

# **Some aspects of H-index and applications**

Thesis submitted in partial fulfillment  
of the requirements for the degree of

*(Master of Science in **Computer Science and Engineering** by Research*

by

AASHAY SINGHAL

201502112

aashay.singhal@research.iiit.ac.in



International Institute of Information Technology  
Hyderabad - 500 032, INDIA  
August 2023

Copyright © AASHAY SINGHAL, 2023  
All Rights Reserved

International Institute of Information Technology  
Hyderabad, India

## **CERTIFICATE**

It is certified that the work contained in this thesis, titled “Generalised aspects of H-index and their applications ” by Aashay Singhal, has been carried out under my supervision and is not submitted elsewhere for a degree.

---

Date

---

Adviser: Prof. Kamalakar Karlapalem

To my advisor and family

## **Acknowledgments**

I would like to take this opportunity to express my sincere thanks and appreciation to all those who have contributed to the successful completion of this thesis.

First and foremost, I would like to express my gratitude to my thesis advisor Prof. Kamal, for his invaluable guidance, support, and patience throughout the entire process. His knowledge, expertise, and insightful comments have been instrumental in shaping the direction of my research and improving the quality of my work. I would have not been able to complete this thesis without his guidance and support.

I would like to acknowledge the support and encouragement of my family and friends, who have been a constant source of motivation and inspiration throughout my academic journey. I am grateful to my father who is always there to support me. Their unwavering support and belief in me have been invaluable and have kept me motivated to pursue my goals. Thank you Goel, Deva and Shubhangi.

Once again, I am deeply grateful to all those who have contributed to the completion of this thesis, and I am honored to have had the opportunity to work with and learn from such exceptional individuals.

## Abstract

The h-index has become a widely used metric for evaluating the research impact of scholars, particularly in academia. It attempts to measure both the productivity and impact of an author's research output by taking into account the number of publications and the number of citations that they have received. However, there are ongoing debates about the accuracy and reliability of the h-index as a measure of research impact. For instance, the h-index can be influenced by factors such as the discipline of the author, the age of their career, and the citation practices of the field. As a result, there have been several attempts to develop alternative metrics that provide a more comprehensive picture of an author's research impact. This thesis seeks to contribute to this ongoing conversation by examining the h-index in detail and exploring its strengths and limitations, as well as its applications in various fields.

There are several variations of h-index being proposed. For example, g-index [18] and h(2)-index[30] which give more weight to highly cited papers, a-index [28], m-index [12] and many more. While these variations provide more nuanced insights into a researcher's impact and productivity, it is important to note that each metric has its own limitations and should be used in conjunction with other measures to provide a comprehensive evaluation of a researcher's output.

The first part of the thesis discusses the generalisation of the h-index. We explore the definition of h-index on a general graph. The aim of this research is to develop a framework for the h-index on a graph. By achieving this, we can apply h-index on number of different domains. Hence, we present a general definition of the h-index on a graph and demonstrate its applicability in various fields. By applying this framework to different contexts, we show that the h-index can reveal underlying semantics and patterns in a network.

The rest of the thesis presents case studies that illustrate the use of the h-index and other derived metrics in various fields. For example, the h-index has been used to evaluate the research impact of single publications, airports, and scientists. These case studies demonstrate the versatility of the h-index as a measure of research impact and highlight the importance of considering the specific context in which it is being used.

Overall, this thesis contributes to a better understanding of the h-index and its applications in different fields. It also highlights the limitations of the h-index and the need for complementary metrics to provide a more comprehensive picture of an author's research impact. This thesis provides a valuable resource for researchers and practitioners who use bibliometric measures to evaluate research impact.

# Contents

Chapter	Page
1 Introduction . . . . .	1
2 Related Work . . . . .	4
2.1 H-index . . . . .	4
2.2 H-index on research papers . . . . .	5
2.3 Airport network analysis . . . . .	6
3 H-index on a graph . . . . .	8
3.1 Introduction . . . . .	8
3.2 Properties of H-index . . . . .	9
3.3 h-index and its variants on a graph . . . . .	14
3.3.1 h-index . . . . .	14
3.3.2 g-index . . . . .	14
3.3.3 h(2)-index . . . . .	15
3.3.4 a-index . . . . .	15
3.3.5 m-index . . . . .	15
3.3.6 r-index . . . . .	16
3.3.7 Wu's w-index . . . . .	16
3.3.8 h(5,2)-index . . . . .	16
3.4 h-index algorithm . . . . .	16
3.4.1 Basic algorithm . . . . .	17
3.4.2 Linear algorithm . . . . .	18
3.5 Summary . . . . .	19
4 H-index and its variants on research papers . . . . .	20
4.1 Introduction . . . . .	20
4.2 Data . . . . .	21
4.2.1 Data Collection . . . . .	22
4.2.2 Dataset Description . . . . .	22
4.2.3 Benchmark Dataset . . . . .	23
4.3 Experiments . . . . .	23
4.3.1 Correlation amongst indices . . . . .	24
4.3.2 Rank Biased Overlap (RBO) . . . . .	24
4.3.3 Performance of indices in predicting trends of awardees . . . . .	25
4.3.4 Performance of indices over time . . . . .	26

4.4	Results . . . . .	26
4.4.1	Correlation amongst indices . . . . .	27
4.4.2	Rank Biased Overlap (RBO) . . . . .	28
4.4.3	Performance of indices in predicting trends of awardees . . . . .	28
4.4.4	Performance of indices over time . . . . .	30
4.5	Conclusion . . . . .	31
5	H-index of authors based on H-index of their papers . . . . .	34
5.1	Introduction . . . . .	34
5.2	H-index of author . . . . .	35
5.2.1	h-index and h-frac-index . . . . .	35
5.2.2	hp-index and hp-frac-index . . . . .	35
5.2.3	Example . . . . .	36
5.3	Citation Data . . . . .	38
5.4	Experiments and results . . . . .	39
5.4.1	Correlation . . . . .	39
5.4.2	RBO . . . . .	39
5.4.3	Awarded researchers . . . . .	39
5.4.4	Further Analysis . . . . .	41
5.5	Conclusion . . . . .	43
6	H-index on airport network . . . . .	47
6.1	Introduction . . . . .	47
6.2	Data and Methods . . . . .	47
6.2.1	Airport network . . . . .	48
6.2.2	Metrics . . . . .	48
6.3	Experiments and Analysis . . . . .	50
6.3.1	out-degree ranks . . . . .	50
6.3.2	h-index ranks . . . . .	51
6.3.3	hh-index ranks . . . . .	51
6.3.4	Correlation . . . . .	53
6.3.5	Rank Biased Overlap . . . . .	55
6.4	Further Analysis . . . . .	56
7	Conclusions . . . . .	57
	Bibliography . . . . .	60



## List of Figures

Figure	Page
3.1 Example of a directed acyclic graph . . . . .	9
3.2 Graph with minimum out-neighbors for node $p$ . . . . .	9
3.3 Smallest graph required for node $p$ to have h-index value as $h$ . . . . .	10
3.4 Graph with maximum nodes having h-index as one . . . . .	11
3.5 Graph with maximum nodes having zero h-index . . . . .	11
3.6 Example graphs for maximum $ S $ case when (a) $k = 1$ (b) $k = 2$ and (c) $k = 3$ . . . . .	12
3.7 Inductive proof of theorem 2 case 2 . . . . .	13
3.8 Resultant graph for Corollary of Theorem 2 . . . . .	13
3.9 Example graph for h-index variants . . . . .	14
4.1 Correlation matrix for SIGMOD conference . . . . .	26
4.2 Correlation matrix for VLDB conference . . . . .	27
4.3 RBO matrix for SIGMOD conference . . . . .	28
4.4 RBO matrix for VLDB conference . . . . .	29
4.5 Index name vs Number of awarded papers ranked in top 5% . . . . .	29
4.6 Index name vs Number of awarded papers ranked in $< 5\%$ , $5\% - 10\%$ , $10\% - 20\%$ and $20\% - 30\%$ . . . . .	30
4.7 Number of years since publishing vs Number of papers in top 5% when ranked on the particular index . . . . .	31
4.8 Number of years since publishing(beyond 10 years) vs Number of papers in top 5% when ranked on the particular index . . . . .	32
4.9 Plot showing number of years since publishing vs H-index and the corresponding rank as per h-index for a few awarded papers. Each of the graphs represents the mentioned curve for a single awarded paper. . . . .	33
5.1 Example graph for h, hp, h-fraction, hp-fraction demonstration . . . . .	36
5.2 Plot showing number of years since publishing vs H-index and the corresponding rank as per h-index for a few awarded papers. . . . .	40
6.1 A sub graph of the airline network . . . . .	48
6.2 Example graph for airport network . . . . .	49
6.3 Correlation of different indices for (a) top 100 airports, (b) top 1000 airports, (c) all airports . . . . .	54
6.4 Rank Biased Overlap of the three indices . . . . .	55

## List of Tables

Table	Page
3.1 Definitions of h-index and its variants . . . . .	15
3.2 Index values for graph in Fig.3.9 . . . . .	16
3.3 Comparison of run times (in seconds) for each algorithm . . . . .	19
4.1 Dataset Description . . . . .	23
4.2 Number of awarded papers . . . . .	23
4.3 Example of set overlap calculation . . . . .	24
5.1 H-index of the papers in given example in Fig. 5.1 . . . . .	37
5.2 Dataset description . . . . .	38
5.3 Percentage of awardees in different ranges of ranked lists as per each index across three fields . . . . .	41
5.4 Percentage of awardees given the highest rank as per each index across the three fields (see Table 5.8, 5.10, and 5.11 . . . . .	42
5.5 Average and maximum values of the calculated metrics . . . . .	42
5.6 Correlation between <i>diff1</i> , <i>diff2</i> , and the average number of co-authors . . . . .	43
5.7 List of awards collected . . . . .	44
5.8 List of award winners with ranks for Biology (highest ranks in bold) . . . . .	44
5.9 List of top 10 authors ranked by hp-frac-index (awarded researchers are in bold) . . . . .	45
5.10 List of award winners with ranks for Computer Science (highest ranks in bold) . . . . .	45
5.11 List of award winners with ranks for Economics (highest ranks in bold) . . . . .	46
6.1 Details about airport network graph . . . . .	48
6.2 The values of out-degree, h-index and hh-index for the example graph . . . . .	50
6.3 Top 10 busiest airport according to Wikipedia . . . . .	51
6.4 Top 10 airports ranked by out-degree . . . . .	51
6.5 Top 10 airports ranked by h-index . . . . .	52
6.6 Top 10 airports ranked by hh-index . . . . .	52
6.7 Connectivity of airports with maximum rank boost by hh-index . . . . .	53
6.8 Percentage of micro hubs connecting to hubs . . . . .	53
6.9 Connectivity of airports with minimum rank boost by hh-index . . . . .	55
6.10 h-index computed by considering only hubs and hubs + micro hubs for top 20 airports. Third column denotes the difference between first two columns. . . . .	56

## *Chapter 1*

### **Introduction**

In the world of academia, the h-index has become a popular measure of an author's research impact. Introduced by Jorge Hirsch in 2005, the h-index reflects both the productivity and impact of a researcher's work by taking into account the number of publications and the number of citations they have received. The h-index is calculated by determining the number of articles published by a researcher and the number of citations each article has received. A researcher has an h-index of  $h$  if  $h$  of their articles have been cited at least  $h$  times. The h-index is a useful tool for both researchers and academic institutions because it provides a quantitative measure of a researcher's scholarly output and impact. The h-index is particularly valuable when comparing researchers from different fields or when evaluating a researcher's entire career, as it accounts for both the number of publications and the quality of those publications.

The h-index has become a popular tool for evaluating the research impact of scholars because it is relatively simple to calculate and provides a single number that can be compared across researchers. However, the h-index has some limitations, and it should not be the sole measure of a researcher's productivity or impact. For example, the h-index does not take into account the context of citations, such as whether they are from highly respected journals or from less reputable sources. Additionally, the h-index can be affected by factors outside of a researcher's control, such as the size of their research community or the time period in which their work was published. Additionally, the h-index is heavily influenced by a researcher's most highly cited papers, which may not necessarily reflect the overall impact of their work. The h-index is discipline-specific and does not account for the varying citation practices across different fields, which may result in unfair comparisons between researchers working in different disciplines. Nevertheless, despite its limitations, the h-index remains a widely used metric in academia for evaluating the research output and impact of scholars.

There are several variations of the h-index that have been proposed to address its limitations. One such variation is the g-index [18], which takes into account the number of highly cited papers a researcher has published. Another variation is the m-index [12], which measures the productivity of a researcher by dividing the h-index by the number of years since their first publication. Other variations include the a-index [28], which measures the average number of citations of papers in the Hirsch

core (which consists of the first  $h$  papers of an author when sorted on number of citations and  $h$  is the h-index of that author). These variations provide more nuanced insights into a researcher’s impact and productivity, but they also have their own limitations and should be used in conjunction with other metrics.

In this thesis, we explore the generalisation of the h-index to different domains such as citation networks of authors, network of papers connected by citations and airport connectivity network. In other words, we apply h-index to different graphs, where nodes represent entities and edges represent relationships between them. Our objective is to develop a comprehensive and consistent framework for the h-index on a graph that can be applied to different domains. To achieve this, we present a general definition of the h-index on a graph and demonstrate its applicability in various fields. We show that the h-index can reveal underlying semantics and patterns in a network of entities and their relationships.

Overall, our work also contributes to the growing body of research on bibliometric analysis. It provides a new perspective on the h-index, which can be used to evaluate the impact and productivity of entities in different domains. We believe that our findings will be useful for researchers and practitioners interested in network analysis, and that they will help to advance our understanding of the structure and dynamics of complex systems.

In particular, we formally define the h-index on a general graph. Given the definition we discuss some properties and inequalities. These include inequalities between number of nodes and h-index. The properties can be useful in application and calculation of h-index. Further, we extend this generalisation to other variants of h-index, namely, g-index [18], h(2)-index [30], a-index [28], m-index [12], r-index [28], w-index [45] and h(5,2)-index [21]. We then discuss three algorithms for h-index calculation and also compare their run times on sample graphs.

After detailing the general definition of h-index we apply it to evaluate individual publications. The h-index is typically used to evaluate the overall productivity and impact of a researcher’s work, but it is also possible to calculate the h-index for individual publications. This can be a useful tool for evaluating the impact of a particular article or book within a field. To calculate the h-index of a single publication, you would need to determine how many times that publication has been cited, and then identify the number  $n$  of other publications that have been cited at least  $n$  times. If the publication in question has an h-index of  $n$ , this means that it has been cited at least  $n$  times, and there are at least  $n$  other publications that have been cited at least  $n$  times. Calculating the h-index of single publications can be useful for researchers who want to understand the impact of their work, or for publishers and editors who want to evaluate the quality of submissions.

In order to understand whether h-index is a good evaluation metric for single publications, we compared it with nine variations of h-index. We collected the papers for VLDB and SIGMOD conference over the years and necessary data required (like their citations, year of publishing, etc.) to calculate the metrics on each of the paper. After calculating each metric value on each paper, we rank them by the metrics. We also gather the awarded papers in each of the two conferences. After comparing ranks

of the awardees, our results show that h-index is the best in ranking awardees at the top. Among our experiments, we also discuss the correlation and overlap of these metrics.

Now that we have shown that h-index is the best metric to evaluate research papers amongst other metrics, the next natural step is to use this h-index to compute the impact of a researcher. Here, we apply four h-index like metrics to evaluate researchers, namely, h-index, h-frac-index, hp-index and hp-frac-index. Amongst these four metrics three of them have been proposed in prior works [24, 29, 20] but hp-frac-index is a newly proposed metric in our work. The traditional h-index uses number of citations of the papers published by an author to calculate the h-index of an author. In hp-frac-index, instead of number of citations we use h-index of the published papers to calculate h-index of an author. Further, we divide the h-index of each paper by the number of authors.

We collect top 1000 researchers in the field of Computer Science, Economics and Biology. Then we gather all the papers published by them and other necessary data required to calculate each of the four metrics. We also gather the list of awarded researchers amongst the list we obtained. Then we calculate the four metrics on each researcher and rank them. On comparing the ranks given by each metric, we find that hp-frac-index is better than the other metrics in ranking the awardees at top. The hp-frac-index is a reliable means of assessing the influence of researchers. hp-frac-index is resistant to manipulation and can capture individual contributions effectively. These combined factors make the hp-frac-index superior in ranking authors.

Lastly, we apply the generalization of the h-index to airport networks, which will allow us to gain valuable insights into the functioning and performance of airports. The use of airport networks has been increasingly relevant due to the growth of air traffic and globalization. Understanding the relationships and dynamics within an airport network is crucial for optimizing resources, improving efficiency, and ultimately enhancing the overall performance of the system.

The proposed generalization of the h-index for airport networks provides a novel approach to evaluating airport performance. By considering both the connectivity and the influence of the airport within the network, we can uncover noteworthy insights that are not captured by traditional measures. These insights could include identifying key airports that are crucial for the functioning of the network, understanding the impact of disruptions, and assessing the overall resilience of the system.

These applications demonstrate the versatility and usefulness of the h-index as a tool for uncovering underlying semantics and evaluating performance in various domain specific graphs. By providing a standardized measure that takes into account both quantity and quality, the h-index has become a valuable tool for researchers and practitioners alike. In this thesis, we explore the application of the h-index to a novel context and demonstrate its potential for uncovering new insights and informing decision-making processes.

The outline of the thesis is as follows: Chapter 2 consists of prior related work, chapter 3 discusses the generalisation of h-index and its properties. In chapter 4, we apply the generalised h-index on papers. In chapter 5, we propose h-index of authors based on h-index of papers. Lastly, chapter 6 describes the application of h-index in airport networks.

## *Chapter 2*

### **Related Work**

#### **2.1 H-index**

The h-index is a widely-used metric for assessing the scientific impact of researchers. It was introduced by Jorge Hirsch in 2005 [24] and has since become a popular way to quantify an individual's research output. The idea behind the h-index originates from scholars' longstanding aspiration to measure the impact of scientific work. In this section, we review some of the related work on the h-index and its applications.

The h-index has been applied in various fields and has become a standard tool for evaluating researchers' scientific output. For example, it is commonly used in academic hiring and promotion decisions, as well as in grant applications and award nominations. Additionally, the h-index has been used to analyze the dynamics of citation networks and to study the relationship between scientific output and career advancement.

Jorge Hirsch's original paper introduced the h-index as a metric that combines the number of publications and the number of citations received by each publication. The advantage of using the h-index is that it measures scholarly impact in a way that recognizes both the quality, and quantity of research by the individual. This simplifies the characterization of researchers' scientific output to great extent. Costas et al. [14] mentioned other good properties of the h-index. For example, it is an objective indicator and therefore, it may play an important role when making decisions about promotions, fund allocation and awarding prizes. Vanclay [43] pays attention to another interesting benefit of the h-index: it is a robust evaluation of impact as it disregards low-citation or no-citation articles, but at the same time does not overvalue high-citation articles in computation. The data needed for computing this index is easily available through different databases like Scopus, and Google Scholar. Furthermore, it does not require tuning thresholds or parameters and it is easily interpretable.

However, the h-index has a number of drawbacks [11, 28], such as the potential influence of self-citation [22], the inclusion of articles whose conclusions are later disproved, and the failure of the metric to credit highly impactful articles that receive far more citations than others. According to Egghe [17], "As the h index is defined now, once an article belongs to the h-defining class, it is totally unimportant

whether or not these papers continue to be cited and, if cited, it is unimportant whether these papers receive 10, 100, or 1000 more citations”. Also, the h-index is highly dependent on the number of years of active research. Moreover, a researcher’s h-index value cannot decline over time. Researchers that do not publish any more papers, or are inactive can maintain the same value of the h-index. Finally, the h-index is time-dependent. It usually takes some amount of time for papers to be fully appreciated before they are cited in other papers. Therefore, older articles and scholars with longer time to accumulate h-index tend to benefit.

Using a single, easy-to-compute indicator poses a risk of indiscriminate use, such as relying solely on it to evaluate scientists. Research performance is a multifaceted and complex endeavor that cannot typically be adequately assessed by a single indicator alone, as noted by [34]. This problem is extended to the evaluation of journals or general research activities using the h-index or similar indicators, as highlighted by [31].

The h-index has been employed by several authors to directly compare the scientific output of researchers. For instance, Hirsch [24] initially employed it to compare well-known physicists, while Schreiber [39] analyzed the h-index for 26 physicists. Imperial et al. [26] utilized the h-index to evaluate research by different authors in multiple areas of Biological Sciences. Oppenheim et al. [15] and Cronin et al. [15] employed the h-index to rank influential scientists. In Bornmann et al. [9, 10] the h-index was examined and employed to evaluate post-doctoral research fellowship applicants. Finally, Salgado and Paez et al. [38] explored scientific productivity in the field of Spanish social psychology.

There are different variations of the h-index that have been proposed to address some of its limitations. For example, the g-index was introduced by Egghe [18] to give more weight to highly cited publications, while the m-index was proposed [12] as the median number of citations received by papers in the Hirsch core. The hg-index [7] was proposed as the geometric mean of a researcher’s h and g indices. The a-index [28] is defined as the average number of citations of papers in the Hirsch core. Hirsch core for a researcher is the set of most cited h papers where h is the h-index of the researcher.

Efforts have also been made to modify the h-index to evaluate research performance across different countries. Guan et al. [23], for instance, compared and assessed the research performance of China in the field of bioinformatics using the h-index and compared it to other countries like the USA, UK, Germany, and others. Similarly, Berhidi et al. [16] created ranked lists of countries worldwide based on their h-index in various scientific fields. Their analysis revealed that EU countries held strong positions in each field but none of them could successfully compete with the USA.

## **2.2 H-index on research papers**

Schubert [40] considered both the direct impact and the indirect impact of citations. He argues that they can reflect unique aspects of the quality of an individual article. He proposed utilizing the paper-level h-index to evaluate a single publication, which would comprehensively reflect both the direct and indirect influence of the said publication. This index is an extension to measure the direct impact of

highly cited publication as well as its indirect influence through the citing papers. Later, Egghe [20] used the single publication h-index to calculate h-index of a researcher. They also discuss formulae for these impact measures in the Lotkaian context. Egghe [19] explained the relationship between a single publication's h-index and its total citations. The resulting relation was a concavely increasing power law derived using the Lotkaian model. Thor et al. [42] built a web application to calculate and show the single publication h-index and related performance measures for publications indexed by Google Scholar. Yan et al. [46] did an empirical study on h-index and various adaptations to evaluate if these indices behave the same in assessing a single publication. They argue that indices which are neither too near to nor too far from the h-index (in terms of correlation) could be much more promising than others. In addition, Bornmann et al. [13] conducted an empirical study to prove the singnificance of the single publication h-index. They collected articles submitted to Angewandte Chemie International Edition (AC-IE) conference which were either accepted or rejected and then accepted at some other journal. The results of their analysis showed that editorial decisions are correlated with the h-index values. This study confirmed that using the h-index for assessing single publications is effective. However, they mention that further studies are required in order to prove the assumption that h-index is useful for the group of the most cited papers.

In our work, we use the work by Schubert [40] to apply h-index on a single publication and compare them with other h-index based metrics. Yan et al. [46] did a similar empirical study to chapter 4 where they compared different h-index based metrics on single publication. They used around 300 papers for this comparison whereas our work covers almost 8000 papers. We also run experiments to compare the ranks of awarded papers. Later in our work, we use [20] to define h-index of a author based on h-index of a paper.

## 2.3 Airport network analysis

In recent years, researchers have begun modeling and analyzing air routes as a complex network [32]. Bagler [8] examined the airport network in India. They found that Airport network of India, despite being small in size, has complex dynamics similar to those of bigger air transportation networks. Jia et al. [27] studied the changes in the US airport network from 1990 to 2000. They presented an argument that the airport network plays a crucial role in US urban and regional development. Wang et al. [25] analyzed the network structure of major airports in China. Their findings suggest that socioeconomic indicators, such as passenger numbers, population, and GDP are highly correlated with the network centrality. Paleari et al. [35] compared the connectivity of air transportation in China, Europe, and the United States in terms of service provision to passengers. They found that China provides the fastest route, Europe has the highest quality level, and US is the most coordinated network.

According to Pere et al. [41], London's Heathrow Airport and regional airports in the United Kingdom play a significant role in national connectivity. The authors noted that while low-cost airlines in provincial cities of the UK have expanded their reach to numerous European cities, they still lack long-



haul routes. Additionally, they observed that Heathrow Airport is losing markets to Amsterdam and Dubai. In a separate study, Oriol et al. [33] evaluated the robustness of three aviation alliances and determined that Star Alliance had the most resilient route network, followed by Sky Team and Oneworld.

To the best of our knowledge, prior work has not applied h-index on airport network analysis. We apply our generalisation of h-index to airport network to reveal noteworthy insights and semantics like categorising hub airports and micro hub airports.

## Chapter 3

### H-index on a graph

#### 3.1 Introduction

Consider a graph  $G(V, E)$  where  $V$  is the set of nodes and  $E$  is the set of directed edges between two nodes. The edge from  $e_i$  to  $e_j$  is called an out edge from  $e_i$ . Also,  $e_j$  is an out-neighbor of  $e_i$ . We will use these terminologies through across the work.

The H-index is defined as, “A scientist has index  $h$  if  $h$  of their  $N_p$  papers have at least  $h$  citations each, and the other  $(N_p - h)$  papers have no more than  $h$  citations each”, according to Hirsch [24]. Here a scientist is a source of knowledge or information, their papers are the outputs produced using the source of information. These output papers of a scientist act as the source of information for the papers citing them. For example, consider an author  $A$  who has published a paper  $p_1$  and  $p_1$  is cited by one paper  $q_1$ . Then  $p_1$  has used  $A$  as the information source and  $q_1$  has further used  $p_1$  as the information source.

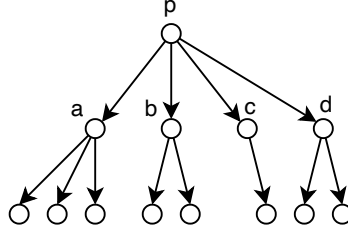
Therefore, we can consider the general definition of h-index as: “A source has index  $h$  if  $h$  of its total  $N_p$  out-neighbors have at least  $h$  further out-edges and the other  $(N_p - h)$  out-neighbors have no more than  $h$  further out-edges each.” Now, this definition can be applied to define the h-index of an author, paper, or any source of information. Thus, it can be applied to any graph where each node is a source of information and the edge from node  $A$  to  $B$  is an out-edge from  $A$  to  $B$ .

In this chapter, we will cover the definitions, properties and algorithms of h-index on a directed acyclic graph.

**Definition 1.** A directed acyclic graph (DAG) is a graph (i.e. a set of objects connected together) where each edge has a direction and there are no cycles in the graph.

In this chapter, we use  $G$  to denote this DAG. Furthermore,  $U$  is the set of edges and  $V$  is the set of vertices of graph  $G$ .

**Definition 2.** The h-index of a node  $p$  in a directed acyclic graph is equal to  $h$  if  $h$  is the largest natural number such that  $p$  has at least  $h$  out-neighbors each with at least  $h$  out-degree. Here, out-neighbors of  $p$  denote the neighbors of  $p$  having an out-edge from the node  $p$ .



**Figure 3.1:** Example of a directed acyclic graph

For example, consider Fig.3.1,  $p$  has four out-neighbors as  $[a, b, c, d]$  with  $[3, 2, 1, 2]$  out-degree respectively. Clearly,  $p$  has a maximum of two out-neighbors with at least 2 out-degree. Therefore, the h-index of  $p$  is two.

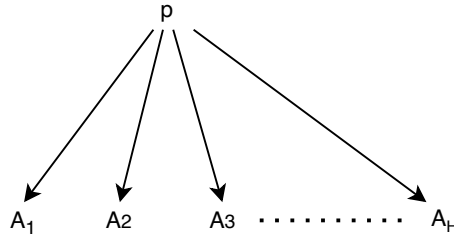
## 3.2 Properties of H-index

**Theorem 1.** *If there exists a node  $p$  with h-index  $h$ , then  $G$  must have*

$$N \geq 2h + 1 \quad (3.1)$$

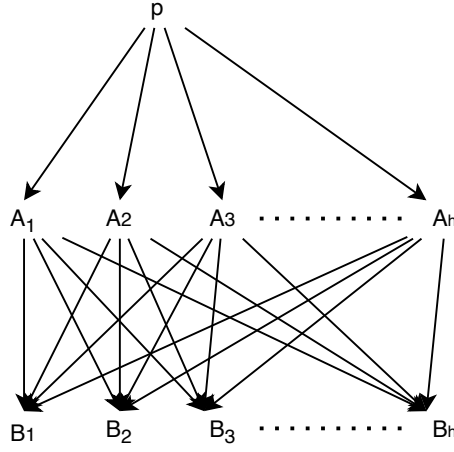
where  $N$  is the number of nodes in the graph

*Proof.* Since node  $p$  has an h-index of  $h$ , it will have at least  $h$  children which in turn will have  $h$  children each. This arises from the definition 2. The graph would look like below:



**Figure 3.2:** Graph with minimum out-neighbors for node  $p$

Each of the nodes  $A_1, A_2, A_3, \dots, A_h$  will have at least  $h$  out-degree. In order to minimise the total number of nodes, the  $h$  children of  $p$  should all have the same set of  $h$  out-neighbors.



**Figure 3.3:** Smallest graph required for node  $p$  to have h-index value as  $h$

The above construction shows that,  $N$  should at least be  $H + H + 1$

□

**Corollary.** If  $G$  has  $N$  nodes then the maximum possible h-index of any node is

$$H \leq (N - 1)/2 \quad (3.2)$$

**Theorem 2.** If  $G$  has  $N$  nodes, let the set of nodes having h-index equal to  $k$  (where  $0 \leq k \leq (N-1)/2$ ) be  $S$ . The following is true:

$$|S| = \begin{cases} [0, N] & k = 0 \\ [2, N - 2k] & k \geq 1 \end{cases}$$

**Case 1.**  $k = 0$

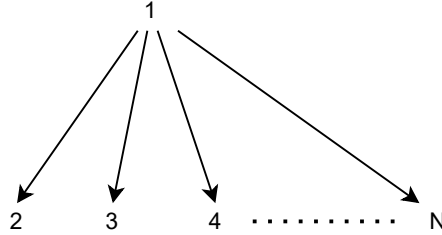
In this case, in order to minimise  $|S|$ , we need to maximise the number of nodes with at least one h-index. Observe that the nodes with no out-neighbors will always have h-index as zero as they have a maximum of zero out-neighbors with at least 0 out-degree. Further their immediate ancestors will also have the h-index as zero. Therefore, we need to minimise the number of leaves and their ancestors. The graph then, should look like this:



**Figure 3.4:** Graph with maximum nodes having h-index as one

The leaves in a graph are the set of nodes with zero out-degree. Ancestors of any node  $n$  are the set of nodes that have an incoming edge towards  $n$ .

Now, in this case, to maximise  $|S|$ , we need to maximise the number of nodes with h-index as zero. As shown below, we can have all the nodes have an h-index of zero. In the minimum case above, we showed that leaf nodes and their immediate ancestor will all have zero h-index. Here, we can construct a graph such that all nodes are either a leaf node or their ancestor.



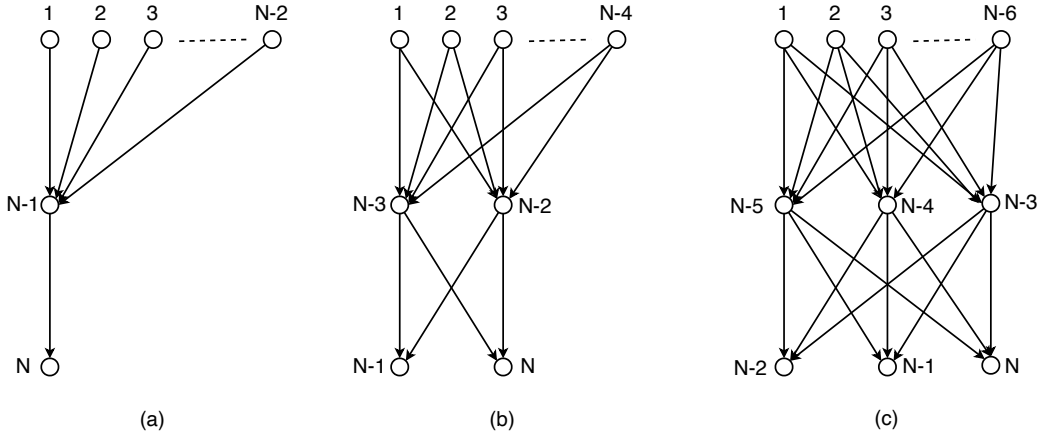
**Figure 3.5:** Graph with maximum nodes having zero h-index

Therefore, the limits for  $k = 0$  are  $[2, N]$

**Case 2.**  $k > 0$

In this case to minimise  $|S|$ , we need to minimise the number of nodes with h-index equal to  $k$  where  $k > 0$ . In other words, we need to maximise the number of nodes with h-index equal to zero. As shown in Case 1, it is always possible for all the nodes to have h-index equal to zero. So the minimum  $|S|$  will be zero.

In this case to maximise  $|S|$ , we need to maximise the number of nodes with the h-index equal to  $k$ . As shown in Theorem 1, we need at least  $2k$  nodes in order to make h-index for one node as  $k$ . Specifically, the node should have  $k$  out-neighbors that each have  $k$  out-neighbors. For the maximum value of  $|S|$ , we need to use the same set of  $2k$  nodes for each node having h-index as  $k$ . Below are some examples for  $k = 1, k = 2$  and  $k = 3$ :



**Figure 3.6:** Example graphs for maximum  $|S|$  case when (a)  $k = 1$  (b)  $k = 2$  and (c)  $k = 3$

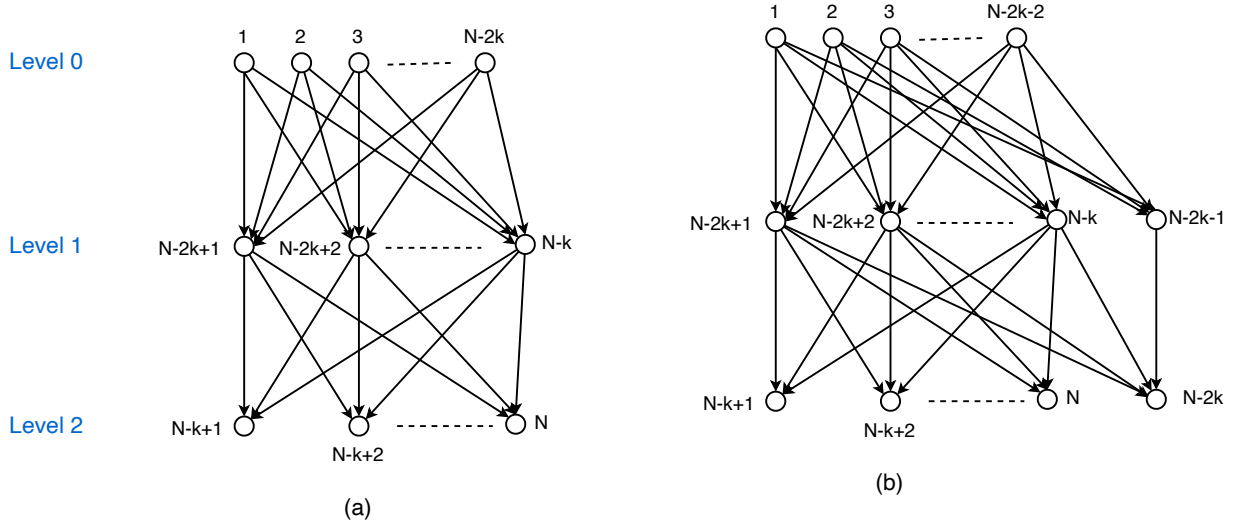
Generalising this pattern, we can say that the maximum value of  $|S|$  can be  $N - 2k$ . We shall prove this via induction.

*Proof by Induction.* Base case:  $k = 1$  As shown in the Fig. 3.6(a), we only need  $N - 2 * 1$  nodes. Therefore, the theorem holds for  $k = 1$

Inductive Step: We will show that the theorem holds for  $k + 1$ , given that it is true for some  $k$ . The graph for the same will look like Fig 3.7(a). The value of  $|S|$  (i.e. the number of nodes with h-index as  $k$ ) will be  $N - 2k$ . To calculate  $|S|$  for the h-index as  $k + 1$ , we will take two nodes from level 0 (node named  $N - 2k - 1$  and  $N - 2k$ ) with  $N - 2k$  nodes and put them on level 1 and 2 as shown in Fig 3.7(b). This way we will make the h-index for the remaining nodes on level 0 as  $k + 1$ . As we can see, the number of nodes on level 0 are  $N - 2k - 2$ . Therefore,

$$|S| = N - 2k - 2 = N - 2(k + 1) \quad (3.3)$$

So the theorem holds true for h-index  $k + 1$ . By the principle of mathematical induction, this theorem is proven.



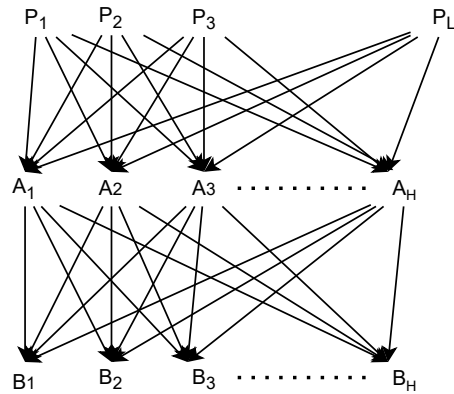
**Figure 3.7:** Inductive proof of theorem 2 case 2

**Corollary.** Suppose there are  $L$  vertices with h-index as  $H$ , then the minimum number of nodes are  $L + 2H$

*Proof.* As shown in Theorem 1, we need at least  $2H$  nodes in order to make the h-index for one node as  $H$ . To minimise the total number of nodes, all the  $L$  nodes need to have the same set of  $2H$  nodes as out-neighbors. Mathematically, the set of nodes can be denoted as follows:

$$[P_1, P_2, \dots, P_L] + [A_1, A_2, \dots, A_H] + [B_1, B_2, \dots, B_H] = L + 2H \quad (3.4)$$

□



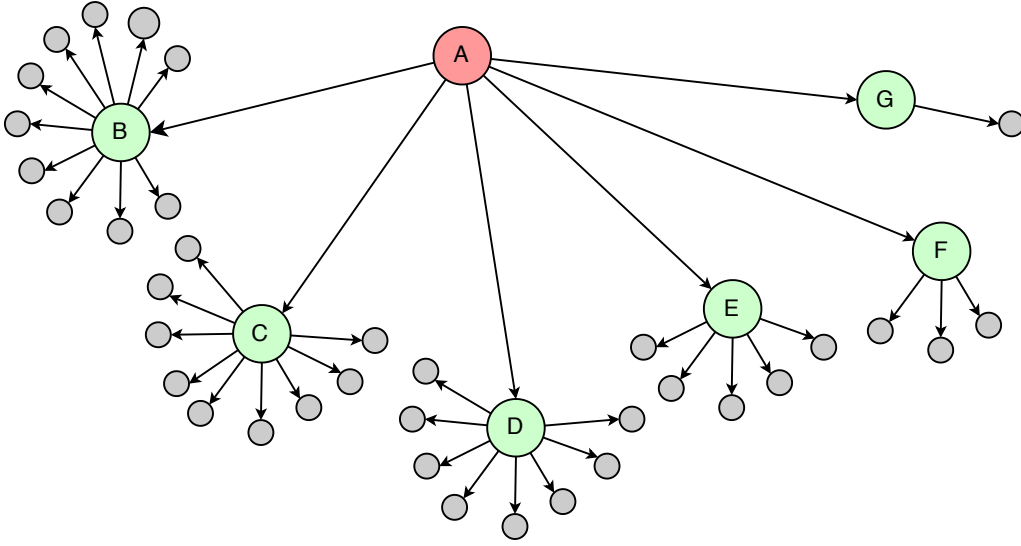
**Figure 3.8:** Resultant graph for Corollary of Theorem 2

### 3.3 h-index and its variants on a graph

We use scientific publication and citation graph as an example in this section, one can consider a paper as a node in the citation graph and an out-edge denotes a citation. Therefore, it can be generalised to any directed acyclic graph.

Rousseau [37] proposed the concepts of first generation and second generation to describe reference networks. For a certain paper  $P$ , Rousseau said that first generation publications are those which reference  $P$ , thereby having a direct influence in its significance. Second generation papers are those which reference a first generation paper but are not themselves first generation, and so on.

In this context, let us assume the paper being evaluated as 0th G-publication (0G-pub), first generation publications as 1th G-publication (1G-pub) and so on. In Fig. 3.9,  $A$  is the 0th G-publication;  $B, C, D, E, F, G$  are 1G-pubs and all the grey colored nodes are 2G-pubs. If the h-index of a node is  $h$ , the Hirsch core denotes the set of top  $h$  papers ranked by decreasing out-degree. In Fig.3.9, the Hirsch core of  $A$  node consists of  $[B, C, D, E]$  as the h-index of  $A$  is four.



**Figure 3.9:** Example graph for h-index variants

#### 3.3.1 h-index

h-index [24] is defined to be the highest number  $h$  such that the 0G-publications has at least  $h$  1G-publications with each of them having  $h$  or more citations from 2G-publications.

#### 3.3.2 g-index

Egghe [18] argued that, “a measure which should indicate the overall quality of a scientist . . . should deal with the performance of the top articles”. Hence, he proposed the g index as a modification of the



h index. This is defined as the highest number such that the top  $g$  1G-publications received together at least  $g^2$  citations from 2G-publications. As compared to the h-index, the g-index gives more weight to highly cited papers. The inflated values of the G-Index give credit to lowly-cited or non-cited papers while giving credit for highly-cited papers.

### 3.3.3 h(2)-index

Similar to the g index, the h(2) [30] index also gives more weight to highly cited articles. This is defined as the highest natural number such that  $h(2)$  most-cited 1G-publications received each at least  $[h(2)]^2$  citations from 2G-publications. Note that, for any 0G-publication, the h(2) index is always lower than the h index.

### 3.3.4 a-index

Rousseau [28] introduced the term Hirsch core. “The Hirsch core can be considered as a group of high-performance publications, with respect to the scientist’s career” (Jin et al., 2007, p. 855). In other words, this is the papers ranking smaller than or equal to h. The a index includes in the calculation only papers that are in the Hirsch core. This is defined as the average number of citations from 2G-publications of papers in the Hirsch core.

Index	Value
h-index	The highest number $h$ of out-neighbors such that each have at least $h$ out-degree
g-index	The highest number $g$ of out-neighbors such that sum of their out-degrees is at least $g^2$
h(2)-index	The highest number $h(2)$ of out-neighbors such that each of them have at least $g^2$ out-degree
a-index	$\frac{1}{h} \sum_{j=1}^h out_j$ , where $out_j$ = out-degree of $j^{th}$ node, and $h$ =h-index
m-index	The median of the out-degree of nodes in Hirsch Core
r-index	$\sqrt{\sum_{j=1}^h out_j}$ , where $out_j$ = out-degree of $j^{th}$ node, and $h$ =h-index
Wu’s w-index	The highest number $w$ of out-neighbors such that each have at least $10w$ out-degree
h(5,2)-index	The highest number $h$ of out-neighbors such that each have at least $5h^2$ out-degree

**Table 3.1:** Definitions of h-index and its variants

### 3.3.5 m-index

The distribution of citation counts is usually skewed, therefore, this modification uses the median as the measure of central tendency. The m-index [12] is defined as the median number of citations from 2G-publications received by papers in the Hirsch core.

### 3.3.6 r-index

[28] observed critically that with the a index, “the better scientist is ‘punished’ for having a higher h-index, as the A-index involves a division by h” (p. 857). The authors suggest computing the index by taking the square root of the total number of citations in the Hirsch core rather than dividing by h. This is defined as the square root of the sum of citations from 2G-publications received by papers in the Hirsch core.

### 3.3.7 Wu’s w-index

Wu [45] proposed this simple way to assess the impact of a work. This is defined to be the highest number  $w$  such that the 0G-publication have at least  $w$  1G-publications with  $10w$  or more citations from 2G-publications. According to their results, there were noticeable differences between the w-index and the h-index, because the w-index plays close attention to the more widely cited papers.

### 3.3.8 h(5,2)-index

The h(5,2)-index is a special form of the h(a,b)-index. Varying a and b yields a large number of combinations. After defining the h(a,b)-index, Ellison [21] selected 12 combinations and conducted empirical research for assessing economists. The results demonstrated that for assessing economists, h(5,2) and h(10,1) are the best. Note that h(10,1) is exactly Wu’s w-index. This is defined to be the highest number  $h$  such that the 0G-publication has at least  $h$  1G-publications with  $5h^2$  or more citations from 2G-publications

Index	Value
h-index	4
g-index	6
h(2)-index	2
a-index	8
m-index	8.5
r-index	5.65
Wu’s w-index	1
h(5,2)-index	1

**Table 3.2:** Index values for graph in Fig.3.9

These different types of indices can be calculated for any DAG.

## 3.4 h-index algorithm

In this section, we will discuss two algorithms for calculation of h-index and one optimisation on the binary search algorithm. The input to each algorithm is the set of out-degrees for each out-neighbor

of the node we are calculating h-index for. Suppose we are calculating h-index for a node  $P$ , the input to each algorithm is the set of out-degrees of the out-neighbors of node  $P$ . The expected output is the h-index of node  $P$ .

### 3.4.1 Basic algorithm

In this algorithm, we first sort the array of out-degrees in a descending order. We traverse the array one by one and keep increasing the value of h-index by 1 (starting with 0). At each step, we will check if the current h-index value is possible or not according to the definition. As soon as we encounter first impossible value, our final h-index will be the last possible h-index value. The psuedo code is shown below:

---

**Algorithm 1:** Basic h-index algorithm

---

**Input:**  $degArr$

**Output:**  $h$

```

1  $h \leftarrow 0$ ;
2  $sort(degArr)$ ;
3  $n \leftarrow size(degArr)$ ;
4 for  $i \in \{1, \dots, n\}$  do
5   if  $i \geq degArr[i]$  then
6      $h \leftarrow i$ ;
7   else
8      $break$ 
9   end
10 end
```

---

The algorithm involves a *for* loop for iterating over all elements in the array. Although this is computationally inefficient, yet it remains the most intuitive and easy way of calculation. The time complexity of this algorithm is  $O(n \log n)$  for sorting the array and additional  $O(n)$  for iterating. Here,  $n$  is the size of  $degArr$ .

This can be improved by using binary search instead of iterating through the array. The algorithm is shown in Algorithm 2.

---

**Algorithm 2:** Basic h-index algorithm using binary search

---

**Input:**  $degArr$ **Output:**  $h$ 

```
1  $h \leftarrow 0$ ;  
2  $sort(degArr)$ ;  
3  $n \leftarrow size(degArr)$ ;  
4  $low \leftarrow 0$ ;  
5  $high \leftarrow n - 1$ ;  
6 while  $low \leq high$  do  
7    $mid \leftarrow (low + high)/2$ ;  
8   if  $degArr[mid] \geq mid + 1$  then  
9      $low = mid + 1$ ;  
10     $h = mid + 1$ ;  
11   else  
12      $high = mid - 1$ ;  
13   end  
14 end
```

---

This algorithm is slightly more optimised than the Algorithm 1. The complexity for sorting the array remains the same but the cost of iteration is substituted by the cost of binary search i.e.  $O(\log n)$ .

### 3.4.2 Linear algorithm

The idea here is based on the bucket sort mechanisms. Suppose,  $n$  is the total number of out-neighbors, if we have  $n + 1$  buckets, numbered from 0 to  $n$ , then for any neighbor with out-degree corresponding to the index of the bucket, we increment the count for that bucket. Note that, for any neighbor with out-degree greater than  $n$ , we put in the  $n^{th}$  bucket.

Then we iterate in the reverse order and keep adding the bucket counts. Whenever the current count exceeds the index of the bucket, meaning that we have the index number of neighbor that has out-degree greater than or equal to the index. This will be our h-index result. We are iterating from the end of the array as we are looking for the greatest h-index.

The time complexity of this algorithm is  $O(n)$  as we are only iterating through the array and not sorting it here. This is the fastest algorithm for h-index calculation.

In Table 3.3, we show the run times (in seconds) for each of the three algorithms using two different graphs with  $10^6$  and  $10^7$  nodes. To calculate the run times, we iterate over each node in the graph and calculate their h-index using one of the algorithms. We repeat this for both of the graphs. Clearly, Algorithm 3 is the fastest one as it has a linear complexity. Only algorithm 3 uses extra memory of the order of  $n$  (size of the array). Thus, there is a run time of memory allocation included for Algorithm 3.

---

**Algorithm 3:** Linear h-index algorithm

---

**Input:**  $degArr$ **Output:**  $h$ 

```
1  $h \leftarrow 0$ ;  
2  $n \leftarrow size(degArr)$ ;  
3  $bucket \leftarrow empty\ array\ of\ size\ (n + 1)$ ;  
4 for  $i \in \{0, \dots, n - 1\}$  do  
5   if  $degArr[i] \geq n$  then  
6      $bucket[n] = bucket[n] + 1$ ;  
7   else  
8      $bucket[degArr[i]] = bucket[degArr[i]] + 1$ ;  
9   end  
10 end  
11  $cnt \leftarrow 0$ ;  
12 for  $i \in \{n, \dots, 0\}$  do  
13    $cnt = cnt + bucket[i]$ ;  
14   if  $cnt \geq i$  then  
15      $h = i$ ;  
16   end  
17 end
```

---

Number of nodes in graph	Run time for Algorithm 1	Run time for Algorithm 2	Run time for Algorithm 3
$10^6$	24s	23s	18s
$10^7$	42s	40s	35s

**Table 3.3:** Comparison of run times (in seconds) for each algorithm

### 3.5 Summary

In this chapter, we first define the generalised h-index. Then we discuss important properties of h-index. We then extend this generalisation to other h-index type metrics like h-index, a-index, m-index and so on. Lastly, we discuss three algorithms and compare their runtimes.

## *Chapter 4*

### **H-index and its variants on research papers**

#### **4.1 Introduction**

Finding the most relevant scientific article from a set of articles may seem to be a simple task at first sight, but the task to rank the articles is specially challenging. Impact of a publication is one of the most important topics in scientometrics. The sheer increase in number of publications per year has made it hard for researchers to keep track of the literature. This problem of inflation in scientific articles makes it a challenging task to find papers that have made significant contributions. This is especially true for the newcomers in the field.

The evaluation of a single publication serves as the foundation for evaluating scientists, organisations, journals, and other aspects of scientific research outputs. Today, citation counts is the most widely used quantitative method to evaluate single publication. However, citation counts can only roughly reflect a publication's impact. Moreover, it cannot effectively reflect a publication's comprehensive influence (i.e. influence beyond just the first level of citations). In recent years, another source of evaluation has emerged which measures the impact of an article in society: alternative metrics (altmetrics [36]). Altmetrics are alternative approaches to measuring the impact of a research article, as demonstrated by users' interest and engagement with it. Altmetric watches social media sites, science blogs, many mainstream media outlets and reference managers for mentions of academic papers. Some of the metrics are as follows: number of views, downloads, clicks, saves, tweets, shares, posts, discussions, and bookmarks. These altmetrics are available on SCOPUS [4] and PLOS [2]. Altmetrics aim to complement traditional research impact measures by showing a more complete picture of how readers engage with and use it.

In the last two decades, the h-index has become a widely used measure of scientific performance. The automatic calculation of h-index has even become a built-in feature of major bibliographic databases such as Web of Science and Scopus, Google Scholar, etc. In the field of academic evaluation, the h-index is a method of quantitative evaluation. The h-index is a single number measuring the cumulative impact of a researcher's output by looking at the amount of citation their work has received. Number of

publications, number of citations, and average number of citations per publication are considered to be the three most important evaluation metrics.

On the other hand, the h-index considers two dimensions combined in a particular manner, namely the number of papers and each paper's citation counts. Hence, the h-index can more effectively reflect a researcher's academic influence. Furthermore, the calculation of h-index is quite simple, fast and easy to implement. The h-index was initially adopted and used to only gauge a researcher's influence on academia (Hirsch [24]). Since then, the h-index has been used to assess groups of individuals, institutions and journals, research themes, and countries. Although, the h-index can be used in different aspects, it is not certain that it is better than traditional methods. Also, large portion of scientific community is not familiar with the variants of h-index. In most cases, variants of h-index are only popular with the bibliometricians. The differences, advantages, and disadvantages of the h-index and its variants are also not clear.

As mentioned in Section 2.2, Schubert [40] proposed to use h-index for assessing single publications. Since the h-index can be used to assess individual articles, we should also use other Hirsch-type indices in assessing single publications. In the current research, eight other variants of h-index (including the original h-index) and one traditional indicator (number of papers). Therefore, nine indicators in total, were chosen to be used for assessing individual articles.

**Contributions:** In this chapter, we answer the following questions:

- Determine which variations of the h-index behaves similar to h-index i.e. which of them follow similar or dissimilar trend. This is achieved by calculating the correlation between the different variants.
- Determine which variations of the h-index have most overlap with each other. For discussing this, we calculate Rank Biased Overlap (RBO).
- Which variant gives the best ranking to national and international award winning authors?
- Does the performance of these indices change over time?

## 4.2 Data

In order to evaluate different citation based indices, data for a large scale of papers is required. It is important to collect suitable data in order to answer the research questions in this work. So we have selected papers from VLDB conference and SIGMOD conference. There are multiple reasons for the selection of these specific conferences. Firstly, their list of papers is easily available on DBLP (dblp.org). Both these conferences are amongst the top conferences in their fields. They have dedicated committees who select the test of time awards each year which we will use later in this chapter. Also, the details for the same are widely available in public domain [6, 5]

### 4.2.1 Data Collection

For the purposes of this research, we collect all the papers in both the conferences over the years. For calculating the indices, we also collect the citation data for these papers.

For collecting these papers and citation data, we implement a three step crawler. The crawler works as follows:

- It queries dblp to retrieve all the paper titles and their year of publishing. This is written using BeautifulSoup library in python.
- Using S2AG API [3], we map these paper titles to the paper IDs in semantic scholar. We have used semantic scholar as it has a easily available API for programmatic retrieval of data. Moreover, it has a vast coverage of papers and citation data. This retrieval process is semi-automatic. Firstly, the script searches for papers with similar title as the given title (from dblp). It then filters papers with the same publishing year. Lastly, it tries to do a fuzzy string match on paper title and if the confidence is very high, we assign this paper ID to the given title. In case the confidence is low, the user is prompted with a question and has to manually select if the two papers match. This last step is the only part where we need human intervention. Out of 8070 (4427+3643; see Table 4.1) papers, we only required the last step for 188 papers.
- Once we have the semantic scholar paper IDs (s2\_id), we can retrieve the citations, year, and other metadata. Using the S2AG API, we retrieve papers that directly cite the collected papers from above. We call this set of papers the 1th Generation publications. Then, we also retrieve papers that directly cite the 1th Generation publications.

At the end of this process, we get a citation graph with the following sets of nodes: (a) all papers from both the conferences (b) their corresponding citations and (c) citations of these citations. We also have the year of publishing for all of these papers. The edges in the graph denote the 'is cited by' relation. In other words, an edge from node A to node B denotes that A is cited by B.

### 4.2.2 Dataset Description

We collected the papers from DBLP for the years 1985-2020 for VLDB and 1988-2020 for SIGMOD. As shown in table 4.1 below, there were 4652 and 3744 papers listed on DBLP respectively. Out of which we were able to map more than 95% of papers to their corresponding semantic scholar paper ID. This was done using our fuzzy match logic on paper titles. After this, we retrieved 174669 direct citations for the 4427 papers from VLDB and 177160 for 3744 papers from SIGMOD. For the next level, we retrieved 1312443 and 1378127 citations of the 1th generation papers. 0th generation to 1th generation is a 100x increase in the number of papers but 1th generation to 2th generation is a 10x increase. Finally, we see total papers in the citation graph as 1316634 and 1382516 respectively.



	VLDB	SIGMOD
DBLP Papers	4652	3744
Papers with <i>s2_id</i>	4427	3643
1th gen papers	174669	177160
2th gen papers	1312443	1378127
Total papers (nodes)	1316634	1382516
Total cites relation (edges)	3135134	3254607

**Table 4.1:** Dataset Description

### 4.2.3 Benchmark Dataset

This specific research problem has no gold standard based on a dataset that could be used to evaluate and assess. A comprehensive and extensive benchmark dataset is required to assess the indices. Hence, in this study, the test of time awards are used as a standard merit or benchmark. In the context of VLDB, a paper is selected from the VLDB Conference from ten to twelve years earlier that best meets the “test of time”. In picking a winner, the committee evaluates the impact of the paper. The committee especially values impact of the paper in practice, e.g., in products and services. Impact on the academic community demonstrated through significant follow-through research by the community is also valued. For SIGMOD, this paper is selected from the conference held exactly ten years ago. Their criterion of identifying the paper is impact (research, products, methodology) over the intervening decade.

	VLDB	SIGMOD
Total awarded papers	29	25
Total awarded papers in our dataset	25	21

**Table 4.2:** Number of awarded papers

In this study, we retrieve all the awarded papers for both the conferences. This data is available on their respective websites. The dataset consists of awardees from 1995 to 2022 for VLDB conference and from 1999 to 2022 for SIGMOD. Total awardees are listed in the table 4.2. The awardees for a few years are not present as either there was no award in a particular year or the corresponding paper did not exist in the crawled dataset. The reason for it missing from the dataset is that it is missing in the semantic scholar database. There are 4 out of 25 missing in SIGMOD dataset and 4 out of 29 missing in VLDB.

## 4.3 Experiments

In this section, we explain three experiments conducted in order to answer the research questions in this work. Firstly, the correlations are evaluated between the h-index and all its variants. We have also

evaluated whether the awarded papers rank on the top by using h-index and its variants. We used the test of time awards that are won by papers for their exceptional impact and performance in a decade to serve as a benchmark (details in the section above). Lastly, we have compared the performance of h-index and its variants by considering change through time.

#### 4.3.1 Correlation amongst indices

The first question to answer is: whether there is some correlation between h-index and its variants for single publication? The purpose of this experiment is to see which indices follow a similar trend and which indices follow a unique trend. The data is prepared as follows: we gather all the papers for both VLDB and SIGMOD. Then, we collect all the citation data required for the calculation of the indices. Lastly, we rank all these papers based on the values of various indices(see Chapter 3), leading to nine different ranked lists for each dataset. Spearman’s correlation coefficient is used to calculate the correlation between all pairs of ranked lists.

#### 4.3.2 Rank Biased Overlap (RBO)

In this experiment, we compare the different variations of the h-index on the basis of their overlap with each other in ranking the papers. Unlike correlation measures, RBO is a similarity measure which denotes how similar are two ranked lists. This will help us determine which variations of the h-index have most overlap or similarity with each other. RBO is based on the simple concept of average set overlap. The idea is to determine the fraction of content overlapping at different depths in the two ranked lists. Suppose we have two ranked lists,  $A : [P_1, P_2, P_3, P_4]$  and  $B : [P_2, P_1, P_4, P_3]$  in order of their ranks. Given below are set intersections at different depths. Set intersection shows the intersection between the two sets of lists at each depth. Fraction denotes the length of intersection set divided by the depth.

Depth(d)	Items in List A@d	Items in List B@d	set intersection	Fraction
1	$P_1$	$P_2$	$\{\}$	$0/1=0$
2	$P_1, P_2$	$P_2, P_1$	$\{P_1, P_2\}$	$2/2=1$
3	$P_1, P_2, P_3$	$P_2, P_1, P_4$	$\{P_1, P_2\}$	$2/3=0.66$
4	$P_1, P_2, P_3, P_4$	$P_2, P_1, P_4, P_3$	$\{P_1, P_2, P_3, P_4\}$	$4/4=1$

**Table 4.3:** Example of set overlap calculation

After calculating the fractions of set overlapping at various depth, one can either plot the distribution to study how similar two lists are or, use the average of the last column (Fraction) to denote the Average overlap. RBO is a further extension of this concept which uses fixed weights for each depth. It uses a geometrically decreasing series for the weights for each depth. This makes the final value to be bound as the sum of indefinite geometric series is finite. RBO also gives higher importance to the top ranks

as compared to the lower ranks due to this geometrically decreasing series. The equation for RBO is denoted by,

$$RBO(S, T, p) = (1 - p) \sum_{d=1}^{\infty} p^{d-1} . A_d$$

The value of RBO lies between 0 and 1 (inclusive) where 0 denotes completely disjoint ranked lists and 1 denotes identical ranked lists.

The key difference between correlation and RBO is that the former is used to evaluate the similarity in the trends of ranking and the latter is used to evaluate the overlap in two ranked lists. RBO also gives more weightage to the top ranks in the ranked lists. In other words, a mismatch in top ranks is given more importance in the final value of RBO.

Similar to correlation, we gather all papers for both VLDB and SIGMOD conference, then calculate all the indices for each paper. Therefore, creating nine different ranked lists. Then these ranked lists are compared pair wise.

### 4.3.3 Performance of indices in predicting trends of awardees

In this experiment, we answer the question: Which variant is the best in ranking the awarded papers on top amongst the possible candidates. The data for the same is prepared as follows:

- Gather the list of all *test of time* awarded papers from VLDB and SIGMOD conferences shown in table 4.2.
- Iterate over all the awarded papers. For each such paper, we retrieve the set of candidates as the papers published from 10 to 12 years ago in the same conference for VLDB. For SIGMOD, the candidates are from the SIGMOD proceedings exactly 10 years ago. For example, while considering an awarded paper in SIGMOD 2020, we will pickup all the papers from SIGMOD 2010 as the candidates for this award.
- Using the crawled citation graph, we then retrieve the citation data for these candidates up to the year in which the award was given. This way we will get the data that the awarding committee uses while selecting. For example, while considering an awarded paper in SIGMOD 2020, we will consider the citation data only up to 2020. In other words, any citation received in 2021 will not be considered.
- Rank all the candidates for this particular award using the nine different indices.
- Retrieve the rank of this particular awarded paper with respect to each index.

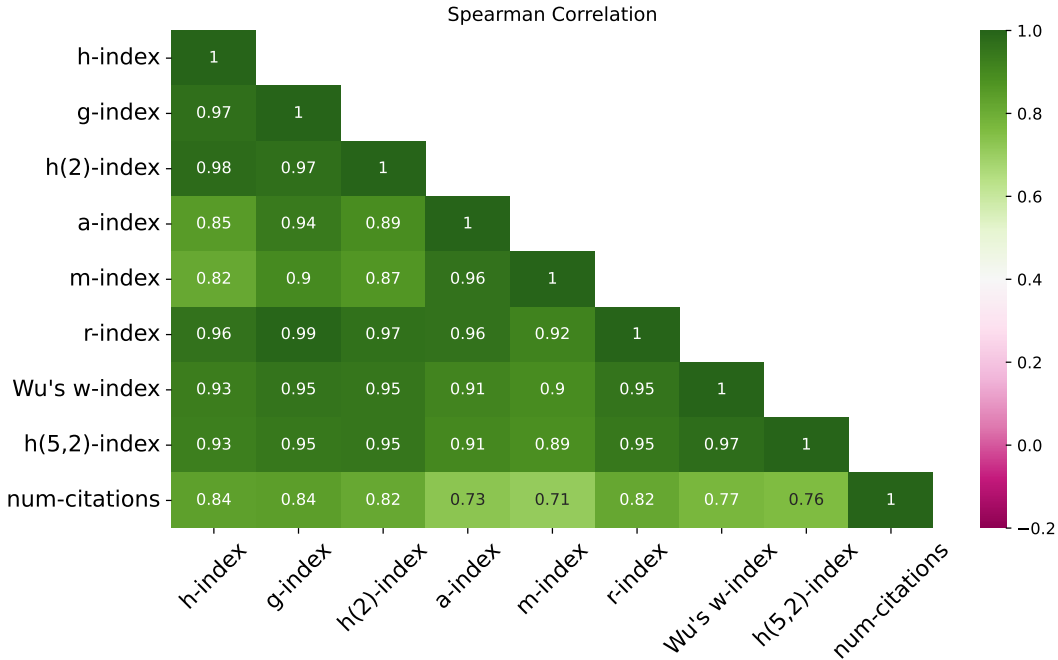
Finally, after calculating these ranks for each awarded paper, we determine how many papers are present in the top 5% of their corresponding list. We also find out the occurrence of awardees in 5–10%, 11–20% up to 31–40%.

For instance, consider an awarded paper in VLDB for year 2018 named  $P_1$ , we take all VLDB papers from 2006, 2007 and 2008 as the candidates for ranking. We then rank these candidates considering the citation data till 2018 only as this is the data available to the awarding committee at the time of selection. Let's say the rank for  $P_1$  as per the h-index is 2 out of 100. Therefore,  $P_1$  ranks in the top 5% of its list according to h-index.

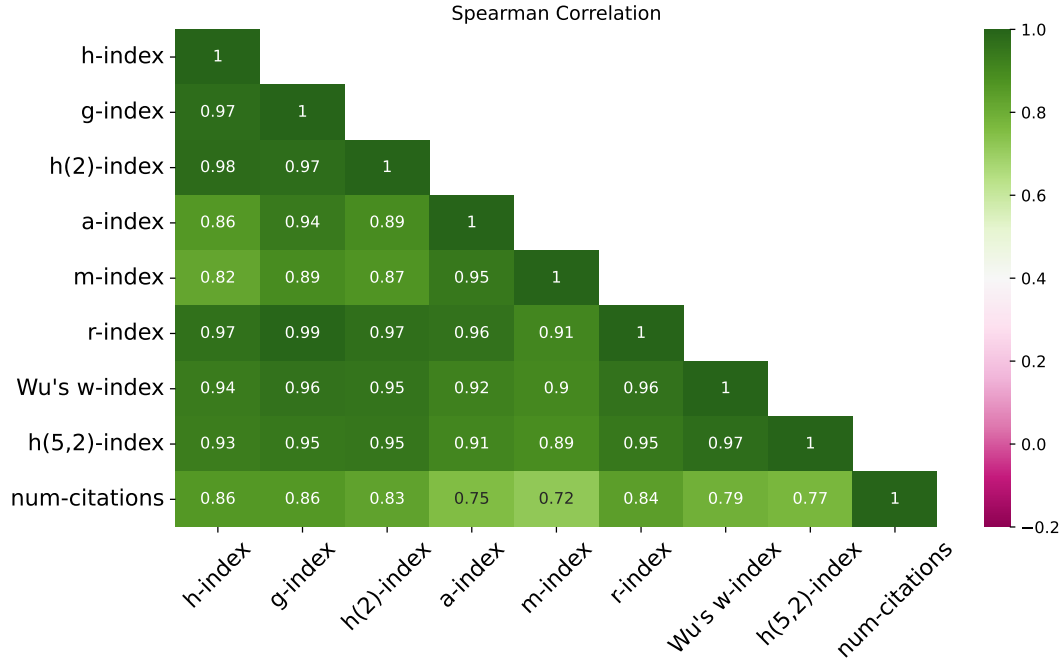
#### 4.3.4 Performance of indices over time

In this experiment, we analyse the performance of each index over time. The measure of performance is the number of papers the index ranks in the top 5% of their corresponding list. The time range is considered from 1st year after publication to 10th year after publication. For example, consider a awarded paper  $P$  published in 2004, we will evaluate the rank for this paper as per each index in the years 2005 to 2014. Let's say this paper is amongst the top 5% in year 2007 as per h-index. We will increment the performance measure of h-index in year 2007 by one. Wherever this rank is in the top 5%, it will be counted towards the performance of that index in that particular year.

## 4.4 Results



**Figure 4.1:** Correlation matrix for SIGMOD conference



**Figure 4.2:** Correlation matrix for Vldb conference

#### 4.4.1 Correlation amongst indices

The purpose of this experiment is to understand the similarities amongst all the nine indices. As mentioned in Section 5.3.1, we have a calculated nine different ranked lists using the indices. Given this, we calculate the Spearman's correlation coefficient for each pair (see Fig 4.1 and Fig 4.2). These correlation values will give us the answer to our first question. There are three possibilities of correlation:

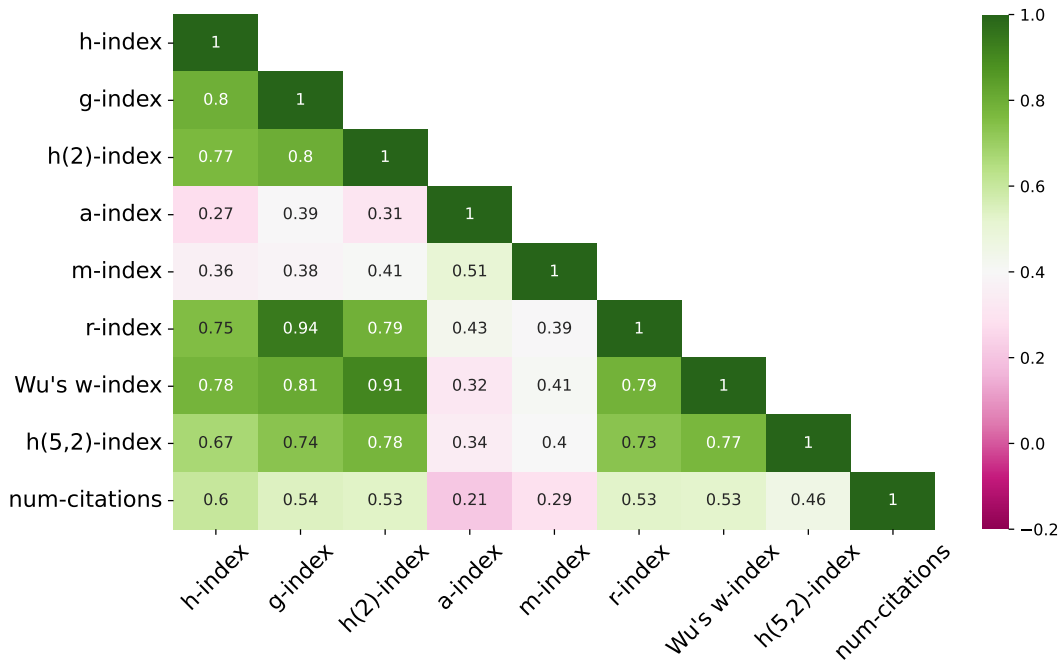
- Positive value denotes that both the ranked list are positively related i.e. if rank for one index increases, the other will also increase.
- Zero value denotes that there is no correlation i.e. they may change independent of each other
- Negative value denotes that if one rank decreases the other will increase.

For all the three cases, the magnitude of the coefficient will determine the strength of correlation. Obviously, the correlation of an index with itself will be 1.

The correlation matrix has been shown in Fig 4.1 and Fig 4.2. Most of the values are higher than 0.8 and some of them even reach more than 0.95. There are no negative values meaning that the indices generally agree with each other and follow a similar trend. It can be seen that h-index is very strongly correlated with 5 indices [g-index, h(2)-index, r-index, w-index, h(5,2)-index]. While number of citations is moderately(around 0.8) correlated with all other indices.

#### 4.4.2 Rank Biased Overlap (RBO)

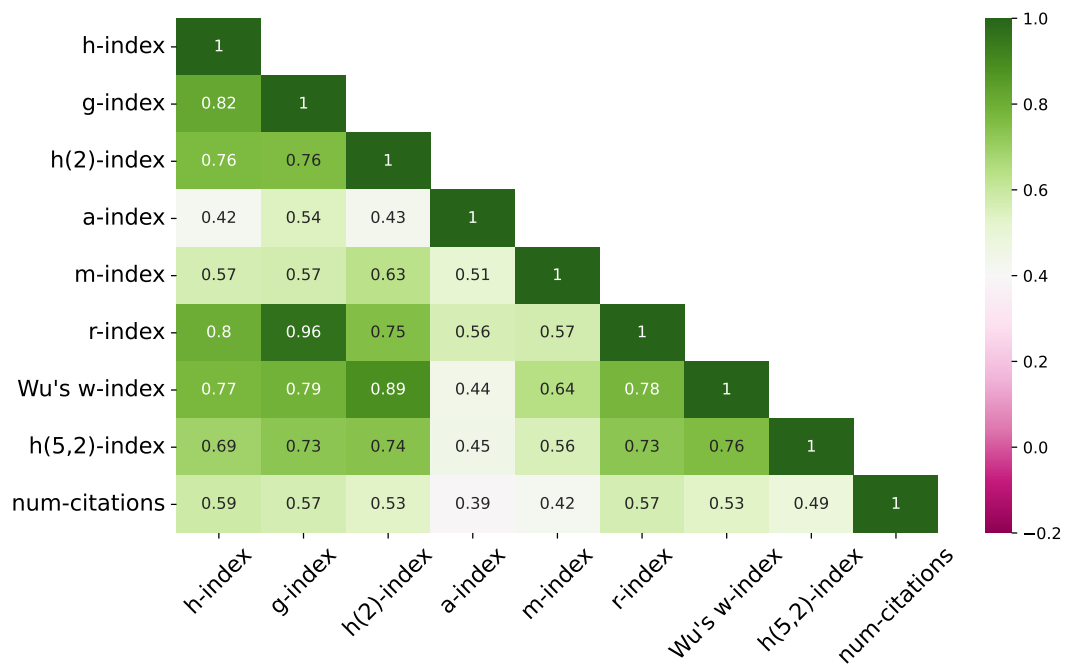
The pairwise RBO values are shown as matrix in Fig 4.3 and 4.4. We can see that a-index and m-index have very low overlap with every other index. They have moderate overlap of 0.51 with each other. From table 3.1, a-index and m-index have similar approach which very different from other indices. Num of citations has the most overlap with h-index (0.6) amongst all indices. Also, both matrices are very similar in trends of overlap but there is difference in magnitude of overlaps. h-index has the most overlap with g-index.



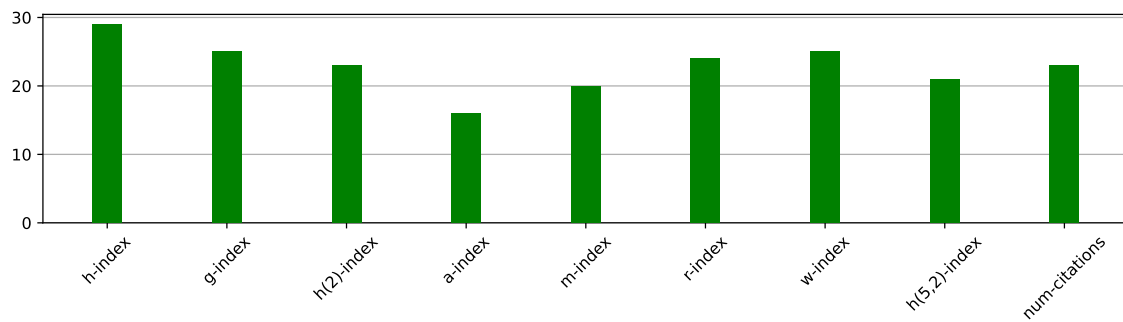
**Figure 4.3:** RBO matrix for SIGMOD conference

#### 4.4.3 Performance of indices in predicting trends of awardees

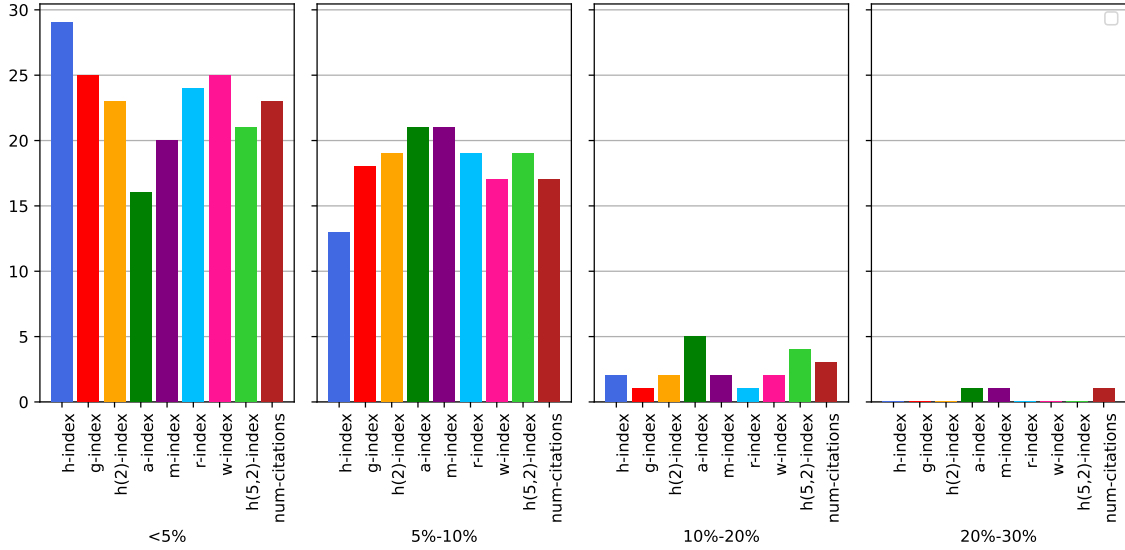
Here we have addressed the second research question, i.e. which variants is the best in ranking the awarded papers at the top. As explained in detail in Section 5.3.2, once we have ranks of all awarded papers in their respective list for all of the indices, we first evaluate how many of the awardees were present in the top 5% of their lists according to each index. From Fig. 4.5, we can see that h-index performs the best with 29 out of 46 (63%) of papers in top 5%. The g-index, h(2)-index, r-index, w-index and number of citations show similar performance of around 52%. The a-index is the worst performing at 34% (16/46).



**Figure 4.4:** RBO matrix for VLDB conference



**Figure 4.5:** Index name vs Number of awarded papers ranked in top 5%



**Figure 4.6:** Index name vs Number of awarded papers ranked in < 5%, 5% – 10%, 10% – 20% and 20% – 30%

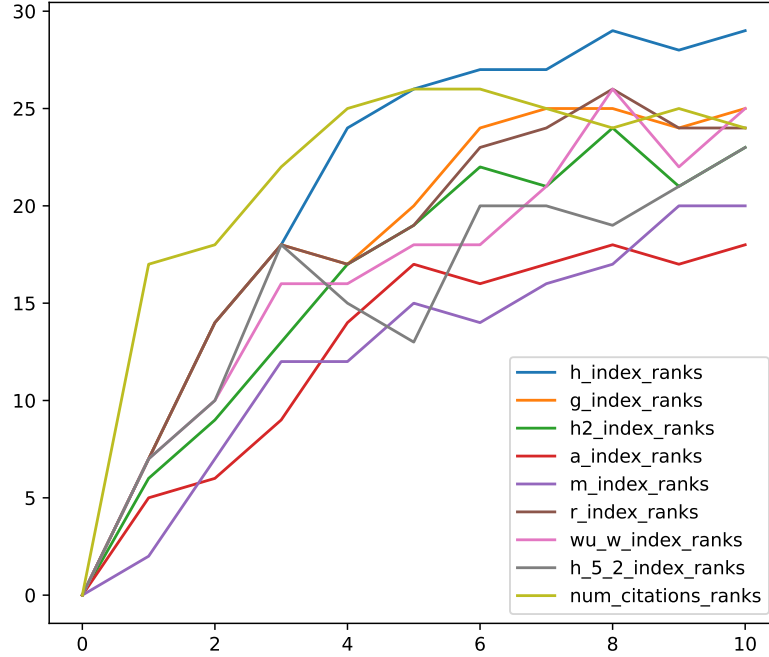
#### 4.4.4 Performance of indices over time

Here, we analyse the change in performance of the different indices over 10 years from publishing of the paper. From the fig. 4.7, we can see that at the end of 10 years when the paper is actually awarded, h-index is best performing index (as discussed in Section 5.4.2). During the initial years (less than 5 years) of the paper, number of citations is the best index. And h-index is the best index in the later years of the paper. This is expected as the number of citations will increase first and then the impact on h-index will be observed. H-index captures a deeper level of impact and hence it needs some amount of time to start seeing an increase.

Fig 4.8 shows the performance of h-index and number of citations beyond 10 years of publishing. We observe that the performance of number of citations start to improve after the 10th year when the papers are awarded publicly. The hypothesis here is that once the papers are awarded they become more popular and it reaches more people. Consequently, it gets more number of citations. Therefore, more awarded papers start to rank in top 5% of the list of papers according to number of citations.

In Fig. 4.9, we show the trends for a few of the awarded research papers. Specifically, we compare the change of h-index and the corresponding rank of the paper over time (starting from the year it was published). We can observe that papers like (a), (d) and (f) have a constantly increasing h-index and higher rank almost every year. Then papers like (b), (e) and (g) see a increase and decrease in the ranks over time. Papers like (c) have an increasing h-index but a lower rank. Lastly, there are papers which have high rank from the first year itself like paper (h). Thus, there is no surety that all awarded papers will have decreasing rank with the increasing h-index like (a), (d) and (f).

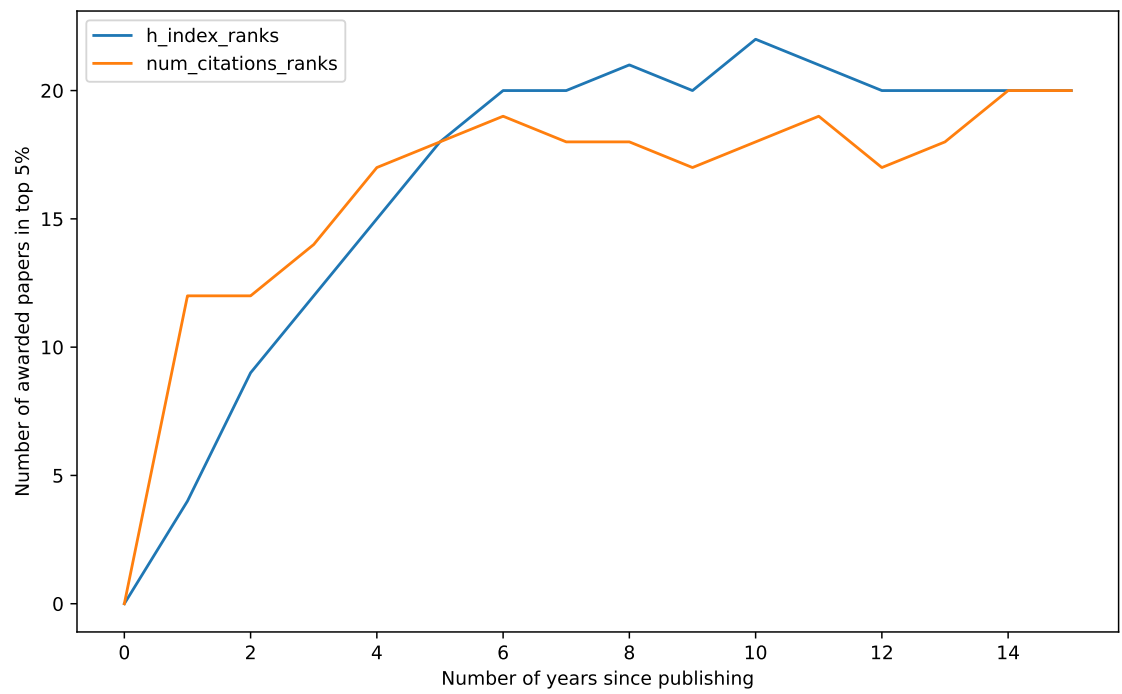




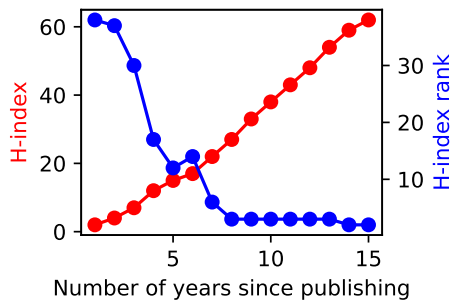
**Figure 4.7:** Number of years since publishing vs Number of papers in top 5% when ranked on the particular index

## 4.5 Conclusion

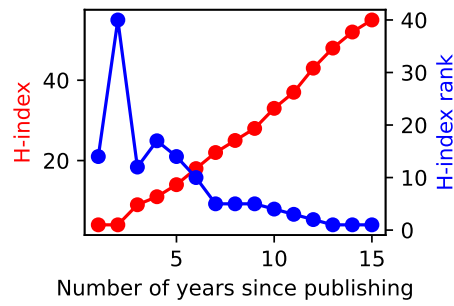
We define and compare h-index with different citation based indexes on research papers. Specifically, we looked at papers from SIGMOD and VLDB conference. We compare the rankings given by these indexes to each paper and also compare the rankings given to awarded papers. Our observations show that h-index is the best performing in ranking awarded papers at the top. That is, VLDB and SIGMOD lay higher recognition to h-index in determining the awards.



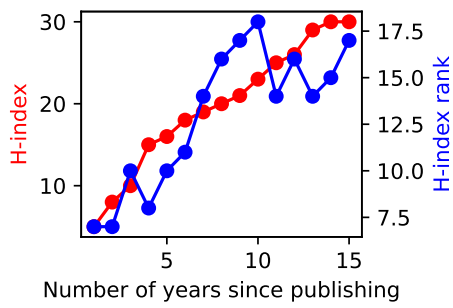
**Figure 4.8:** Number of years since publishing(beyond 10 years) vs Number of papers in top 5% when ranked on the particular index



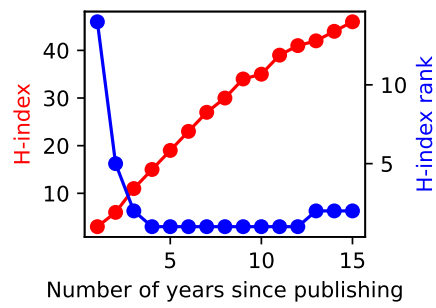
(a)



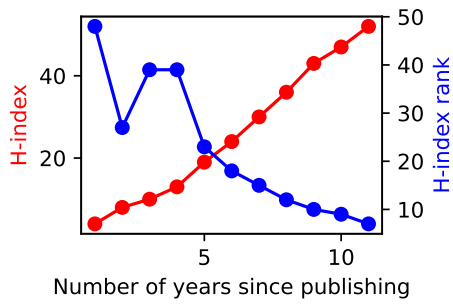
(b)



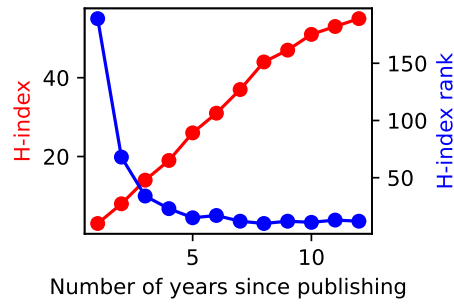
(c)



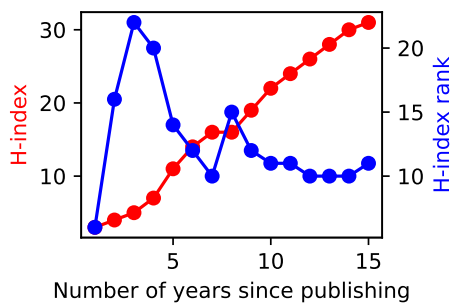
(d)



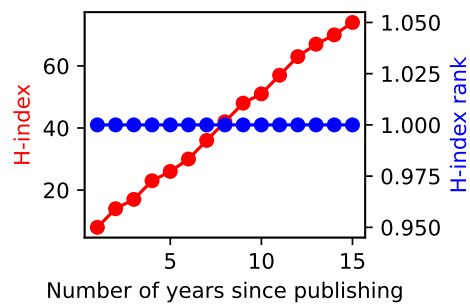
(e)



(f)



(g)



(h)

**Figure 4.9:** Plot showing number of years since publishing vs H-index and the corresponding rank as per h-index for a few awarded papers. Each of the graphs represents the mentioned curve for a single awarded paper.

## *Chapter 5*

### **H-index of authors based on H-index of their papers**

#### **5.1 Introduction**

Researchers contribute to advancing the horizons of knowledge in the world by establishing facts and reaching new conclusions through methodical analysis; and thereafter publishing the results of their findings in the form of research articles. Bibliometrics are key factors in the assessment and comparison of the research productivity of individuals or groups. Quantifying the impact of their work can help researchers not only recognise significant contributors in their field of research, but also provide a measure of an author's perceived value - by demonstrating the citation patterns of one's work.

Measuring the qualitative value of a researcher is much easier as compared to making a quantitative analysis. According to an individual's personal opinion, one could simply state that a researcher is good if they publish many good papers. But quantitatively measuring the value proposition of these papers is much more complicated, since it can be measured in several distinct ways. In the past few years, multiple different metrics have been put forth to determine a researcher's scientific merit based on the quantity and quality of their peer-reviewed publications.

The merit of an individual researcher is commonly quantified by using citation-based metrics. Amongst the varied set of citation-based metrics, the most common is the h-index. The h-index of a researcher influences decisions about financing, promotion, and employment, thus shaping the researcher's career. As a result, it influences how the scientific community develops, and how research advances. The h-index, proposed by Hirsch in 2005, has emerged as the most prominent metric for calculating the impact of a scientist's published work. The h-index is readily available in various citation databases, for example, Scopus and Google Scholar.

In this work, we discuss four modifications of the traditional h-index. First is the traditional h-index of a researcher. Secondly, we use h-frac-index proposed in [29]. They argue that h-index is not a good measure for scientific impact. This is due to changing authorship patterns, including a higher prevalence of hyper authorship. The major finding is that fractional allocation of citations among co-authors can mitigate the issues with h-index.

In chapter 4, we argue that the h-index of a paper is a better quantifier of paper impact when compared to many other metrics. Hence for the third metric, we use the h-index of a paper to calculate the h-index of an author. This metric was proposed in [20]. Finally, we propose hp-frac-index, as the h-index of an author using h-index of a paper divided by total authors of the paper (details in next section).

We retrieve the top 1000 researchers ranked on decreasing h-index in three different fields of research, namely, Computer Science, Economics, and Biology. Then we cross reference our dataset with the scientific award winners in each field. The award lists used are Turing award winners for Computer Science, Nobel Prize in Economics, and Nobel prize in Chemistry, Physiology and Medicine for biology. The traditional and proposed metrics are calculated for each researcher. Our experiments show that hp-index and hp-frac-index outperform traditional indices by giving better ranks to the awarded researchers in all three fields. We also compute the correlation amongst the metrics across the three fields.

## 5.2 H-index of author

In this section, we will cover all the metrics being used in our experiments. We are considering four metrics.

### 5.2.1 h-index and h-frac-index

#### **h-index**

We use the h-index as the first traditional metric for quantifying an author's research impact. H-index of an author is the largest number  $h$  such that the given author has published at least  $h$  papers that have each been cited at least  $h$  times.

#### **h-frac-index**

Secondly, we use a recent extension to h-index called h-frac [29]. This is a variant of the h-index that allocates citations fractionally among co-authors. In other words, when using the number of citations to calculate the h-index of an author, they divide each paper's number of citations by its total number of authors. This mitigates the cluttering of the ranking by hyper authors.

### 5.2.2 hp-index and hp-frac-index

We formally define the h-index of a *paper*. Then we discuss the two metrics being used to evaluate researchers, namely, hp-index and hp-frac-index.

Consider a paper  $p$  and the set of papers citing  $p$  be the set  $C = [c_1, c_2, c_3, \dots, c_n]$ . The h-index of  $p$  is equal to the largest number  $h$  such that at least  $h$  papers from  $C$  have at least  $h$  citations each.

#### **hp-index**

In this metric, we first calculate the h-index of all papers of an *author*  $X$  using the definition above.

Suppose, the set of papers published by  $X$  is  $[P_1, P_2, P_3, \dots, P_n]$  and the corresponding h-index values of these papers be  $[h_1, h_2, h_3, \dots, h_n]$ . We compute the hp-index of the author  $X$  as follows:

$$hp(X) = H([h_1, h_2, h_3, \dots, h_n]) \quad (5.1)$$

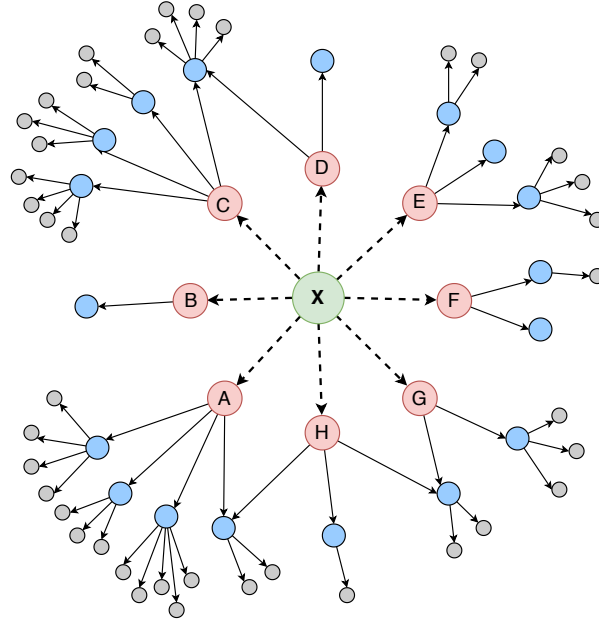
where  $H$  is the function to calculate the h-index of any given set of values. To sum up, the hp-index of an author is the h-index of the h-index of all the author's papers.

### hp-frac-index

Similar to hp-index, we calculate the h-index of all the papers  $[P_1, P_2, P_3, \dots, P_n]$  as  $[h_1, h_2, h_3, \dots, h_n]$ . Let the number of authors for each paper be  $[a_1, a_2, a_3, \dots, a_n]$ . Note that for all  $i, a_i \geq 1$ . We then compute the hp-frac-index of an author  $X$  as follows:

$$hp-frac(X) = H([\frac{h_1}{a_1}, \frac{h_2}{a_2}, \frac{h_3}{a_3}, \dots, \frac{h_n}{a_n}]) \quad (5.2)$$

The main difference between hp-index and hp-frac-index is that we use fractional h-index values of a paper in the latter.



**Figure 5.1:** Example graph for h, hp, h-frac, hp-frac demonstration

### 5.2.3 Example

Consider the graph in Fig. 5.1, node  $X$  (in green) is the author,  $A$  to  $H$  are the papers co-authored by  $X$ . A dotted edge from  $X$  to  $A$  denotes that  $X$  has co-authored the paper  $A$ . All the blue papers are the ones that cite the red papers and grey papers cite the blue ones. A solid edge from  $A$  to  $B$  denotes

that  $A$  is cited by  $B$ . We can see that  $X$  has written 8 papers and paper  $A$  has been cited 4 by 4 different papers.

Node	# authors	# citations	list of citation of citation	h-index
A	2	4	[4, 3, 5, 2]	3
B	1	1	[0]	0
C	1	4	[4, 3, 2, 4]	3
D	2	2	[0, 4]	1
E	2	3	[0,2,3]	2
F	1	2	[1, 0]	1
G	3	2	[2, 3]	2
H	2	3	[1,2,2]	2

**Table 5.1:** H-index of the papers in given example in Fig. 5.1

In the table 5.1 shown above, *#authors* denotes the number of authors of a paper, *#citations* denotes the number of citations and *list of citation of citation* denotes the number of citations of each blue paper that has been cited by the nodes  $A$  to  $H$ . Lastly, *h-index* is the h-index of each paper. Using the values of h-index for the nodes  $A$  to  $H$  from the table above, values of the four indices (defined in section 7.2.1 and 7.2.2) for the author can be calculated as follows:

$$h\text{-index}(X) = H(4, 4, 3, 3, 2, 2, 2, 1) = 3$$

$$h\text{-frac-index}(X) = H(\frac{4}{2}, \frac{4}{1}, \frac{3}{2}, \frac{3}{2}, \frac{2}{2}, \frac{2}{1}, \frac{2}{3}, \frac{1}{1}) = 2$$

$$hp\text{-index}(X) = H(3, 3, 2, 2, 2, 1, 1, 0) = 2$$

$$hp\text{-frac-index}(X) = H(\frac{3}{2}, \frac{3}{1}, \frac{2}{2}, \frac{2}{2}, \frac{1}{2}, \frac{1}{1}, \frac{2}{3}, \frac{0}{1}) = 1$$

### 5.3 Citation Data

In order to compare the four indices, we crawled the list of top 1000 researchers in Computer science, Economics and Biology field in the order of decreasing number of citations from google scholar [1]. The steps followed to complete the data collection are:

- Crawl the list of names of top 1000 authors from google scholar for each field.
- Match these author names to *author\_ids* in Semantic Scholar [3]. We use Semantic Scholar as it has a highly accessible database of scientific literature with author and paper details readily available.
- Once we have *author\_ids*, we retrieve the papers published by them in one set of API calls. Let this set of papers for a researcher be called  $P_w$ .
- Next, we run another set of API calls to get the papers citing any paper in  $P_w$ . Let us call this set of papers as  $P_c^1$ .
- Lastly, we retrieve the papers citing any paper in  $P_c^1$ . Let us call this set of paper  $P_c^2$ .

At the end of this process, we have a graph for each author. One graph consists of the author and all the other papers from the sets  $P_w, P_c^1, P_c^2$  as nodes. The author connects to the nodes in  $P_w$  with a 'written by' edge. The nodes from  $P_w$  connect to  $P_c^1$  and  $P_c^1$  connect to  $P_c^2$  with a 'cited by' edge. This graph looks similar to the example in Fig. 5.1

	CS	Economics	Biology
Crawled authors	1000	1000	1000
Matched authors	803	856	842
Total papers published ( $P_w$ )	285622	159445	468683
Total Citations ( $P_c^1$ )	6232959	3251169	12054908
Total Citations of citations ( $P_c^2$ )	21053913	11497621	43194158

**Table 5.2:** Dataset description

From the table above, we can see that Biology has the most number of papers published per author followed by Computer Science and then Economics. Subsequently, the two sets of citations follow the same trend. All three fields have more than 800 authors matching with their corresponding IDs in semantic scholar dataset. The total citation to total papers published ratio ( $P_c^1:P_w$ ) is around 20 for all the fields. Whereas, the second level of citations to first level of citations ( $P_c^2:P_c^1$ ) is around 3.5 for the three datasets.



## 5.4 Experiments and results

This section explains all the experiments done in order to compare the proposed indices with the traditional ones. We took the data collected, and calculated the h-index of each *paper* published (i.e. all papers from set  $P_w$ ). We use the definition discussed in Section 2.2 for the same. Using these values we calculated the four indices for an *author* as explained in Section 2. Lastly, we ranked the authors on the decreasing order of each index (h, hp, h-frac, hp-frac) resulting in four different ranked lists of authors. Note that, if the value of an index is same for two different authors, then the one with more citations is given a higher rank.

### 5.4.1 Correlation

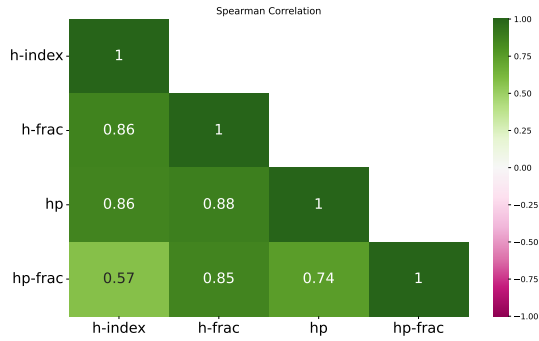
The aim of this experiment is to gauge how unique are the trends followed by these indices are. For this, we calculated the Spearman's correlation between the four ranked lists obtained. As shown in Fig. 5.2, h-index is only moderately related to hp-frac with 0.54, 0.7 and 0.58 correlation for Computer Science, Economics and Biology respectively. This shows that the information presented by hp-frac is unique as compared to h-index. The h-index is highly related with hp-index across the three fields. Also, hp-frac and h-frac are highly correlated. h-index and h-frac-index are also highly correlated for Computer Science and Economics with a value of 0.86 for both. But in the case of Biology, this falls down to 0.72. Interestingly, hp-frac and hp are moderately related with a value of around 0.7. This is interesting because the only difference between both of them is that hp-frac has a division by number of authors for each paper.

### 5.4.2 RBO

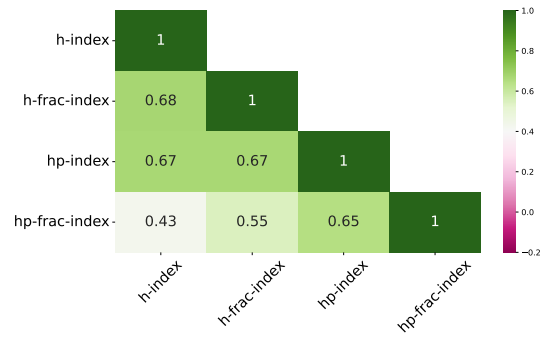
In this section, we compare the pairwise overlap of the four ranked lists. We use Rank Biased Overlap as proposed by Webber et al. [44]. From Fig 5.2, the overlap between h-index and hp-frac-index is the lowest with the values of 0.43, 0.51 and 0.36. The overlap is the highest for h-index and h-frac-index. The h-index, hp-index and hp-frac-index have high correlation (greater than 0.65) amongst each of them. Whereas, the hp-frac-index has low overlap of 0.43, 0.55, 0.65 with h-index, h-frac-index and hp-index respectively. Hence, the hp-frac-index is ranking differently from the other indices. This stimulates our next experiment to evaluate how each index ranks the awarded researchers and *which one can rank awarded researchers higher*.

### 5.4.3 Awarded researchers

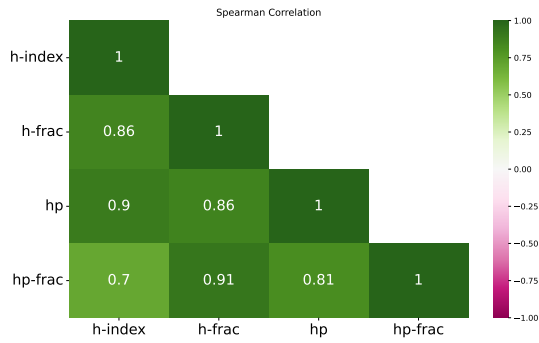
In this experiment, we compare the position of each award winning researcher in the four different ranked lists obtained for each of the three fields. We compiled the list of awardees for the awards listed in Table 5.7 and cross referenced them to our list of top 1000 researchers in each field. We found eighteen



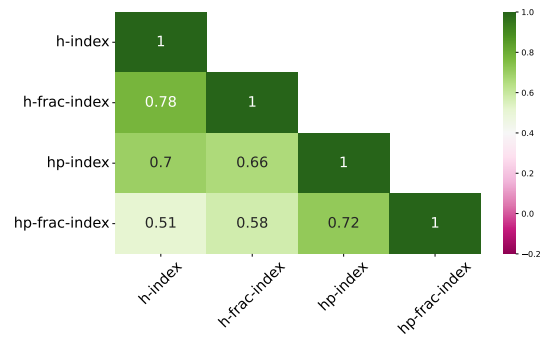
(a) Correlation for Computer Science



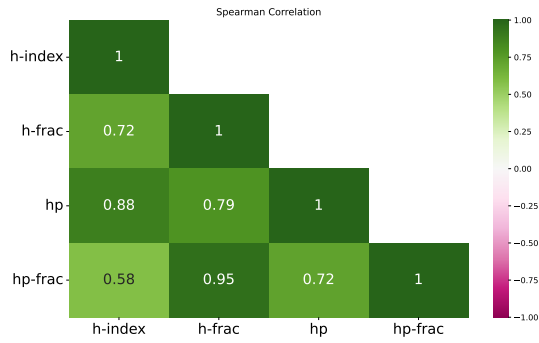
(b) RBO for Computer Science



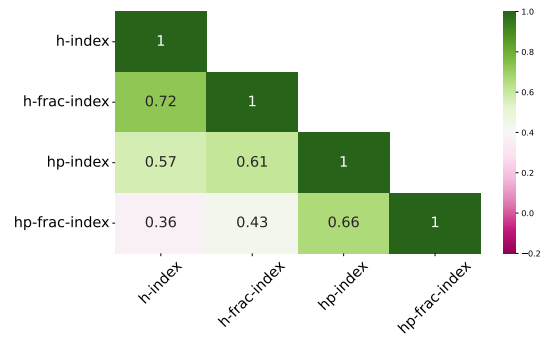
(c) Correlation for Economics



(d) RBO for Economics



(e) Correlation for Biology



(f) RBO for Biology

**Figure 5.2:** Plot showing number of years since publishing vs H-index and the corresponding rank as per h-index for a few awarded papers.

	Biology				Computer Science				Economics			
	h	h-frac	hp	hp-frac	h	h-frac	hp	hp-frac	h	h-frac	hp	hp-frac
Top 5%	23.08	38.4	38.4	<b>61.5</b>	27.7	38.8	27.7	<b>50</b>	26.4	29.4	29.4	<b>44.1</b>
5% - 10%	<b>23.1</b>	15.3	15.3	7.6	<b>16.6</b>	5.5	<b>16.6</b>	<b>16.6</b>	20.5	<b>26.4</b>	17.6	20.5
10% - 15%	<b>15.3</b>	7.6	7.6	<b>15.3</b>	<b>16.6</b>	11.1	<b>16.6</b>	11.1	11.7	<b>20.5</b>	11.7	8.8
15% - 20%	<b>7.6</b>	<b>7.6</b>	<b>7.6</b>	0	<b>11.1</b>	5.5	0	5.5	<b>14.7</b>	5.8	8.8	8.8
Total $\leq 20\%$	69.08	68.9	68.9	<b>84.4</b>	72	60.9	60.9	<b>77.7</b>	73.3	<b>82.1</b>	67.5	<b>82.2</b>

**Table 5.3:** Percentage of awardees in different ranges of ranked lists as per each index across three fields

such awarded researchers in Computer Science, thirty four in Economics and thirteen in Biology. Then, we extracted the values and ranks as per the indices for all awardees (see Table 5.8, 5.10, 5.11).

Table 5.3 shows the percentage of awardees ranked amongst top 5%, 5%-10%, 10%-15% and 15%-20% for each index across the three fields. We can observe that hp-frac-index is the best performing index with around 61% (Biology), 50% (Computer Science) and 44% (Economics) of awardees being ranked in the top 5% of the list. The other three indices perform 10% poorer than hp-frac-index in ranking awarded researchers in top 5%. For Biology, h-index performs the best in all ranges except top 5%. Although other indices perform equally in some ranges. We observe similar trend in Computer Science as well. However, in Economics, the next best indicator is h-frac in 5%-10% and 10%-15% range. Overall, in the range of  $< 20\%$ , we see that hp-frac-index performs better than the other indices for Biology and Computer Science. It performs at par with h-frac index for Economics.

Table 5.8, 5.10, and 5.11 show the ranks and values as per each index for all the awarded researchers in our data set. In Table 5.4, we show the percentage of awardees (for a particular field) that receive the best or highest rank as per a given index. We do this for all the three fields and all the four indices. To elaborate, consider an awardee  $R$ , they have four ranks as per each index. Suppose, out of these four ranks the rank given by h-index is the highest, then increment the count for h-index by 1. For example in Table 5.4, under Biology, 15.3% for h-index means that 15.3% of all Biology awardees had the best rank as per h-index (among the four indices). We can observe that hp-frac-index far outperforms the other indices across all three fields. Note, that if more than one index for an awardee yields equal best rank, it is considered towards all those indices, hence the sum of each column may surpass 100%.

The outcome here is that the hp-frac-index perform better than all the other methods. The hp-index and hp-frac-index capture a deeper level of impact by taking into account one extra level of papers as compared to h-index and h-frac-index. This extra depth of information helps in recognising the sustained research impact of the author better than h-index.

#### 5.4.4 Further Analysis

The table 5.5, displays a number of statistics about the calculated indices and their average and maximum values. We can see that the average h-index for Biology is the highest with a big margin.

	Biology	Computer Science	Economics
h-index	15.3	16.6	8.8
h-frac-index	0	0	29.4
hp-index	30.7	22.2	17.6
hp-frac-index	<b>61.5</b>	<b>66.6</b>	<b>47</b>

**Table 5.4:** Percentage of awardees given the highest rank as per each index across the three fields (see Table 5.8, 5.10, and 5.11)

	CS	Economics	Biology
Average h-index	70.3	47.8	111.7
Max h-index	184	149	285
Average h-frac-index	41.6	34.6	44.5
Max h-frac-index	107	114	148
Average hp-index	32.9	26.2	53.09
Max hp-index	72	69	122
hp-index of 100th author	30	32	53
Average hp-frac-index	16.7	16.9	18.08
Max hp-frac-index	42	47	51
hp-frac of 100th author	23	24	29
Average number of publications	383.9	201	622.09

**Table 5.5:** Average and maximum values of the calculated metrics

The average number of publications is the highest for Biology. This also explains the high maximum h-index value for the same.

The average hp-frac-index does not have a big margin for Biology. This shows that the range of hp-frac values are more tightly packed. The 100th author as per hp-frac index has a hp-frac value of 23 for Computer Science, 24 for Economics and 29 for Biology. For hp-index, we can observe that the average value and the value for 100th author are very close. This means that there is sharp fall in the values of hp-index from 1st to 100th author. Whereas for hp-frac-index, there is a consistent difference between average value and 100th author value. This shows that hp-frac-index has slower decline in the values (when we move down the rank list) as compared to hp-index.

Both hp-frac-index and h-frac-index include a division by number of authors in their calculation. One might wonder if these two indices only punish authors who have high number of co-authors? To evaluate this, we calculated the difference of h-index and h-frac-index, and hp-index and hp-frac-index for each author. Let us call them *diff1* and *diff2* respectively. We ranked the authors on descending order of *diff1* and *diff2*. Then we ranked the authors on decreasing order of average co-authors. Finally, we calculate the Spearman's correlation of *diff1* and *diff2* ranked lists with the ranked list ordered by average co-authors. Note that we did this for all three fields. As shown in table 5.6, the average number of co-authors has a very low correlation with the difference of h-index and h-frac-index, and hp-index and hp-frac-index values. This shows that hp-frac-index and h-frac-index are not antithetical to collaboration

with others. The fact that they have very low correlation shows that these two rank lists do not follow the same trend.

	CS	Economics	Biology
Correlation between <i>diff1</i> and average co-authors	0.08	0.15	0.07
Correlation between <i>diff2</i> and average co-authors	0.17	0.22	0.039

**Table 5.6:** Correlation between *diff1*, *diff2*, and the average number of co-authors

Furthermore, authors like Gregg L. Semenza with a average co-authorship of 14.17 has the highest rank in Biology. We also notice that highly collaborative authors like Yoshua Bengio and Michael I. Jordan rank amongst the top 5 authors for Computer Science.

## 5.5 Conclusion

In this chapter, we collected large-scale data for evaluation of author from three fields, namely, Computer Science, Economics and Biology. We used four different metrics: two traditional metrics, one re-applied metric (hp) and one proposed metric (hp-frac). Our experimental analysis show that hp-frac-index gives a unique ranking order to authors and outperforms all the other metrics in ranking the awarded researchers higher. The hp-frac-index is a robust way to evaluate the impact of researchers. Its ability to capture individual contributions and resist manipulation makes it a valuable tool for assessing the impact of researcher. It takes into account the importance of a paper’s impact by using the paper’s h-index, therefore, capturing a second level of research impact. These factors together make hp-frac-index better at ranking the authors. One of the problems to address in further work is the ability to predict future award winners using these metrics more accurately.

Award	Total awardees	Matched awardees	Acronym
Turing award winners (for Computer Science)	70	8	CS1
ACM Prize in Computing	13	10	CS2
Nobel Prize in Economics	84	15	EC1
Fellows of the American Finance Association	66	19	EC2
Nobel Prize in Chemistry	184	2	B1
Nobel Prize in Physiology or Medicine	219	2	B2
Breakthrough Prize in Life Sciences	48	9	B3

**Table 5.7:** List of awards collected

Author name	Award	Avg. co-authors	h-index		h-frac-index		hp-index		hp-frac-index	
			value	rank	value	rank	value	rank	value	rank
Gregg L. Semenza	B2	14.18	177	44	103	10	84	32	<b>51</b>	<b>1</b>
Robert A. Weinberg	B3	<b>5.2</b>	177	43	101	13	92	15	<b>41</b>	<b>9</b>
Lewis C. Cantley	B3	9.25	175	49	77	44	75	53	<b>40</b>	<b>14</b>
David Botstein	B3	7.53	159	86	76	47	77	46	<b>38</b>	<b>21</b>
Eric S. Lander	B3	<b>26.54</b>	<b>285</b>	<b>1</b>	95	19	<b>122</b>	<b>1</b>	37	25
Bert Vogelstein	B3	10.98	255	5	103	8	<b>111</b>	<b>3</b>	36	32
Robert J. Lefkowitz	B1	5.48	<b>214</b>	<b>16</b>	91	26	85	28	36	33
James P. Allison	B2, B3	13.85	140	151	63	125	66	136	<b>35</b>	<b>40</b>
Gary B. Ruvkun	B3	5.42	101	500	55	214	59	260	<b>32</b>	<b>60</b>
Karl Deisseroth	B3	9.54	158	89	62	134	<b>71</b>	<b>89</b>	29	94
Aaron Ciechanover	B1	6.92	106	447	58	180	55	355	<b>28</b>	<b>117</b>
Xiaowei Zhuang	B3	9.63	89	618	47	340	47	546	<b>20</b>	<b>321</b>
Masashi Yanagisawa	B3	9.12	129	203	50	297	<b>63</b>	<b>181</b>	17	418

**Table 5.8:** List of award winners with ranks for Biology (highest ranks in bold)

Biology		Computer Science		Economics	
Author	Awarded?	Author	Awarded?	Author	Awarded?
<b>Gregg L. Semenza</b>	<b>Yes</b>	<b>Geoffrey E. Hinton</b>	<b>Yes</b>	Cass R. Sunstein	No
Michael Karin	No	Ronald R. Yager	No	<b>James J. Heckman</b>	<b>Yes</b>
Edmund T. Rolls	No	<b>Judea Pearl</b>	<b>Yes</b>	<b>Richard H. Thaler</b>	<b>Yes</b>
Joan Massagué	No	<b>Yoshua Bengio</b>	<b>Yes</b>	Dani Rodrik	No
K. J. Friston	No	Andrew P. Zisserman	No	<b>William D. Nordhaus</b>	<b>Yes</b>
Douglas G. Altman	No	Michael I. Jordan	No	Colin F. Camerer	No
Joseph E. LeDoux	No	<b>Yann Le Lecun</b>	<b>Yes</b>	<b>Paul A. Samuelson</b>	<b>Yes</b>
Solomon H. Snyder	No	Tomaso A. Poggio	No	Gary S. Becker	No
<b>Robert A. Weinberg</b>	<b>Yes</b>	Lotfi A. Zadeh	No	Robert W. McGee	No
Mark P. Mattson	No	<b>Jon M. Kleinberg</b>	<b>Yes</b>	<b>Jean Tirole</b>	<b>Yes</b>
% of awardees	20%	% of awardees	50%	% of awardees	50%

**Table 5.9:** List of top 10 authors ranked by hp-frac-index (awarded researchers are in bold)

Author name	Awards	Average co-authors	h-index		h-frac-index		hp-index		hp-frac-index	
			value	rank	value	rank	value	rank	value	rank
Geoffrey E. Hinton	CS1	3.14	142	9	106	2	65	3	<b>42</b>	<b>1</b>
Judea Pearl	CS1	<b>1.68</b>	104	69	89	8	42	89	<b>37</b>	<b>3</b>
Yoshua Bengio	CS1	5.03	<b>184</b>	<b>1</b>	105	4	<b>72</b>	<b>1</b>	34	4
Yann Le Lecun	CS1	6.74	117	40	73	27	<b>62</b>	<b>4</b>	33	7
Jon M. Kleinberg	CS2	3.64	108	61	72	30	49	27	<b>30</b>	<b>10</b>
Daphne L. Koller	CS2	7.07	129	23	68	36	<b>52</b>	<b>19</b>	26	28
Jeffrey David Ullman	CS1	3.51	99	89	66	44	46	42	<b>26</b>	<b>26</b>
Dan Boneh	CS2	5.44	117	42	72	31	46	44	<b>26</b>	<b>29</b>
Ronald L. Rivest	CS1	4.35	79	234	50	163	37	193	<b>26</b>	<b>25</b>
Pat M. Hanrahan	CS1	4.58	86	160	52	131	41	115	<b>25</b>	<b>54</b>
Stefan Savage	CS2	5.41	87	152	43	327	41	118	<b>24</b>	<b>78</b>
David M. Blei	CS2	3.48	92	115	56	92	38	165	<b>24</b>	<b>65</b>
M. Frans Kaashoek	CS2	4.09	77	254	42	350	38	167	<b>23</b>	<b>95</b>
Pieter Abbeel	CS2	5.59	<b>129</b>	<b>24</b>	56	97	45	54	22	119
David A. Patterson	CS1	4.58	<b>92</b>	<b>114</b>	50	164	38	164	21	136
John Leroy Hennessy	CS2	3.76	67	415	41	376	33	372	<b>17</b>	<b>339</b>
David Silver	CS2	<b>7.86</b>	66	434	40	408	<b>35</b>	<b>262</b>	15	462
Jeffrey Dean	CS2	7.5	35	758	30	658	29	541	<b>14</b>	<b>533</b>

**Table 5.10:** List of award winners with ranks for Computer Science (highest ranks in bold)

Author name	Awards	Avg. co-authors	h-index		h-frac-index		hp-index		hp-frac-index	
			value	rank	value	rank	value	rank	value	rank
James J. Heckman	EC1	2.91	<b>149</b>	<b>1</b>	<b>114</b>	<b>1</b>	<b>69</b>	<b>1</b>	45	2
Richard H. Thaler	EC1, EC2	3.15	89	23	73	16	53	10	<b>39</b>	<b>3</b>
William D. Nordhaus	EC1	2.27	86	30	76	12	43	28	<b>39</b>	<b>5</b>
Paul A. Samuelson	EC2	2.13	86	29	<b>84</b>	<b>6</b>	39	50	38	7
Jean Tirole	EC1, EC2	3.02	121	4	<b>98</b>	2	58	4	36	10
Jeremy C. Stein	EC2	3.91	74	66	61	46	46	18	<b>33</b>	<b>17</b>
Ben S. Bernanke	EC1, EC2	<b>1.38</b>	68	94	64	32	40	42	<b>33</b>	<b>16</b>
Alvin E E. Roth	EC1	3.25	90	22	69	23	39	52	<b>33</b>	<b>21</b>
Christopher A. Sims	EC1	1.57	61	148	56	74	34	117	<b>33</b>	<b>19</b>
René M. Stulz	EC2	3.59	96	18	<b>76</b>	<b>11</b>	47	14	32	23
Raghuram G. Rajan	EC2	3.26	73	70	63	35	<b>46</b>	<b>16</b>	32	22
Joshua D. Angrist	EC1	4.21	75	64	56	72	<b>45</b>	<b>23</b>	32	24
John Y. Campbell	EC2	3.61	66	106	54	80	40	45	<b>31</b>	<b>29</b>
G. William Schwert	EC2	2.28	60	158	47	112	35	103	<b>31</b>	<b>31</b>
John H. Cochrane	EC2	3.38	52	250	49	102	30	200	<b>31</b>	<b>33</b>
Esther Dufo	EC1	<b>12.06</b>	87	27	61	45	<b>45</b>	<b>22</b>	27	57
Luigi Zingales	EC2	3.02	71	80	54	79	<b>40</b>	<b>43</b>	27	58
Guido W. Imbens	EC1	3.68	<b>81</b>	<b>43</b>	56	73	40	44	27	60
Franklin Allen	EC2	2.86	80	51	60	50	<b>40</b>	<b>49</b>	27	68
Campbell R. Harvey	EC2	2.79	80	49	<b>63</b>	<b>36</b>	35	99	27	62
Abhijit V. Banerjee	EC1	4.36	<b>89</b>	<b>24</b>	61	47	35	100	27	63
Lars Peter Hansen	EC1, EC2	3.05	66	109	52	89	32	140	<b>27</b>	<b>66</b>
Lloyd S. Shapley	EC1	2.05	53	234	45	126	32	139	<b>24</b>	<b>96</b>
Eduardo S. Schwartz	EC2	2.89	66	110	<b>51</b>	<b>91</b>	30	193	23	108
John R. Graham	EC2	2.63	60	159	44	142	30	195	<b>23</b>	<b>113</b>
José A. Scheinkman	EC2	3.17	59	169	<b>46</b>	<b>117</b>	32	146	21	162
David Hirshleifer	EC2	3.44	64	130	<b>50</b>	<b>97</b>	30	192	21	159
Franco Modigliani	EC2	2.95	53	235	<b>43</b>	<b>163</b>	29	229	21	166
Robert B. Wilson	EC1	1.39	33	726	31	435	22	547	<b>21</b>	<b>183</b>
David S. Scharfstein	EC2	4.39	42	446	34	337	<b>30</b>	<b>197</b>	19	230
Laura T. Starks	EC2	3.1	55	212	<b>42</b>	<b>176</b>	29	235	19	234
William F. Sharpe	EC2	2.35	38	567	35	305	25	385	<b>18</b>	<b>263</b>
Robert H. Litzenberger	EC2	2.47	31	769	25	664	22	546	<b>16</b>	<b>409</b>
Philip H. Dybvig	EC1	4.31	30	786	25	666	18	775	<b>13</b>	<b>614</b>

**Table 5.11:** List of award winners with ranks for Economics (highest ranks in bold)



## *Chapter 6*

### **H-index on airport network**

#### **6.1 Introduction**

The aviation industry plays a significant role in global transportation and trade, making airports an essential part of the modern society. Airport networks have evolved over the years, becoming increasingly complex and interconnected. Understanding the structure and dynamics of these networks is essential to ensure efficient operations, strategic planning, and infrastructure development.

The h-index is a widely used bibliometric index that provides a measure of an individual's research productivity and impact. It has been applied to various fields of research, including science, engineering, and medicine. However, its application to the analysis of airport performance is a relatively new area of research. This chapter aims to investigate the application of the h-index to the analysis of airport performance, using a large dataset of airport network and routes between them. Specifically, we will explore how the h-index can be used to rank and compare airports based on their productivity and impact.

This chapter will conduct an analysis that includes calculating the out degree, h-index, and hh-index for every airport, as detailed in the upcoming section. We will present the top 10 ranked airports based on each metric, and investigate the correlation between the h-index and other relevant metrics. The results of this study will demonstrate the effectiveness of applying the h-index to airport networks, revealing noteworthy insights.

#### **6.2 Data and Methods**

In this section, we will first explain the data collection and processing. Secondly, the metrics being used to evaluate the airports are defined.

### 6.2.1 Airport network

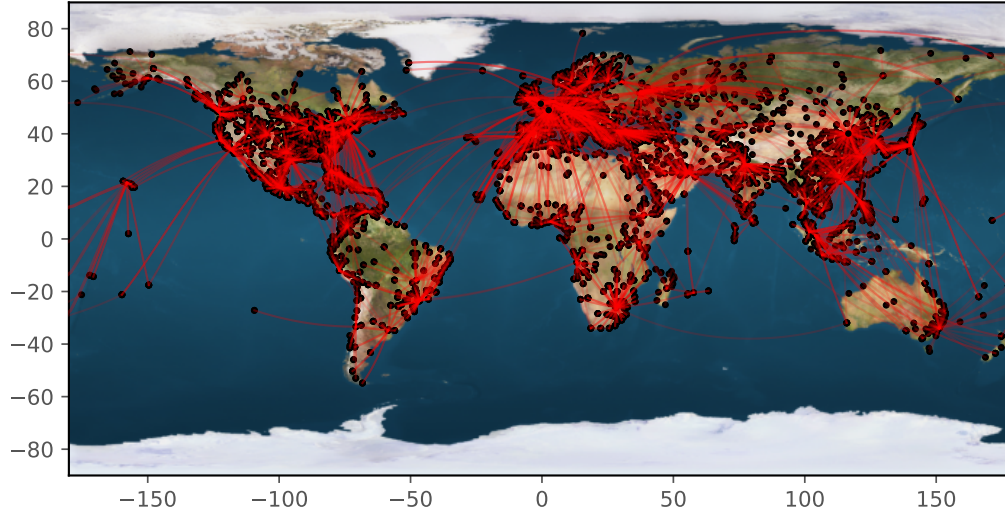
In order to calculate h-index and other metrics on airports, we collect the airport network consisting of all routes from one airport to another. The data is retrieved from OpenFlights.org. The data processing is as follows:

- We first retrieve two sets of data:
  - List of all the routes where each line corresponds to a route from airport  $A$  to airport  $B$ . A route exists when there is at least one flight from  $A$  to  $B$ .
  - Metadata for each airport like coordinates, city, country, airport code, etc.
- Create a directed graph  $G$ . For each route  $A$  to  $B$  from the list of routes, add an edge to  $G$  from  $A$  to  $B$ .

After processing, we get the graph of airport network and the metrics explained in the next section are evaluated. The details about the graph are shown in table 6.1. In Fig. 6.1, we show a subset of airports and routes connecting them. The black dots are the airports and red lines are the routes.

Total number of airports	3425
Total number of routes	67663

**Table 6.1:** Details about airport network graph



**Figure 6.1:** A sub graph of the airline network

### 6.2.2 Metrics

We use three metrics on the airport network to reveal interesting patterns.

### Out-degree

This denotes simply the out-degree of an airport i.e., the number of routes leading outwards from the given airport.

### h-index

We use the generalised h-index method here, explained in Section 3.3.1. Specifically, h-index of an airport is  $h$  if  $h$  of its connected airports each have  $h$  connected airports.

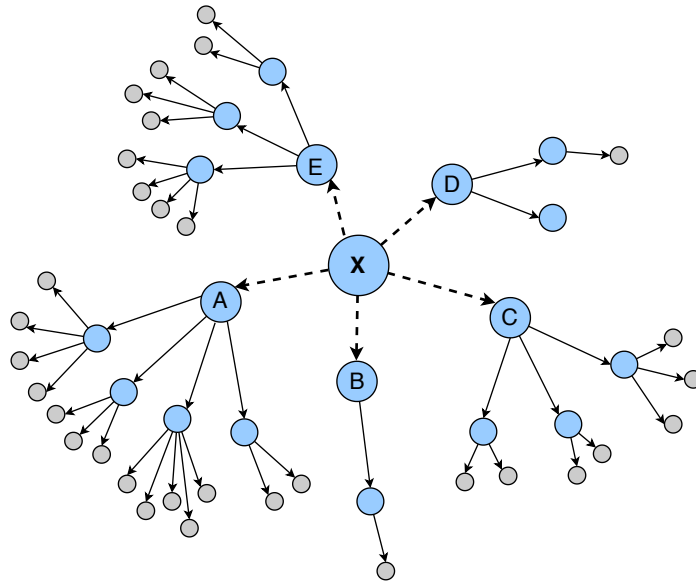
### hh-index

To calculate hh-index of an airport, we first calculate the h-index of each airport. Then, we apply the h-index method once again but using h-index values calculated previously instead of out-degree. For example, an airport has hh-index as  $w$ , if  $w$  of its connected airports each have their h-index as  $w$ . Formally, let us assume an airport  $X$  has a list of its  $k$  neighbors as  $[n_1, n_2, \dots, n_k]$ . We calculated the h-index of each node. Let the list of h-index values of its  $k$  neighbors be  $[h_1, h_2, \dots, h_k]$ . Now, hh-index is represented as:

$$hh-index = h-index([h_1, h_2, \dots, h_k]) \quad (6.1)$$

In the graph in Fig. 6.2, the nodes  $A, B, C, D, E$  and  $X$  represent airports and each edge represents a route from one node to another. The table 6.2 shows the values of the three indices discussed above. Note that, hh-index is only shown for node  $X$  as we would need another level of airport routes to calculate hh-index for the other nodes. The calculation of out-degree and h-index is straight forward from the definition. For hh-index of node  $X$ , we used h-index values of nodes  $A, B, C, D, E$ .

$$hh-index(X) = h-index([4, 1, 3, 2, 4]) = 3 \quad (6.2)$$



**Figure 6.2:** Example graph for airport network

Node	out-degree	h-index	hh-index
A	4	3	-
B	1	1	-
C	3	2	-
D	2	1	-
E	3	2	-
X	5	3	3

**Table 6.2:** The values of out-degree, h-index and hh-index for the example graph

### 6.3 Experiments and Analysis

In this section, we demonstrate several experiments to show that the h-index when applied to airport network, yields interesting patterns.

As discussed above, we ranked different airports based on three metrics: out-degree of the airport, h-index of the airport and hh-index of the airport. For each metric we have discussed the ranks and the values of each metric for top 10 airports ranked on the given metric of the subsection. The ranks are denoted in brackets.

We also define the concept of out-degree connectivity, h-index connectivity and hh-index connectivity. We only explain h-index connectivity below as the other two can be defined similarly.

**Definition 3.** The *h-index connectivity<sub>k</sub>* of an airport  $X$  is defined as: Given the list of top  $k$  airports ranked on decreasing h-index is  $T_k$ . The list of airports directly connecting to airport  $X$  be  $L_k$ . The h-index connectivity of  $X$ :

$$h\text{-index connectivity}_k(X) = \frac{|L \cap T|}{|T|} \times 100 \quad (6.3)$$

For example, if Chicago connects to 10 out of the top 20 airports (ranked on h-index) then *h-index connectivity<sub>k</sub>*( $X$ ) = 50

Similar to definition 3, we can define out-degree connectivity and hh-index connectivity. Intuitively, a connectivity score denotes how connected an airport is to the top airports of any ranked list.

#### 6.3.1 out-degree ranks

Table 6.4 shows the top 10 airports ranked on the decreasing order of their out-degree. 7 out of 10 of these airports are also on the top 10 busiest airports list (see Table 6.3). This is expected as out-degree can be seen as how busy an airport is. Higher the number of flights coming in and out, the busier the airport.

Airport	Passengers
Atlanta	96,178,899
Beijing	86,128,270
London	73,408,489
Tokyo	72,826,565
Los Angeles	70,663,265
Dubai	70,475,636
Chicago	69,999,010
Paris	63,813,756
Dallas-Fort Worth	63,554,402
Hong Kong	63,121,786

**Table 6.3:** Top 10 busiest airport according to Wikipedia

Airport	out-degree	h-index	hh-index
Atlanta	915(1)	55(39)	40(52)
Chicago	558(2)	59(28)	43(40)
Beijing	535(3)	65(14)	44(37)
London	527(4)	79(2)	51(1)
Paris	524(5)	79(3)	51(2)
Frankfurt	497(6)	80(1)	51(3)
Los Angeles	492(7)	60(25)	43(41)
Dallas-Fort Worth	469(8)	49(60)	37(83)
New York	456(9)	72(8)	47(15)
Amsterdam	453(10)	79(4)	51(4)

**Table 6.4:** Top 10 airports ranked by out-degree

### 6.3.2 h-index ranks

Table 6.5 shows the top 10 airports ranked on the decreasing order of their h-index. 9 out of 10 of these airports are located around the European region. These airports serve as hubs connecting the Eastern and the Western world. As shown in the table 6.5, these airports have lower out-degree ranks as they do not have the highest traffic. They are highly ranked in h-index because of their high connectivity. This is expected as h-index takes into account a deeper level of connection.

### 6.3.3 hh-index ranks

Table 6.6 shows the top 10 airports ranked on the decreasing order of their hh-index. This list is very similar to h-index ranked list with the exception of Copenhagen being included.

The Table 6.7 shows the airports which got the highest rank increase in hh-index as compared to out-degree. To elaborate, for each airport we calculated the difference of ranks in out-degree and hh-index. Then ranked all airports on the decreasing order of this difference. As we can see the airports getting the maximum boost have high h-index connectivity i.e. these airports are directly connected to a large

Airport	out-degree	h-index	hh-index
Frankfurt	497(6)	80(1)	51(3)
London	527(4)	79(2)	51(1)
Paris	524(5)	79(3)	51(2)
Amsterdam	453(10)	79(4)	51(4)
Munich	368(15)	76(5)	51(5)
Rome	331(23)	73(6)	50(8)
Zurich	247(47)	73(7)	51(6)
New York	456(9)	72(8)	47(15)
Madrid	330(24)	71(9)	49(10)
Istanbul	358(18)	69(10)	50(7)

**Table 6.5:** Top 10 airports ranked by h-index

Airport	out-degree	h-index	hh-index
London	527(4)	79(2)	51(1)
Paris	524(5)	79(3)	51(2)
Frankfurt	497(6)	80(1)	51(3)
Amsterdam	453(10)	79(4)	51(4)
Munich	368(15)	76(5)	51(5)
Zurich	247(47)	73(7)	51(6)
Istanbul	358(18)	69(10)	50(7)
Rome	331(23)	73(6)	50(8)
Copenhagen	229(55)	66(13)	50(9)
Madrid	330(24)	71(9)	49(10)

**Table 6.6:** Top 10 airports ranked by hh-index

number of top 20 h-index airports. Therefore, we conclude that the airports that get a higher rank in hh-index are highly connected to the hub airports but only moderately connected to the busiest airports. We call these airports micro hubs. These micro hubs play an integral role in making an airport a hub airport as they feed traffic to them.

Consider micro hubs as airports having a rank boost of more than 35 and an h-index of more than 20. Rank boost is defined as the difference between rank as per out-degree and rank as per hh-index. A positive value for rank boost signifies that the hh-index rank is better than the out-degree for that airport. The purpose of introducing a minimum value on h-index is to remove the airports that have very low amount of traffic and hence may not feed the major hubs. As per the above mentioned definition there are 160 micro hubs. The table 6.8, shows for each hub airport (discussed in Section 6.3.2), the percentage of micro hubs connecting to it. For example, Frankfurt has connection from 60.625% of micro hubs, i.e. 60.625% of the 160 micro hubs connect to Frankfurt. As we go down the table, the percentage drops, which shows that the top hub airports are better connected with micro hubs.

Moreover, the airports with minimum boost to the rank have very low h-index connectivity as shown Table 6.9. This means that the airports that have the lower rank as per hh-index have very low connec-

Airport	Rank boost	out-degree connectivity	h-index connectivity	hh-index connectivity
Naples	148	35.0	60.0	70.0
Keflavik	144	40.0	60.0	55.0
Luxemburg	132	40.0	65.0	80.0
Budapest	132	40.0	75.0	90.0
Sofia	122	40.0	65.0	70.0
Mulhouse	118	40.0	65.0	75.0
Bologna	113	40.0	65.0	75.0
Riga	110	40.0	70.0	85.0
Hannover	108	35.0	70.0	75.0
Toulouse	97	40.0	65.0	65.0

**Table 6.7:** Connectivity of airports with maximum rank boost by hh-index

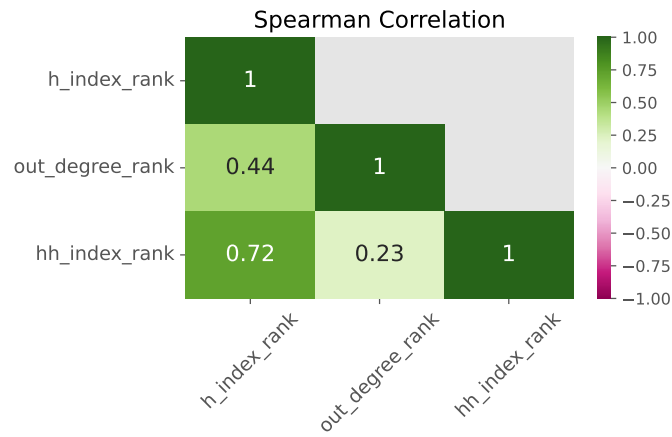
Airport	percentage of micro hubs connecting
Frankfurt	60.625
London	60.625
Paris	59.375
Amsterdam	67.5
Munich	60.625
Rome	51.25
Zurich	45.625
New York	48.75
Madrid	43.75
Istanbul	41.25

**Table 6.8:** Percentage of micro hubs connecting to hubs

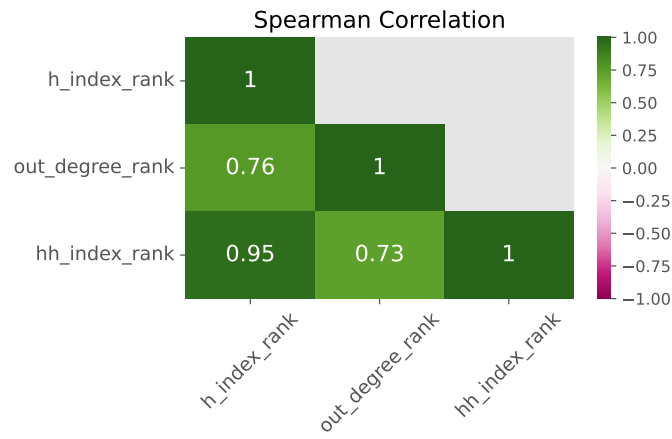
tivity to hub airports. In other words, we can say there is a direct correlation between rank boost and h-index connectivity i.e. when the rank boost is high the h-index connectivity is high and vice versa. To reiterate, having high h-index connectivity means that the airport is highly connected to the hub airports.

#### 6.3.4 Correlation

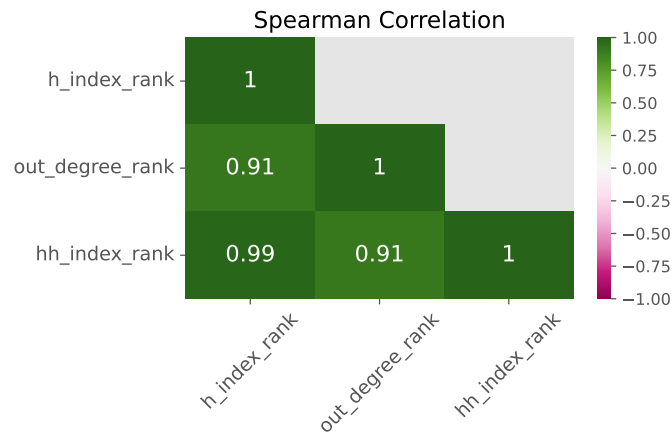
As shown in the figures 6.3, we calculated the Spearman correlation for the ranked lists by considering top 100, 1000 and lastly all airports. These three figures show that the correlation between h-index and hh-index is high in all three cases. The correlation of out-degree with h-index grows across the three subsets. This shows that the difference is created mainly for the highly ranked airports. As we go down the list of airports, all three ranked lists are similar in ranking them. In Fig. 6.3(c), we can see that the correlation amongst each pair is very high meaning, after the top airports the ranks are very similar to each other.



**(a)** Correlation of top 100 airports



**(b)** Correlation of top 1000 airports



**(c)** Correlation of all airports

**Figure 6.3:** Correlation of different indices for (a) top 100 airports, (b) top 1000 airports, (c) all airports

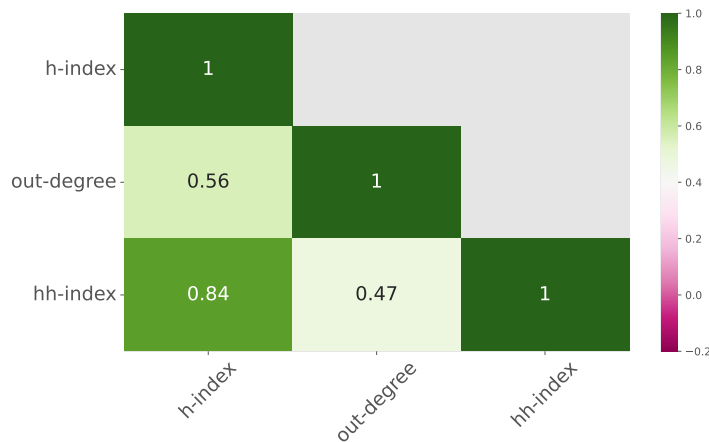


Airport	Rank boost	out-degree connectivity	h-index connectivity	hh-index connectivity
Brisbane	-316	20.0	5.0	5.0
Melbourne	-193	20.0	5.0	5.0
Nairobi	-192	30.0	35.0	35.0
Sydney	-181	35.0	10.0	5.0
Johannesburg	-180	50.0	45.0	45.0
Bogota	-172	35.0	35.0	25.0
Denpasar	-167	15.0	0.0	0.0
Jakarta	-162	25.0	10.0	5.0
Bangalore	-162	25.0	20.0	20.0
Tokyo	-161	45.0	30.0	25.0

**Table 6.9:** Connectivity of airports with minimum rank boost by hh-index

### 6.3.5 Rank Biased Overlap

In this experiment, we compare the overlap of ranked lists pairwise. We use Rank Biased Overlap as explained in Section 4.3.2. From Fig 6.4, the overlap between h-index and hh-index is the highest with a value of 0.84. The other two pairs (h-index and outdegree, and hh-index and out-degree) have similar overlap of around 0.5. When we compare this with the correlations from Section 6.3.4, we can see that RBO shows greater difference in magnitude of overlap as compared to difference in magnitude of correlation. The reason being RBO gives more importance to top ranks as compared to lower ranks.



**Figure 6.4:** Rank Biased Overlap of the three indices

## 6.4 Further Analysis

Consider hub airports as the list of top 25 airports as per h-index. As explained in section 6.3.3, micro hub airports are the ones with rank boost of at least 35 and h-index of at least 20. Now, consider the subgraph of airport network having only hubs airports called as hub graph. Also, the subgraph of airports having only hubs and micro hubs is called hub + micro hub graph. We calculate the h-index of top 20 airports by only considering these two graphs. The values are shown in Table 6.10. For example, Frankfurt has a h-index of 20 in hub graph and 47 in hub + micro hub graph.

Airport	h-index in hub graph(A)	h-index in hub + microhub graph(B)	Column (B) - (A)
Amsterdam	20	49	29
Frankfurt	20	47	27
Paris	20	47	27
Munich	20	47	27
Rome	20	47	27
Zurich	20	46	26
London	20	44	24
Istanbul	20	43	23
Madrid	20	42	22
Barcelona	20	40	20
Dubai	19	38	19
Manchester	15	34	19
Brussels	19	37	18
Copenhagen	20	36	16
Milano	19	35	16
New York	20	33	13
Newark	18	30	12
Toronto	18	29	11
Beijing	19	29	10
Dublin	17	27	10

**Table 6.10:** h-index computed by considering only hubs and hubs + micro hubs for top 20 airports. Third column denotes the difference between first two columns.

## *Chapter 7*

### **Conclusions**

In this work, we explored the concept of the h-index, which is a metric used to evaluate the research output of scholars. The h-index is defined based on the number of publications and their citation count. First, a generalised h-index is defined that can be applied to any graph. The properties of this general h-index were discussed. We also discuss the algorithms for h-index calculation and compare their runtimes on sample graphs.

In Chapter 4, we present a comparison of nine different variants of the h-index and evaluate their effectiveness in ranking research papers. Our analysis focused on papers presented at the SIGMOD and VLDB conferences, and we examined how each index ranked individual papers as well as award-winning papers. Our findings revealed that the h-index was the most effective in ranking award-winning papers at the top. We also compare the performance of each index over time and h-index comes out on top. This demonstrates that h-index is the best metric (amongst the list of metrics considered) to evaluate papers.

In Chapter 5, we apply 4 metrics to evaluate researchers, namely, h-index, h-frac-index, hp-index and hp-frac-index. Amongst these four metrics three of them have been proposed in prior works but hp-frac-index is a newly proposed metric in our work. In our experiments, the hp-frac-index was found to be a robust and effective tool for evaluating the impact of researchers. Various experiments conducted in this work demonstrate that the hp-frac-index outperforms other metrics in ranking awarded researchers higher. The ability to capture individual contributions and resist manipulation makes it a valuable tool for assessing the impact of researcher. It takes into account the importance of a paper's impact by using the paper's h-index, therefore, capturing a second level of research impact. These factors together make hp-frac-index better at ranking the authors. This finding can have significant implications for universities and funding agencies that rely on metrics to evaluate the research output of scholars.

Furthermore in Chapter 6, we apply the h-index concept to the airport network and demonstrates its usefulness in categorising airports and explaining the nature of traffic on a particular airport. We also explain hh-index that is a derivative of h-index. While h-index ranks hub airports (i.e. highly connected airports) on top, the hh-index index uncovers the presence of micro hubs (i.e. airports that feed traffic to hub airports) in the airport network.

Finally, this work highlights the potential of the h-index and its derivatives to be applied in numerous other domains, such as social networks, road networks, Wikipedia links, etc. The h-index and its variations are pragmatic in nature. Moreover, the h-index also exposes the underlying semantics of the domain on which it is applied. For instance, in the academic world, the h-index reflects the scholarly impact of researchers and their contribution to their respective fields. Similarly, in the airport network, the h-index can reflect the connectivity and importance of airports in a given region. This demonstrates the versatility of the h-index concept and its ability to capture the nuances of different domains.

## **Related Publications**

1. Aashay Singhal, Kamalakar Karlapalem. hp-fraction: An index to determine Awarded Researchers. In Companion Proceedings of the ACM Web Conference 2023 (WWW '23 Companion), April 30-May 4, 2023, Austin, TX, USA.

## Bibliography

- [1] Google Scholar. <https://scholar.google.com/>, 2023. [Online; accessed 31-Jan-2023].
- [2] PLOS. <https://journals.plos.org/>, 2023. [Online; accessed 31-Jan-2023].
- [3] S2AG API. <https://www.semanticscholar.org/product/api>, 2023. [Online; accessed 31-Jan-2023].
- [4] Scopus. <https://www.scopus.com/>, 2023. [Online; accessed 31-Jan-2023].
- [5] SIGMOD test of time awards. <https://sigmod.org/sigmod-awards/sigmod-test-of-time-award/>, 2023. [Online; accessed 31-Jan-2023].
- [6] VLDB test of time awards. [https://www.vldb.org/awards\\_10year.html](https://www.vldb.org/awards_10year.html), 2023. [Online; accessed 31-Jan-2023].
- [7] S. Alonso, F. Cabrerizo, E. Herrera-Viedma, and F. Herrera. hg-index: A new index to characterize the scientific output of researchers based on the h-and g-indices. *Scientometrics*, 82(2):391–400, 2010.
- [8] G. Bagler. Analysis of the airport network of india as a complex weighted network. *Physica A: Statistical Mechanics and its Applications*, 387(12):2972–2980, 2008.
- [9] L. Bornmann and H.-D. Daniel. Does the h-index for ranking of scientists really work? *Scientometrics*, 65:391–392, 2005.
- [10] L. Bornmann and H.-D. Daniel. Convergent validation of peer review decisions using the h index: extent of and reasons for type i and type ii errors. *Journal of Informetrics*, 1(3):204–213, 2007.
- [11] L. Bornmann and H.-D. Daniel. What do we know about the h index? *Journal of the American Society for Information Science and technology*, 58(9):1381–1385, 2007.
- [12] L. Bornmann, R. Mutz, and H.-D. Daniel. Are there better indices for evaluation purposes than the h index? a comparison of nine different variants of the h index using data from biomedicine. *Journal of the American Society for Information Science and technology*, 59(5):830–837, 2008.
- [13] L. Bornmann, H. Schier, W. Marx, and H.-D. Daniel. Does the h index for assessing single publications really work? a case study on papers published in chemistry. *Scientometrics*, 89(3):835–843, 2011.
- [14] R. Costas and M. Bordons. The h-index: Advantages, limitations and its relation with other bibliometric indicators at the micro level. *Journal of informetrics*, 1(3):193–203, 2007.
- [15] B. Cronin and L. Meho. Using the h-index to rank influential information scientistss. *Journal of the American Society for Information Science and technology*, 57(9):1275–1278, 2006.

- [16] E. Csajbók, A. Berhidi, L. Vasas, and A. Schubert. Hirsch-index for countries based on essential science indicators data. *Scientometrics*, 73(1):91–117, 2007.
- [17] L. Egghe. How to improve the h-index. *The scientist*, 20(3):15–16, 2006.
- [18] L. Egghe. Theory and practise of the g-index. *Scientometrics*, 69(1):131–152, 2006.
- [19] L. Egghe. On the relation between schubert’s h-index of a single paper and its total number of received citations. *Scientometrics*, 84(1):115–117, 2010.
- [20] L. Egghe. The single publication h-index of papers in the hirsch-core of a researcher and the indirect h-index. *Scientometrics*, 89(3):727–739, 2011.
- [21] G. Ellison. How does the market use citation data? the hirsch index in economics. *American Economic Journal: Applied Economics*, 5(3):63–90, 2013.
- [22] L. Engqvist and J. G. Frommen. The h-index and self-citations. *Trends in ecology & evolution*, 23(5):250–252, 2008.
- [23] J. Guan and X. Gao. Comparison and evaluation of chinese research performance in the field of bioinformatics. *Scientometrics*, 75:357–379, 2008.
- [24] J. E. Hirsch. An index to quantify an individual’s scientific research output. *Proceedings of the National academy of Sciences*, 102(46):16569–16572, 2005.
- [25] C.-C. Hwang and G.-C. Shiao. Analyzing air cargo flows of international routes: an empirical study of taiwan taoyuan international airport. *Journal of Transport Geography*, 19(4):738–744, 2011.
- [26] J. Imperial and A. Rodríguez-Navarro. Usefulness of hirsch’s h-index to evaluate scientific research in spain. *Scientometrics*, 71(2):271–282, 2007.
- [27] T. Jia, K. Qin, and J. Shan. An exploratory analysis on the evolution of the us airport network. *Physica A: Statistical Mechanics and its Applications*, 413:266–279, 2014.
- [28] B. Jin, L. Liang, R. Rousseau, and L. Egghe. The r-and ar-indices: Complementing the h-index. *Chinese science bulletin*, 52(6):855–863, 2007.
- [29] V. Koltun and D. Hafner. The h-index is no longer an effective correlate of scientific reputation. *PLoS One*, 16(6):e0253397, 2021.
- [30] M. Kosmulski et al. A new hirsch-type index saves time and works equally well as the original h-index. *ISSI newsletter*, 2(3):4–6, 2006.
- [31] L. Leydesdorff. Caveats for the use of citation indicators in research and journal evaluations. *Journal of the American Society for Information Science and Technology*, 59(2):278–287, 2008.
- [32] O. Lordan, J. M. Sallan, and P. Simo. Study of the topology and robustness of airline route networks from the complex network approach: a survey and research agenda. *Journal of Transport Geography*, 37:112–120, 2014.
- [33] O. Lordan, J. M. Sallan, P. Simo, and D. Gonzalez-Prieto. Robustness of airline alliance route networks. *Communications in Nonlinear Science and Numerical Simulation*, 22(1-3):587–595, 2015.

- [34] B. Martin. The use of multiple indicators in the assessment of basic research. *Scientometrics*, 36(3):343–362, 1996.
- [35] S. Paleari, R. Redondi, and P. Malighetti. A comparative study of airport connectivity in china, europe and us: Which network provides the best service to passengers? *Transportation Research Part E: Logistics and Transportation Review*, 46(2):198–210, 2010.
- [36] J. Priem, D. Taraborelli, P. Groth, and C. Neylon. Altmetrics: A manifesto. 2011.
- [37] R. Rousseau. The gozinto theorem: Using citations to determine influences on a scientific publication. *Scientometrics*, 11(3-4):217–229, 1987.
- [38] J. F. Salgado and D. Paez. Scientific productivity and hirsch’s h index of spanish social psychology: convergence between productivity indexes and comparison with other areas. *Psicothema*, 19(2):179–189, 2007.
- [39] M. Schreiber. A case study of the hirsch index for 26 non-prominent physicists. *Annalen der Physik*, 16(9):640–652, 2007.
- [40] A. Schubert. Using the h-index for assessing single publications. *Scientometrics*, 78(3):559–565, 2009.
- [41] P. Suau-Sanchez, A. Voltes-Dorta, and H. Rodríguez-Déniz. The role of london airports in providing connectivity for the uk: regional dependence on foreign hubs. *Journal of Transport Geography*, 50:94–104, 2016.
- [42] A. Thor and L. Bornmann. The calculation of the single publication h index and related performance measures: A web application based on google scholar data. *Online Information Review*, 35(2):291–300, 2011.
- [43] J. K. Vanclay. On the robustness of the h-index. *Journal of the American Society for information Science and Technology*, 58(10):1547–1550, 2007.
- [44] W. Webber, A. Moffat, and J. Zobel. A similarity measure for indefinite rankings. *ACM Transactions on Information Systems (TOIS)*, 28(4):1–38, 2010.
- [45] Q. Wu. The w-index: A measure to assess scientific impact by focusing on widely cited papers. *Journal of the American Society for Information Science and Technology*, 61(3):609–614, 2010.
- [46] Z. Yan, Q. Wu, and X. Li. Do hirsch-type indices behave the same in assessing single publications? an empirical study of 29 bibliometric indicators. *Scientometrics*, 109:1815–1833, 2016.