the clusters was required, in order to obtain a semantic description of the clusters that can help the user in his/her decision making tasks.

In task T3 the goal is to evaluate the deployment of the methods designed in the previous tasks in a particular case study: a personalized recommendation system of touristic destinations. A Web application will be designed to offer this kind of recommendations to any user. The tool will be focused on searching touristic destinations in the different types of touristic resources available in Internet using the tools developed in task T1. The clustering methods defined in T2 will then be applied to obtain a classification of touristic destinations based on the domain knowledge and the user preferences, in order to be able to recommend the set of places that match better with the user's interests. In task T3, the experience and touristic knowledge of the EPOs of the project is very valuable (EPO=organism or industry interested in the results of the project). They are particularly very helpful in the construction of the domain ontology and the evaluation of the recommendations provided by the system. Particularly, this document (Internal project report T3.1) is the result of the design and creation of the aforementioned ontology (Task 3.1). As this work is focused on the touristic domain, it has been implemented in order to model the most important aspects which represent a touristic destination such as monuments, landmarks, etc. To do so, three existent ontologies modelled in this domain were studied and merged.

# 2 Ontology

This section states some ontologies related in the area of tourism (section 2.1) and presents the proposed ontology for DAMASK (section 2.2).

#### 2.1 Background ontologies

#### TourismOWL.owl ontology

This ontology models touristic points of interest for different kinds of tourist profiles. It was designed in a final year project [4] based on information extracted through Wikipedia articles. It consists of 315 classes and a depth of 5 hierarchical levels. Its main classes represent concepts related with administrative divisions (borough, city, country, village, etc.), buildings (commercial buildings, cultural buildings, religious buildings, sport buildings, etc.), festivals (art festivals, music festivals, carnival, etc.), landmarks (commemorate landmarks, geographical landmarks, memorial landmarks, etc.), museums (archaeology museum, history museum, science museum, etc.) and sports (football, basketball, hockey, formula one, etc.). See Annex I.





#### Space.owl ontology

This ontology was found by looking up ontologies using the Web search engine SWOOGLE<sup>1</sup> that is specialized in ontologies. The Space.owl<sup>2</sup> ontology consists of 188 classes and a depth of 6 hierarchical levels. It contains concepts related with three main topics: geographical features (i.e., archipelago, beach, river, forest, etc.), geopolitical entities (i.e., country, capital, city, district, street, etc.) and places, which includes business places (factory, convention centre, etc.), private places (residential structure, home, etc.) and public places (educational and medical structures, entertainment places, shopping facilities, transportation connections, etc.). See Annex II.

#### **PCTTO.owl ontology**

This ontology was created with the collaboration of the project EPOs: the Scientific and Technological Park for Tourism and Leisure of Vila-Seca (Tarragona), and the Tourism Observatory for Costa Daurada. It is focused on tourist activities. The ontology represents up to 203 connected concepts in 5 hierarchy levels. It is structured around eight main concepts, that constitute the first level of the hierarchy: "Events", "Nature", "Culture", "Leisure", "Sports", "Towns", "Routes" and "ViewPoints". The last three classes are considered transversal concepts, since they share children nodes with other main classes, e.g., "Routes" and "Nature" are both superclasses of the "NatureRoutes" class. The rest of the concepts in the ontology are connected via is-a (subclass) relationships with these main classes. The ontology is not a pure taxonomy, as it contains multi-inheritance between concepts, e.g., EthnographicMuseum is a subclass of both Museum and Tradition-alCulturalElement. This ontology was developed using the Thesaurus of the World Tourism Organization as a reference guide to represent the touristic and leisure activities in the "Costa Daurada and Terres de l'Ebre" region. See Annex III

#### 2.2 Damask Ontology

This ontology<sup>3</sup> is the result of merging and combining the aforementioned ontologies. It consists of 538 classes connected in 9 hierarchy levels. It is structured around 4 main concepts that constitute the first level of the hierarchy: "geopolitical division", "activity", "point of interest" and "geographical feature". The Damask Ontology is not a pure taxonomy, as it contains multi-inheritance between concepts. Annex IV depicts the taxonomy of those main concepts from the first level. The whole ontology is available at the DAMASK project Web page<sup>4</sup> (section "ontologies") in "png", "jpg" and "owl" formats.

<sup>&</sup>lt;sup>4</sup> http://deim.urv.cat/~itaka/CMS2/index.php?option=com\_content&task=view&id=29&Itemid=51



<sup>&</sup>lt;sup>1</sup> http://swoogle.umbc.edu/

<sup>&</sup>lt;sup>2</sup> http://deim.urv.cat/~itaka/CMS2/images/ontologies/space.owl

<sup>&</sup>lt;sup>3</sup> http://deim.urv.cat/~itaka/CMS2/images/ontologies/damaskontology.owl



## **3** References

[1] C. Vicient, D. Sánchez, and A. Moreno, "Ontology-Based Feature Extraction," in *Proceedings* of the IEEE/WIC/ACM International Conferences on Web Intelligence and Intelligent Agent Technology, NLPOE, Lyon, France, 2011, pp. 189-192.

[2] C. Vicient, D. Sánchez, and A. Moreno, "A methodology to discover semantic features from textual resources," in *Proceedings of the 6th International Workshop on Semantic media adaptation and personalization (SMAP 2011)*, Vigo, Spain, 2011.

[3] M.Batet, "Ontology-based semantic clustering". PhD Thesis, Univ. Rovira i Virgili, Tarragona, February 2011.

[4] C. Vicient, "Extracció basada en ontologies d'informació de destinacions turístiques a partir de la Wikipedia," . Computer Science Final Year Project, Universitat Rovira i Virgili, Tarragona, 2009.

# 4 Annex I – TourismOWL

Following it is shown the taxonomy of the TourismOWL.owl ontology. Each sub-tree depicts the main classes of the ontology and its hierarchy.







Outlet

Fashion\_Shop

is-a





# 5 Annex II – Space.owl

Following it is shown the taxonomy of the Space.owl ontology. Each sub-tree depicts the main classes of the ontology and its hierarchy.







## 6 Annex III – PCTTO.owl

Following it is shown the taxonomy of the PCTTO.owl ontology. Each sub-tree depicts the main classes of the ontology and its hierarchy.







# 7 Annex IV – damaskontology.owl

Following it is shown the taxonomy of the damaskontology.owl ontology. Each sub-tree depicts the main classes of the ontology and its hierarchy.







