

On the Complexity of Shared Conceptualizations

Gonzalo A. Aranda-Corral¹, Joaquín Borrego-Díaz², and Jesús Giráldez-Cru³

¹ Universidad de Huelva, Department of Information Technology,
Crta. Palos de La Frontera s/n, 21819 Palos de La Frontera, Spain

² Universidad de Sevilla, Department of Computer Science and Artificial Intelligence,
Avda. Reina Mercedes s/n. 41012 Sevilla, Spain

³ Artificial Intelligence Research Institute (IIIA-CSIC)
Campus Universidad Autónoma de Barcelona, Barcelona, Spain

Abstract. In the Social Web, folksonomies and other similar knowledge organization techniques may suffer limitations due to both different users' tagging behaviours and semantic heterogeneity. In order to estimate how a social tagging network organizes its resources, focusing on sharing (implicit) conceptual schemes, we apply an agent-based reconciliation knowledge system based on Formal Concept Analysis. This article describes various experiments that focus on conceptual structures of the reconciliation process as applied to Delicious bookmarking service. Results will show the prevalence of sharing tagged resources in order to be used by other users as recommendations.

1 Introduction

The availability of powerful technologies for sharing information among users (social network members) empowers the organization of social resources. Among them, collaborative tagging represents a very useful process for users that aim to add metadata to documents, objects or, even, urls.

As with other social behaviours, tagging shows advantages but also deficiencies, e.g. semantic heterogeneity. Projects like *Faviki* (<http://www.faviki.com>) or *CommonTag* (<http://commontag.org>) attempt to resolve these deficiencies. Within the network, and also based on user preferences, different tagging behaviours exist that actually obstruct automated interoperability. Although solutions exist that assist the user's folksonomy (tag clouds, tools based on related tag ideas, collective intelligence methods, data mining, etc.), personal organization of information leads to implicit logical conditions that often differ from the global interpretation of these conditions. Tagging provides a manner of weak organization for information that, although useful, is mediated by the individual user's behaviour. In order to make the concept of semantic heterogeneity explicit,

* Supported by TIN2009-09492 project of Spanish Ministry of Science and Innovation and *Excellence project* TIC-6064 of *Junta de Andalucía* cofinanced with FEDER funds.

we use Formal Concept Analysis (FCA) [5]. FCA is a mathematical theory that, applied to tagging systems, results in explicit sets of concepts that users manage by tagging, thereby organizing information into structured relationships.

As is argued in [6], tagging is essentially about sensemaking, a process where information is categorized, labeled and, most importantly, through which meaning emerges [8]. Even in a personal tagging structure, concept boundaries and categories are vague, so some items can be doubtfully labeled. Finally, users also use tagging task for their own benefit, but nevertheless they contribute usefully to the public good [6]. Therefore, it seems interesting to apply concept mining technologies to facilitate semantic interoperability. Since the users' tagging reflects their own set of concepts about documents, tag-driven navigation among different resources could be insufficient due to semantic heterogeneity. Thus, to ensure an efficient use of another user's tag sets, some thought must be given to tags in order to achieve some consensus (also using FCA based tools), which allows us to navigate between different conceptual structures. In this scenario, it could be very important to attempt to delegate these tasks to intelligent agents. In [2], an agent-based knowledge conciliation method is presented.

The aim of this paper is to show how a Multiagent System (MAS) can be applied to shape the complexity of users' conceptual structures into a social bookmarking service, by comparing the *resource sharing* relationship among users against the *tagging sharing* relationship between users. The first relationship comprises a complex network where semantic similarities could be weak, while one expects that the second allows us some understanding about semantic interoperability based on tags and achieved by conciliation. The paper aims to show the prevalence of semantic similarity (knowledge conciliation) in *tagging sharing* relation.

The following paper is organized as follows. Section 2 is devoted to the introduction of FCA. Section 3 reviews original agent-based reconciliation, which is applied in this paper. Section 4 describes the relational structure of tagging in Delicious. Sect. 5 provides a specific implementation of knowledge reconciliation. Section 6 presents the experiments and some results. Finally, Sect. 7 discusses some conclusions.

2 Formal Concept Analysis

Convergence between Mobile Web 2.0 and Semantic Web will depend on the specific management of ontologies. Ontologies and tags/folksonomies must be reconciled in these kinds of projects. A useful bridge between these two kinds of knowledge representation could be *Formal Concept Analysis* [5]. According to Wille, FCA mathematizes the philosophical understanding of a concept as a unit of thought, composed by the *extent* and the *intent*. The extent covers all objects belonging to the concept, while the intent comprises all of the common attributes valid for all the objects under consideration. FCA also allows us to compute concept hierarchies from data tables.

The process of transforming data into structured information by means of FCA starts from an entity called *Formal Context*. This formal context is a tuple



Fig. 1. Formal context and associated concept lattice and Stem Basis

$M = (O, A, I)$ composed of two sets, O (objects) and A (attributes), and a relation $I \subseteq O \times A$. Given $X \subseteq O$ and $Y \subseteq A$, a derivative operator can be defined such as

$$X' := \{a \in A \mid oIa \text{ for all } o \in X\}, \quad Y' := \{o \in O \mid oIa \text{ for all } a \in Y\}$$

From this, a definition of (formal) concept can be obtained as a pair (X, Y) which holds $X' = Y$ and $Y' = X$. If we define the subconcept relation, $C_1 \subseteq C_2$ if $O_1 \subseteq O_2$, a hierarchy among concepts can be obtained and represented as a lattice.

Finally, logical expressions in FCA are *implications between attributes*, a pair of sets of attributes, written as $Y_1 \rightarrow Y_2$. This expression holds in M if for all $o \in O$, its derivative set, $\{o\}'$, models $Y_1 \rightarrow Y_2$, and it is said that $Y_1 \rightarrow Y_2$ is an *implication* of M . A set \mathcal{L} of implications is a (implication) basis, for M , if \mathcal{L} is complete and non-redundant. Also, FCA defines a method to calculate an implication basis [5], which is called Stem Basis. It is important to note that the Stem Basis is only a particular case of implication basis, any other implication basis could be used as well. SB will be used as set of rules in production systems for reasoning (as in [2]). This rules (implication) support can be defined as the number of objects that contain all attributes Y_1 and hold the implication. Based on this property, a variant of implicational basis is defined, called Stem Kernel basis (SKB), the SB's subset where support of each rule is greater than zero.

To illustrate these three entities -formal context, concept lattice, and Stem Basis- an example based on a living being is depicted in fig. 1, left, center, and right, respectively.

2.1 Tagging, Contexts and Concepts

There are several limitations to collaborative tagging in sites such as Delicious. The first is that a tag can be used to refer to different concepts, i.e. there is a context dependent feature of the tag associated with the user. This dependence -called "Context Dependent Knowledge Heterogeneity" (CDKH)- limits both the effectiveness and adequacy of collaborative tagging. The second is the Classical Ambiguity (CA) of terms, inherited from natural language and/or the consideration of different "basic levels" among users [6]. CA would not be critical when users work with urls (content of url induces, in fact, a disambiguation of terms because of its specific topic). In this case, the contextualization of tags in a graph structure (by means of clustering analysis) distinguishes the different terms associated with the same tag [4]. However, CDKH is associated with

concept structures that users do not represent in the system, but that FCA can extract. Thus, navigation among concept structures of different users faced with CDKH. So, the use of tagged resources for automatic recommendation is not advisable without some kind of semantic analysis. More interesting is the idea of deciphering the knowledge that is hidden in user tagging to understand their tagging behaviour and its implied meaning. In sites such as Delicious, CDKH is the main problem, because tags perform several functions as bookmarks [6].

3 Agent-Based Reconciliation

Users's Knowledge Conciliation aims to exploit an important benefit of the Web 2.0, namely information and knowledge sharing. A potential threat is that semantic techniques are adapted to each user. Over time, the user's knowledge can vary a great deal, and this difference could create knowledge incompatibility issues. In order to navigate through the set of tags and documents from different users, SinNet¹ has delegated this process to agents in order to make these different conceptualizations compatible. A agent-based conciliation algorithm was presented in [2]. It is based on the idea that conceptual structure associated with tags gives more information about the user's tagging. The algorithm runs in six steps:

- 1. Agent Creation:** It starts creating two Jade² agents, passing through agent names and SinNet data as parameters.
- 2. Each Agent Then Builds Its Own Formal Contexts and Stem Basis**
- 3. Initializing Dialogue Step:** The agent executes tasks related to communications: It sends its own language (attribute set) to the other agent, and also prepares itself to receive the same kind of messages from the other agent.
- 4. Restrictions of Formal Contexts:** After this brief communication, each agent creates a new (reduced) set of common attributes, and with them a new context to which are added all of the objects from the original context, along with the values and attributes of the common language.
- 5. Extraction of the Production System.** (Stem Basis) for the new contexts.
- 6. Knowledge Negotiation between Agents:** Agents establish a conversation based on objects, accepting them (or not) according to their tag set and their own Stem Kernel Basis: if the object matches the rules, it is accepted, if not the production system is applied, considering the object's tags as facts, getting the answer (new facts which should be added in order to be accepted as a valid object) that is added to the object and re-sent to the other agent to be accepted.

¹ <http://www.semanticville.org/sinnet/>

² <http://jade.tilab.com>

Once this process is completed, the agents will achieve a common context. So, they can extract new concepts and suggestions from a common context, and therefore, a shared conceptualization.

4 Delicious Bookmarking Service

We have chosen the bookmarking service Delicious (<http://www.delicious.com/>) due to its large volume of data. In Delicious, objects are web links (urls), and attributes are tags. Users save their personal web links tagged with their personal tags. But several users may share common objects (with different attributes for each one), or common attributes (tagged in different links). The structure and dynamics of tagging with Delicious have been extensively analyzed [6]. Because of limited computing capacity, certain reduction operations must be performed in order to ensure the normal functioning of the solution presented in this paper. Therefore, a subset of public Delicious data has been extracted, in which all the links are tagged with the tag *haskell*, and saved in a private database (DB) used to drive experiments.

The process of obtaining this data is achieved through a query by tag (*haskell*), and the extraction of the associated results content: link, user, and others tags, which have been saved in the DB. Thereafter, we optimize this data. For example, one of the optimization operations achieved consisted of simplifying equal and equivalent links that have different registers in DB. Our DB is composed of 4259 users, 3028 links, 2427 tags, and 45079 tuples of {user, link, tag}. Data extraction was performed on March 1st, 2011. This data set has a volume large enough to expect significant results. However, this set of data does not encompass all the links related to the *haskell* tag, instead only the first query results.

4.1 The Relational Structure of Tags

In order to estimate the complexity of the relationships among tags of data source, a graph was generated, in which nodes appear as tags, which were interconnected by weighted edges, whose weight represents the amount of links commonly shared according to a Delicious user. To understand the structure of the graph and the number of relevant tags, some simplifications have to be made.

Fig. 2 shows data resulting from *semantic communities* computing (using the method [3]), which is a simplified graph. This graph shows 5 different communities, demonstrating that tags of a same community are very interconnected, unlike tags of different communities, which display little connectivity. In the graph, each node is characterized by its color (determining the community it belongs to), its size (scaled according to its degree), and by the width of its edges (scaled according to the weight of the edge). Finally, only the most relevant nodes (27) and edges (138) are shown - accordingly measured by their importance in terms of degree and weight, respectively.

attributes number, in descending order. Additional methods are equally needed to verify that a pair of users is only referenced in one of their request queues. Further experiments use the number of common objects of a pair of users as the threshold in order to compare results from both executions.

2. Negotiation: User agents must execute a dual behaviour in order to perform the negotiation process: sending and receiving requests. This negotiation establishes a very simple method to decide when a pair of users starts the reconciliation process. Each user is only allowed to perform one reconciliation process at a time. Furthermore, received requests have priority over the sent ones. Two possible states for each user are defined: *free*, if it is not performing any reconciliation at the moment, and otherwise, *busy*. As such, only free users may send or receive requests. On one hand, every user sends proposals to the user having the highest priority in its request queue. If it receives a response, the reconciliation process with the addressee starts. Should this not be the case, it reiterates with the user having the next highest priority. On the other hand, every free user accepts any incoming proposals, even if it has already sent another proposal, which will be cancelled by timeout. The following conditions ensure that all of the conciliations will be processed: their number is finite, and there is always free users ready to accept new conciliations, reducing the number of unsolved processes. When starting a reconciliation, user's state switches from free to busy.

3. Reconciliation: The algorithm presented in section 3 is used to calculate the common knowledge between two users. The steps 1 and 2 (user's concept lattice and SB) are executed only once, when the user runs it for the first time. The rest of the steps (3-6), are executed each time the user runs the algorithm. The obtained common knowledge, a formal context with objects and common attributes, is stored in the DB. Both users switch from busy to a free state.

4. Finalization: When a user's request queue becomes empty, its behaviour is limited to receiving incoming proposals. However, if all the users' request queues are empty, no proposal is received by any of them. Therefore, this situation requires that the execution stops. The *control* agent is used to manage it. It is informed by every user when its request queue becomes empty. When all the users have completed this action, the control agent stops the MAS execution.

6 Experiments

Different experiments have been conducted with data described in section 4 using several criteria. The first criterion is setting a threshold of common attributes (tags) between users. The second criterion is setting a different threshold of common objects (urls). In both cases, the threshold is a necessary condition of a minimum number of attributes or objects that two users must have in common in order to execute the reconciliation algorithm. For each executed reconciliation process, a common knowledge is obtained. This knowledge is a formal context where the attributes are common to both users, and objects belong either to one of them, or both. In this way, the global result is a set of reconciled contexts.

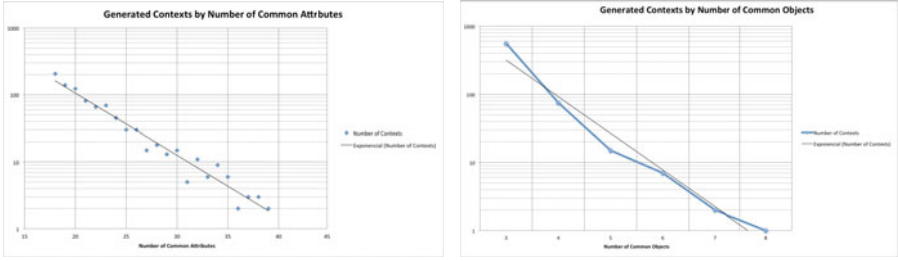


Fig. 3. Contexts generated by number of common attributes or objects

The results obtained for both experiments using numerical and graphic representations are presented below. In order to do so, the results have been measured with five parameters for a fixed value of the threshold. They are the number of contexts obtained (there are as many contexts as number of executed reconciliation processes), and the average values of objects, attributes, concepts and implications per context. Finally, both experiments are compared.

Reconciliation from Common Attributes: In this experiment, the threshold value is set to 18. It implicates that two users having a language³ size greater or equal to 18, reconcile their knowledge. It is assumed that users having a number of common attributes less than 18 do not share a relevant amount of information. In fig. 3 (left), the graphics are plotted in logarithmic scale.

A total of 908 contexts were obtained, with an average value of 44.18 objects, 18 attributes, 6.91 concepts, and 17.20 implications per context. As the threshold value increases, the number of generated contexts decreases exponentially. However, the four average values tend to increase. Although the number of contexts is smaller, they are semantically better, since the two users generating these contexts share more information. In this DB, the maximum number of common attributes is 64. Such a threshold value results in one matching context. It is concluded that one pair of users share a minimum of 64 common attributes. In this context, 138 objects, 38 concepts, and 114 implications are obtained.

Reconciliation from Common Objects: In the second experiment, the threshold value is set to 3. The implication is that two users having a set of common objects with size greater or equal than 3 reconcile their knowledge. As previously mentioned, it is assumed that sharing less than 3 objects is not relevant for the purpose of this study. In fig. 3 (right), the results are represented. This case shows a total of 663 contexts, with an average value of 33.55 objects, 9.41 attributes, 6.09 concepts, and 79.75 implications per context. As the threshold value increases, the number of obtained contexts also decreases, but in this case, more than exponentially. The maximum number of common objects is 11, which is very small: We obtain 98 objects, 37 attributes, 29 concepts, and 56 implications.

³ The language between two users is the set of common tags that both of them use, independent of whether or not these tags have been used in different urls or not.

6.1 Results

The results draw the conclusion that common attributes criterion is better than common objects criterion. On one hand, the decrease in generated contexts is higher when using common objects rather than common attributes. In the first case, this decrease is higher than exponential (more curved than an exponential line). On the contrary, the second case shows an exponential progression. On the other hand, the *semantic* validity of the generated contexts, measured along with their average values, is higher using attributes rather than objects. In the first case, average values increase linearly. It is then thought that the higher number of common attributes, the more reconciled context. Unlike the first case, the second shows a constant function from a certain value of the number of common objects. It seems that the validity of the generated contexts does not depend on the number of common objects.

In conclusion, previous results lead us to think that the common attributes criterion separates more effectively the sample of generated contexts. Indeed, despite the fact that it returns a smaller amount of contexts, increasing the threshold value leads to results *semantically* better. Therefore, it is a good measurement of the semantic similarity of two users.

7 Conclusions and Future Work

The experiments described in this paper show the prevalence of semantic techniques (tags) in resource sharing when users aim to exploit knowledge organization from other users in Delicious as a recommendation source. Although this result seems evident, Web 2.0 shows several examples where url sharing by social networks represent a powerful method for information diffusion (e.g. Twitter).

Therefore, we have empirical evidence that semantic similarity between users is better supported by using the method of reconciling the knowledge among users that have a large set of common attributes, rather than any other method. One of our lines of research is the intensive application of definability methods based on completion [1] in order to enrich the bookmarking system and to facilitate the reconciliation.

References

1. Alonso-Jiménez, J.A., Aranda-Corral, G.A., Borrego-Díaz, J., Fernández-Lebrón, M.M. and Hidalgo-Doblado, M.J. Extending Attribute Exploration by Means of Boolean Derivatives. In: Proc. 6th Int. Conf. on Concept Lattices and Their Applications. CEUR Workshops Proc., vol. 433 (2008)
2. Aranda-Corral, G.A., Borrego-Díaz, J.: Reconciling Knowledge in Social Tagging Web Services. In: Corchado, E., Graña Romay, M., Manhaes Savio, A. (eds.) HAIS 2010. LNCS (LNAI), vol. 6077, pp. 383–390. Springer, Heidelberg (2010)
3. Blondel, V.D., Guillaume, J.L., Lambiotte, R., Lefebvre, E.: Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment* (10) (2008)

4. Au Yeung, C.M., Gibbins, N., Shadbolt, N.: Contextualising Tags in Collaborative Tagging Systems. In: Proceedings of the 20th ACM Conference on Hypertext and Hypermedia (2009)
5. Ganter, B., Wille, R.: Formal Concept Analysis - Mathematical Foundations. Springer, Heidelberg (1999)
6. Golder, S., Huberman, B.A.: The structure of collaborative tagging systems. *Journal of Information Science* 32(2), 98–208 (2006)
7. Jäschke, R., Hotho, A., Schmitz, C., Ganter, B., Stumme, G.: Discovering shared conceptualizations in folksonomies. *Journal of Web Semantics* 6(1), 38–53 (2008)
8. Weick, K.E., Sutcliffe, K.M., Obstfeld, D.: Organizing and the Process of Sensemaking. *Organization Science* 16(4), 409–421 (2005)