

Cooperative Authorship Social Network

Giseli Rabello Lopes¹, Mirella M. Moro², Leandro Krug Wives¹, and
José Palazzo Moreira de Oliveira^{1*}

¹ Universidade Federal do Rio Grande do Sul - UFRGS
Porto Alegre, Brazil

`{grlopes,wives,palazzo}@inf.ufrgs.br`

² Universidade Federal de Minas Gerais - UFMG
Belo Horizonte, Brazil
`mirella@dcc.ufmg.br`

Abstract. This paper introduces a set of challenges for developing a dissemination service over a Web collaborative network. We define specific metrics for working on a co-authorship research social network. As a case study, we build such a network using those metrics and compare it to a manually built one. Specifically, once we build a collaborative network and verify its quality, the overall effectiveness of the dissemination services will also be improved.

Key words: Social Networks, Dissemination Systems.

1 Introduction

Web 2.0 is the second generation of communities and services characterized by providing techniques for the personal publication, sharing, collaboration, and organization of information on the World Wide Web. In this perspective, not only the technological and content aspects but also the social interactions and its relational aspects must be considered. In this context, the web-based communities, hosted services, and web applications emerged, including Social Networks.

The *Social Network Analysis* (SNA) is based on the assumption that the relationship's importance between interaction units is a central point to the evaluation and analysis of social interaction. Some fundamental concepts used on SNA include actors and relational ties [1]. *Actors* are social entities that have social linkages modeled by the Social Network (SN). Actors are linked to other actors by *relational ties*.

The increasing interest in researching in SNA was encouraged by the popularization of online social networks, which are very interesting Web applications. Another example of such concepts application is a co-authorship social network representing a scientific collaboration network. In this network, actors represent authors and relational ties represent the relationships between pairs of authors. The presence of at least one co-authored paper between two authors determines a

* This research is partially supported by CNPq (Brazil), and is part of the InWeb research project.

relational tie between them. Some examples of data sources for the construction of this kind of networks are DBLP, Google Scholar, CiteSeer, among others.

The relational tie between authors may help to identify long term collaborations, common research interests, preferred conferences, research groups under formation, among others. Furthermore, as the social ties evolve, new research interests and new collaborations will be identified. Any person who wants to keep updated about such an evolution can be notified of such novel aspects by adding a dissemination service to the social network.

A dissemination service is formed by data producers and consumers. Specifically, consumers subscribe to the service by defining a profile, which is usually composed of different queries. As the producers inject the system with new data, usually through messages, the dissemination service evaluates each message against the profiles. Once there is a match between a profile and a message, the service sends that message to the profile's consumer [2].

The contributions of this paper are twofold. First, we introduce a set of challenges for developing a dissemination service over a Web collaborative network. Then, we tackle the challenges from the SN perspective. Specifically, we present an architecture for such a dissemination service over a collaborative network. The architecture is formed by different layers, from the Web to digital libraries, social network, and the dissemination service. Based on the architecture, we were able to identify research challenges that are innovative to the SN area. We define specific metrics for working on a co-authorship research SN. Then, we build a network using those metrics and compare it to a manually built one. Specifically, once we build a collaborative network and verify its quality, the effectiveness of the dissemination services will also be improved. Therefore, based on such an evaluation, the dissemination service can identify (and recommend) the more pertinent publications as well as identify possible hidden collaboration nets.

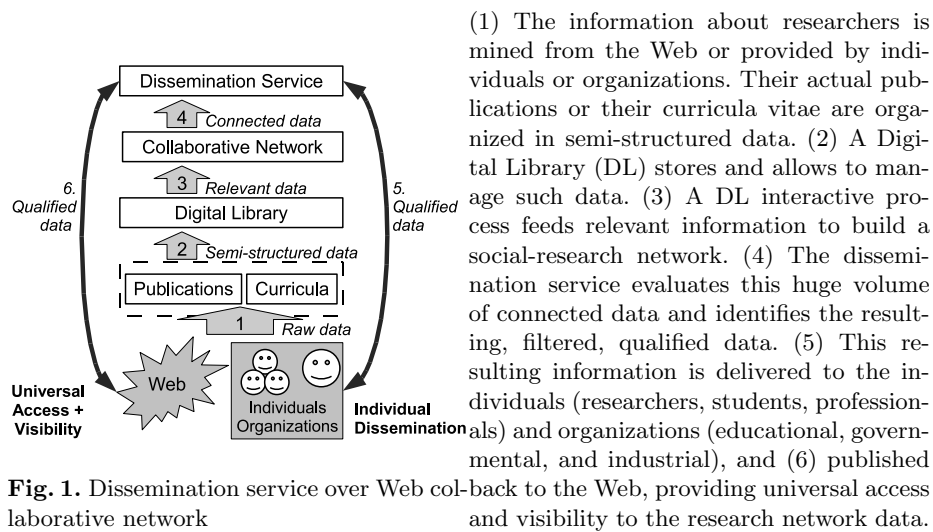
The paper is organized as follows. Section 2 describes the general context of dissemination services and defines the base architecture. Section 3 introduces the metrics to determine the weights of relational ties of a co-authorship Social Network. Section 4 presents a case study that shows the construction of collaboration Social Network. It also evaluates the metrics employed to analyze the SN. Section 5 presents some related work. Section 6 concludes this paper.

2 Dissemination Service in Social Network Context

Content-based dissemination is a form of data delivery that differs from traditional communications since the messages are delivered according to their content rather than the IP address of their destination. There is a continuous stream of messages from data producers to consumers, without any of the human parties having knowledge of the other [2, 3]. This form of communication is widely employed by dissemination services, which may be employed within publish/subscribe systems (pub/sub for short).

In order to clarify how a dissemination service can work on a Web collaborative network, we present a case study based on the academic field. It exemplifies

a service that disseminates new publications and research connections. Specifically, individuals (or organizations) can subscribe to research topics or researcher names, for example. Once a new publication or a new collaboration is detected, this information is disseminated to those subscribers whose keywords match such new data. It is important to notice that not only publications are recommended but also (and more important) new possible cooperations among researchers are identified and suggested. The whole process is composed by six phases, as illustrated in Figure 1. Each step of this process works as follows.



The dissemination service from Figure 1 illustrates tasks with challenges to different Computer Science areas. Specifically, Information Retrieval techniques may be employed along with Data Mining algorithms in order to recover the researchers' data from the Web (1). Moreover, Web Management issues become critical when considering that the data will be extracted from the Web (for example privacy, security, provenance, and credibility). The Digital Library maintenance presents new challenges due to the interactive nature of the framework (2), where individuals and organization will access the data through the dissemination service, and not through the Digital Library interface as usual. Social Network's mechanisms are necessary for defining the collaborative network (3). Then, the challenges appear on the Dissemination service level, which also include Network Management (4). Finally, the actual dissemination and evaluation of data involve Document Management, Distributed Systems, Parallel Computing, Security and Networks as well (5, 6).

It is important to notice that each of those disciplines is complex by nature. Instead of discussing each of such areas, the focus of this paper is on the social networks challenges. Specifically, with the increasing interest in Social Networks, the interaction of the parties (data producers and consumers) within

the dissemination service will soon conquer the spotlight. In social networks, it is important to qualify and quantify how individuals (people and organizations) are connected, how tightly (or loosely) they interact, and what their common interests are. Due to the large volume of data involved and the high complexity of those connections, the development of an automatic mechanism capable of efficiently identifying and analyzing such interactions is imperative.

3 Social Scientific Networks Analysis

Social Networks are based on the assumption of the relationship's importance between interaction units. The weights of the relational ties in a social network aim to measure the importance of the ties between actors. It is necessary to establish approaches to automatically determine these weights based on information available about the actor's relationships.

In this paper, we employ a scientific collaboration network as base example. We present approaches to determine two types of associations namely *Collaboration in Co-authorship* and *Collaboration in Research Areas*. These associations were chosen because they cover certain facets of the relational ties of the collaboration network. According to Newman [4], that studied scientific collaboration networks in which two scientists are considered connected if they have coauthored a paper, this seems a reasonable definition of scientific acquaintance.

3.1 Collaboration-based association - Co-authorship (Ca)

Formally, a Social Network SN of a co-author relationship a is a pair: $SN_a = (N, E)$ where N and E are the set of Nodes and Edges. Each edge $e \in E$ is a tuple of the form $\langle a_i, t, w, a_j \rangle$, where the edge is directed from a_i to a_j , t denotes the type of association between a_i and a_j , and w denotes the weight affected to the association. This weight is a numerical value between 0 and 1. In our approach, the equation 1 determines the *Collaboration in Co-authorship* weight.

$$w_{t_{Ca(a_i \rightarrow a_j)}} = \frac{|a_j co_authorship|}{|a_i author|} \quad (1)$$

where:

- $w_{t_{Ca(a_i \rightarrow a_j)}}$ corresponds to the weight of the recommendation based on the co-author relationship. The weight is different according to the relation direction (the weight in the direction $a_i \rightarrow a_j$ is different than in $a_j \rightarrow a_i$);
- $|a_j co_authorship|$ corresponds to the number of times that the author a_j was a co-author of a paper with author a_i ;
- $|a_i author|$ corresponds to the total number of papers of the author a_i .

In other words, the higher this weight is, the more relevant is the relationship with author a_j to the author a_i . The use of Ca metric implies that there is a graph with 0 or 2 links between two authors. The weights represent the degree of collaboration in co-authorship between the authors. This metric is an asymmetric variant of the *Jaccard Coefficient* and it was applied in the context of Social Networks by other works as [5, 6].

3.2 Collaboration-based association - Research Areas (*Ra*)

In this case, we consider the same definition of Social Network *SN* of co-author relationship (as defined in the previous section). However, each edge $e \in E$ is a tuple of the form $\langle a_i, t, r, w, a_j \rangle$, where the edge is directed from a_i to a_j , t denotes the type of association between a_i and a_j , r denotes the research area associated to the relationship represented, and w denotes the weight affected to the association. This weight is a numerical value between 0 and 1. The equation 2 provides the *Collaboration in Research Areas* weight.

$$w_{t_{Ra}(a_i \rightarrow a_j)} = \frac{|Cr_{research_areas}(a_i, a_j)|}{|research_areas_{a_i}|} \times \frac{|co_authorship_{research_area_rx}(a_i, a_j)|}{|co_authorship_{research_areas}(a_i, a_j)|} \quad (2)$$

where:

- $w_{t_{Ra}(a_i \rightarrow a_j)}$ corresponds to the weight of the recommendation based on the co-author relationship according to research areas. Again, the weight is different according to the relation direction;
- $|Cr_{research_areas}(a_i, a_j)|$ corresponds to the number of research areas in which the authors a_i and a_j published co-authored papers;
- $|research_areas_{a_i}|$ corresponds to the total number of research areas in which author a_i published;
- $|co_authorship_{research_area_rx}(a_i, a_j)|$ corresponds to the number of co-author relationship between authors a_i and a_j in the x area;
- $|co_authorship_{research_areas}(a_i, a_j)|$ is the total number of co-author relationship between authors a_i and a_j in every research areas in which they published together.

The use of *Ra* metric implies that there are $2n$ links between two authors, being that n indicates the number of research areas in which the authors published together. Each link has a direction, a research area and a weight associated. The higher this weight is, the more relevant is the relationship with author a_j to the author a_i in the research area x . In such an approach, we have the idea of collaboration in research areas.

4 Case Study

This paper proposes an approach to construct a social network for collaborative research. The complete work is under development as research project of the In-Web (MCT/CNPq Grant Number 573871/2008-6), the Brazilian National Institute of Science and Technology for the Web. In fact, we have built a collaborative social network based on the publications of the researchers associated to INWeb. The Institute is formed by 27 researchers and their students. All researchers are professors in a major education institution (namely UFMG, UFRGS, UFAM, and CEFET-MG) with graduate program in Computer Science.

4.1 Building the Social Network: Manually and Automatically

Initially, this group of researchers was *manually* analyzed by a specialist. The resulting network can be visualized in Figure 2. This network is used as *baseline*.

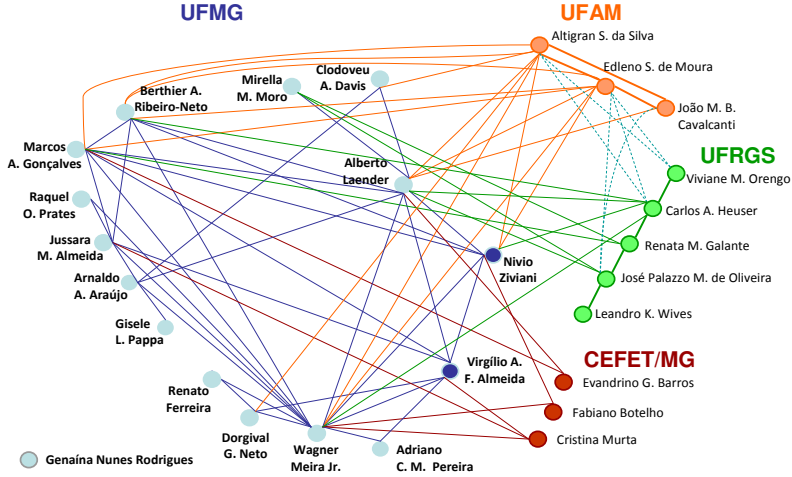


Fig. 2. Manual INWeb Social Network

For validating our metrics, we have implemented a tool to automatically generate a Social Network. This SN was built using information about authors provided by the DBLP digital library. It is important to notice that this library is exported as an XML document. Instead of using the whole dataset, we extracted from the library just the papers written by the considered researchers and published in conferences proceedings and in journals (as elements *inproceedings* and *article*). Such a subset was chosen because this information is significantly important for representing the co-author relationship between authors and, consequently, to determine the research collaborations among them.

The actors of the SN can be chosen and they are a subset of authors with scientific papers indexed by the DBLP. The relational ties between actors are the relationships between pairs of authors. These social ties represent the co-author relationships. The weights of the linkages are determined by equation 1. In that equation, $|a_i \text{ author}|$ corresponds to the total number of papers of the author a_i , and it considers all papers to this author a_i indexed at DBLP, including papers that are not co-authored by authors in the SN who will be graphically presented.

The resultant INWeb Social Network constructed automatically is presented in Figure 3. The data used in this case was collected from the DBLP repository on January 21, 2009. This data gathering process summed up 677,345 authors; 692,431 conference proceedings papers and 432,663 journal articles.

After building them, we compared the two Social Networks: the manually constructed SN (called Manual INWeb SN) and the automatically generated one

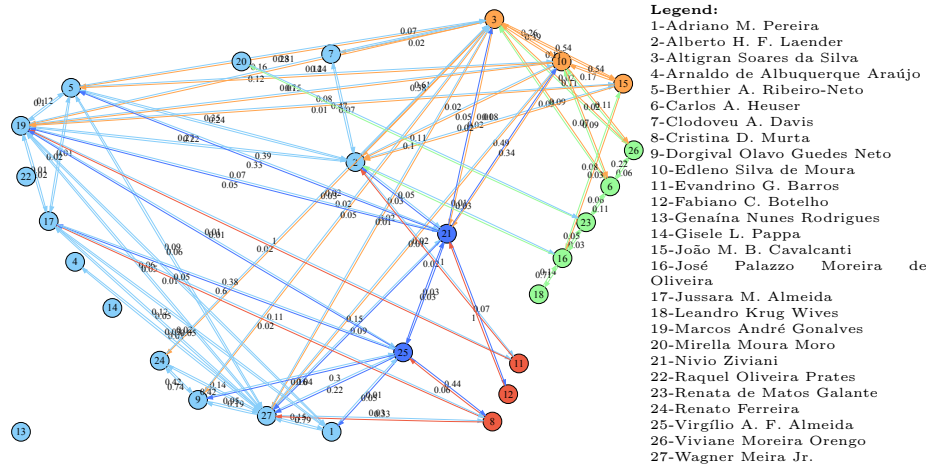


Fig. 3. Automatic INWeb Social Network

(called Automatic INWeb SN). Comparing them against each other, we observed that the Manual INWeb SN covers 93.44% of the Automatic INWeb SN. The Automatic INWeb SN covers 83.82% of the Manual INWeb SN. Furthermore, if we consider that the ideal network (144 edges) is the union between the edges of the Manual INWeb SN (136 edges, considering that each linkage was reciprocal) and the edges of the Automatic INWeb SN (122 edges), we have the following results. The Manual INWeb SN recall is 94.44% and the Automatic INWeb SN recall is 84.72%. The ideal network was considered the union because the Manual INWeb SN was carefully developed by a specialist and the Automatic INWeb SN was based on an occurrence of a co-authorship between two authors for the establishment of the relational ties.

The main goal of this comparative analysis between the two networks was to validate the Social Network constructed automatically by our system using the DBLP dataset. The results obtained demonstrate that the DBLP digital library is a good data source that considerably covers the co-authorship relations in Computer Science, more specifically in Information Systems research area.

4.2 Analysis of the Automatic Co-authorship Network

In this section, we further analyze the Automatic INWeb Social Network. The goal is to use other metrics to understand the properties of the Social Network on this case study. In the next subsections, we present the metrics considered and discuss the results obtained (observation: the results of the metrics were plotted in decreasing order of the values obtained in all graphics and the authors were represented by numbers in the range of 1 to 27 into accordance to the ascending order of the full names (see Legend of Figure 3)).

Clustering Metrics. Clustering is a process that aims to identify subsets or clusters of “similar” elements (or data items). The goal of clustering algorithms

is to create groups that are coherent internally, but clearly different from each other. Thus, elements within a cluster should be as similar as possible; and elements in one cluster should be as dissimilar as possible from elements in other clusters [7]. In order to evaluate the clusters generated by those algorithms, we can employ internal quality measures that require no human intervention, such as cohesion and coupling [8]. *Cohesion* is the average pairwise similarity of elements within the cluster. *Coupling* is the average pairwise similarity of elements in which one element belongs to cluster C and the other does not.

The clustering metrics were adapted for evaluating our case study. We considered each group constituted by an author and all his co-authors as a cluster. For each cluster (each author), we calculated the respective cluster metrics. The similarity values for the metrics calculation are the weights of the relational ties between authors. In our case, the best results will be that whose cohesion and coupling measure high values. Such result is important because each cluster is a subnet of the social network being analyzed.

The cohesion metric was adapted to consider two similarity values between each pair of authors. This was necessary because our SN is represented by a directional graph. The new equation is defined as follows (Equation 3).

$$cohesion(C) = \frac{\sum_{i=1}^{m-1} \sum_{j=i+1}^{m-1} w_{t(a_i \rightarrow a_j)} + w_{t(a_j \rightarrow a_i)}}{m(m-1)} \quad (3)$$

where, m corresponds to the total number of authors in the group considered ($m=1(\text{author})+n(\text{total number of his/her INWeb co-authors})$).

In this case, the similarity values used (w_t) in the calculation were the weights $w_{t_{C_a}}$. Figure 4 presents the cohesion results obtained to each cluster formed by one author and all his INWeb co-authors. The results obtained show the average of importance between all pairs of authors in each cluster considered. The more cohesive groups are those formed by authors with high number of collaborations whose weights indicate a high importance in these co-authorships.

As Figure 4 illustrates, some clusters formed by few authors have the best results. This probably happened because these clusters are formed by young authors whose importance weights in relation to their co-authors are high. Some senior authors formed clusters with low cohesion values. This probably happened because those worked with many co-authors over time and/or have a much larger collaboration (cooperation) network that the one formed by INWeb authors.

Figure 5 presents the results for coupling metric. This graphic plots the authors in x axis and the coupling values obtained for each cluster (formed by the author and his co-authors) in y axis. Equation 4 was used. This metric was evaluated by using the output weights to the author a_i whose cluster C is being analyzed as similarity value. Indeed, C is the cluster formed by an author and his co-authors; m is the number of elements in the cluster C ; and n is the number of elements outside the cluster C belonging to a cluster Q formed by the co-authors of a_i and all co-authors of these co-authors of a_i (including a_i). In this case, the similarity values used in the calculation are the weights $w_{t_{Ca(a_i \rightarrow a_j)}}$ where a_i was

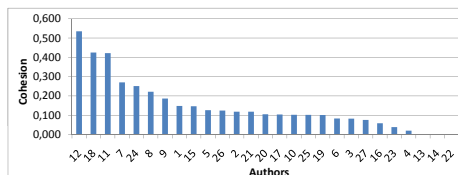


Fig. 4. Cohesion results for INWeb

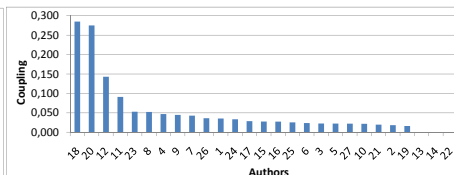


Fig. 5. Coupling results for INWeb

the author been analyzed and a_j varies among each author of the cluster Q .

$$coupling(C) = \frac{\sum_{i,j} sim(c_i, q_j)}{m \times n} \quad (4)$$

Note that the nonzero similarity values are between a_i and his co-authors, and between a_i and a_i himself. On the equation, the weight between the author and himself was considered 1. This shows the coupling among the group of researchers formed by each author and his co-authors. The results show that some young researchers that have “good” publications present high coupling. This probably occurred because such researchers work in more “condensed” groups while the others have a larger network and/or work in several groups.

Complementary Analysis. This subsection presents other analysis performed on the Automatic INWeb Social Network.

First, Figure 6 presents the percentage of INWeb Co-authors in relation of the total Co-authors indexed by DBLP, for each author. This metric prioritizes authors that have high number of his total co-authors within the INWeb Social Network. The results show higher values to the authors that have his co-author relationships represented more significantly by the INWeb partnerships.

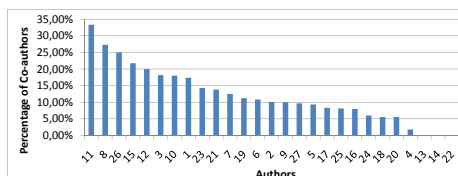


Fig. 6. Percentage of Co-authors

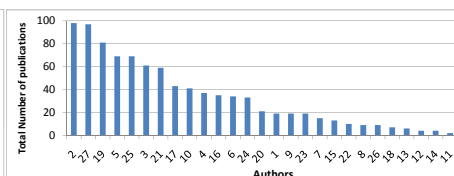


Fig. 7. Total number of publications

Figure 7 presents the total number of publications by author. This metric is presented in order to help to understanding the results. The INWeb Social Network shows that some authors do not have co-author relationship with any INWeb author. However, Figure 7 shows that all authors of INWeb Social Network have at least one publication indexed by DBLP.

Figure 8 shows the average importance of each author to his INWeb co-authors. This metric was calculated according to the equation 5.

$$In_Avg_Imp(a_i) = \frac{\sum_{j=1}^n w_{t_{a_j \rightarrow a_i}}}{n} \quad (5) \quad Out_Avg_Imp(a_i) = \frac{\sum_{j=1}^n w_{t_{a_i \rightarrow a_j}}}{n} \quad (6)$$

where a_i corresponds to the author being analysed, a_j varies among the co-authors of a_i , and n corresponds to the total number of co-authors of a_i in the Social Network being considered.

The graph in Figure 8 plots the authors in x axis and the input average importance values obtained for each author in y axis. For calculating the importance ($w_{t_{Ca(a_i \rightarrow a_j)}}$), it considered the DBLP Social Network (i.e., all publications indexed by DBLP were considered, whether they are co-authored by an INWeb author or not). However, the co-authors considered were only those belonging to the INWeb Network. Figure 8 also illustrates the relative importance of each author to the others. The result shows that the equation prioritizes authors who have a high average importance value to his co-authors. Some authors that have few co-authors but have a meaningful importance value to his co-authors overcame other authors that have a high number of co-authors.

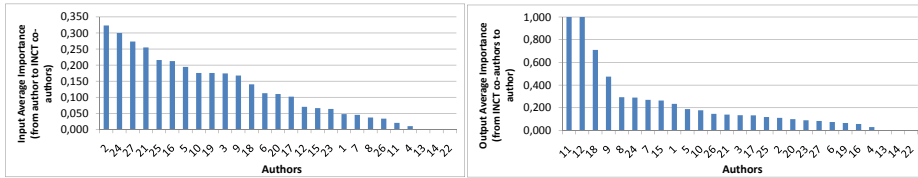


Fig. 8. Input Average Importance (from author to INWeb co-authors) **Fig. 9.** Output Average Importance (from INWeb co-authors to author)

Figure 9 shows the average importance of all INWeb co-authors to each author. This metric was calculated according to the equation 6. This graph plots the authors in x axis and the output average importance values obtained for each author in y axis. This graph illustrates the importance of the other INWeb authors in relation to all collaboration network represented by DBLP SN. The result shows that authors that have a group of co-authors more “condensed” and, sometimes, without interaction with other people outside of this group, will often have higher values of output average importance.

5 Related Work

This section overviews some work related to recommender systems (a type of dissemination system) and social networks.

Weng and Chang [9] propose a recommender method that employs ontologies and the *spreading activation model*. The ontologies are employed for defining user profiles, being the basis to reason about the users’ interests. The spreading activation model is used to search for other influential users in a Social Network

Golbeck et al. [10] present a website that integrates Social Networks on the Semantic Web context and the *trust* concept for the generation of movies' recommendations. The Social Networks then indicate the trust ratings between users by considering the path length between them.

Aleman-Meza et al. [5] define a solution for the *Conflict of Interest* (COI) problem using Social Networks. The goal is to detect COI relationships among authors of scientific papers and potential reviewers of these papers. Moreover, rules are established to determine a possible degree of COI among the authors based on the Social Networks built and the relationship's weights between them.

Jeh et al. [11] propose a measure of structural-context similarity, called *Sim-Rank*. The recommender systems were used as motivation. The base idea of the model is that two objects are similar if they are related to similar objects.

Zaiane et al. [12] explore a Social Network coded within the DBLP database. It considers a new random walk approach to reveal interesting knowledge about the research community and even to recommend collaborations.

Menezes et al. [13] developed a geographical analysis of knowledge production in Computer Science. They analyzed co-authorship Social Networks of the Computer Science area.

Ganev et al. [14] developed a set of tools for building, exploring and querying academic Social Networks. They proposed a measure reputation called *visibility* as an adjusted PageRank applied on the Social Network context.

Our paper is related to all those since it focuses on solutions for Social Networks. However, we presented a case study to clarify how a dissemination service can work on top of a Web collaborative network. We presented an approach to construct a Social Network for collaborative research that considers new metrics. Our paper also adapts evaluation metrics to analyze the quality of the social network obtained using the proposed approach.

6 Concluding Remarks

The section 4 analyzed the Automatic INWeb Social Network. In the future, we plan to analyze the evolution of these results. We will also be able to compare them against new analysis from other Social Networks. Regarding the dissemination service, these results will also be useful. Specifically, once we build a collaborative network and verify its quality (using the aforementioned metrics), the quality of the dissemination services will also be improved. In other words, the evaluation of the relational ties among the researchers (authors) ensures better quality to the dissemination service. Therefore, based on such an evaluation, the dissemination service can identify (and recommend) the more pertinent publications as well as identify possible hidden collaboration nets.

As dissemination systems have recently grown from topic-based systems to XML-enabled systems, we believe that the next step is for them to follow the data technology and support any type of data uniformly (e.g. relational and XML). Moreover, considering all the aspects involved from the other research areas, we believe that the database technology must evolve to consider uniformly and

seamlessly any type of data there exist with *extensible* and *Web-scalable features*. This complex scenario brings new and exciting issues to be handled by many different Computer Science communities. Our final goal is to have a working system that integrates our research groups. The results will be evaluated, at the end of a four year period, by the access patterns and users evaluation of the quality of the disseminated papers and, more important, by the increase in the cooperation pattern among inter-institutional researchers. From the social point of view, those features are the fundamental element to the integration to the access of the content available at INWeb.

References

1. Wasserman, S., Faust, K.: Social Network Analysis: methods and applications. Cambridge University Press (1994)
2. Diao, Y., Rizvi, S., Franklin, M.J.: Towards an internet-scale xml dissemination service. In: VLDB. (2004) 612–623
3. Moro, M.M., Vagena, Z., Tsotras, V.J.: Recent Advances and Challenges in XML Document Routing. In: Open and Novel Issues in XML Database Applications: Future Directions and Advanced Technologies. IGI Global (2009) 136–150
4. Newman, M.E.J.: The structure and function of complex networks. SIAM Review **45** (2003) 167–256
5. Aleman-Meza, B., Nagarajan, M., Ding, L., Sheth, A.P., Arpinar, I.B., Joshi, A., Finin, T.W.: Scalable semantic analytics on social networks for addressing the problem of conflict of interest detection. TWEB **2**(1) (2008)
6. Mika, P.: Social networks and the semantic web. In: WI '04, Washington, DC, USA, IEEE Computer Society (2004) 285–291
7. Manning, C.D., Raghavan, P., Schütze, H.: Introduction to Information Retrieval. Cambridge University Press (July 2008)
8. Kunz, T., Black, J.P.: Using automatic process clustering for design recovery and distributed debugging. IEEE Trans. Softw. Eng. **21**(6) (1995) 515–527
9. Weng, S.S., Chang, H.L.: Using ontology network analysis for research document recommendation. Expert Syst. Appl. **34**(3) (2008) 1857–1869
10. Golbeck, J., Hendler, J.: Filmtrust: movie recommendations using trust in web-based social networks. In: IEEE CCNC - Consumer Communications and Networking Conference. Volume 1. (2006) 282–286
11. Jeh, G., Widom, J.: Simrank: a measure of structural-context similarity. In: ACM SIGKDD. (2002) 538–543
12. Zaiane, O.R., Chen, J., Goebel, R.: Dbconnect: mining research community on dblp data. In: WebKDD/SNA - Workshop on Web Mining and Social Network Analysis. (2007) 74–81
13. Menezes, G.V., Ziviani, N., Laender, A.H., Almeida, V.: A geographical analysis of knowledge production in computer science. In: WWW. (2009) 1041–1050
14. Ganev, V., Guo, Z., Serrano, D., Tansey, B., Barbosa, D., Stroulia, E.: An environment for building, exploring and querying academic social networks. In: MEDES '09, New York, NY, USA, ACM (2009) 282–289