# Towards a Soft Evaluation and Refinement of Tagging in Digital Humanities

**Gonzalo A. Aranda-Corral, Joaquín Borrego Díaz**
**and Juan Galán Páez**

**Abstract** In this paper we estimate the soundness of tagging in digital repositories within the field of Digital Humanities by studying the (semantic) conceptual structure behind the folksnonomy. The use of association rules associated to this conceptual structure (Stem and Luxenburger basis) allows to faithfully (from a semantic point of view) complete the tagging (or suggest such a completion).

## 1 Introduction

According to Wikipedia, Digital Humanities (DH) is an area of research and teaching at the intersection between computer science and humanities. DH embraces a variety of topics, from on line collections curation to data mining on large cultural data sets, where researchers use tools from Computing as Knowledge Extraction (KE), Machine Learning, Agent-Based Modeling techniques, as well as solutions from the Social Web. In order to bridge the gap between Humanities and Computing methodologies for Knowledge Organization, it is usual to provide humanists with services to self-organize digital content. Often resources are indexed and classified by categories. Also it is interesting to create tagging services for the community of researchers in order to enrich the content and provide a better navigation.

The reason of the success of tagging in the social web is that it does not have any kind of limitation. Tagging has several use cases in the social web [14]: personal information management (navigate through our selected and tagged resources), digital objects tagging helps to share and spread them, or even to improve user experience

G.A. Aranda-Corral
Departamento de Tecnologías de la Información, Universidad de Huelva,
Crta. Palos de La Frontera S/n, 21819 Palos de La Frontera, Spain

J.B. Díaz · J.G. Páez (✉)
Departamento de Ciencias de la Computación e Inteligencia Artificial,
Universidad de Sevilla, Avda. Reina Mercedes S/n, 41012 Sevilla, Spain

in e-commerce platforms. This allows the user take advantage of its *personomy* as well as other users personomies: an user can access resources uploaded and/or tagged by other users [11]. In the case of collaborative tagging, the full set of resources and tags represents a folksonomy which weakly represents an implicit ontology on community's knowledge. Tagging is an activity which produces folksonomies (inducing consensual vocabularies for the community), that can be understood as a kind of *emergent ontology* which facilitates the organization and navigation.

In general terms, these ontologies suffer of a number of limitations and deficiencies. On the one hand, since the vocabulary has user-dependent intentionality, semantic heterogeneity occurs: some tags represent distinct features for distinct users. Semantic heterogeneity is an intrinsic problem of tagging which prevents the user from exploiting other user's tagging with reliability. Another major drawback of existing social tagging systems is that social tags are used as keywords in keyword-based search. They focus on keywords and their interpretation by humans rather than on computer interpretable semantic knowledge [10]. In the case of DH it is usual to take into account that users share tag interpretation (to a certain extent). Despite these limitations, shared tagging is a potential solution to provide the community with semantically organized knowledge. However, this knowledge is not machine processable and the semantic heterogeneity is a common problem in multi-topic tagging services.

In [8] tagging is described as a task providing resources with sense and aims to categorize resources producing emergent meaning [15]. As a consequence of this, an individual tagging will not be really useful as a public one. Objects tagging made by the community can show the same problem, although in a different scale. However, this problem can be solved by means of Collective Intelligence: when the community tags collectively, they tend to unify the use of tags. Thus the most common tags set associated to an object provides a collective description of a certain concept [6]. In fact, these collective tagging are useful to build recommender systems [5].

The aim of this paper is to propose an (soft) estimate of the soundness of existent tagged digital objects in repositories, as well as to propose a rule set for the automatic refinement of existent tagging (semantics-based). The idea consists in estimating the topological structure of a conceptual network extracted from the tagging system by using Formal Concept Analysis, a mathematical theory which also provides reasoning tools useful for this second goal. The proposed methodology is applied to two tagged repositories relevant in Digital Humanities.

## 2 Formal Concept Analysis for Tagging Services

Formal Concept Analysis (FCA) [7] provides powerful semantic tools for classification, data mining and KE and Discovery (KD). Among these tools particularly interesting are concepts extraction and organization, and implication basis. The last one, represents a sound approach to rule extraction for classification. This task is a

significant issue in KD where FCA applications in the Soft Computing field have been implemented (see for example [16]).
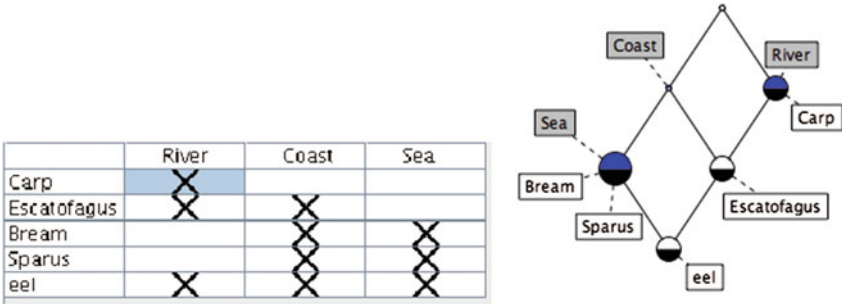


| | River | Coast | Sea |
|---|---|---|---|
| Carp | X | | |
| Escatofagus | X | X | |
| Bream | | X | X |
| Sparus | | X | X |
| eel | X | X | X |

**Fig. 1** Formal context on fish, and its associated concept lattice

A *formal context* $M = (O, A, I)$ consists of two sets, $O$ (objects) and $A$ (attributes), and a relation $I \subseteq O \times A$. Finite contexts can be represented by a 1-0-table (identifying $I$ with a boolean function on $O \times A$). Given $X \subseteq O$ and $Y \subseteq A$, it defines $X' = \{a \in A \mid oIa$ for all $o \in X\}$ and $Y' = \{o \in O \mid oIa$ for all $a \in Y\}$

The classical method for defining a "concept" is actually twofold: The concept is defined *extensionally* by some set of objects that are instances of that concept. The concept is defined *intensionally* by a property that all the instances have in common but that is not possessed by any of the remaining objects. FCA mathematizes this philosophical understanding as a unit of thoughts composed of two parts: the extent and the intent. The extent covers all objects belonging to the concept, while the intent comprises all common attributes valid for all the objects under consideration [7]. The main goal of FCA is the computation of the set of concepts associated with the context.

A (formal) concept is a pair $(X, Y)$ such that $X' = Y$ and $Y' = X$. The set of concepts of a context given $M$, $CL(M)$, can be endowed with the lattice structure by means of the "subconcept" relationship [7]. For example, the concept lattice from the formal context of fishes of Fig. 1, left (attributes are understood as "live in") is shown in Fig. 1, right. Each node is a concept, and its intension (or extension) is formed by the set of attributes (or objects) included along the path to the top (or bottom). For example, the bottom concept ($\{eel\}, \{Coast, Sea, River\}$) is the concept *euryhaline fish*.

Knowledge Bases (KB) in FCA are formed by *implications between attributes*. An implication is a pair of sets of attributes, written as $Y_1 \rightarrow Y_2$. We say that the implication is true with respect to $M = (O, A, I)$ according to the following definition: A subset $T \subseteq A$ *respects* $Y_1 \rightarrow Y_2$ if $Y_1 \not\subseteq T$ or $Y_2 \subseteq T$. $Y_1 \rightarrow Y_2$ is said to hold in $M$ ($M \vDash Y_1 \rightarrow Y_2$ or $Y_1 \rightarrow Y_2$ is an implication of $M$) if for all $o \in O$, the set $\{o\}'$ respects $Y_1 \rightarrow Y_2$.

**Definition 1** *Let $\mathcal{L}$ be a set of implications and L be an implication.*

1. *L follows from $\mathcal{L}$ ($\mathcal{L} \vDash L$) if each subset of A respecting $\mathcal{L}$ also respects L.*
2. *$\mathcal{L}$ is complete if every implication L*

$$M \vDash L \Rightarrow \mathcal{L} \vDash L$$

3. *$\mathcal{L}$ is non-redundant if for each $L \in \mathcal{L}$, $\mathcal{L} \setminus \{L\} \nvDash L$.*
4. *$\mathcal{L}$ is a (implication) basis for M if $\mathcal{L}$ is complete and non-redundant.*

A particular basis is the *Duquenne-Guigues* or so called *Stem* Basis (SB) [9]. In order to work with formal contexts, stem basis and association rules, the Con-exp[1]software has been selected. The reasoning system we use is a production system described in [3]. Initially it works with SB, and the entailment is based on the following result (see [3] for details):

**Theorem 1** *Let $S$ be a basis for M and $\{A_1, \dots, A_n\} \cup Y \subseteq A$. The following statements are equivalent:*

1. *$S \cup \{A_1, \dots A_n\} \vdash_p Y$ ($\vdash_p$ is the entailment by means of a production system).*
2. *$S \vDash \{A_1, \dots A_n\} \rightarrow Y$*
3. *$M \vDash \{A_1, \dots A_n\} \rightarrow Y$.*

In conditions of above definition, let define

$$S[\{A_1, \dots, A_n\}] := \{a \in A \ : \ S \cup \{A_1, \dots A_n\} \vdash_p a\}$$

In FCA, association rules are also implications between sets of attributes. Confidence and support are defined as usual in data mining. The analogous to Stem Basis for association rules is the Luxenburger basis [12]. The reasoning system for SB can be adapted for reasoning with Luxemburger basis [3]. Recall that $Y$ is closed if $Y'' = Y$.

**Definition 2** *Let be $M = (O, A, I)$ a formal context and $Y, Y_1, Y_2 \subset A$.*

- *Given $Y_1, Y_2$ closed, we denote $Y_1 \prec Y_2$ if there is not Y closed such that $Y_1 \subset Y \subset Y_2$.*
- *The support of an attribute set $Y \subseteq A$ is $supp(Y) = |Y'|$.*
- *The support of an implication $L = Y_1 \rightarrow Y_2$ is $supp(L) = |(Y_1 \cup Y_2)'|$*
- *The confidence of L is $conf(L) = \dfrac{supp(Y_1 \cup Y_2)}{supp(Y_1)}$*

**Definition 3** *Given $\gamma$ and $\delta$, the Luxenburger basis of a context M with confidence $\gamma$ and support $\delta$, denoted by $\mathcal{L}(M, \gamma, \delta)$, is*

$$\mathcal{L}(M, \gamma, \delta) := \{L : Y_1 \rightarrow Y_2 \mid Y_1, Y_2 \text{ closed, } Y_1 \prec Y_2, \ conf(L) \geq \gamma, \ sup(L) \geq \delta\}$$

---

[1]http://sourceforge.net/projects/conexp/.

Implications from the Luxenburger basis can be interpreted as association rules from classic data mining, and therefore they allow reasoning under uncertainty. Conexp software provides association rules (and their confidence) associated to formal contexts. The subset of implications from the Luxenburger basis having confidence equal to one (those which are always true within the context) are the same than in the Stem Basis.

For the example from Fig. 1, the basis $\mathcal{L}(M, 0.8, 5)$ and $\mathcal{L}(M, 0.5, 2)$ are[2]

```
1 < 3 > Sea =[100%]=> < 3 > Coast;        1 < 3 > Sea =[100%]=> < 3 > Coast;
2 < 5 > { } =[80%]=> < 4 > Coast;         2 < 5 > { } =[80%]=> < 4 > Coast;
                                          3 < 4 > Coast =[75%]=> < 3 > Sea;
                                          4 < 3 > River =[67%]=> < 2 > Coast;
                                          5 < 5 > { } =[60%]=> < 3 > River;
                                          6 < 2 > River Coast =[50%]=> < 1 > Sea;
```

In general terms, when applying FCA to tagging systems, it is necessary to adapt its environment (formed by resources, tags, and users) to the format required by formal contexts. In this case our aim is to analyze the global structure of the whole tagged repository, thus it is not necessary to take into account which user tagged what resource. The general methodology to apply FCA on tagging is to consider tagged items as objects of the formal context and its tags as attributes on those objects. In this way, a formal context is associated to a folksonomy without taking users into account. Once the context is built, the associated concept lattice $CL(M_{\mathbb{F}})$ can be extracted. This concept lattice represents a concept hierarchy on the universe of the given folksonomy.

## 3 Meaning-Free Tagging Evaluation

In this section it is shown how to evaluate the semantic suitability of folksonomies. It should be noted that this analysis must be independent with respect to the topics of the repository and the field of study it belongs to. Therefore, any methodology used with this aim should take into account, only from a structural point of view, the structure of concepts $CL(M_{\mathbb{F}})$ obtained by means of the process formerly described.

An important feature in semantic networks is the degree distribution given by the connectivity of its nodes (concepts in this case, related by $\prec$), which have been deeply studied. It is expected that concept networks sharing a similar structure could share as well other properties, for instance, those related with its semantics and its suitability as knowledge representation in a certain domain. Therefore it is expected that the topological analysis of the $CL(M)$ shows a big picture of the semantics implicit in the folksonomy itself.

It should be noted that the $CL(M)$ is a complex network of semantic relationships that is not bounded by the self language, as in other semantic networks [13]. That is to say, there are concepts that are not represented by a single language term nor a

---

[2]the format of $L = Y_1 \rightarrow Y_2$ is $<$ supp($Y_1$) $>$ $Y_1$ = conf(L) => $<$ supp($Y_2$) $>$ $Y_2$.

**Table 1** Features of $CL(M_\mathbb{F})$ for case studies. *Density* is $|I|/|O \times A|$ and $< k >$ if the mean degree of the nodes (concepts) of $CL(M_\mathbb{F})$

| | $|\mathbb{O}|$ | $|\mathbb{A}|$ | $|\mathbb{I}|$ | Density | $|CL|$ | $\langle k \rangle$ |
|---|---|---|---|---|---|---|
| Baroque Art | 11.062 | 221 | 74.993 | 3,067 % | 17.817 | 7,949 |
| Gothic Past | 3.246 | 1.781 | 66.432 | 1,149 % | 416.896 | 9.834 |

intelligible definition by the observer. Thus, is a task of the field specialist to interpret such concepts. This feature produces complex networks with extreme structural topology.

A scale-free network is one whose degree distribution follows a power law, at least asymptotically: the fraction $P(k)$ of nodes in the network having $k$ connections to other nodes goes for large values of k as $P(k) \sim ck^{-\gamma}$ where $c$ is a normalization constant and $\gamma$ is a parameter whose value is typically in the range $2 < \gamma < 3$, although occasionally it may lie outside these bounds (as we will see below). It is more common for this behavior to appear from a certain threshold $x_{min}$. The *scale-free* residue of a $CL(M)$ is the set of its nodes whose degree is greater than $x_{min}$ (Table 1).

The analysis of the topology of Concept Lattices is a promising method for addressing the issue raised in the introduction, namely, whether sound qualitative modelizations (in our case, the Concept Lattices) share a similar structure. In [2] the following working hypothesis, called *Scale-Free Conceptualization Hypothesis* (SFCH) is stated, analyzed and experimentally validated:

*Only if the attribute set selected to observe the System is computable, objective, and induces a Concept Lattice that provides a sound analysis of the CS, then its degree-distribution is scale-free.*

This hypothesis (SFCH) has been tested in different experiments, In one of these experiments it was shown that random formal contexts do not respect the SFCH [1, 2]. In the case of the present work, regarding the analysis of folksonomies representing cultural complexity, the statement of the SFCH would be as follows: *A tags set is a suitable knowledge representation for a repository if the conceptual structure that it induces is an scale free network.*

## 4 Analyzing Tagging in DH Repositories

Two DH digital repositories have been chosen as example case study for the proposed methodology: *Baroque Art* from CulturePlex lab[3] and *Gothic Past*.[4]

*The Hispanic Baroque: Complexity in the first Atlantic culture*[5] is a multidisciplinary project carried out by a group of researchers from different universities and financed by the *Social Sciences and Humanities Research Council of Canada*. The

[3]http://baroqueart.cultureplex.ca/ in http://www.cultureplex.ca/ .

[4]http://www.gothicpast.com/.
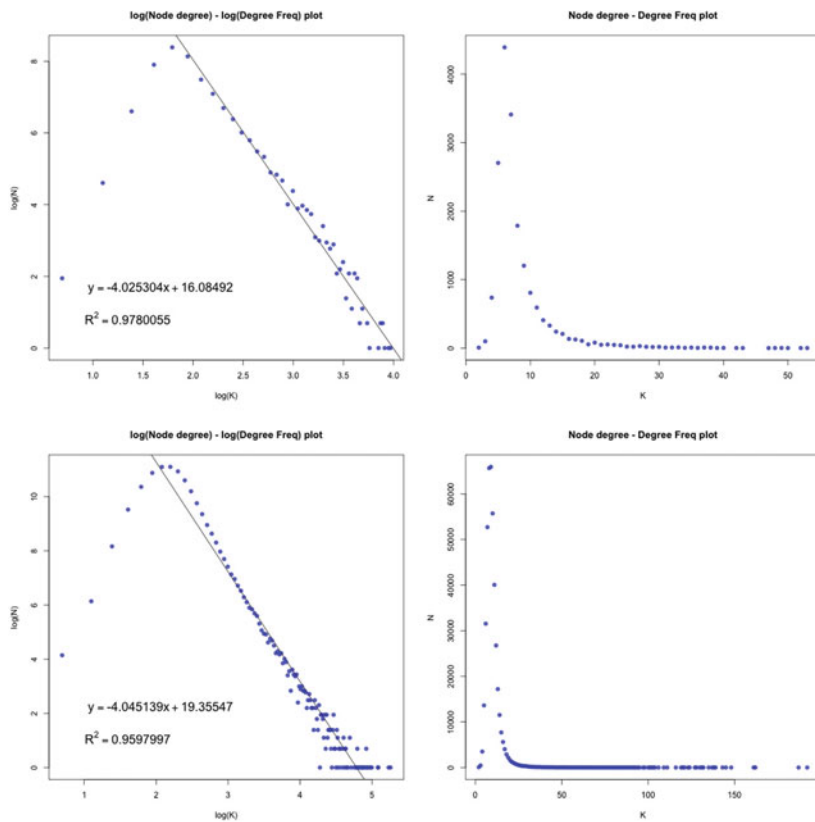
[5]http://www.hispanicbaroque.ca/.

**Fig. 2** Degree distribution of Concept Lattices associated to *Baroque Art* (up) and *Gothic Past* (down)

artwork repository selected, *Baroque Art*, belongs to this project. *Baroque Art* is a repository of artworks tagged by small and closed group of people and following common tagging rules and using as tags a preset common vocabulary (an ontology). The dataset consists of 11.000 artworks and 200 tags approximately.

Figure 2 (up) shows the degree distribution of the conceptual structure extracted from the tagged artworks, which presents an scale free distribution. According to the SFCH that means that the tags set used provides a sound and consistent knowledge representation for the artworks.

*Gothic Past* is a public on line repository for the study of the medieval architecture in Ireland. The repository provides for each element, different information items as pictures, tags, detailed descriptions, etc. In this case, the system allows other users to add new elements to the repository or to modify existing ones. Therefore the number of people involved in this tagging process is higher than in the former repository, possibly leading to more heterogeneous tagging criteria. Figure 2 (down) shows the degree distribution of the conceptual structure associated to the tagged repository on Irish Gothic monuments, which also presents a scale free distribution.

# 5 Luxenburger Basis for Automated Tagging Completion

By considering folksnomies as formal contexts it is possible to use (Luxenburger) Implication basis for suggesting new tags:

**Definition 4** *Let $M_{\mathbb{F}}$ be the context associated to a folksonomy $\mathbb{F}$, $S$ be a basis and $r$ be a resource of the folksonomy (that is, an object of $M_{\mathbb{F}}$)*
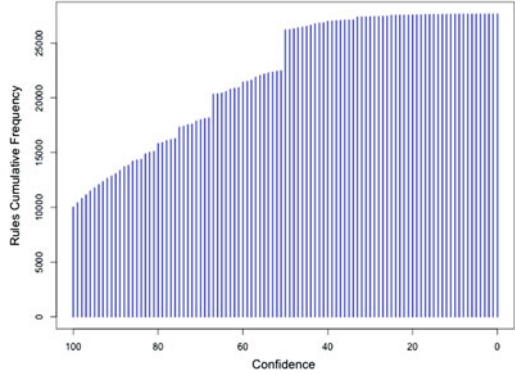
- *The completion tagging of $r$, is $c(r, \mathbb{F}) := S[\{r\}']$*
- *The suggested tagging for $r$ is $s(r, \mathbb{F}) := c(r, \mathbb{F}) \setminus \{r\}'$*
- *The Luxenburger tagging with confidence $\gamma$ and weight $\delta$ is $\mathcal{L}(M, \gamma, \delta)$ is*

$$c(r, \mathbb{F}, \gamma, \delta) := \mathcal{L}(M_{\mathbb{F}}, \gamma, \delta)[\{r\}']$$

- *The suggested Luxenburger tagging with respect to $\mathcal{L}(M, \gamma, \delta)$ is*

$$s(r, \mathbb{F}, \gamma, \delta) := s(r, \mathbb{F}, \gamma, \delta) \setminus \{r\}'$$

**Fig. 3** Distribution of $L(M_{\mathbb{F}}, \gamma, \delta)$ according to rules' confidence, associated to *Baroque Art* repository



**Proposition 1** *The completion tagging does not depend on the basis selected.*

*Proof* Let $r \in O$ and $S_1, S_2$ be two basis. If $a \in S_1[\{r\}']$ then $S_1 \cup \{r\}' \vDash a$ so $S_1 \vDash \{r\}' \to \{a\}$ and thus $M \vDash \{r\}' \to \{a\}$. Therefore $S_2 \vDash \{r\}' \to \{a\}$

Moreover, if the intent of all objects in the context is augmented by applying a Luxenburguer basis, then the implications of this turns to be true within the new context (thus they belong to a basis formed by implications with confidence 1 of the new context):

**Proposition 2** *Let be* $M_{\mathbb{F}}^{\gamma,\delta} = (O, A, I^{\gamma,\delta})$ *where*

$$(o, a) \in I^{\gamma,\delta} \iff a \in \mathcal{L}(M_{\mathbb{F}}, \gamma, \delta)[\{o\}']$$

*Then* $\mathcal{L}(M_{\mathbb{F}}, \gamma, \delta) \subseteq \mathcal{L}(M_{\mathbb{F}}^{\gamma,\delta}, 1, 0)$

## Fray Pedro Machado



**Creator:** Zurbarán, Francisco de
**Dated In:** 1630 - 1632
**Original Location:** Convento de la Merced Calzada
**Current Location:** Real Academia de Bellas Artes de San Fernando (Madrid, Spain)
**It belongs to series:** Retratos de personajes ilustres de la Orden Mercedaria-Zurbarán

**General Description**

[ Objeto ] → Clasificación → Género → Retrato
[ Objeto ] → Clasificación → Temática → Religioso
[ Objeto ] → Clasificación → Tipo → Pintura
[ Objeto ] → Propiedades Físicas → Color → Blanco
[ Objeto ] → Propiedades Físicas → Color → Negro
[ Objeto ] → Propiedades Físicas → Color → Rojo
[ Objeto ] → Propiedades Físicas → Material → Lienzo
[ Objeto ] → Propiedades Físicas → Tamaño: 193 x 122 cm
[ Objeto ] → Propiedades Físicas → Técnica → Óleo

**Suggested tags**

[ Objeto ] → Propiedades Físicas → Color → Café
[ Objeto ] → Clasificación → Género → Santos

**Fig. 4** Suggested tagging for the object $o353$, Zurbaran's artwork *Fray Pedro Machado*, by means of $\mathcal{L}(M_{\mathbb{F}}, 0.5, 30)[\{o353\}']$

Given a collaborative tagging service inducing a folksonomy $\mathbb{F}$, and a resource $r$, the tag set $c(r, \mathbb{F})$ extends the tagging $\{r\}'$ in order to allocate the object (the resource) in the most specific concept (as possible), according to its original tagging. However, $c(r, \mathbb{F}, \gamma, \delta)$ provides suggested tagging with a certain confidence degree. Thus, the user acceptability (or community of users) is important. Figure 3 shows the distribution of $|\mathcal{L}(M_{\mathbb{F}}, \gamma, \delta)|$ for *Baroque Art* repository. It is worthy to note that the tags set for a resource is very small with respect to the set of all tags. Therefore the computing of $\vdash_p$ (with confidence propagation [3]) is very fast. Particularly in the case of *Baroque*, the ontology-assisted tagging makes the basis $\mathcal{L}(M_{\mathbb{F}}, 1, 0)$ to have a relevant size: $|\mathcal{L}(M_{\mathbb{F}}, 1, 0)| = 10.007$ whilst $|\mathcal{L}(M_{\mathbb{F}}, 0.5, 0)| = 22.457$. As example, Fig. 4 shows some suggested tags (in red) for an artwork $o353$ (http://baroqueart.cultureplex.ca/artworks/353/) from *Baroque Art*. The tags belong to $\mathcal{L}(M_{\mathbb{F}}, 0.5, 30)[\{o353\}']$.

# 6 Conclusions and Related Work

Two uses of FCA within DH projects are described: the evaluation of the soundness of Knowledge Organization in tagging services and reasoning with implication basis to augment its tagging. Future work is focused on the use of *attribute exploration* [7] as a web service for accepting new tags (offered as plug-in). This idea (suggested in [4]) could be useful in cases where the repository is complete enough to extract useful knowledge from it, in the form of expert system.

In [11] authors study folksonomies by means of using triadic concepts, by considering the user as responsible of the tag. In our case, tagging in collaborative platforms, those are anonymized. It is also possible to exploit domain ontologies for suggesting tags (see for example [10]). In the first of the presented case studies (Baroque), the main tag vocabulary is provided by an ontology, thus it would be possible to expand or refine suggested tagging. In the second one this is not possible because, to the best of our knowledge, there is not a similar ontology.

The consensus a community can reach on collaborative tagging on a specific topic is different from personal information organization systems as for example Delicious (http://delicious.com/) or Diigo (https://www.diigo.com/). In this case it is necessary to reconcile their knowledge with other users to leverage their information (as in [4] by using FCA also).

## References

1. Aranda-Corral GA, Borrego-Díaz J, Galán-Páez J (2013) Complex concept lattices for simulating human prediction. Sport J Syst Sci Complex 26(1):117–136
2. Aranda-Corral GA, Borrego-Díaz J, Galán-Páez J (2013) On the phenomenological reconstruction of complex systems-the scale-free conceptualization hypothesis. Syst Res Behav Sci 30(6):716–734
3. Aranda-Corral GA, Borrego-Díaz J, Galán J (2011) Confidence-based reasoning with local temporal formal contexts. In Proceedings 11th international conference artificial neural networks conference on advances in computational intelligence - volume part II (IWANN'11). Lecture notes in computer science 6692, pp 461–468 Springer
4. Aranda-Corral GA, Borrego-Díaz J, Giráldez-Cru J (2012) Agent-mediated shared conceptualizations in tagging services. Multimedia Tools Appl 65(1):5–28
5. Carmel D, Roitman H, Yom-Tov E (2010) Social bookmark weighting for search and recommendation. VLDB J 19(6):761–775
6. Halpin H, Robu V, Shepherd H (2007) The complex dynamics of collaborative tagging. In: Proceedings 16th international Conference WWW '07, pp 211–220
7. Ganter B, Wille R (1999) Formal concept analysis - mathematical foundations. Springer, Berling
8. Golder S, Huberman BA (2006) The structure of collaborative tagging systems. J Inf Sci 32(2):98–208
9. Guigues J-L, Duquenne V (1986) Familles minimales d' implications informatives resultant d'un tableau de donnees binaires. Math Sci Humaines 95:5–18
10. Hsu I-C (2013) Integrating ontology technology with folksonomies for personalized social tag recommendation. Appl Soft Comput 13(8):3745–3750

11. Jäschke R, Hotho A, Schmitz C, Ganter B, Stumme G (2008) Discovering shared conceptualizations in folksonomies. J Web Semant 6(1):38–53
12. Luxenburger M (1991) Implications partielles dans un contexte. Math Inf Sci Hum 11:335–55
13. Motter AE, de Moura APS, Lai Y, Dasgupta P (2002) Topology of the conceptual network of language. Phys Rev E 65
14. Smith G (2007) Tagging: people-powered metadata for the social web. First New Riders Publishing, Berkeley
15. Weick K-E, Sutcliffe K-M Obstfeld D (2005) Organizing and the process of sensemaking. Organ Sci 16(4):409–421
16. Jianping Y, Chen L, Wenxue H, Shaoxiong L, Deming M (2015) A new approach of rules extraction for word sense disambiguation by features of attributes. Appl Soft Comput 27:411–419