

Dynamic Taxonomies for Intelligent Information Access



Giovanni Maria Sacco
Università di Torino, Italy

INTRODUCTION

End-user interactive access to complex information is a key requirement in most applications, from knowledge management, to e-commerce, to portals. Traditionally, only access paradigms based on the retrieval of data on the basis of precise specifications have been supported. Examples include queries on structured databases and information retrieval. There is now a growing perception that this type of paradigm does not model a large number of search tasks, such as product selection in e-commerce sites among many others, that are imprecise and require exploration, weighting of alternatives and information thinning. The widespread feeling that “search does not work” and “information is too hard to find” shows evidence of the crisis of traditional access paradigms.

New access paradigms supporting exploration are needed. Since the goal is end-user interactive access, a holistic approach in which modeling, interface and interaction issues are considered together, must be used and will be discussed in the following.

BACKGROUND

Four retrieval techniques are commonly used: a) information retrieval (IR) systems (van Rijsbergen, 1979; Baeza-Yates & Ribeiro-Neto, 2011) also recently known as search engines; b) queries on structured databases; c) hypertext/hypermedia links and d) static taxonomies, such as Yahoo!.

IR systems exhibit an extremely wide semantic gap between the user model (concepts) and the model used by commercial retrieval systems (words). This leads to a significant loss of relevant information (Blair & Maron, 1985), and to poor user interaction because query formulation is difficult and no or very little assistance is given. In addition, since results are

normally presented as a flat list with no systematic organization, no exploration is possible. Database queries require structured data and are not applicable to situations in which information are textual and not structured or loosely structured. Exploration is usually limited to sorting flat result lists according to different ordering criteria.

Hypermedia techniques (Groenbaek & Trigg, 1994) have become pervasive and support exploration. However, they do not support abstraction so that exploration is performed one-document-at-a-time, which is quite time consuming. Building and maintaining non-trivial hypermedia networks is very expensive.

Traditional taxonomies are based on a hierarchy of concepts that can be used to select areas of interest and restrict the portion of the infobase to be retrieved. They are easily understood by end-user, but they are not scalable for large information bases (Sacco, 2006a), so that the average number of documents retrieved becomes rapidly too large for manual inspection.

A more recent approach is the Semantic Web (Berners-Lee et al., 2001). Although one of the driving forces behind it is retrieval, the general semantic schemata proposed are intended for programmatic access and are known to be difficult to understand and manipulate by the casual user. User interaction must be mediated by specialized agents, which increases costs, time to market and decreases the transparency and flexibility of user access.

DYNAMIC TAXONOMIES

Dynamic taxonomies (Sacco, 2000, later also called *faceted search systems*) are a general knowledge management model based on a multidimensional classification of heterogeneous data items and are used to explore/browse complex information bases in a guided yet unconstrained way through a visual interface. The

DOI: 10.4018/978-1-4666-5888-2.ch382

reader is addressed to Sacco and Tzitzikas, 2009, for the most complete and up-to-date book on this model.

The intension of a dynamic taxonomy is a taxonomy designed by an expert. This taxonomy is a concept hierarchy going from the most general to the most specific concepts. A dynamic taxonomy does not require any other relationships in addition to *subsumptions* (e.g., IS-A and PART-OF relationships). Directed acyclic graph taxonomies modeling multiple inheritance are supported but rarely required.

In the extension, items can be freely classified under n ($n > 1$) concepts at any level of abstraction (i.e. at any level in the conceptual tree). The multidimensional classification required by dynamic taxonomies is a generalization of the monodimensional classification scheme used in conventional taxonomies and models common real-life situations. First, items are very often about different concepts: for example, a news item on September 11th, 2001 can be classified under “terrorism,” “airlines,” “USA,” etc. Second, items to be classified usually have different features, “perspectives” or facets (e.g. Time, Location, etc.), each of which can be described by an independent taxonomy.

In dynamic taxonomies, a concept C is just a label that identifies all the items classified under C . Because of the subsumption relationship between a concept and its descendants, the items classified under C ($\text{items}(C)$) are all those items in the *deep extension* of C , i.e. the set of items identified by C includes the *shallow extension* of C (i.e. all the items directly classified under C) union the deep extension of C 's sons. By construction, the shallow and the deep extension for a terminal concept are the same. This set-oriented approach implies that logical operations on concepts can be performed by the corresponding set operations on their extension, and therefore the user is able to restrict the information base (and to create derived concepts) by combining concepts through all the standard logical operations (and, or, not).

A fundamental feature of this model is that dynamic taxonomies can find all the concepts related to a given concept C : these concepts represent the conceptual summary of C . Concept relationships other than subsumptions are inferred on the basis of empirical evidence through the extension only, according to the following *extensional inference rule*: two concepts A and B are related iff there is at least one item d in the knowledge base which is classified at the same time under A or under one of A 's descendants and under B

or under one of B 's descendants. For example, we can infer an unnamed relationship between *terrorism* and *New York*, if an item classified under *terrorism* and *New York* exists. At the same time, since *New York* is a descendant of *USA*, also a relationship between *terrorism* and *USA* can be inferred.

The extensional inference rule can be easily extended to cover the relationship between a given concept C and a concept expressed by an arbitrary subset S of the universe: C is related to S iff there is at least one item d in S which is also in $\text{items}(C)$. Hence, the extensional inference rule can produce conceptual summaries not only for base concepts, but also for any logical combination of concepts. In addition, since it is immaterial how S is produced, dynamic taxonomies can produce summaries for sets of items produced by other retrieval methods such as database queries, shape retrieval, etc. and therefore access through dynamic taxonomies can be easily combined with any other retrieval method.

Dynamic taxonomies are defined in terms of conceptual descriptions of items, so that heterogeneous items of any type and format can be managed in a single, coherent framework. Finally, since concept C is just a label that identifies the set of the items classified under C , concepts are language-invariant, and multilingual access can be easily supported by maintaining different language directories, holding language-specific labels for each concept in the taxonomy.

EXPLORATION

The user is initially presented with a tree representation of the initial taxonomy for the entire knowledge base. The initial user focus F is the universe, i.e. all the items in the information base. In the simplest case, the user selects a concept C in the taxonomy and zooms over it. The *zoom* operation changes the current state in the following way:

1. Concept C is used to refine the current *user focus* F , which becomes $F \cap \text{items}(C)$. Items not in the focus are discarded.
2. The tree representation of the taxonomy is modified in order to summarize the new focus. All and only the concepts related to F are retained and the count for each retained concept C' is updated

8 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:

www.igi-global.com/chapter/dynamic-taxonomies-for-intelligent-information-access/112829

Related Content

Reasoning on vague ontologies using rough set theory

(). *International Journal of Rough Sets and Data Analysis* (pp. 0-0).

www.irma-international.org/article//288522

Nth Order Binary Encoding with Split-Protocol

Bharat S. Rawal, Songjie Liang, Shiva Gautam, Harsha Kumara Kalutarage and P Vijayakumar (2018).

International Journal of Rough Sets and Data Analysis (pp. 95-118).

www.irma-international.org/article/nth-order-binary-encoding-with-split-protocol/197382

Technology Integration in a Southern Inner-City School

Molly Y. Zhou and William F. Lawless (2015). *Encyclopedia of Information Science and Technology, Third Edition* (pp. 2609-2617).

www.irma-international.org/chapter/technology-integration-in-a-southern-inner-city-school/112677

Glycoinformatics and Glycosciences

Anita Sarkar and Serge Pérez (2015). *Encyclopedia of Information Science and Technology, Third Edition* (pp. 414-425).

www.irma-international.org/chapter/glycoinformatics-and-glycosciences/112352

The Consequences of New Information Infrastructures

(2012). *Perspectives and Implications for the Development of Information Infrastructures* (pp. 175-195).

www.irma-international.org/chapter/consequences-new-information-infrastructures/66262