

## **FINDING REFERENCES OF REPUTATION**



ALBERTO HIDEKI UEDA

## FINDING REFERENCES OF REPUTATION

Dissertation proposal presented to the  
Graduate Program in Computer Science of  
the Federal University of Minas Gerais in  
partial fulfillment of the requirements for  
the degree of Master in Computer Science.

ADVISOR: BERTHIER RIBEIRO-NETO

CO-ADVISORS: NÍVIO ZIVIANI

RODRYGO L. T. SANTOS

Belo Horizonte

April 2016



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	What is Reputation? . . . . .	1
1.2	An Academic Example . . . . .	2
1.3	What Do We Propose? . . . . .	3
<b>2</b>	<b>Related Work</b>	<b>5</b>
2.1	Ranking in Academic Search . . . . .	5
2.2	Ranking with Random Walks . . . . .	7
2.3	Random Walks in Academic Search . . . . .	8
<b>3</b>	<b>Methodology</b>	<b>11</b>
3.1	Datasets . . . . .	11
3.2	Experiments . . . . .	12
3.2.1	Markov Cluster Algorithm . . . . .	12
3.2.2	Preliminary Results . . . . .	13
3.2.3	Next Steps . . . . .	14
3.3	Tasks and Dates . . . . .	16
	<b>Bibliography</b>	<b>19</b>



# Chapter 1

## Introduction

### 1.1 What is Reputation?

Reputation is a widespread notion in society, albeit an arguably ill-defined one. In general, the reputation of an entity reflects the public perception about this entity developed over time [32]. This public perception may be either good or bad, and touches a variety of aspects that may impact the identity of the entity before the public, such as its competence, integrity and trustworthiness. Moreover, the reputation of an entity can change rapidly following an event in which the entity is involved, by means of word-of-mouth dissemination – whether traditional or electronic. As a result, reputation has been subject of professional management by public relations departments as well as of collective management by members of online communities, such as question-answering forums and online marketplaces [15].

The identification of reputable entities is an important task in many fields. Indeed, more reputable entities are presumably a better fit for most purposes. However, the subjective nature of reputation makes its quantification – and hence the identification of reputable entities – challenging. As a result, existing attempts to quantify the reputation of an entity rely on either manual assessments or on a restrictive definition of reputation, e.g., in terms of authority [18, 28], influence [3], or expertise [4].

In this work, we take an agnostic view of reputation. In particular, instead of relying on a single, precise definition of reputation, we exploit the *transference* of reputation among entities in order to identify the most reputable ones.

## 1.2 An Academic Example

For instance, evaluating a group of researchers is a permanent problem within research and academic institutions, laboratories and funding agencies [22]. Usually, this process involves forming highly qualified committees that must meet, define evaluation criteria and perform the evaluation. Moreover, it is very costly in terms of time, because evaluating numerous researchers (their curricula and publications) is not a simple process.

The decision of which researchers should be at the top – for hiring, promoting, funding, or distributing grants, scholarships, awards and so on – is typically based on criteria such as number of publications, impact of publications, number of undergraduate and graduate students under supervision, number of advised MSc and PhD theses, and participation in committees (conferences, journal editorial boards, technical committees, etc). Clearly, the effectiveness of the resulting ranking depends on how each criterion is assessed and the period of time covered in the assessment. Several indices have become widely used to measure the productivity of researchers. Examples include the raw number of citations, h-index [14], g-index [10] and citation z-score [23]. Likewise, most academic search platforms, such as Google Scholar<sup>1</sup>, Microsoft Academic Search<sup>2</sup> and ArnetMiner<sup>3</sup>, use some of such indices to rank researchers.

The same problem – finding the top academic entities – applies to the context of research groups, departments, universities and publication venues (conferences and journals). The general nature of the problem and the different relationships between the entities motivate us to use network analysis techniques as random walks and clustering.

For illustrative purposes, let us show the distribution of Brazilian researchers according to CNPq productivity levels in Computer Science, in Table 1.1. This distribution is the result of an evaluation process quite similar to that described in the first paragraphs of this section. The table shows that there are just 390 CS researchers in all of the Brazil, out of many thousands, who receive an individual research grant from CNPq. To compute this distribution, an important rule is applied: regardless the metric used to rank researchers, there is a percentage limit for the number of authors in each productivity level, except the last one (level 2). This limitation applies to the whole area of Computer Science, without considering the specificities of each of its sub-areas. One can see that independently of the academic metrics used, some sub-areas of Computer Science have intrinsic disadvantages in this classification. This problem motivates the current work and it is best addressed in next section.

---

<sup>1</sup><http://scholar.google.com>

<sup>2</sup><http://academic.research.microsoft.com>

<sup>3</sup><http://arnetminer.org>



**Table 1.1.** Distribution of researchers in CNPq productivity levels in Computer Science (CS) and in all areas of science

Level	Researchers	
	CS	All Areas
1A	23 (5.9%)	1320 (9.2%)
1B	22 (5.6%)	1308 (9.1%)
1C	31 (7.9%)	1376 (9.6%)
1D	70 (17.9%)	2386 (16.7%)
2	244 (62.6%)	7933 (55.4%)
Total	390	14,323

### 1.3 What Do We Propose?

Ranking researchers without regarding the specificity of their areas or sub-fields of research is arguably unfair and potentially error-prone [22]. For instance, consider the area of Human-Computer Interaction within the broad field of Computer Science (CS) [2]. Experimental evaluation in this area usually takes more time than in other CS areas when arranging and assessing users' feedback is necessary. On the other hand, CS areas such as Databases and Computer Graphics do not usually face the same problem because their experimental evaluations depend on assessing the outcome of an automatic process, such as a query evaluation or a graphics rendering engine. Likewise, researchers from some areas may have fewer publications than others, but with a potentially higher impact in their community.

In order to have a better understanding of this current problem in academic rankings, see the scenario of Computer Science in Table 1.2. It contains a list of 24 sub-areas of Computer Science according to Microsoft Academic Search.

In several academic communities, Brazilian researchers included, the distribution of grants, scholarships, awards, productivity levels, among others, are computed from the perspective of the whole area of Computer Science. Further, funding agencies as the *Conselho Nacional de Pesquisa (CNPq)* impose numerical limits on the members of grants issued to any given broad area of knowledge such as Computer Science. It so, one can realize that the most popular sub-areas in Computer Science have large advantages in this classification process over the other sub-areas in CS. For instance, the sub-area of computational theory cannot be equally compared to data mining, since the latter has more researchers, publications, and venues available today than the former one. Therefore, there is a necessity of a closer look to these academic scenarios from a perspective of sub-areas.

**Table 1.2.** Sub-areas of Computer Science according to Microsoft Academic Search

#	Subarea
1	Algorithms & Theory
2	Artificial Intelligence
3	Bioinformatics & Computational Biology
4	Computer Vision
5	Data Mining
6	Databases
7	Distributed & Parallel Computing
8	Graphics
9	Hardware & Architecture
10	Human-Computer Interaction
11	Information Retrieval
12	Machine Learning & Pattern Recognition
13	Multimedia
14	Natural Language & Speech
15	Networks & Communications
17	Operating Systems
18	Programming Languages
19	Real-Time & Embedded Systems
20	Scientific Computing
21	Security & Privacy
22	Simulation
23	Software Engineering
24	World Wide Web

One of the goals of this work is to find a fair and trustworthy way to rebalance these academic rankings in the context of sub-areas. To do that, we intend to use different approaches and modern metrics of relevance, including a generic model of reputation called P-score, which we describe in Chapter 2.

Some of the research questions we want to answer are: i) How to be fair with small sub-areas of knowledge or sub-areas with a low frequency of publications considering scientific impact categorization? ii) How well Brazilian research groups are internationally doing from the perspective of sub-areas? iii) How to characterize sub-areas in terms of scientific impact? iv) Can we find groups with the highest reputation in any topic of knowledge only adjusting the input data, but using the same method for sub-areas?

# Chapter 2

## Related Work

In this section, we review the related literature on ranking based on random walks, as well as approaches devoted to generating rankings in an academic search setting.

### 2.1 Ranking in Academic Search

Ranking has traditionally played an important role in academic search, particularly for tasks related to assessing the scientific productivity of academic entities. In particular, one of the earliest metrics proposed to quantify academic impact was Garfield’s Impact Factor [11]. Despite its wide usage since it was proposed in 1955, it has been largely criticized [33]. As a result, many alternatives have been proposed in the literature, such as other citation-based metrics like the H-Index [14], download-based metrics [6], and PageRank-like metrics [39]. As argued by Leydesdorff [21], each metric has its own bias and there are both advantages and disadvantages associated with each one.

Citation-based metrics have been applied to rank computer and information science journals [17, 26]. Also, several citation-based metrics have been proposed to measure the quality of a small set of conferences and journals in the database field [30], and to rank documents retrieved from a digital library [19]. Mann et al. [24] introduced topic modeling to further complement the citation-based bibliometric indicators, producing more fine-grained impact measures. Yan and Lee [39] proposed two measures for ranking the impact of academic venues which aim at efficiency and at mimicking the results of the widely accepted Impact Factor. An alternative method was presented by Zhuang et al. [42], who proposed a set of heuristics to automatically discover prestigious and low-quality conferences by mining the characteristics of program committee members.

Piwowar [29] recently claimed that citation-based metrics are useful, but not sufficient to evaluate research. In particular, he observed that metrics like the H-Index are slow. Indeed, the first citation of a scientific article can take years. As a result, he argued for the development of alternative metrics to complement citation analysis. In a similar vein, Lima et al. [22] argued that productivity indices should account for the singularities of the publication patterns of different research areas, in order to produce an unbiased assessment of the impact of academic output. Accordingly, they proposed to assess a researcher’s productivity by aggregating his or her impact indicators across multiple areas. Finally, Gonçalves et al. [13] investigated the importance of various academic features to scholar popularity and concluded that only two features are needed to explain all the variation in popularity across different scholars: (i) the number of publications and (ii) the average quality of the scholar’s publication venues. The metric P-score – one of the methods used in this work – exploits exactly these two features to rank different venues and different researchers.

The idea of reputation, instead of citations, was discussed by Nelakuditi et al. [25]. In particular, they proposed a metric called peers’ reputation, which measures the selectivity of a publication venue based upon the reputation of its authors’ institutions. The proposed metric was shown to be a better indicator of selectivity than the acceptance ratio. In addition, the authors observed that many conferences have similar or better peers’ reputation than journals. Another approach related to ours was proposed by Cormode et al. [7], who attempted to rank authors according to their similarity with respect to a reference author.

Alves et al. [1] and Lima et al. [5] analyze the structure of the communities formed by the flagship conferences of ACM SIGs. Their findings show that most of the ACM SIGs are able to connect their main authors in large and visually well-structured communities. However, they note that a few conferences, such as the ACM Symposium on Applied Computing, flagship conference of SIGAPP, and the ACM Conference on Design of Communications, flagship conference of SIGDOC, do not form a strong research community, presenting a structure with several disconnected components. They have opened their results to the research community as an interactive visualization tool<sup>1</sup> that allows one to browse the scientific communities, visualizing their structures and the contribution of each specific researcher to connect its coauthorship graph. The concept of scientific community and its natural development has enforced the clustering approach we adopt in this work.

---

<sup>1</sup><http://acmsig-communities.dcc.ufmg.br>

## 2.2 Ranking with Random Walks

Random walks have been widely used in information retrieval. They allow us to produce rankings of interconnected entities in a network, highlighting the most important nodes from a structural perspective. Page and Brin [28] designed the PageRank algorithm to calculate the importance of pages on the Web. PageRank simulates a web surfer's behavior. In particular, with probability  $p < 1$ , the surfer randomly chooses one of the hyperlinks of the current page and jumps to the page it links to; otherwise, with probability  $1 - p$ , the user jumps to a web page chosen uniformly at random from the collection. This defines a Markov chain on the web graph, where each probability of the stationary distribution corresponds to the rank of a web page, referred to as its *pagerank*.

Kleinberg [18] divided the notion of "importance" of a webpage into two related attributes: *hub*, measured by the authority score of other pages that the page links to, and *authority*, measured by the hub score of the pages that link to the page. These attributes are calculated in his Hyperlinked-Induced Topic Search (HITS) algorithm. Both algorithms, PageRank and HITS, have been successfully applied to rank the importance of different web pages through analyzing the link structure of the web graph.

Extensions of the random walk model were also studied for scoring several types of objects – e.g., products, people, and organizations – in different applications. For instance, Nie et al. [27] presented PopRank, a domain-independent object-level link analysis model to rank objects within a specific domain, by assigning a popularity propagation factor to each type of object relationship. Different popularity propagation factors for these heterogeneous relationships were assessed with respect to their impact on the global popularity ranking. Xi et al. [37] proposed a unified link analysis framework, called Link Fusion, which considers two different categories of links: intra-type links, which represent the relationship of data objects of a homogeneous data type (e.g., web pages), and inter-type links, which represent the relationship of data objects of different data types (e.g., between users and web pages). Regarding the recommendation of generic types of object, Jamali and Ester [16] proposed TrustWalker, a random walk method that combines trust-based and item-based recommendation, considering not only ratings of the target item, but also those of similar items.

Under the context of social networking systems, social friendship and random walks have been shown to be beneficial for collaborative filtering-based recommendation systems. These works argue that social friends – for instance, in Facebook or Twitter – tend to share common interests and thus their relationships should be considered

in the process of collaborative filtering [40]. In this context, a random walk sees a social network as a graph with probabilistically weighted links that represent social relations and thus is able to accurately predict users' preferences to items and their social influence with respect to other users. Backstrom and Leskovec [2] proposed an algorithm based on supervised random walks that combines the information from the network structure with node and edge level attributes, using these attributes to guide the random walk on the graph.

## 2.3 Random Walks in Academic Search

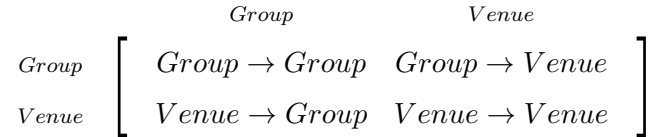
Earlier works have studied the application of random walks for ranking authors, papers, and venues in an academic setting. For instance, Sun and Giles [34] proposed a popularity weighted ranking algorithm for academic digital libraries that uses the popularity factor of a publication venue. Their approach overcomes some limitations of the Impact Factor and performs better than PageRank, citation counting, and HITS. Relatedly, Zhou et al. [41] proposed a method for co-ranking authors and their publications using several networks. Similarly, Yan et al. [38] presented a new informetric indicator, P-Rank, for measuring prestige in heterogeneous scholarly networks containing articles, authors, and journals. P-Rank differentiates the weight of each citation based on its citing papers, citing journals and citing authors.

In a narrower perspective, random walks have also been used for the task of expert finding in academic search collections. For instance, Deng et al. [8] proposed a joint regularization framework to enhance expertise retrieval in academia by modeling heterogeneous networks as regularization constraints on top of a document-centric model [4]. Relatedly, Wu et al. [36] proposed to model authors and publications as nodes of a publication network, with additional edges representing co-authorship information (author-author edges). In a similar vein, Tang et al. [35] proposed a probabilistic topic modeling approach to enrich a heterogeneous graph comprising multiple academic entities as nodes, including authors, papers, and publication venues, with directed edges representing a variety of relationships such as “written by” and “published in”. The stationary distribution computed after a random walk on this graph was then used to rank these entities with respect to an input query. A very similar approach was proposed by Gollapalli et al. [12], by assigning topics to nodes and then computing the unique stationary distribution of the associated Markov chain.

## Reputation Flows

In contrast to the aforementioned works, we use random walks to model the *transference* of reputation from *multiple* reference sources to selected targets in a reputation graph, as discussed in our previous work [32]. In order to validate our model, we instantiated it in an academic search setting by using research groups as reputation sources and publication venues as reputation targets. Moreover, while previous approaches have exploited multiple ranking signals, we demonstrated the power of the notion of reputation transfer by relying on publishing behavior as the only reputation signal.

Here we briefly discuss the instantiation of our conceptual framework of *reputation flows* in the academic context to model the transference of reputation between authors, papers, research groups and publication venues. The relations between these scientific entities may be captured through distinct metrics and, as far as we know, the most important ones (including citation-based metrics) fit well in our conceptual framework. In our experiments, we study how the reputation of a reference set of research groups is propagated to the venues they publish in and to other individual researchers by applying the concept of reputation flows. In this conceptual framework, publication venues are aggregations of papers and research groups are aggregations of authors, as shown in Figure 2.1. These aggregations are sufficient to establish core relations that allow ranking these entities.



**Figure 2.1.** Reputation flows between groups and venues.

The basic idea of the P-score metric is to associate a reputation with publication venues based on the publication patterns of a *reference set* of research groups in a given area or sub-area of knowledge. Given a pre-selected set of reference research groups, P-score associates weights with the publication venues the researchers in the reference groups publish in. Further, these weights can be used to rank other research groups or authors.

The reputation of a research group is strongly influenced by the reputation of its members, which is largely dependent on their publication records. We assume that i) A research group conveys reputation to a publication venue proportionally to its own reputation, and ii) A publication venue conveys reputation to a research group

proportionally to its own reputation. Once a reference group is selected, the reputation of its members is transferred to the venues. Recursively, since the reputation of research groups is correlated with the reputation of the venues in which they published, the venues transfer reputation to the groups. We use these hypotheses to build a Markov chain in which the nodes are research groups and publication venues. This chain can be solved by a stochastic computation which associates steady state probabilities to each node in the chain. These probabilities are taken as weights associated with venues which we refer to as P-scores.

In summary, the main contributions of this previous work are i) a novel random walk model for ranking entities according to the reputation collectively transferred to them by other entities in a reputation graph; b) an empirical validation of the effectiveness and robustness of the proposed model for two academic search tasks, namely, venue and researcher ranking; c) a preliminary investigation of the suitability of automatically choosing effective reputation sources. One of the goals of our current work aims to study and propose improvements for the latter item, i.e., the selection of academic entities – venues, researchers or research groups – as reputation sources [31].



# Chapter 3

## Methodology

In this chapter we describe the datasets we will use in this work, some initial experiments we have already done, the next studies we intend to perform and a time schedule for the main tasks of this research.

### 3.1 Datasets

#### DBLP

We have a collection of scientific publications extracted from DBLP [20], a digital library containing more than 3 million publications from more than 1.5 million authors that provides bibliographic information on major computer science conference proceedings and journals<sup>1</sup>. Each publication is accompanied by its title, list of authors, year of publication, and publication venue, i.e., conference or journal.

#### Microsoft Academic Graph

The Microsoft Academic Graph<sup>2</sup> (MAG) is a heterogeneous graph containing scientific publication records, citation relationships between publications – as well as authors, institutions, journals, conferences, and fields of study. It contains individual information about more than 120 million papers.

In Table 3.1, we present some statistics of the MAG dataset that are relevant to this work. Notice that a large fraction (59%) of the papers in this dataset have no information on citations, that is, the paper can be represented as a node in the citation

---

<sup>1</sup><http://dblp.uni-trier.de>

<sup>2</sup><http://research.microsoft.com/en-us/projects/mag/>

graph that has neither inlinks nor outlinks. There are two possible reasons for this to happen. The first alternative is the zero degree (i.e., both in- and out-degree) is a true representation of reality, it is a paper that in fact does not receive any citation yet and does not cite any other paper. The second alternative is due to the fact that any big repository offers an approximation of the reality, which also happens in the MAG dataset. Collecting and organizing a real-world dataset of such size is not a trivial task. In fact, it is a process that involves the treatment of huge amounts of semi-structured data, which usually causes inconsistencies. This lack of information reinforces the necessity of alternative strategies to citation-based methods.

**Table 3.1.** Relevant statistics on the Microsoft Academic Graph

Papers with citation information	49,870,036
Papers without citation information	71,017,797
Total number of papers	120,887,833

## 3.2 Experiments

Initially, we explored some clustering methods in order to help us characterize research communities. To do that, we built a graph of coauthorships, using the data from the DBLP dataset. In this graph, each node is a distinct researcher and each edge represents a cowork of the authors – an article written by the both authors. The higher the number of collaborations between the authors, the greater the weight of the edge connecting them. We use this weight as a measure of collaboration among the authors. This allows us to run clustering algorithms and analyze the resulting communities, as we now discuss.

### 3.2.1 Markov Cluster Algorithm

One of the main algorithms used in our initial studies, the Markov Cluster Algorithm (MCL) discussed in [9] is a fast and scalable unsupervised cluster algorithm for graphs based on simulation of stochastic flow in graphs. It finds cluster structures in graphs by a mathematical bootstrapping procedure. The process deterministically computes (the probabilities of) random walks through the graph, and uses two operators transforming one set of probabilities into another. It does so using the language of stochastic matrices (also called Markov matrices) which capture the mathematical concept of random walks on a graph. The MCL algorithm simulates random walks within a graph by an alternation of two operators called expansion and inflation. The expansion coincides

with taking the power of a stochastic matrix using the normal matrix product (i.e. matrix squaring). Inflation corresponds with taking the Hadamard power of a matrix (taking powers entrywise), followed by a scaling step, such that the resulting matrix is stochastic again, i.e. the matrix elements (on each column) correspond to probability values.

### 3.2.2 Preliminary Results

The first of our preliminary experiments has been the application of the MCL algorithm to the sub-area of Computer Networking. We did so by taking Infocom as a single source of reputation. We then selected the most productive authors (in number of papers published) of Infocom. Our goal in this step was to take a small set of researchers that represents the academic production of the venue Infocom, which we will call *reference set*. But, as we were interested in studying the whole sub-area of Computer Networking instead of a single venue, we added to the reference set all the coauthors of the researchers from the original set. A coauthor of author  $a$  is any researcher who has at least one published work together with  $a$ .

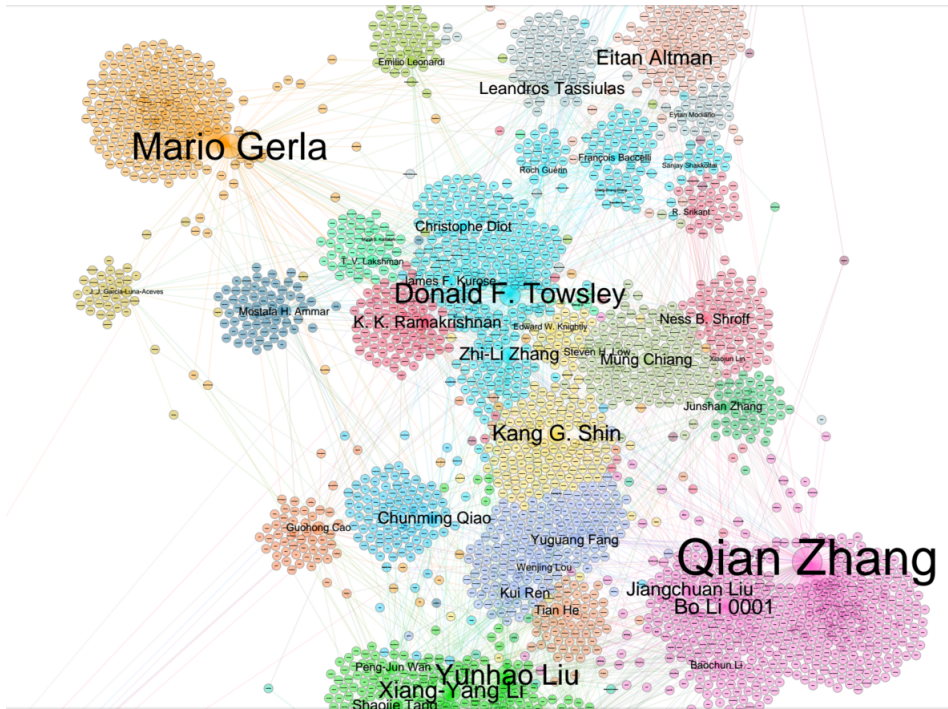
Using the researchers as nodes and the coauthorships as edges, we then obtained a graph of coauthorships. The weights of the edges in this graph are the number of papers two researchers have published together. In order to generate clusters in the graph of coauthorships, we used the number of papers which two researchers published together as a metric of similarity (or proximity) between these two researchers. We then ran clustering methods over this graph of coauthorships (e.g. the MCL algorithm, discussed in next session) and could also analyze the resulting academic communities with graph visualization tools. In our preliminary experiments we used the software Gephi<sup>3</sup> to visualize the graphs of coauthorships.

The result of this initial experiment is illustrated in Figure 3.1. It shows that the clusters produced are good indicators of productive authors in Computer Networking.

Additionally, in Figures 3.2, 3.3, 3.4, and 3.5 we show other experiments generated with the same approach described above but changing the initial single venue as input – respectively, Transaction on Networks (TON), Computer Networks (CN), SIGIR, and Web Search and Data Mining (WSDM). The venues TON and CN are considered as publication venues of the sub-area of Computer Networking while SIGIR and WSDM are considered venues of the sub-area of Information Retrieval.

---

<sup>3</sup><http://gephi.org/>



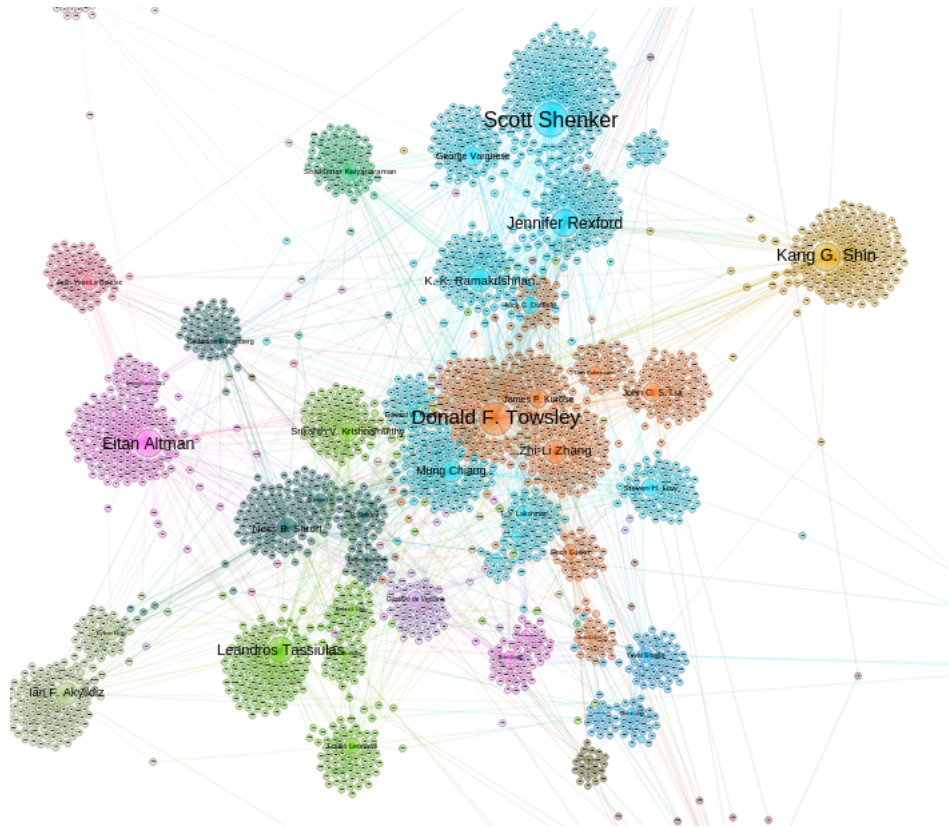
**Figure 3.1.** Graph of couauthorships with researchers from the sub-area of Computer Networking, using the venue Infocom as source of reputation. This visualization was generated by Gephi, an open-source framework for manipulating graphs.

### 3.2.3 Next Steps

The use of clustering methods is just an initial step to generate reputable sources in a sub-area. The goal is to rank venues, authors, and departments in a given sub-area by computing P-scores based on the reputable sources identified.

The framework we want to propose has multiple components between the initial input (reputation reference) and the final output (a ranking of venues, research groups, or individual researchers). This implies that there are several decision points within the framework. Each of these points must be studied and tested before we choose a given strategy as a standard approach. Examples of questions which we should answer in this framework are: i) How to select representative authors of a sub-area based only on a single publication venue of this sub-area? ii) Could we use the information provided by the resulting clusters as relevance feedback to the process (e.g. most centered clusters, largest clusters)? iii) Does the final result improve in quality if we use more than one publication venue as reputation source? If so, how many publication venues lead to the optimal result?

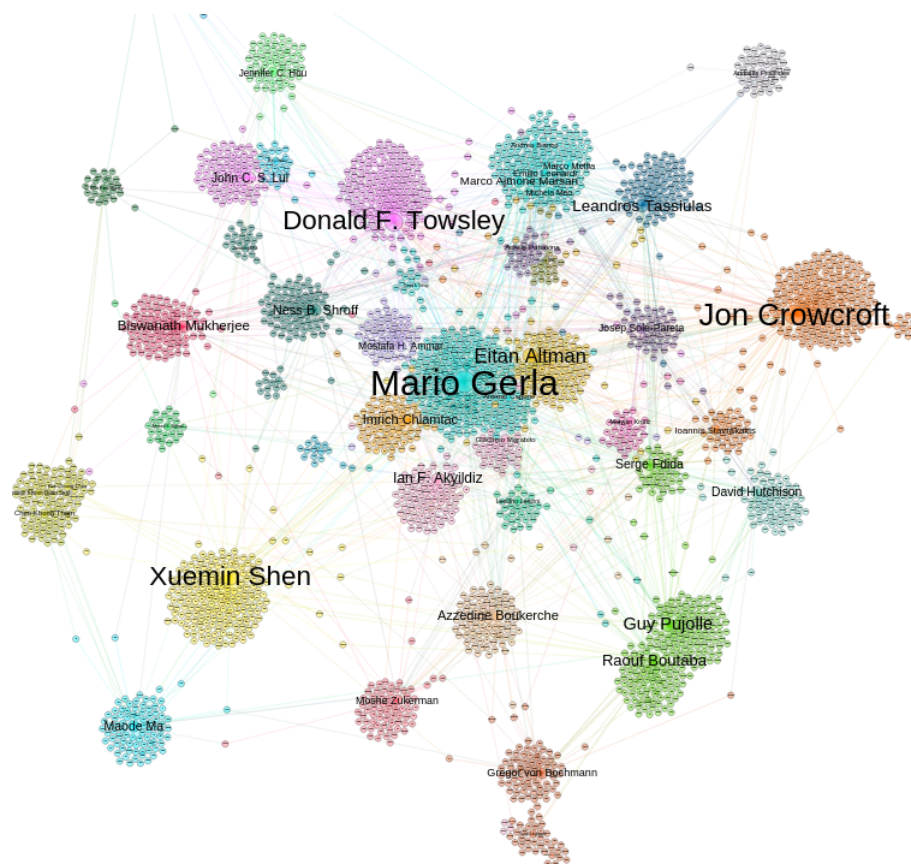
Also, we intend to test our strategies of ranking academic entities varying the



**Figure 3.2.** Graph of couauthorships with researchers from the sub-area of Computer Networking, using the venue TON as source of reputation.

method parameters of the algorithms. In this way, we want to check if our approach is robust enough to generate rankings of good quality – even with perturbations on the data input – and which are the more indicated values for the parameters and the reasons for that. Specifically, the reputation flow model allows us to increase or decrease the quantity of reputation that is transferred in each of the quadrants shown in section 2.3.

With a stable ranking algorithm defined, we can make a characterization of sub-areas of Computer Science in Brazil and compare the Brazilian departments to other international departments. We could also analyze which Brazilian sub-areas should be highlighted and which must be updated in terms of scientific impact. Besides that, following the research questions we want to answer with this work, we can make initial experiments for the problem to find good references of reputation in any topic of study in any sub-area (e.g., “land policies”, “infectious diseases”, “deep learning”, among others).

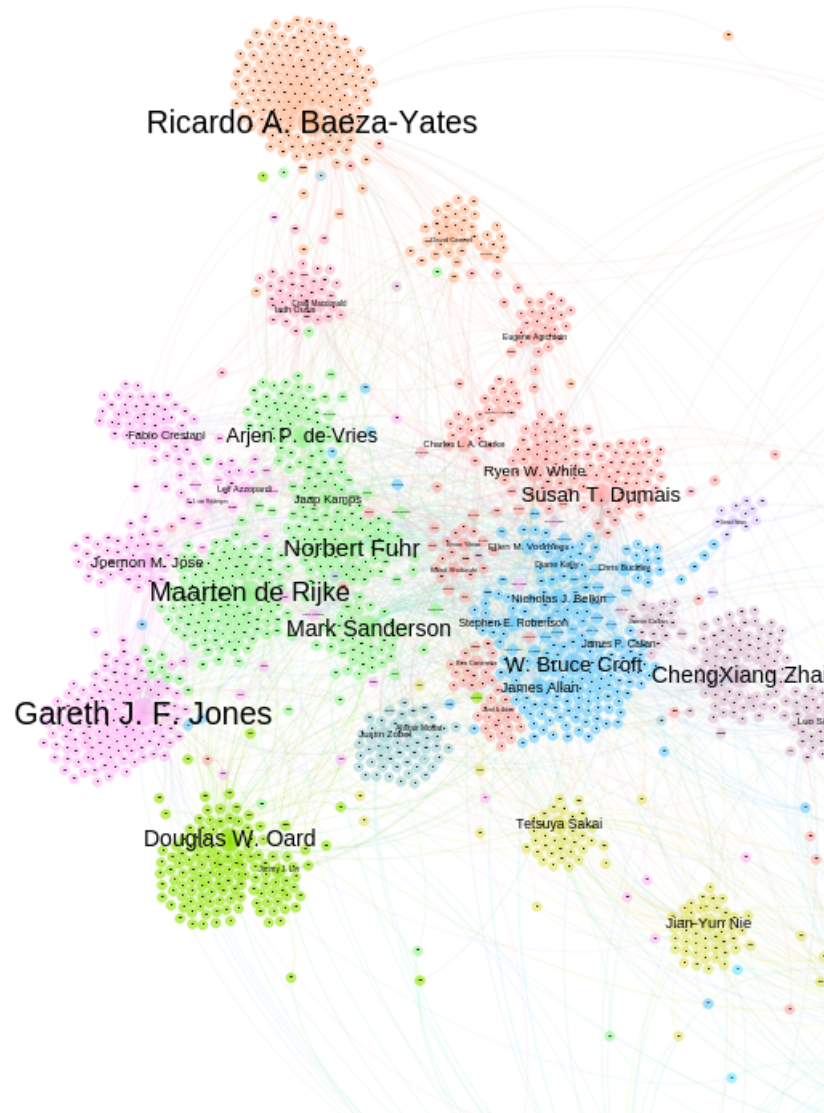


**Figure 3.3.** Graph of couauthorships with researchers from the sub-area of Computer Networking, using the venue Computer Networks as source of reputation.

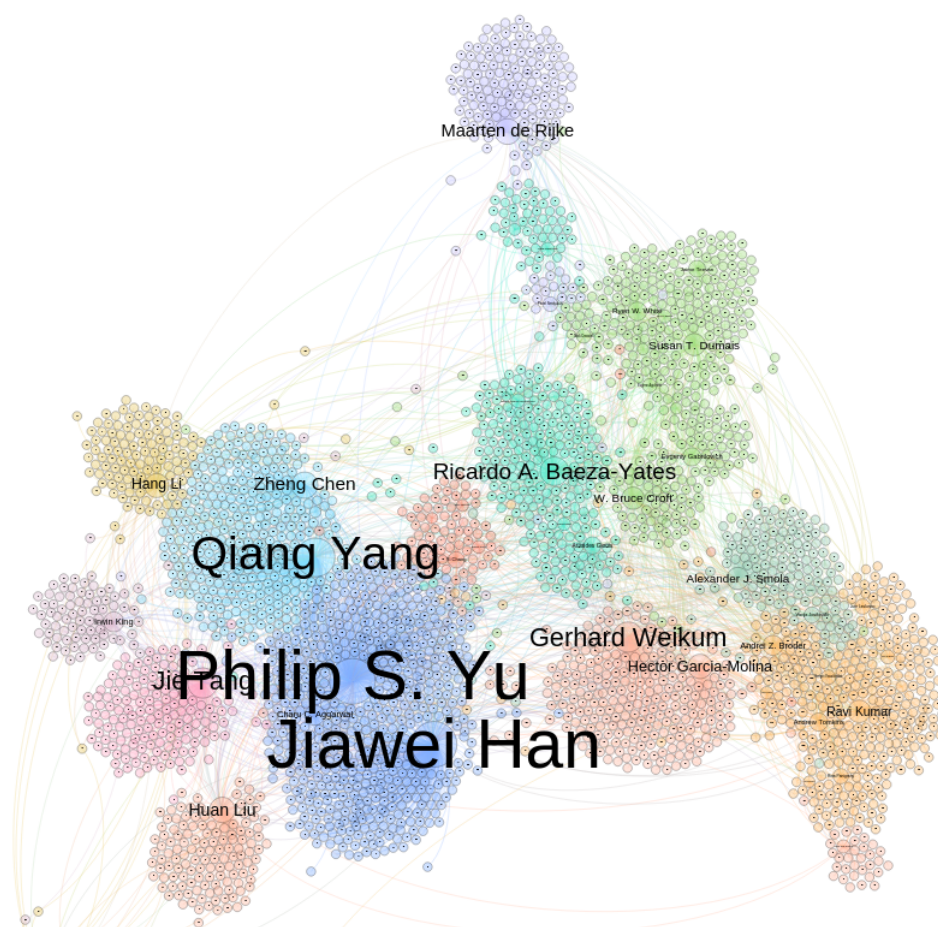
### 3.3 Tasks and Dates

2016 Schedule								
Task	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec
Thesis' project proposal	X							
Advisers assistant	X	X	X	X	X	X	X	X
Experiments	X	X	X	X	X			
Writing			X	X	X	X		
Defense						X		





**Figure 3.4.** Graph of couauthorships with researchers from the sub-area of Information Retrieval, using the venue SIGIR as source of reputation.



**Figure 3.5.** Graph of coauthorships with researchers from the sub-area of Information Retrieval, using the venue WSDM as source of reputation.



# Bibliography

- [1] Alves, B. L., Benevenuto, F., and Laender, A. H. (2013). The role of research leaders on the evolution of scientific communities. In *Proc. of WWW*, pages 649--656.
- [2] Backstrom, L. and Leskovec, J. (2011). Supervised random walks: Predicting and recommending links in social networks. In *Proc. of WSDM*, pages 635--644.
- [3] Bakshy, E., Rosenn, I., Marlow, C., and Adamic, L. (2012). The role of social networks in information diffusion. In *Proc. of WWW*, pages 519--528.
- [4] Balog, K. (2012). Expertise retrieval. *Found. Trends Inf. Retr.*, 6(2-3):127--256.
- [5] Benevenuto, F., Laender, A. H., and Alves, B. L. (2015). How Connected are the ACM SIG Communities? *SIGMOD Record*, 44(4).
- [6] Bollen, J., van de Sompel, H., Smith, J., and Luce, R. (2005). Toward alternative metrics of journal impact: A comparison of download and citation data. *Information Processing & Management*, 41(6):1419--1440.
- [7] Cormode, G., Muthukrishnan, S., and Yan, J. (2014). People like us: mining scholarly data for comparable researchers. In *Proc. of WWW*, pages 1227--1232.
- [8] Deng, H., Han, J., Lyu, M. R., and King, I. (2012). Modeling and exploiting heterogeneous bibliographic networks for expertise ranking. In *Proc. of JCDL*, pages 71--80.
- [9] Dongen, S. (2000). A cluster algorithm for graphs. Technical report, Amsterdam, The Netherlands, The Netherlands.
- [10] Egghe, L. (2006). Theory and practise of the g-index. *Scientometrics*, 69(1):131--152.
- [11] Garfield, E. (1955). Citation indexes for science. *Science*, 122(3159):108--111.

- [12] Gollapalli, S. D., Mitra, P., and Giles, C. L. (2011). Ranking authors in digital libraries. In *Proc. of JCDL*, pages 251--254.
- [13] Gonçalves, G. D., Figueiredo, F., Almeida, J. M., and Gonçalves, M. A. (2014). Characterizing scholar popularity: A case study in the computer science research community. In *Proc. of JCDL*, pages 57--66.
- [14] Hirsch, J. (2005). An index to quantify an individual's scientific research output. *Proc. Nat. Acad. Sciences*, pages 16569--16572.
- [15] Hutton, J. G., Goodman, M. B., Alexander, J. B., and Genest, C. M. (2001). Reputation management: the new face of corporate public relations? *Pub. Rel. Rev.*, 27(3):247--261.
- [16] Jamali, M. and Ester, M. (2009). Trustwalker: A random walk model for combining trust-based and item-based recommendation. In *Proc. of SIGKDD*, pages 397--406.
- [17] Katerattanakul, P., Han, B., and Hong, S. (2003). Objective quality rankings of computing journals. *Commun. ACM*, 45.
- [18] Kleinberg, J. M. (1999). Authoritative sources in a hyperlinked environment. *J. ACM*, 46(5):604--632.
- [19] Larsen, B. and Ingwersen, P. (2006). Using citations for ranking in digital libraries. In *Proc. of JCDL*, pages 370--370.
- [20] Ley, M. (2009). Dblp: Some lessons learned. *Proc. VLDB Endow.*, 2:1493--1500.
- [21] Leydesdorff, L. (2009). How are new citation-based journal indicators adding to the bibliometric toolbox? *J. Am. Soc. Inf. Sci. Technol.*, 60(7):1327--1336.
- [22] Lima, H., Silva, T. H. P., Moro, M. M., Santos, R. L. T., Jr., W. M., and Laender, A. H. F. (2013). Aggregating productivity indices for ranking researchers across multiple areas. In *Proc. of JCDL*, pages 97--106.
- [23] Lundberg, J. (2007). Lifting the crown-citation z-score. *Journal of Informetrics*, 1.
- [24] Mann, G., Mimno, D., and McCallum, A. (2006). Bibliometric impact measures leveraging topic analysis. In *Proc. of JCDL*, pages 65--74.

- [25] Nelakuditi, S., Gray, C., and Choudhury, R. R. (2011). Snap judgement of publication quality: how to convince a dean that you are a good researcher. *Mobile Computing and Commun. Review*, 15(2):20–23.
- [26] Nerur, S., Sikora, R., Mangalaraj, G., and Balijepally, V. (2005). Assessing the relative influence of journals in a citation network. *Commun. ACM*, 48(11):71–74.
- [27] Nie, Z., Zhang, Y., Wen, J.-R., and Ma, W.-Y. (2005). Object-level ranking: Bringing order to web objects. In *Proc. of WWW*, pages 567–574.
- [28] Page, L., Brin, S., Motwani, R., and Winograd, T. (1998). The pagerank citation ranking: Bringing order to the web. In *Proc. of WWW*, pages 161–172.
- [29] Piwowar, H. (2013). Altmetrics: Value all research products. *Nature*, 493(7431):159–159.
- [30] Rahm, E. and Thor, A. (2005). Citation analysis of database publications. *ACM Sigmod Record*, 34(4):48–53.
- [31] Ribas, S., Ribeiro-Neto, B., de Souza e Silva, E., Ueda, A. H., and Ziviani, N. (2015a). Using reference groups to assess academic productivity in computer science. In *Proc. of WWW*, pages 603–608.
- [32] Ribas, S., Ribeiro-Neto, B., Santos, R. L., de Souza e Silva, E., Ueda, A., and Ziviani, N. (2015b). Random walks on the reputation graph. In *Proceedings of the 2015 International Conference on The Theory of Information Retrieval*, pages 181–190. ACM.
- [33] Saha, S., Saint, S., and Christakis, D. (2003). Impact factor: a valid measure of journal quality? *J. Med. Lib. Assoc.*, 91(1):42–46.
- [34] Sun, Y. and Giles, C. L. (2007). *Popularity weighted ranking for academic digital libraries*. Springer.
- [35] Tang, J., Jin, R., and Zhang, J. (2008). A topic modeling approach and its integration into the random walk framework for academic search. In *Proc. of ICDM*, pages 1055–1060.
- [36] Wu, H., Pei, Y., and Yu, J. (2009). Detecting academic experts by topic-sensitive link analysis. *Front. Comp. Science in China*, 3(4):445–456.

- [37] Xi, W., Zhang, B., Chen, Z., Lu, Y., Yan, S., Ma, W.-Y., and Fox, E. A. (2004). Link fusion: A unified link analysis framework for multi-type interrelated data objects. In *Proc. of WWW*, pages 319–327.
- [38] Yan, E., Ding, Y., and Sugimoto, C. R. (2011). P-rank: An indicator measuring prestige in heterogeneous scholarly networks. *J. Am. Soc. Inf. Sci. Technol.*, 62(3):467–477.
- [39] Yan, S. and Lee, D. (2007). Toward alternative measures for ranking venues: a case of database research community. In *Proc. of JCDL*, pages 235–244.
- [40] Ye, M., Yin, P., Lee, W.-C., and Lee, D.-L. (2011). Exploiting geographical influence for collaborative point-of-interest recommendation. In *Proc. of SIGIR*, pages 325–334.
- [41] Zhou, D., Orshanskiy, S. A., Zha, H., and Giles, C. L. (2007). Co-ranking authors and documents in a heterogeneous network. In *Proc. of ICDM*, pages 739–744.
- [42] Zhuang, Z., Elmacioglu, E., Lee, D., and Giles, C. (2007). Measuring conference quality by mining program committee characteristics. In *Proc. of JCDL*, pages 225–234.