



**HAL**  
open science

## It is not What but Who you Know

Mario Cataldi, Luigi Di Caro, Myriam Lamolle, Claudio Schifanella

► **To cite this version:**

Mario Cataldi, Luigi Di Caro, Myriam Lamolle, Claudio Schifanella. It is not What but Who you Know. 24th International Conference on World Wide Web, May 2015, Florence, Italy. 10.1145/2740908.2742023 . hal-03580993

**HAL Id: hal-03580993**

**<https://hal.science/hal-03580993>**

Submitted on 18 Feb 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/311491273>

# It is not What but Who you Know: A Time-Sensitive Collaboration Impact Measure of Researchers in Surrounding Communities

Conference Paper · May 2015

DOI: 10.1145/2740908.2742023

CITATION

1

READS

137

4 authors, including:



**Luigi Di Caro**

Università degli Studi di Torino

97 PUBLICATIONS 1,079 CITATIONS

[SEE PROFILE](#)



**Myriam Lamolle**

Université de Vincennes - Paris 8

109 PUBLICATIONS 451 CITATIONS

[SEE PROFILE](#)



**Claudio Schifanella**

Università degli Studi di Torino

72 PUBLICATIONS 853 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Learning Café [View project](#)



ConTraffic [View project](#)

# It is not What but Who you Know: a Time-Sensitive Collaboration Impact Measure of Researchers in Surrounding Communities

Mario Cataldi · Luigi Di Caro · Myriam  
Lamolle · Claudio Schifanella

Received: date / Accepted: date

**Abstract** In the last decades, many measures and metrics have been proposed with the goal of automatically providing quantitative rather than qualitative indications over researchers' academic productions. Such interest is due to actual needs related to promotions, allocation of funds, and employment in general. Scientific profiles are often analyzed by means of measures calculated over publication lists, like h-index, supporting human efforts in making crucial decisions.

However, when evaluating a researcher, most of the commonly-applied measures do not consider one of the key aspect of every research work: the collaborations between researchers and, more specifically, the impact that each co-author has on the scientific production of another. In fact, in an evaluation process, some co-authored works can unconditionally favor researchers working in competitive research environments surrounded by experts able to lead high-quality research projects, where state-of-the-art measures usually fail in trying to distinguish co-authors from their pure publication history.

In the light of this, instead of focusing on a pure quantitative/qualitative evaluation of curricula, we propose a novel temporal model for formalizing and estimating the *dependence* of a researcher on individual *collaborations*, *over time*, in surrounding communities. We then evaluate our model with a set of experiments on real case scenarios and through an extensive user study.

## 1 Introduction

Considering the research community, one of the most commonly adopted method to evaluate the career of a researcher is to consider his/her authored

---

M.Cataldi · M. Lamolle  
Université Paris 8, Paris, France, E-mail: {m.cataldi,m.lamolle}@iut.univ-paris8.fr  
L. Di Caro · C. Schifanella  
Università degli Studi di Torino, Torino, Italy, E-mail: {dicaro,schi}@di.unito.it

papers and evaluate their “impact” on the surrounding research community. But how to evaluate this impact is still debated: most of the existing methods rely on counting the number of papers co-authored by the researcher and/or estimating, by applying different approaches, their citations number and quality.

Although these measures represent valuable tools for analyzing researchers outputs, they usually assume the co-authorship to be a proportional collaboration among the involved parts, missing out their relationships and their relative scientific impact on the resulting work. Moreover, considering that many of these measures are also used for recruitment purposes, it could be crucial to analyze the scientific relationships among authors in order to estimate the capacity of an author to work and produce research outcomes *without* the people that assisted his/her work until that time. A research collaboration can be indeed defined as a two-way process where individuals and/or organizations share learning, ideas and experiences to produce together scientific outcomes. Collaborations are necessary because of the evident difficulty for individual scientists to conduct several groundbreaking research on their own. For this, one of the key aspect (and more demanded on recruitment processes) of a successful researcher is the development of a large, active, network of collaborators that can help the researcher to bring new solutions and propose, continuously, novel ideas and approaches to the research community. On the other hand, evaluation of individuals needs a sort of inverse process with the primary goal of understanding the role of each researcher, and his/her specific impact on the research community, in this collaborative environment.

Let us consider for example a case of a young researcher who finishes a Ph.D. program supervised by an expert scientist; considering this situation, we can suppose his/her publication history as a set of research papers in which the supervisor is deeply involved as co-author. When this young researcher applies for some research position, the evaluators will be probably interested not only to a pure numerical productivity quantification, which mainly reflects the scientific effectiveness of his/her research group, but they will also analyze the curriculum with the goal of understanding how much the young scientist would be effective *without* his/her past research environment. In some sense, based on his/her publication record, the evaluator could be interested in analyzing the impact of each collaboration on the career of the young scientist in order to empirically estimate his/her potential capacity of maintaining the same research level without the environment that assisted his/her work until that time.

With this goal, we propose a novel temporal model that aims at evaluating the scientific collaborations of an author, over time, and their impacts on his/her entire research production (intended as the set of papers co-authored by him/her). Moreover, based on the DBLP bibliographic database<sup>1</sup>, we also developed a web environment (<http://d-index.di.unito.it/>) that implements the presented model and proposes a set of visualization tools to permit

---

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

to analyze, study and compare the careers of all the indexed authors based on their entire bibliographic records. Finally, relying on this platform, we present case and user studies that test both the validity and the reliability of the proposed evaluation measure.

The remainder of this paper is organized as follows: Section 2 surveys related research, focusing on the works that study the structure of the co-authorship networks and estimate the links connecting the authors/communities from different points of view. Section 3 presents our novel temporal metric for analyzing the impact of each co-authorship and its evolution along time. Section 4 presents a set of real-world scenarios and experiments applying our approach to the whole DBLP data corpus. Empirical evidence collected through experiments supports the effectiveness of our approach. Finally, Section 5 draws conclusions about the work presented in this paper, discussing limitations and future research directions.

## 2 Related Work

Bibliometric indicators are increasingly used to evaluate scientific careers based on personal publication records. The simple number of papers published by an author rather than the received citations are still common ways to capture both the quantity and the impact of an author's set of works.

However, such methods are still far from deducing the actual contribution of a researcher in the community. In this respect, it has been much discussed whether co-authors should have all the same value in quantifying the impact of a paper. In [26], for example, the author first pointed out the problem of *undeserved coauthorship*. In various works like [13, 23, 29] it has been stated that further efforts have to be done in this direction. However, the simple analysis of the position of an author in the list is not enough [14]. Indeed, this generalizes over something that is actually unknown. Which are the rules governing the position of a person in the authors list? An objective and universally-recognized point of view on that simply does not exist. For instance, a research group can have specific traditions when compiling the order of the authors.

Another interesting point of analysis is given by the following case: an author with one single highly-cited paper and a set of poorly-cited works. Surely, he/she could be advantaged when measuring the number of citations (or even the average). The same could happen when focusing on the number of papers without taking into account their impact. Hybrid approaches that use both citations and number of papers are ineffective since they need to rely on numeric thresholds, which are commonly not easy to choose and refine. An example is the number of *significant papers* (as in [11]), defined as the number of papers with more than a specific number of citations. In [11] and later in [12], the author introduced the *h-index*, a well-known metric for evaluations of academic careers, impact of journals, and research communities. An author has index  $h$  if he/she published  $h$  papers having at least  $h$  citations. As it has been fully demonstrated in [3, 5, 15, 25], it indirectly measures both quality and

quantity. Despite this, it has some disadvantages given the fact that it is not able to capture the quality of the papers that got more than  $h$  citations. In fact, [7] and [9] presented cases of authors with similar  $h$ -index but different citation patterns. The main problem of the  $h$ -index is that once a paper belongs to the  $h$ -core (i.e., the set of papers that have more than  $h$  citations), it does not matter how many more citations this paper will receive [9].

Since the  $h$ -index has been introduced, several extensions have been studied to avoid its drawbacks. The  $g(m)$  index [24], for instance, counts the papers equally fractionally according to the number of authors. In [8] the author gives a credit to a co-authorship based also on the number of received citations.

Finally, other works presented interesting ideas and insights on such a complex and multi-faceted domain. The authors of [1], for instance, stated that an author's scientific relevance should not be based on the number of citations of her/his papers, but is about how much co-workers she/he has been able to connect to in order to produce (joint) scientific publications. In [2] the authors started from the same motivation and proposed the *independence indicator* made up of three different dimensions of independence: the ability of developing own co-author networks, novel thematic directions, and strong quality of the research focus. However, considering that these measures are eigenvector analysis, the obtained results are not easily perceivable (and understandable) through search and navigation processes. In [17] the authors analyzed the temporal trend of the  $h$ -index to give some hints on how to evaluate an author with respect to a specific career phase. In [31], the authors shown how the eigenfactor score can be adapted to rank the scholarly output of authors. Again, eigenfactor is not a direct measure of quality, as the same authors claimed in their paper. Our approach, instead, is able to produce both analytical rankings as well as navigable explanations.

Many works tried also to take into account implicit and/or explicit edges of the collaboration network for detection/evaluation purposes. In [19] and [22], the authors aim at detecting characteristics like academic department, position, and country of origin from socio-academic networks, while [28] focus on the evolution of research teams. In [33], the authors use four centrality measures (closeness, betweenness, degree and PageRank) within a restricted collaboration network, showing that they are significantly correlated with citation counts. The work introduced by [10] compares the ways in which people seek advice and support from women in their networks. In [30], the authors learn advisor-advisee relationships from research publication networks starting from labeled data. Our approach is instead completely unsupervised. [21] aims at discovering the diffusion of scientific credits in the community relying on a citation network. Conversely, our technique does not measure some impact of an author in the whole scientific community but analyzes it at a local level. [4] apply the PageRank algorithm [18] to collaboration networks as a more effective way to compute scientific influences of papers with respect to the simple citation counts. [20] evaluate scientific careers taking into account different journals and citations over time.

In this paper, we focus on a different aspect of the problem, which is the analysis of the relationships in the co-authoring network to find out the impact and the effectiveness of such collaborations in his/her career. This permits to eventually support decisions when comparing  $h$ -index-related measures.

### 3 Analyzing Scientific Co-Authorship Relationships in Collaborative Environments

In this section we analyze the problem of evaluating research profiles by proposing a novel temporal metric for estimating the impact, over time, of a specific scientific collaboration on the production of researcher. Then we leverage this measure to summarize the overall dependence, over time, of an author, on the collaboration with the surrounding community. This measure will permit to estimate the independence of an author from the collaboration with other scientists and his/her potential capacity to maintain the same production in a different collaborative environment.

Please notice that, in this paper, we then define a *collaboration* as the common activity of two (or more) researchers intended to achieve the goal of producing new scientific papers.

For this we will first introduce the formalizations we use along the paper and explain the assumptions and the ideas that guided our work. We will then expose our approach providing different experiments to test its validity and reliability.

#### 3.1 Quantification of the Scientific Productivity of Authors over Time

In literature there is plenty of methods for evaluating the output of an author (either called scientist or researcher in the paper). Most of them consider their publication records as the basis for their scientific evaluation. In our paper, given an author  $a_i$ , we formalize his/her set of research outputs (also called papers, works or outcomes from now on)  $O_{a_i}^t$ , published until the time  $t$ , as<sup>2</sup>

$$O_{a_i}^t = \{o_{a_i,1}^t, o_{a_i,2}^t, \dots, o_{a_i,n}^t\}, \quad (1)$$

where  $o_{a_i,k}^t$  is the  $k$ -th research output authored, or co-authored, by him/her at the time  $t$  (for example, if  $O_{a_i}^t$ , with  $t=2000$ , contains all papers authored by  $a_i$  from the beginning of his/her career until 2000). Considering this information, it is possible to quantify the “productivity” of  $a_i$ ,  $p_{a_i}^t$ , at the time  $t$ , as

$$p_{a_i}^t = |O_{a_i}^t|, \quad (2)$$

where  $|O_{a_i}^t|$  is the cardinality of  $O_{a_i}^t$ .

---

<sup>2</sup> In this paper, the considered time intervals represent publication years.

In the same way, we can define the common outcome  $O_{a_i, a_j}^t$ , at the time  $t$ , of two authors,  $a_i$  and  $a_j$ , as

$$O_{a_i, a_j}^t = O_{a_i}^t \cap O_{a_j}^t = \{o_{(a_i, a_j), 1}^t, o_{(a_i, a_j), 2}^t, \dots, o_{(a_i, a_j), m}^t\}, \quad (3)$$

where  $o_{(a_i, a_j), k}^t$  is the  $k$ -th research output co-authored by both  $a_i$  and  $a_j$  at the time  $t$ , and  $m$  is the total number of papers co-authored by both of them at  $t$ . It is then possible to quantify their productivity  $p_{a_i, a_j}^t$  at the time  $t$  as

$$p_{a_i, a_j}^t = |O_{a_i, a_j}^t|. \quad (4)$$

Notice that this approach is extendable to any set of authors with any cardinality.

Moreover, given an author  $a_i$ , we formalize the scientific network in which he/she produced research outcomes as

$$Net_{a_i}^t = \{a_1^t, a_2^t, \dots, a_h^t\} \quad (5)$$

where  $h$  is the total number of co-authors, at the time  $t$ , of  $a_i$ . In the same way, given two authors  $a_i$  and  $a_j$ , we formalize their common scientific co-authorship network  $Net_{a_i, a_j}^t$ , at the time  $t$ , as the set of authors who co-authored at least one paper with both  $a_i$  and  $a_j$ , in the same output (i.e.  $a_k \in Net_{a_i, a_j}^t \Rightarrow \exists o_x^t \in O_{a_i, a_j}^t$  s.t.  $o_x^t$  is also co-authored by  $a_k$ ).

In the next sections we will leverage these formalizations to introduce our temporal model.

### 3.2 The $d$ -index: Analyzing Dependences on Collaborative Environments

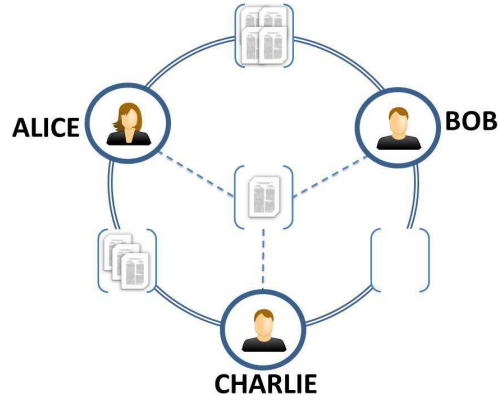
Given the entire set of scientific authors, theoretically, it is now possible to model a temporal co-authorship network,  $N^t$ , as a directed graph that expresses the dependence of each author, at the time  $t$ , on the scientific collaboration with a co-author. Formally we define  $N^t$  as

$$N^t = \{V^t, E^t, d\}, \quad (6)$$

where

- $V^t = \{a_1^t, a_2^t, \dots, a_n^t\}$  is the complete set of  $n$  scientific authors at the time  $t$  (i.e. researcher having published at least one outcome at the time  $t$ );
- $E^t$  is the set of undirected edges, where each  $e_{i,j} \in E$  represents an existing collaboration at the time  $t$  between  $a_i$  and  $a_j$  (where  $a_i, a_j \in V$ ) motivated by at least one research output co-authored by both of them at the time  $t$ ;
- $d$  is the weighting function ( $d : E \rightarrow [0; +1]$ ) representing the dependence, at the time  $t$  of  $a_i$  on the scientific collaboration with  $a_j$ .





**Fig. 1** A simplified example of co-authorship relations: three scientists, *Alice*, *Bob* and *Charlie* published, all together, one paper. However, *Bob* and *Charlie* did not publish any scientific outcome without *Alice*.

Following this formalization, in order to measure this dependence value, called *d*-index, we first aim to study each scientific collaboration of an author and estimate his/her autonomy from the surrounding scientific communities. Then, we will try to quantify the overall dependence of a considered researcher on the scientific collaboration with a specific co-author by analyzing how much each collaboration of the author were autonomous from the contribution of the considered co-author. In a sense, we aim at quantifying the impact of an author on the career of another by analyzing his/her average impact on all the his/her scientific collaborations.

In order to better understand the problem, let us consider a simplified situation as the one shown in Figure 1. Three authors, *Alice*, *Bob* and *Charlie* collaborated by publishing several scientific works. In particular, the collaboration between *Alice* and *Bob* (without *Charlie*) resulted in many research outputs as for the collaboration between *Alice* and *Charlie* (without *Bob*). On the other hand, the scientific relationships without *Alice* did not result in any published work. This situation can be summarized as follows: *Bob* and *Charlie* can be thought as young researchers who are supervised by *Alice*. In this case, *Alice* is leading this research group and most of the necessary expertise can be easily credited to her. Notice that, this fact does not reduce the merits or the contribution of *Bob* and *Charlie* in the considered research outputs. We just state that this situation suggests that the scientific production of *Bob* and *Charlie* results highly dependent on the scientific collaboration with *Alice*, which is also confirmed by the fact that, for each of them, any collaboration without *Alice* resulted poorly productive.

Considering this example, the scientific dependence has been highlighted from the analysis of their *co-authorship network* that models the environment (and, therefore, the relationships existing among authors) in which they work. In a sense, through this model, we analyze the productivity and the autonomy

of each collaboration with respect to all their co-authors and understand the impact of each author on the scientific production of each scientist in this collaborative environment.

Based on these assumptions, we first introduce a method to measure the ‘‘autonomy’’ of a collaboration, by taking into account the common scientific production, at the time  $t$ , of the involved authors. At this point, given two authors  $a_i, a_j$  and their common scientific network  $Net_{a_i, a_j}^t$ , the autonomy of their collaboration  $w_{a_i, a_j}^t$  is calculated as

$$w_{a_i, a_j}^t = \begin{cases} 0 & \text{if } Net_{a_i, a_j}^t = \{\} \\ \frac{1}{\sum_{a_k \in Net_{a_i, a_j}^t} \left( \frac{collab(a_k, O_{a_i, a_j}^t)}{\sum_{x=1}^x} \right)} & \text{if } Net_{a_i, a_j}^t \neq \{\} \end{cases} \quad (7)$$

where the function  $collab(a_k, O_{a_i, a_j}^t)$  returns the number of times the author  $a_k$  co-authored a paper with both  $a_i$  and  $a_j$  at the time  $t$ . Intuitively, this formula permits to measure the independence, at the time  $t$ , of the collaboration between  $a_i$  and  $a_j$  from the collaboration with any other author of the common scientific environment, expressed by  $Net_{a_i, a_j}^t$ . In this way, we take into account number and frequency of each collaboration; from one side, we count how many external co-authors, along their collaboration history (until the time  $t$ ), have been involved in the collaboration between  $a_i$  and  $a_j$ . From the other side, we also aim at evaluating the frequency of each contribution on their collaborations. In a sense, the autonomy of the collaboration will be lower when a high number of external co-authors are repetitively involved in the scientific outputs of the collaboration. Intuitively, the higher the autonomy, the more independent the work of  $a_i$  and  $a_j$  from the collaboration with any other co-author (and the other way around).

From this, given an author  $a_i$ , we aim at calculating his/her overall dependence on the collaboration with a co-author  $a_j$  by taking into account the capacity of  $a_i$  of working in his/her scientific environment without the scientific support of  $a_j$ . For this, given an author  $a_i$  and his/her scientific environment  $Net_{a_i}^t$ , at the time  $t$ , we define the dependence value,  $d$ -index, of the co-author  $a_i$  on the collaboration with  $a_j$  as  $d_{a_i \rightarrow a_j}^t$

$$d_{a_i \rightarrow a_j}^t = \frac{p_{a_i, a_j}^t}{p_{a_i}^t} \times \frac{w_{a_i, a_j, Net_{a_i}^t}^t + w_{a_j, \neg a_i, Net_{a_i}^t}^t}{w_{a_i, a_j, Net_{a_i}^t}^t + w_{a_j, \neg a_i, Net_{a_i}^t}^t + w_{a_i, \neg a_j, Net_{a_i}^t}^t}, \quad (8)$$

where

- $p_{a_i}^t$  returns the productivity of  $a_i$  at the time  $t$ ;
- $p_{a_i, a_j}^t$  is the productivity of the collaboration between  $a_i, a_j$  at the time  $t$ ;
- $w_{a_i, a_j, Net_{a_i}^t}^t$  is the autonomy of the collaboration, at the time  $t$ , among  $a_i, a_j$  and  $Net_{a_i}^t$  (i.e. the autonomy score of the collaboration between  $a_i$  and  $a_j$ , and at least one author  $a_k$  in  $Net_{a_i}^t$ );

- $w_{a_i, -a_j, Net_{a_i}^t}^t$  is the autonomy score of the collaboration between at least one author in  $Net_{a_i}^t$  and  $a_i$  without the contribution of  $a_j$  (i.e., excluding the research outputs in which  $a_j$  is also involved);
- $w_{a_j, -a_i, Net_{a_i}^t}^t$  is the autonomy score of of the collaboration between a least one author in  $Net_{a_i}^t$  and  $a_j$  without the contribution of  $a_i$  (i.e., excluding the research outputs in which  $a_i$  is also involved).

It is important to notice that the  $d$ -index value  $d_{a_j \rightarrow a_i}^t$  ranges from 0 to 1; in particular,  $d_{a_i \rightarrow a_j}^t \approx 0$  indicates that the dependence of  $a_i$  on  $a_j$ , at the time  $t$ , is negligible, while a  $d_{a_i \rightarrow a_j}^t \approx 1$  highlights the contrary. In fact the second term of the formula increases when the autonomy score of  $a_i$  and  $Net_{a_i}^t$ , without the contribution of  $a_j$ , is negligible ( $w_{a_i, -a_j, Net_{a_i}^t}^t \approx 0$ ) and the other collaborations are significantly autonomous ( $w_{a_i, a_j, Net_{a_i}^t}^t > 0$  and  $w_{a_i, -a_j, Net_{a_i}^t}^t > 0$ ). On the other hand, the higher the  $w_{a_i, -a_j, Net_{a_i}^t}^t$ , the lower the relative dependence.

Please also notice that  $d_{a_i \rightarrow a_j}^t \neq d_{a_j \rightarrow a_i}^t$ ; in fact their mutual dependences can significantly differ, since they are also based on their personal collaborations (which are obviously not the same, even when they share the same co-authors).

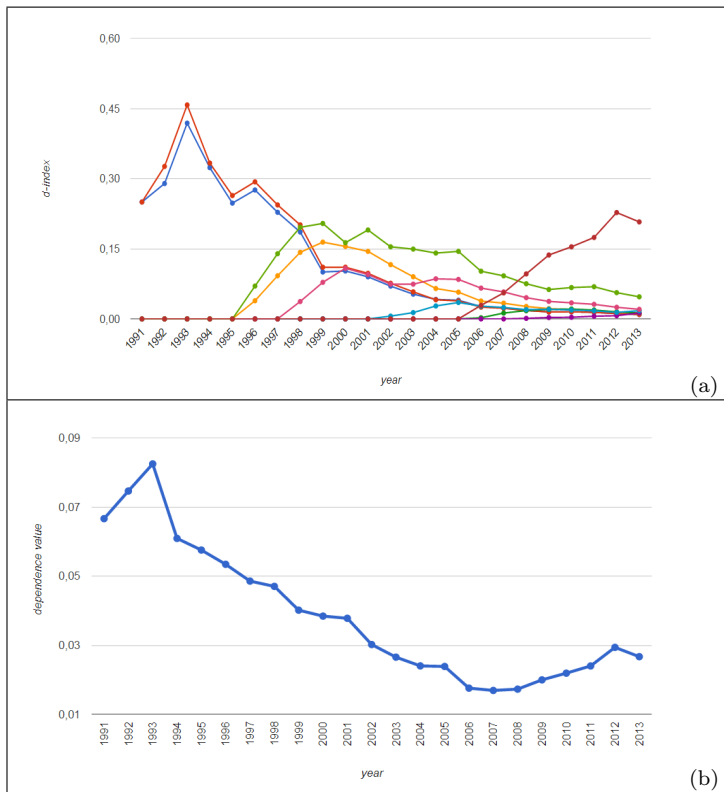
### 3.3 Dependence Trajectory: Leveraging the $d$ -index Values to Estimate the Evolution of the Dependence over Time

In Section 3.2, we introduced a novel way to estimate the dependence, at a specific time, of a given author on the scientific collaboration with a co-author, based on their scientific network and the productivity of each collaboration within this network. These values can now be leveraged to graphically map the scientific dependences of an author, along his/her career, on the collaboration with each co-author, as a set of curves that plots the relative  $d$ -index values. For this, we define the dependence curve of an author  $a_i$  with respect to a co-author  $a_j$  as

$$\overrightarrow{d_{a_i \rightarrow a_j}} = \{d_{a_i \rightarrow a_j}^t, d_{a_i \rightarrow a_j}^{t+1}, \dots, d_{a_i \rightarrow a_j}^{t+n}\}, \quad (9)$$

where  $t$  is the year of the first publication of  $a_i$ , and  $n$  expresses the arithmetical difference between the last and the first year of publication of  $a_i$ . Thus, given an author  $a_i$ , and the complete set of his/her coauthors expressed by  $Net_{a_i}$ , it is now possible to graphically represent, in the same chart, his/her dependence on each co-author  $a_k \in Net_{a_i}$ , along the career of  $a_i$ , to obtain a first sight on this mined knowledge. An example is shown in Figure 2 (a).

Each of these curves can graphically highlight the evolution of the collaboration with a specific co-author along the time and understand how much the considered author became independent (or dependent) from him/her with the years. Considering the example in 2 (a), nine dependence curves are provided.



**Fig. 2** (a) The dependence curves of an author  $a_i$  and his/her dependence trajectory (b).

Eight of them are visibly decreasing (highlighting the fact that the author becomes increasingly independent from the collaboration with the eight related co-authors, along the considered intervals) while the last one, in red, significantly increases after 2005 (thus, even if the number of co-authors increases, the author seems becoming dependent on the collaboration with the related co-author).

In order to better evaluate this situation, considering that many authors can have hundreds of co-authors on a career of multiple decades, in this section we also aim at obtaining a one-curve evaluation system to summarize, at best, the overall independence of the author. This curve can permit to graphically highlight how the career of an author evolved with respect to the dependence on the collaboration with the surrounding scientific community.

Thus, given the complete set of dependence-curves, we calculate the author's *dependence trajectory*, by calculating the standard deviation, along the time, of each  $d$ -index value, for each co-author, from the optimal attended value of 0 (which would mean a dependence score of 0; i.e., the production of the considered author is independent from the collaboration with the considered co-author). In a sense, we aim at evaluating the overall independence of

an author from the surrounding community. More formally, given an author  $a_i$ , we define his/her dependence trajectory  $\vec{d}_{a_i}$  as

$$\vec{d}_{a_i} = \{sd_{a_i}^t, sd_{a_i}^{t+1}, \dots, sd_{a_i}^{t+n}\}, \quad (10)$$

where  $sd_{a_i}^t$  is calculated as

$$sd_{a_i}^t = \sqrt{\frac{\sum_{a_k \in \text{Net}_{a_i}} (d_{a_i \rightarrow a_k}^t)^2}{|\text{Net}_{a_i}|}}. \quad (11)$$

The meaning of this formula is evident: calculate the average standard deviation of the previously calculated  $d$ -index values from the optimal value of 0. The higher the  $sd_{a_i}^t$ , the more dependent the work of  $a_i$  on the collaboration with any of his/her co-authors at the time  $t$ .

In Figure 2 (b) we show an example of dependence trajectory (calculated based on the dependence curves shown in Figure 2 (a)). In this example, we can easily see an overall increment in the dependence trajectory; this is mainly due to the significant increase in the dependence values of the considered author related to a specific co-author (visualized through the red line in Figure 2 (a)). The reason of this behavior is evident: the system tries to detect anomalies in the collaboration patterns with respect to some expected values. Authors in fact are expected, along the career, to increment their collaboration network and, therefore, become independent from the collaboration with each single co-author. Equation 7 in fact leverages the number of collaborations to estimate the autonomy of a scientific relationship. However, in the considered example, even in presence of a general (expected) increase in the number of scientific collaborations along the career, the increment in the dependence on a single co-author is so significant to lead the system to a visible boost in the dependence trajectory (which, however, is expected to constantly decrease).

Please also notice that, theoretically, it is not sufficient a single, increasing, dependence curve to obtain an overall increment in the trajectory curve. In fact, we obtain this result only when the dependence on a specific co-author (or a subset of co-authors) increases so much to counterbalance (in the Equation 11) the overall expected drop on the dependence values related to the rest of his/her research community.

In a sense, through the dependence trajectory, we analyze the impact of the surrounding environment, along the time, on the publication record of the considered researcher. In fact, this curve permits to estimate how much the average dependence of a researcher evolved, over time, with respect to his/her scientific environment. Moreover, it can be also used to analyze the collaboration patterns of scientific authors. Using the dependence trajectory, we can now compare researchers focusing on their average collaboration behavior.

In the next sections we will show how to use this evaluation system for analysis and comparison purposes.

Position	Avg. # of papers	Avg. # co-authors	Cardinality
Ph.D Stud/Post Doc.	19.28	26.21	24
Researchers/Assist. Professors	31.57	29.04	29
Professors	64.43	45.93	28

**Table 1** The number of users that replied to our questionnaire, grouped based on their academic position.

## 4 Evaluation

In this paper we have presented a novel temporal model that permits to focus the evaluation of an author on the analysis of his/her scientific collaborations. The model has been implemented and all the data are freely accessible at <http://d-index.di.unito.it>. This allows complementary evaluations to make easier real tasks like recruiting, prize-giving, and so forth.

Moreover, in order to better analyze the introduced measures, in this section, we provide real, high-level, scenarios and user studies that can help understand the introduced evaluation scheme and motivate the benefits of using the proposed approach. For our experiments, we considered a data-set extracted from the DBLP bibliographic database<sup>3</sup> containing information about 1,342,723 authors and 2,446,236 scientific papers<sup>4</sup>.

### 4.1 User Study and Experiments

In this section, we illustrate the results of a user study that we conducted to evaluate the approach and the developed system in actual and qualitative analysis of collaboration dependences. In particular, we made available a web application where all people from universities and research institutes could answer specific questions on the impact of some of their collaborators in different time ranges. The aim was to validate the ability of our approach in the identification of the key actors within individual and evolving scientific careers over time.

The web site was left open to answers for an entire week, after having published the news on different networks and with a word-of-mouth strategy. To make the study valid, we used an email-checking control, avoiding possible abuses. This way, we are able to state that all participants only evaluated their own scientific profile. We received a total of 81 answers, among which 78 users finished the entire questionnaire (we did not consider incomplete questionnaires). A detailed view of the users set is shown in Table 1.

<sup>3</sup> <http://www.informatik.uni-trier.de/ley/db/>

<sup>4</sup> Information updated on March 2014. Please notice that, within our system, we rely on the disambiguated authors name provided by [27]. This name disambiguation system is based on a probabilistic model and aims at finding, extracting, and fusing the semantic-based profiling information of a researcher from the Web.

Going through the details of the study, we asked users to answer a questionnaire in which they had to order three randomly-picked co-authors according to who they felt being more important in her/his career. We repeated this question 6 times, i.e. 2 times for 3 different time frames (randomly picked within the beginning, middle and final time intervals of the career<sup>5</sup>). This way, we aimed at analyzing the performance of the introduced measure over time. We also added two check-boxes called “*difficult choice*” to let the users express a possible doubt between the first-ranked co-author and the second-ranked co-author, and/or between the second-ranked and the third-ranked respectively. This option ensured the possibility for the user to express doubtful responses.

Then, in order to assess the validity of the approach, we compared the users rankings with the ones provided by our system through the  $d$ -index values. This comparison has been first performed by means of the Pearson product-moment correlation coefficient  $r$ , which is the covariance of the two variables divided by the product of their standard deviations. The correlation coefficient ranges from -1 to 1 (where 0 means independence), and in our test we achieved a total  $r$  score of 0.76. Usually, values greater 0.7 indicate a strong correlation [6].

Since this coefficient is only able to capture linear correlations, and it is known to be not robust especially in case of outliers [32], we further evaluated our  $d$ -index values making use of the Kendall’s Tau coefficient [16]. Using this measure it is possible to capture the rank correlation between different orderings. In detail, this measure works as follows: given a set of non-ordered items  $I = \{i_1, i_2, \dots, i_n\}$  (in our case, representing co-authors), two ordered sets  $R_a$  and  $R_b$  of all the elements in  $I$  (i.e., the orderings provided by the user and the system), and the set  $T$  of all the ordered pairs of elements in  $I$ ,  $\langle i_p, i_q \rangle$  (where  $p > q$ ), the Kendall’s Tau coefficient calculates the distance between the two orderings by relying on the number of pairs in  $T$  that are concordant and discordant in the considered rankings. More formally, the Kendall’s Tau coefficient is calculated as:

$$kendall(R_a, R_b) = \frac{conc - disc}{|T|} \quad (12)$$

where  $|T|$  is the number of the pairs,  $conc$  is the number of pairs in  $T$  that are equally ordered within  $R_a$  and  $R_b$ , and  $disc$  is the number of pairs in  $T$  that are differently ordered within  $R_a$  and  $R_b$ . Please notice that Kendall’s Tau coefficient ranges from -1 (one ordering is the contrary of the other) to +1 (the considered orderings are exactly equals).

The results are shown in Table 2, aggregated and averaged by different time frames. The overall average Kendall’s Tau coefficient is 0.55 (within the interval  $[-1, 1]$ ), which proves a positive correlation among the orderings provided by the users and the ones elaborated through the  $d$ -index measure. These

---

<sup>5</sup> We define the beginning of a career as the period included within the first third of his/her publication history, the middle as the second third and the end as the third third of the career of an author. The decision of selecting these periods aimed to analyze the performances of the presented metric along the whole career of an author, even in presence of a very large discard between the first and last year publication year.

Time Range	KT	Weighted KT	KT ( <i>DC</i> )	Weighted KT ( <i>DC</i> )
Early career	0.561	0.648	0.645	<b>0.704</b>
Mid-career	0.556	0.685	0.634	<b>0.741</b>
End career	0.536	0.668	0.552	<b>0.662</b>
Overall	0.551	0.667	0.610	<b>0.702</b>

**Table 2** The results of the user study in terms of both original and weighted Kendall’s Tau coefficient. The term (*DC*) indicates that the difficult choices marked by the users have not been considered.

experiments clearly highlight both the validity and the reliability of our approach: the resulting Kendall’s Tau coefficients result very coherent for all the aggregations and the considered time intervals. The Kendall’s Tau coefficients are indeed similar for all the considered academic intervals (0.561 for Early career, 0.556 for mid-career and 0.536 for the end of the career), therefore proving the capacity of the proposed approach in positively capturing collaboration dependences over time.

Moreover, in order to better evaluate the approach, we slightly modified the reported Kendall’s Tau definition in order to take into account the weighted distances provided by the ordered  $d$ -index values. We believe that a comparison among users and system orderings should take into account also the relative distances, in terms of  $d$ -index values, among the considered co-authors. In other words, in the Kendall’s Tau computation, the presence of a discordant pair with a high  $d$ -index distance between the two items should have an higher negative impact with respect to a discordant ordering with similar  $d$ -index values, and the other way around. For this reason, we adapted the Kendall’s Tau coefficient to fit with our range of dependence values by weighting accordingly both the concordant and the discordant pairs. More formally, we then computed these values as

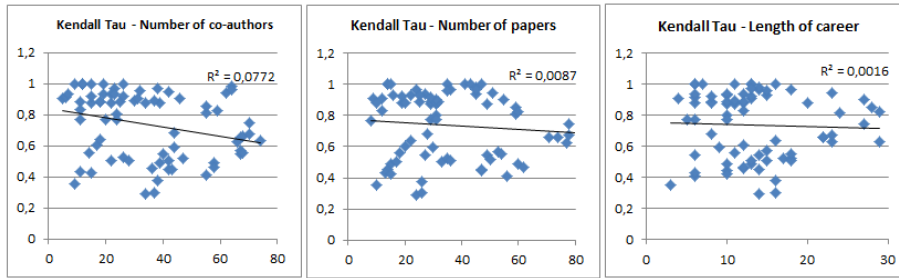
$$conc = \sum_{\langle i_p, i_q \rangle \in T} \left( 1 - \frac{dist(i_p, i_q)}{max(I) - min(I)} \right) \quad (13)$$

and

$$disc = \sum_{\langle i_p, i_q \rangle \in T} \left( \frac{dist(i_p, i_q)}{max(I) - min(I)} \right) \quad (14)$$

where  $dist(i_p, i_q)$  is the distance between the  $d$ -index value of  $i_p$  and the one of  $i_q$  (calculated as the absolute value of the difference), while  $max(I)$  and  $min(I)$  are respectively the highest and the lowest  $d$ -index values related to some item in  $I$ . This way, we are able to weight the differences between system and users orderings accordingly to the  $d$ -index values (i.e., the higher the difference between two  $d$ -index values the more significant the correct matching between the ranks, and vice-versa). Notice that, even with this weighted normalization, the Kendall’s Tau coefficient still ranges between -1 and +1. Again, the results shown in the “Weighted Kendall’s Tau” column in Table 2 demonstrate that





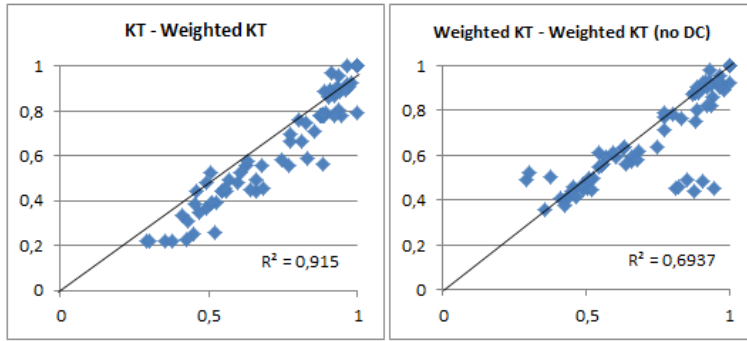
**Fig. 3** The correlation between the Kendall's Tau scores and information like number of papers, number of co-authors, and length of career. Notice that there is no linear correlation (the value of  $R^2$  is very low), demonstrating the reliability of the approach on distinct contexts of analysis.

all the Kendall's Tau coefficients increase when the  $d$ -index distances are taken into account, highlighting that the system is able to capture differences among scientific dependencies where they matter while it can fail mostly when they are minimal.

Finally, we thought that making decisions between collaborators in terms of their scientific dependence could be difficult in some cases. For this reason, as already reported, the users had the possibility to mark those choices that they felt difficult to take (check-box “*difficult choice*” in the User Study). Table 2 shows the results of the evaluation where these doubtful answers are not taken into account (*DC* columns). At this point, considering users' difficult choices and weighting the rankings with respect to the obtained  $d$ -index values, we reached an overall positive Weighted Kendall's Tau score of 0.702, which highlights the capacity of the proposed model to capture the dependences where they have an evident impact on a research curriculum. To sum up, the  $d$ -index resulted to be able to achieve high positive correlation values with the users feelings. Moreover, we also demonstrated that the approach is also able to capture the concept of scientific impact of collaboration over time, i.e., in different time frames. An interesting insight is about the light decrease of the Kendall's Tau scores in the last years of the authors' career, probably due to the incremental complexity of the overall collaboration networks.

As additional test, we then evaluated the relationship between the Kendall's Tau scores and information like the number of published papers, the number of co-authors, and the length of the career. This way, we tested the reliability of our measure by evaluating it on data samples with distinct characteristics. Figure 3 shows the results of this analysis.

Finally, we studied the correlation between the Kendall's Tau and its weighted version. As it can be seen in Figure 4 (plot on the left), there is a high linear correlation between them, and the weighted scores are always in favor of the  $d$ -index approach. The impact of the difficult choices in the evaluation of the  $d$ -index is shown on the right plot of Figure 4, demonstrating that



**Fig. 4** The correlation between the Kendall's Tau and its weighted version (on the left), and the analysis of the difficult choices impact (on the right). Notice that the diagonal lines, in this case, are not tendency lines, and they only split the plot to highlight the difference in values between KT and weighted KT.

choices marked as difficult by the users make the effectiveness of the  $d$ -index lower than what it actually is.

## 5 Conclusions

The problem of evaluating the quality of researchers' outputs has been broadly studied whereas few works attempted to discover collaboration dependences among researchers in case of co-authored papers. In this sense, we proposed a novel temporal model that aims at uncovering dependences among the authors over time according to their research environment and their publication history. We then evaluated the presented model through several examples and user studies that validated the model under different points of view. We also introduced the freely available web platform <http://d-index.di.unito.it/> that implements the presented ideas.

However, some limitations emerged from the study of this challenging task; for this reason, we plan to extend our work in order to analyze collaborations among research groups and evaluate their relative dependences based on their co-authored works. In particular, we also aim to study how these collaborations evolve and how much they are dependent on the collaboration among individual researchers (for example, the leaders of the groups). We also plan to extend the work in order to include alternative productivity measures (as for example the number of citations, or the h-index of the shared research outputs) and evaluate the impact of existing social rules that can condition their quantification. Finally, the impact of collaborations might be also analyzed according to research areas. In this sense, by incorporating textual contents (titles, abstracts, entire papers) in the data, we will study topic-based dependences and impact values within surrounding communities.

## References

1. Ausloos, M.: A scientometrics law about co-authors and their ranking: the co-author core. *Scientometrics* **95**(3), 895–909 (2013)
2. Van den Besselaar, P., Sandström, U., Van der Weijden, I.: The independence indicator: Towards bibliometric quality indicators at the individual level (2012)
3. Bornmann, L., Daniel, H.: Does the h-index for ranking of scientists really work? *Scientometrics* **65**(3), 391–392 (2005)
4. Chen, P., Xie, H., Maslov, S., Redner, S.: Finding scientific gems with google’s pagerank algorithm. *Journal of Informetrics* **1**(1), 8–15 (2007)
5. Cronin, B., Meho, L.: Using the h-index to rank influential information scientists. *Journal of the American Society for Information Science and Technology* **57**(9), 1275–1278 (2006)
6. Dancey, C.P., Reidy, J.: *Statistics without maths for psychology*. Pearson Education (2007)
7. Egghe, L.: An improvement of the h-index: the g-index. *ISSI Newsletter* **2**(1), 8–9 (2006)
8. Egghe, L.: Mathematical theory of the h-and g-index in case of fractional counting of authorship. *Journal of the American Society for Information Science and Technology* **59**(10), 1608–1616 (2008)
9. Egghe, L.: The hirsch-index and related impact measures. *Annual Review of Information Science and Technology* **44**, 65–114 (2010)
10. Feeney, M., Bernal, M.: Women in stem networks: who seeks advice and support from women scientists? *Scientometrics* **85**(3), 767–790 (2010)
11. Hirsch, J.: An index to quantify an individual’s scientific research output. *Proceedings of the National Academy of Sciences of the United States of America* **102**(46), 16,569 (2005)
12. Hirsch, J.: Does the h index have predictive power? *Proceedings of the National Academy of Sciences* **104**(49), 19,193 (2007)
13. Hunt, R.: Trying an authorship index. *Nature* **352**(6332), 187–187 (1991)
14. Imperial, J., Rodríguez-Navarro, A.: Usefulness of Hirschs h-index to evaluate scientific research in Spain. *Scientometrics* **71**(2), 271–282 (2007)
15. Kelly, C., Jennions, M.: The h index and career assessment by numbers. *Trends in Ecology & Evolution* **21**(4), 167–170 (2006)
16. Kendall, M.: A new measure of rank correlation. *Biometrika* **30**(1/2), 81–93 (1938)
17. Mannella, R., Rossi, P.: On the time dependence of the h-index. *Journal of Informetrics* **7**(1), 176 – 182 (2013). DOI <http://dx.doi.org/10.1016/j.joi.2012.10.003>. URL <http://www.sciencedirect.com/science/article/pii/S1751157712000855>
18. Page, L., Brin, S., Motwani, R., Winograd, T.: The pagerank citation ranking: Bringing order to the web. In: *Proceedings of the 7th International World Wide Web Conference*. Stanford InfoLab (1998)
19. Pepe, A., Rodriguez, M.: An in-depth longitudinal analysis of mixing patterns in a small scientific collaboration network. *Scientometrics* **85**(3) (2010)
20. Petersen, A., Wang, F., Stanley, H.: Methods for measuring the citations and productivity of scientists across time and discipline. *Physical Review E* **81**(3), 036,114 (2010)
21. Radicchi, F., Markines, B., Vespignani, A.: Diffusion of scientific credits and the ranking of scientists. *Physical Review E* **80**(5), 056,103 (2009)
22. Rodriguez, M., Pepe, A.: On the relationship between the structural and socioacademic communities of a coauthorship network. *Journal of Informetrics* **2**(3), 195–201 (2008)
23. Schmidt, R.: A worksheet for authorship of scientific articles. *Bulletin of the Ecological Society of America* **68**(1), 8–10 (1987)
24. Schreiber, M.: To share the fame in a fair way, hm modifies h for multi-authored manuscripts. *New Journal of Physics* **10**(4), 040,201 (2008)
25. Sidiropoulos, A., Katsaros, D., Manolopoulos, Y.: Generalized Hirsch h-index for disclosing latent facts in citation networks. *Scientometrics* **72**(2), 253–280 (2007)
26. Slone, R.M.: Coauthors’ contributions to major papers published in the *ajr*: frequency of undeserved coauthorship. *AJR. American journal of roentgenology* **167**(3), 571–579 (1996)

27. Tang, J., Zhang, D., Yao, L.: Social network extraction of academic researchers. In: Proceedings of the 2007 Seventh IEEE International Conference on Data Mining, pp. 292–301. IEEE Computer Society, Washington, DC, USA (2007)
28. Taramasco, C., Cointet, J., Roth, C.: Academic team formation as evolving hypergraphs. *Scientometrics* **85**(3), 721–740 (2010)
29. Verhagen, J., Wallace, K., Collins, S., Scott, T.: QUAD system offers fair shares to all authors. *Nature* **426**(6967), 602–602 (2003)
30. Wang, C., Han, J., Jia, Y., Tang, J., Zhang, D., Yu, Y., Guo, J.: Mining advisor-advisee relationships from research publication networks. In: Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining, pp. 203–212. ACM (2010)
31. West, J.D., Jensen, M.C., Dandrea, R.J., Gordon, G.J., Bergstrom, C.T.: Author-level eigenfactor metrics: Evaluating the influence of authors, institutions, and countries within the social science research network community. *Journal of the American Society for Information Science and Technology* **64**(4), 787–801 (2013)
32. Wilcox, R.R.: Introduction to robust estimation and hypothesis testing. Academic Press (2012)
33. Yan, E., Ding, Y.: Applying centrality measures to impact analysis: A coauthorship network analysis. *Journal of the American Society for Information Science and Technology* **60**(10), 2107–2118 (2009)