

The H-index Paradox: Your Coauthors Have a Higher H-index than You Do

Accepted for publication in *Scientometrics*

Fabício Benevenuto ·
Alberto H. F. Laender ·
Bruno L. Alves

the date of receipt and acceptance should be inserted later

Abstract One interesting phenomenon that emerges from the typical structure of social networks is the *friendship paradox*. It states that your friends have on average more friends than you do. Recent efforts have explored variations of it, with numerous implications for the dynamics of social networks. However, the friendship paradox and its variations consider only the topological structure of the networks and neglect many other characteristics that are correlated with node degree. In this article, we take the case of scientific collaborations to investigate whether a similar paradox also arises in terms of a researcher's scientific productivity as measured by her H-index. The H-index is a widely used metric in academia to capture both the quality and the quantity of a researcher's scientific output. It is likely that a researcher may use her coauthors' H-indexes as a way to infer whether her own H-index is adequate in her research area. Nevertheless, in this article, we show that the average H-index of a researcher's coauthors is usually higher than her own H-index. We present empirical evidence of this paradox and discuss some of its potential consequences.

Introduction

One interesting phenomenon that emerges from the typical structure of social networks is the *friendship paradox* [6]. It states that, on average, your friends have more friends than you do. This paradox basically exists because of the discrepancy on node degree values in typical social networks [2], in which

Computer Science Department
Federal University of Minas Gerais
E-mail: {fabricio, laender, bruno.leite}@dcc.ufmg.br
Tel.: +5531-3409-5860
Fax: +5531-3409-5858

individuals with a high number of friends are over-represented when averaging over them [10]. As a consequence, the friendship paradox can dramatically skew an individual's local observation, making such an observation appear far more common than it is in reality [4, 14].

In this context, identifying variations of this paradox in different ecosystems has been the topic of some important recent research efforts [5, 9, 12]. For instance, two new paradoxes have been verified on Twitter [9]: (1) the *virality paradox* that states that your friends receive more viral content than you do, and (2) the *activity paradox* that states that your friends post more frequently than you do. More recently, the friendship paradox was generalized to any complex network [5] and its origins are highly correlated with the skewed distribution of node degree (i.e., the number of network friends) [10]. In a nutshell, these efforts suggest that any attribute that is highly correlated with node degree is likely to produce this kind of paradox [5]. Thus, in this article we take the case of scientific collaborations to investigate whether a similar paradox also arises in terms of a researcher's scientific productivity as measured by her H-index.

The H-index [8] is a metric originally proposed to measure a researcher's scientific output. Its calculation is quite simple as it is based on the researcher's set of most cited publications and the number of citations they have received. More specifically, a researcher has an H-index h if she has at least h publications that have received at least h citations. Thus, if a researcher has at least 10 publications with at least 10 citations, her H-index is 10.

Like any metric that attempts to summarize a complex and subjective evaluation in a single number, the H-index has its limitations, including being biased towards the researchers' scientific lifetime, not accounting for the number of coauthors in the publications and ignoring the distinct citation patterns across different areas [3]. Nevertheless, the H-index became popular as it provides a notion of both quality and quantity of a researcher's scientific output in a simple and easy-to-compute metric. As a consequence, researchers are often tempted to evaluate themselves based on the H-index. Systems like Google Scholar¹ and ArnetMiner² help researchers track their publication impact and coauthors, as well as to maintain their profiles, where the H-index is clearly stamped. Thus, it is natural to assume that researchers may use their coauthors' H-indexes as a way to estimate whether their own H-index is adequate in their respective research areas or within a department or university.

Despite recent efforts to generalize the friendship paradox [5], it is still unclear whether a similar paradox actually happens when we consider the H-index in a coauthorship network. However, we have been able to show that the average H-index of a researcher's coauthors is usually higher than her own H-index.

¹ <http://scholar.google.com/intl/en/scholar/citations.html>

² <http://arnetminer.org>

Next, we briefly discuss how we have estimated the H-index for researchers from distinct Computer Science research communities, and then provide empirical results that corroborates the existence of the *H-index paradox*.

Estimating H-index

In order to provide evidence of the H-index paradox, we need to be able to (1) identify the coauthors of a large set of researchers and (2) estimate the H-index of these researchers as well as of their respective coauthors.

We focus on constructing the coauthorship network of Computer Science researchers from different areas. To do that, we gathered data from DBLP³, as it offers its entire database in XML format for download. We gathered this data for those researchers who published in the flagship conferences of 10 major ACM SIGs (Special Interest Groups)⁴: SIGCHI, SIGCOMM, SIGCSE, SIGDOC, SIGGRAPH, SIGIR, SIGKDD, SIGMETRICS, SIGMOD and SIGPLAN.

There are several tools that measure the H-index of researchers, of which Google Scholar is today the most prominent one. However, in order to have a profile in this system, a researcher needs to sign up and explicitly create it. In a preliminary collection of part of the profiles of the DBLP authors, we found that less than 30% of these authors had a profile at Google Scholar [1]. Thus, this strategy would largely reduce our dataset.

To overcome this limitation, we used data from the SHINE (Simple HINdex Estimation) project⁵ to estimate the researchers' H-index. SHINE provides a website that shows the H-index of almost 1,800 Computer Science conferences. It was created based on a large scale crawl of Google Scholar. Its strategy consisted of searching for the title of all papers published in such conferences, thus effectively estimating their H-index based on the citations computed by Google Scholar. Although SHINE only allows searching for the H-index of conferences, their developers kindly allowed us to use its dataset to infer the H-index of researchers based on the citations received by their conference papers.

However, besides covering only conferences, SHINE does not track all existing conferences in Computer Science, which might cause the researchers' H-index to be underestimated when computed with this data. To investigate this issue, we compared the H-index of researchers with a profile on Google Scholar with their estimated H-index based on the SHINE data. For this, we randomly selected 10 researchers for each of the ACM SIG's flagship conferences and extracted their H-indexes from their respective Google Scholar profiles. In comparison with the H-index we estimated from SHINE, the Google Scholar values are, on average, 50% higher. Figure 1 shows the scatterplot for the two H-index measures. We can note, however, that although the SHINE H-index is lower, the two measures are highly correlated. The Pearson's correlation

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

⁴ <http://www.acm.org/sigs>

⁵ <http://shine.icomp.ufam.edu.br>

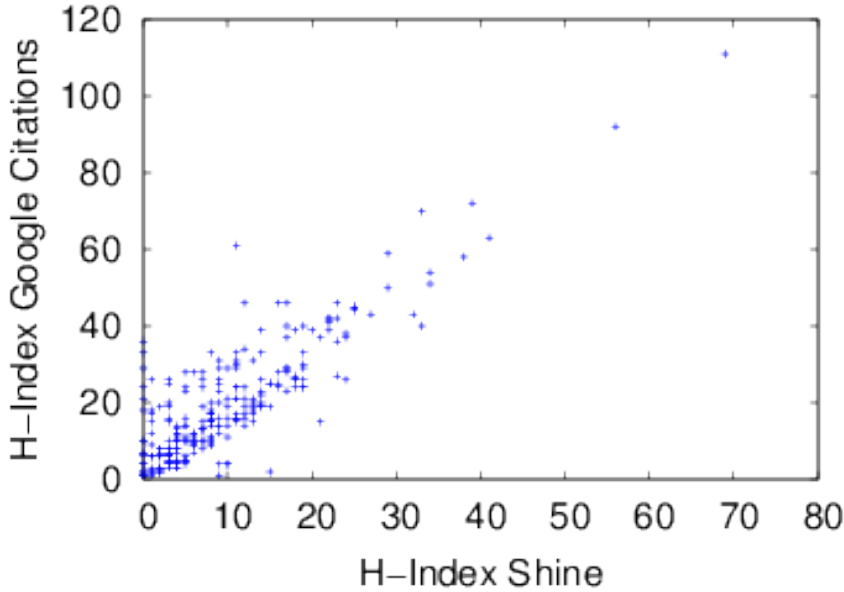


Fig. 1 Correlation between the inferred H-index and Google Citations one.

Table 1 The DBLP data of the 10 ACM SIG flagship conferences.

SIG	Conference	Period	H-Index	Authors	Publications	Editions
SIGDOC	SIGDOC	1989-2010	23	1071	810	22
SIGCHI	CHI	1994-2012	144	5095	2819	19
SIGIR	SIGIR	1978-2011	116	3624	2687	34
SIGKDD	KDD	1995-2011	124	3078	1699	17
SIGCOMM	SIGCOMM	1988-2011	140	1593	796	24
SIGCSE	SIGCSE	1986-2012	51	3923	2801	27
SIGGRAPH	SIGGRAPH	1985-2003	119	1920	1108	19
SIGMETRICS	SIGMETRICS	1981-2011	71	2083	1174	31
SIGPLAN	POPL	1975-2012	85	1527	1217	38
SIGMOD	SIGMOD	1975-2012	147	4202	2669	38

coefficient is 0.85, indicating that the H-index estimations are proportional in both systems.

Table 1 summarizes the collected data, including the SIG, the conference acronym, the period considered (some conferences had their period reduced to avoid a hiatus in the data), the conference’s SHINE H-index and the total number of authors, publications and editions. This dataset is useful to our purposes, since it allows us to investigate the H-index paradox on real Computer Science communities, in which researchers might tend to compare themselves with their peers.

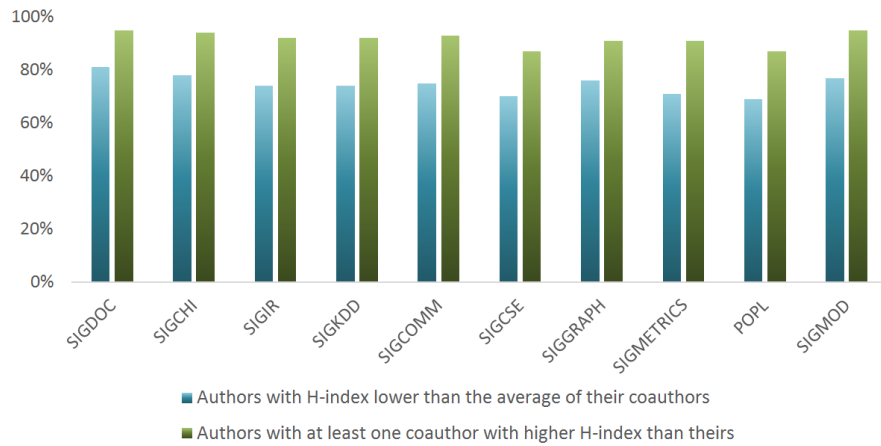


Fig. 2 Comparison results of a researcher H-index with her coauthors.

Comparing the H-index of a Researcher with her Coauthors'

Having estimated the H-index of each researcher, we can compare it with her coauthors'. Figure 2 shows the fraction of authors with an H-index that is lower than the average of their coauthors for the 10 conferences we have considered. We note that even focusing on authors that have published in flagship conferences of ACM SIGs, the fraction of authors that are below average is quite high for all research communities analyzed, varying from 69% (POPL) to 81% (SIGDOC). When we look at the percentage of authors with at least one coauthor with a higher H-index than theirs, the numbers are higher than 90% for most of the conferences.

These results confirm the H-index paradox since one's coauthors in a research community have, on average, a higher H-index than hers. The reasons behind the H-index paradox might be explained by the high correlation between node degree and H-index in a research community. Usually, high degree nodes tend to be senior researchers that not only advise a large number of students but also establish more collaborations, often with different groups along their career [1]. To further investigate this issue, Figure 3 shows the distribution of the number of authors as a function of the H-index. It clearly resembles a long tail distribution, thus suggesting that some authors disproportionately contribute to the average H-index. This disproportion on the average H-index might be even sharper with the typical structural properties of coauthorship networks, which are similar to many social networks [11,13], i.e., they have a long tail degree distribution, in which highly connected authors create bridges across multiple highly connected components, leading to the properties of high clustering coefficient and short diameter. Finally, we measured the Pearson's coefficient correlation between a researcher's H-index and her degree. Such correlation is 0.36 a value that, although not very high, is positive, thus sug-

gesting that a small number of researchers simultaneously have a high H-index and a large number of connections in the network.

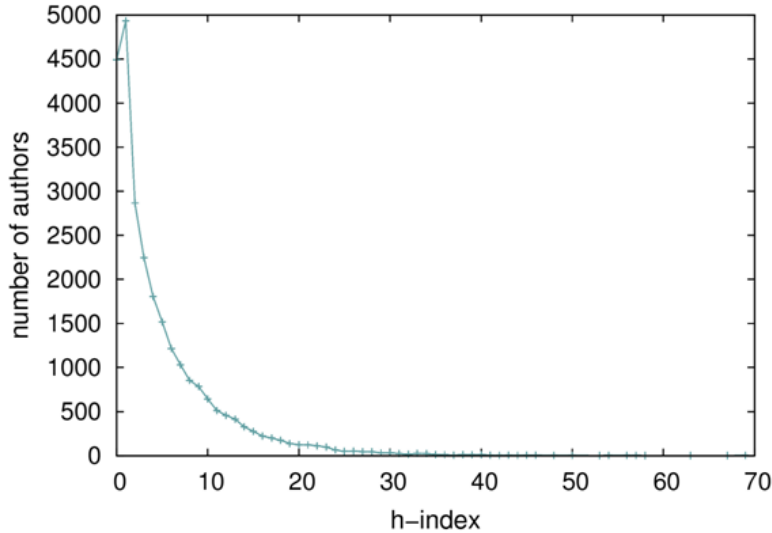


Fig. 3 Distribution of authors according to their H-indexes.

Conclusions

In this article we have analyzed a variation of the well-known friendship paradox. By analyzing the average H-index of a researcher’s coauthors for different Computer Science research communities, we show that the H-index paradox arises because the H-index is positively correlated with node degree. One of the implications of the friendship paradox is the fact that it leads to systematic biases in our perceptions. Thus, similarly, the H-index paradox induces researchers to feel that they rank below average in comparison with their coauthors. This phenomenon is an instantiation of a sensation that occurs in different scenarios and is popularly captured by an expression that is common to many languages and cultures: *the grass is always greener on the other side of the fence* [7].

Acknowledgments

This research was partially funded by InWeb - The Brazilian National Institute of Science and Technology for the Web (MCT/CNPq/FAPEMIG grant 573871/2008-6), and by the authors’ individual grants from CAPES, CNPq and FAPEMIG.

References

1. Alves, B.L., Benevenuto, F., Laender, A.H.F.: The Role of Research Leaders on the Evolution of Scientific Communities. In: Proceedings of the 22nd International Conference on World Wide Web (companion volume), pp. 649–656 (2013)
2. Barabási, A.L., Albert, R.: Emergence of scaling in random networks. *Science* **286**(5439), 509–512 (1999)
3. Bornmann, L., Daniel, H.D.: Does the h-index for ranking of scientists really work? *Scientometrics* **65**(3), 391–392 (2005)
4. Centola, D.: The spread of behavior in an online social network experiment. *Science* **329**(5996), 1194–1197 (2010)
5. Eom, Y.H., Jo, H.H.: Generalized friendship paradox in complex networks: The case of scientific collaboration. *Scientific Reports* **4** (2014)
6. Feld, S.L.: Why your friends have more friends than you do. *American Journal of Sociology* **96**(6), 1464–1477 (1991)
7. Giansante, S., Kirman, A., Markose, S., Pin, P.: The grass is always greener on the other side of the fence: The effect of misperceived signalling in a network formation process. In: *Artificial Markets Modeling*, pp. 223–234. Springer (2007)
8. Hirsch, J.E.: 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), 16,569–16,572 (2005)
9. Hodas, N.O., Kooti, F., Lerman, K.: Friendship Paradox Redux: Your Friends Are More Interesting Than You. In: *Proceedings of the International Conference on Web and Social Media*, pp. 8–10 (2013)
10. Hodas, N.O., Kooti, F., Lerman, K.: Network Weirdness: Exploring the Origins of Network Paradoxes. In: *Proceedings of the International Conference on Web and Social Media*, pp. 8–10 (2014)
11. Huang, J., Zhuang, Z., Li, J., Giles, C.L.: Collaboration Over Time: Characterizing and Modeling Network Evolution. In: *Proceedings of the 2008 International Conference on Web Search and Data Mining*, pp. 107–116 (2008)
12. Lerman, K., Yan, X., Wu, X.Z.: The Majority Illusion in Social Networks. *arXiv preprint arXiv:1506.03022* (2015)
13. Mislove, A., Marcon, M., Gummadi, K.P., Druschel, P., Bhattacharjee, B.: Measurement and Analysis of Online Social Networks. In: *Proceedings of the 7th ACM SIGCOMM Conference on Internet Measurement*, pp. 29–42 (2007)
14. Salganik, M.J., Dodds, P.S., Watts, D.J.: Experimental study of inequality and unpredictability in an artificial cultural market. *Science* **311**(5762), 854–856 (2006)