

# CITEX: A new citation index to measure the relative importance of authors and papers in scientific publications

Arindam Pal † and Sushmita Ruj ‡

† TCS Innovation Labs, Kolkata, India. Email address: arindamp@gmail.com

‡ Indian Statistical Institute, Kolkata, India. Email address: sush@isical.ac.in

**Abstract**—Evaluating the performance of researchers and measuring the impact of papers written by scientists is the main objective of citation analysis. Various indices and metrics have been proposed for this. In this paper, we propose a new citation index CITEX, which gives normalized scores to authors and papers to determine their rankings. To the best of our knowledge, this is the first citation index which simultaneously assigns scores to both authors and papers. Using these scores, we can get an objective measure of the reputation of an author and the impact of a paper.

We model this problem as an iterative computation on a publication graph, whose vertices are authors and papers, and whose edges indicate which author has written which paper. We prove that this iterative computation converges in the limit, by using a powerful theorem from linear algebra. We run this algorithm on several examples, and find that the author and paper scores match closely with what is suggested by our intuition. The algorithm is theoretically sound and runs very fast in practice. We compare this index with several existing metrics and find that CITEX gives far more accurate scores compared to the traditional metrics.

**Keywords:** CITATION ANALYSIS, GRAPH ALGORITHMS, MATRIX COMPUTATIONS, EIGENVALUES AND EIGENVECTORS, INFORMATION RETRIEVAL.

## I. INTRODUCTION

In today’s world, numerous papers are written by authors in many journals and conferences. It is difficult for people to judge the quality and impact of an author or a paper, even if they are experts, just by reading a few papers. Thus, measuring the relative importance of authors and papers published in scientific conferences and journals is very important. More importantly, there is a need for an index giving accurate results, which can be computed easily for a large collection of authors and papers.

Many metrics are available to evaluate the importance of journals like Impact factor [12], Immediacy index, citation page rank [18], and Y-factor [4]. There are many metrics which give scores to authors. Most famous of these are Hirsch’s h-index [13], Individual h-index [1], Egghe’s g-index [10], and Zhang’s e-index [21]. However, till now, the only way to evaluate the impact of a paper, is to count the number of citations. It has been observed that surveys and review articles receive more citations than high quality

original research papers. Self citations also increase the number of citations of a paper.

This is why we propose a new metric for evaluating papers for the first time. This new citation index CITEX, gives normalized scores to authors and papers to determine their rankings. Using these scores, we can get an objective measure of the reputation of an author and the impact of a paper. Paper scores are calculated not solely based on the number of citations. Apart from giving scores to papers, we give scores to authors. The author scores and paper scores reinforce each other. Thus, an influential author will increase the score of the paper (s)he writes. An author’s score increases if (s)he writes a good paper. Since the author score increases the score of a paper written by the author, it will be a general tendency to write a paper with an influential author. To prevent this, we assign scores to authors in such a way that the score of a paper gets uniformly divided by the number of people who have co-authored the paper. This basic model can be very easily extended to weighted distribution of scores, where a first author who has the highest contribution receives more weightage than an author who has less contribution.

Existing literature rate an author based on the number of papers that he/she has written, the total number of citations received, average number of citations per paper etc. In Section II, we discuss each of these metrics and state their advantages and disadvantages. None of the existing techniques take into account how the paper scores of an author influence the author’s standing in the academic community, because the paper score is calculated solely based on the number of citations. Our CITEX index is inspired by the ideas of PageRank [16], [6], [5] and HITS [15] algorithms for ranking web pages. In this scheme, paper scores are updated depending on author scores and author scores are updated based on paper scores. This is often referred to as the *Principle of Repeated Improvement* [9]. We prove that the scores asymptotically converge in the limit when the number of iterations is large. In practice, the scores converge within a few iterations.

The main idea is to consider the authors and papers as a network with a disjoint set of nodes, with the set of authors and the set of papers as vertices in the two sets. An edge exists between an author and a paper, if the author has written

the corresponding paper. Apart from these, there is a citation subgraph, consisting of papers as the vertices. A directed edge exists from node  $i$  to node  $j$ , if paper  $i$  cites paper  $j$ . At each step, a paper score gets uniformly divided amongst all its authors. The paper score is the sum of the scores of the authors who have written the paper and the scores due to citation.

The currently available scores like h-index count the number of citations, but not the impact of these papers which have cited this paper. This can be easily manipulated by self-citations. Moreover, since many indices like h-index and number of citations are integers, it is difficult to distinguish between two authors or papers with the same value of the index. Our index has a higher *discriminatory power*, since it is a real number between 0 and 1. It is highly unlikely that two authors or papers will have the same value for CITE<sub>X</sub>. There is also a non-uniformity across various disciplines. In sciences, there is a general tendency to work in large groups. These papers also receive large number of citations compared to papers in computer science. Since the score is divided uniformly among multiple authors, each author score will be reduced, thus affecting paper scores as well. Our index also gives credit to authors who write single-author papers. This should not deter people to do collaborative research.

The problem can be fine-tuned in a number of ways. We can take into consideration the recommendation of authors by other authors and consider weighted distribution. The problem also has a number of ramifications. A similar index can be designed for product-customer recommendation system, where customers can recommend each other depending upon the reputations (similar to citation of papers), and a customer can recommend a product (similar to writing papers). The difference here is that, the score of a product is not uniformly divided amongst customers, and information cascades have to be taken into account when calculating the product scores. Other interdependent networks, which reinforce each other, can be treated similarly.

The rest of the paper is organized as follows. In section II, we discuss the related work on citation analysis, define the metrics and compare them. In section III, we define the problem and propose the model to analyze it. We give an informal description of the algorithm and present the rules to iteratively compute the author and paper scores in section IV. In section V, we mathematically analyze the iterative procedure and prove that it converges to the eigenvectors of certain matrices. In section VI, we execute the algorithm on some illustrative examples, and show that the author scores and paper scores give good indication of their importance, as can be seen from the underlying graph structure. In section VII, we discuss some extensions of the basic algorithm and future direction to work on. We conclude the paper in section VIII with some future directions to work on.

## II. RELATED WORK

### A. Different metrics used in citation analysis

Previously, there have been several attempts to measure the impact of authors and papers. We list here some of them along with their definition.

- 1) **Number of papers ( $N_p$ ):** Total number of papers written by an author.
- 2) **Number of citations ( $N_c$ ):** Total number of citations for all papers written by an author.
- 3) **Average number of citations per paper:** Ratio of total number of citations and total number of papers, *i.e.*,  $\frac{N_c}{N_p}$ . This is sometimes also called the *impact factor*.
- 4) **Average number of citations per author:** For each paper, its citation count is divided by the number of authors for that paper to give the normalized citation count for the paper. The normalized citation counts are then summed across all papers to give the average number of citations per author.
- 5) **Average number of papers per author:** For each paper, the inverse of the number of authors gives the normalized author count for the paper. The normalized author counts are then summed across all papers to give the average number of papers per author.
- 6) **Average number of authors per paper:** The sum of the author counts across all papers, divided by the total number of papers.
- 7)  **$h$ -index [13]:** An author has index  $h$ , if  $h$  of his  $N_p$  papers have at least  $h$  citations each, and the rest of the  $N_p - h$  papers have no more than  $h$  citations each.
- 8)  **$g$ -index [10]:** Given a set of articles ranked in decreasing order of the number of citations that they received, the  $g$ -index is the (unique) largest number such that the top  $g$  articles together received at least  $g^2$  citations.
- 9)  **$e$ -index [21]:** It is the square root of surplus citations in the  $h$ -set beyond the theoretical minimum ( $h^2$ ) required to obtain a  $h$ -index of  $h$ . It is useful for highly cited scientists and for comparing those with the same  $h$ -index but different citation patterns.
- 10) **Number of significant papers:** Total number of papers with more than  $c$  citations for some integer  $c$ .
- 11) **Number of citations to the most cited papers:** Total number of citations to the  $k$  most cited papers for some integer  $k$ .
- 12) **Eigenfactor:** The Eigenfactor score [3] is a rating of the total importance of a scientific journal. Journals are rated according to the number of incoming citations, with citations from highly ranked journals weighted to make a larger contribution to the Eigenfactor than those from poorly ranked journals [2]. As a measure of importance, the Eigenfactor score scales with the total impact of a journal. Journals generating higher impact to the field have larger Eigenfactor scores. However, it is not clear whether Eigenfactor gives better estimate than raw citation count. [8]

TABLE I  
COMPARISON BETWEEN DIFFERENT CITATION METRICS

Citation Metric	Advantage	Disadvantage
<i>Number of papers</i>	Measure of productivity.	Importance of papers not considered.
<i>Number of citations</i>	Measures impact of an author.	A few highly cited papers increase the total. Survey and review articles are cited more than original research papers. Favors established authors.
<i>Average number of citations per paper</i>	Allows comparison of scientists of different ages.	Rewards low productivity.
<i>Average number of citations per author</i>	Measures impact of an author.	Difficult to distinguish between authors whose average is same, but citation patterns are different.
<i>Average number of papers per author</i>	Measures average productivity	Does not measure impact of papers.
<i>Average number of authors per paper</i>	Measures collaboration between authors.	Does not consider importance of authors and papers.
<i>h-index</i>	Measures both the quality and quantity of scientific output.	Does not account for the number of authors of a paper.  Different fields with different number of citations will have different h-index. Can be manipulated through self-citations.
<i>g-index</i>	Gives more weight to highly-cited articles.	Unlike the h-index, the g-index saturates whenever the average number of citations for all published papers exceeds the total number of published papers.
<i>e-index</i>	Differentiates between scientists with identical h-indices but different citations.	Can't be used independently. Must be used together with the h-index.
<i>Number of papers with at least c citations</i>	Measures the broad and sustained impact of an author.	Difficult to find the right value of $c$ . Different values of $c$ favors different authors.
<i>Number of citations to the k most cited papers</i>	Identifies the most influential authors.	Not a single number, so difficult to compare. Different values of $k$ favors different authors.
<i>Eigenfactor</i>	Takes into account impact of the citing papers in addition to the number of citations.	Does not give author scores. Importance of authors have to inferred indirectly from the papers (s)he has written.

### B. Comparison between different citation metrics

In this section, we compare the different citation indices and state their advantages and disadvantages. The comparison is presented in Table I.

### C. Comparison with similar works

There are a number of previous attempts to rank authors and papers based on importance. The SIMRANK algorithm by Jeh and Widom [14] gives a measure of the similarity between two objects based on their relationships with other objects. Their basic idea is that two objects are similar if they are related to similar objects. Note that this only measures the similarity of two objects, not their relative ranking, so this is different from what we are trying to do. Zhou et. al. [22] proposed a method for co-ranking authors and their publications using several networks associated with authors and papers. Although there is some similarity between our algorithm and their approach, there are fundamental differences between the two. Their co-ranking framework is based on coupling two random walks that separately rank authors and documents using the PAGERANK algorithm. Our algorithm is designed from scratch and does not use the PAGERANK algorithm. Moreover, our algorithm is much simpler and the computations required is also far lesser than what is required in their method. Walker et. al. [20] gave a new algorithm called CITERANK. The ranking of papers is based on a network traffic model, which uses a variation of the PAGERANK algorithm. A paper is selected randomly from the set of all papers with a probability that decays

exponentially with the age of the paper. Chen et. al. [7] uses a PAGERANK based algorithm to assess the relative importance of all publications. Their goal is to find some exceptional papers or “gems” that are universally familiar to physicists. Sun and Giles [19] propose a popularity weighted ranking algorithm for academic digital libraries. They use the popularity of a publication venue and compare their method with the PAGERANK algorithm, citation counts and the HITS algorithm.

### D. Some structures in citation analysis

- **Collaboration graph:** This is a graph associated with the authors. The nodes of the graph are the authors. There is an undirected edge between two nodes, if the corresponding authors have written a paper together.
- **Citation graph:** This is a graph associated with the papers. The nodes of the graph are the papers. There is a directed edge from a paper to another paper, if the first paper has cited the second paper.
- **Publication graph:** This is a graph relating the authors with the papers. The nodes of the graph are the authors and the papers. There is an undirected edge between two nodes, if the author has written the paper.

## III. PROBLEM DEFINITION AND MODEL

We have a set of  $m$  authors  $\mathcal{A} = \{a_1, \dots, a_m\}$  and a set of  $n$  papers  $\mathcal{P} = \{p_1, \dots, p_n\}$ . We represent this by a *publication graph*  $G_P = (V_P, E_P)$ , whose vertices are the set of authors and papers, i.e.,  $V_P = \mathcal{A} \cup \mathcal{P}$ . There is an

undirected edge between author  $a_i$  and paper  $p_j$ , if author  $a_i$  has written paper  $p_j$ . Note that this is a *symmetric* relation, so the edges are *undirected*. Since there are only edges between authors and papers, the publication graph is a *bipartite graph*. Associated with this, there is an  $m \times n$  *publication matrix*  $M$ , whose rows and columns are  $a_1, \dots, a_m$  and  $p_1, \dots, p_n$  respectively, and whose  $(i, j)^{th}$  entry  $m_{ij} = 1$ , if and only if author  $a_i$  has written paper  $p_j$ .

Moreover, there is a *citation graph*  $G_C = (V_C, E_C)$  associated with the papers, whose vertices are the set of papers, *i.e.*,  $V_C = \mathcal{P}$ . There is a directed edge from paper  $p_j$  to paper  $p_k$ , if paper  $p_j$  has cited paper  $p_k$ . Note that this is an *asymmetric* relation, so the edges are *directed*. Associated with this, there is an  $n \times n$  *citation matrix*  $C$ , whose both rows and columns are  $p_1, \dots, p_n$ , and whose  $(j, k)^{th}$  entry  $c_{jk} = 1$ , if and only if paper  $p_j$  has cited paper  $p_k$ . Note that the citation graph can't have any *directed cycle*. This is because a paper can only cite a previously published paper, so they are *totally ordered* in time. This also means that if the papers are numbered in *decreasing order of time* (newer first), the resulting citation matrix will be *upper-triangular*. An example of a publication graph and a citation graph is given in Figure 1.

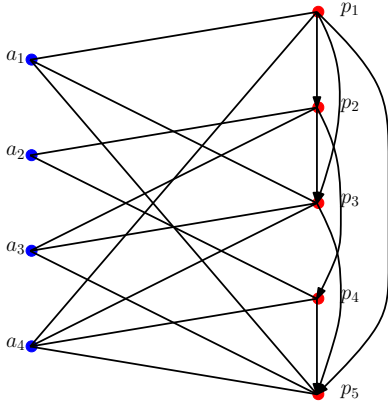


Fig. 1. The publication and citation graphs for Example 1 showing authors, papers and citations.

The following sets are important for further development.

- 1) For an author  $a$ ,  $PAPERS(a)$  is defined as the set of papers written by author  $a$ . In other words,  $PAPERS(a) = \{p \in \mathcal{P} : (a, p) \in E_P\}$ .
- 2) For a paper  $p$ ,  $AUTHORS(p)$  is defined as the set of authors who have written paper  $p$ . In other words,  $AUTHORS(p) = \{a \in \mathcal{A} : (a, p) \in E_P\}$ .
- 3) For a paper  $p$ ,  $CITE(p)$  is defined as the set of papers who have cited paper  $p$ . In other words,  $CITE(p) = \{q \in \mathcal{P} : (q, p) \in E_C\}$ .
- 4) For a paper  $p$ ,  $REF(p)$  is defined as the set of papers which have been given as reference (cited) by paper  $p$ . In other words,  $REF(p) = \{q \in \mathcal{P} : (p, q) \in E_C\}$ .

Our goal is to assign scores to authors and papers using

the structure of the publication and citation graphs, so that important authors and papers get higher scores.

## IV. DESCRIPTION OF THE CITEX INDEX

### A. Informal description of the algorithm

In this section, we give an overview of our proposed algorithm. The algorithm maintains a set of author scores and paper scores, which are initially set to 1. This initial choice of scores is arbitrary, and the scores can be set to any nonzero value. Then we update the scores considering the relationship between authors and papers (who has authored which paper) and relationship between papers (which paper has cited which paper). This critically uses the publication graph and the citation graph. We use the *Principle of Repeated Improvement* [9] to iteratively compute the new scores based on the previous scores. More specifically, the author scores for the next iteration is computed from the paper scores for the current iteration. The paper scores for the next iteration is computed from the author scores and the paper scores for the current iteration. In every iteration, we normalize the scores by dividing them by the sum of the individual scores, so that each of them lies between 0 and 1, and they add up to 1. We continue to do this till the author scores and the paper scores converge or a specified number of iterations have been completed. The *Principle of Repeated Improvement* states that each improvement of author scores will lead to a further improvement of paper scores, and vice versa. The final author scores and paper scores are the measure of importance of the authors and the papers. The higher the score is, the higher is the impact of an author and a paper.

### B. Computing author and paper scores

For each author  $a_i$ , we have an *author score* (*a-score*)  $x_i$ , and for each paper  $p_j$ , we have a *paper score* (*p-score*)  $y_j$ . We represent the set of author scores as a column vector  $\mathbf{x} = (x_1, \dots, x_m)^T$  and the set of paper scores as a column vector  $\mathbf{y} = (y_1, \dots, y_n)^T$ . We initialize all author and paper scores to one, *i.e.*,  $\mathbf{x} = \mathbf{y} = \mathbf{1}$ . Then, we iteratively update the *a-scores* and *p-scores* using the following rules.

- 1) For each paper  $p_j$ , its *adjusted p-score*  $\bar{y}_j$  is given by the *p-score*  $y_j$  divided by the number of authors who have written the paper. In other words,  $\bar{y}_j = \frac{y_j}{k}$ , for  $j = 1, \dots, n$ , where  $k = |AUTHORS(p_j)|$  is the number of co-authors of the paper  $p_j$ .
- 2) For each author  $a_i$ , set his *a-score*  $x_i$  to be the sum of the adjusted *p-scores* of all the papers that he has authored. In other words,  $x_i = \sum_{j \in PAPERS(i)} \bar{y}_j$ , for  $i = 1, \dots, m$ .
- 3) For each paper  $p_j$ , set its *p-score*  $y_j$  to be the sum of the *a-scores* of all the authors who have co-authored the paper  $p_j$ . In other words,  $y_j = \sum_{i \in AUTHORS(j)} x_i$ , for  $j = 1, \dots, n$ .
- 4) For each paper  $p_j$ , add to its *p-score*  $y_j$ , the sum of the *p-scores* of all the papers who have cited the paper  $p_j$ . In other words,  $y_j = y_j + \sum_{k \in CITE(j)} y_k$ , for  $j = 1, \dots, n$ .

We normalize the scores by dividing the author (paper) scores by the sum of the author (paper) scores, so that each score lies between 0 and 1, and the sum of the scores is 1.

## V. MATHEMATICAL ANALYSIS

### A. Analysis of author scores and paper scores

We observe that the rule  $y_j = \sum_{i \in AUTHORS(j)} x_i$  can be rewritten as  $y_j = \sum_{i=1}^m m_{ij} x_i$ , since  $m_{ij} = 1$  if and only if  $i \in AUTHORS(j)$ . Consider the matrix-vector equation  $\mathbf{y} \leftarrow M^T \mathbf{x}$ . The  $j^{\text{th}}$  row of this equation is  $y_j = \sum_{i=1}^m m_{ij} x_i$ . Hence, this matrix-vector equation succinctly encodes all  $n$  scalar equations for  $j = 1, \dots, n$ .

The corresponding equation for  $\mathbf{x}$  is similar, but a little more involved. The equation  $x_i = \sum_{j \in PAPERS(i)} \bar{y}_j$  can be written as  $x_i = \sum_{j=1}^n m_{ij} \bar{y}_j$ , since  $m_{ij} = 1$  if and only if  $j \in PAPERS(i)$ . Let  $\bar{\mathbf{y}} = (\bar{y}_1, \dots, \bar{y}_n)^T$ . Consider the matrix-vector equation  $\mathbf{x} \leftarrow M \bar{\mathbf{y}}$ . The  $i^{\text{th}}$  row of this equation is  $x_i = \sum_{j=1}^n m_{ij} \bar{y}_j$ . Hence, this matrix-vector equation succinctly encodes all  $m$  scalar equations for  $i = 1, \dots, m$ . Further note that,  $\bar{y}_j = \frac{y_j}{|AUTHORS(p_j)|} = \frac{y_j}{\sum_{i=1}^m m_{ij}}$ . Hence,  $x_i = \sum_{j=1}^n \left( \frac{m_{ij}}{\sum_{i=1}^m m_{ij}} \right) y_j = \sum_{j=1}^n w_{ij} y_j$ , where  $w_{ij} = \frac{m_{ij}}{\sum_{i=1}^m m_{ij}}$  is the weight associated with the paper  $p_j$ . Now,  $\mathbf{x}$  can be written as  $\mathbf{x} \leftarrow W \mathbf{y}$ , where  $W$  is the  $m \times n$  weight matrix whose  $(i, j)^{\text{th}}$  entry is  $w_{ij}$ .

The equation  $y_j = y_j + \sum_{k \in CITE(j)} y_k$  can be rewritten as  $y_j = y_j + \sum_{k=1}^n c_{kj} y_k$ , since  $c_{kj} = 1$  if and only if  $k \in CITE(j)$ . This can be written as the matrix-vector equation  $\mathbf{y} \leftarrow (I + C^T) \mathbf{y}$ , where  $I$  is the  $n \times n$  identity matrix.

Let the initial author vector and paper vector be  $\mathbf{x}^{(0)}$  and  $\mathbf{y}^{(0)}$  respectively. If we start with the equation  $\mathbf{x} \leftarrow W \mathbf{y}$ , the successive iterations proceed as below.

$$\mathbf{x}^{(1)} = W \mathbf{y}^{(0)}, \quad (1)$$

$$\mathbf{y}^{(1)} = M^T \mathbf{x}^{(1)} = M^T W \mathbf{y}^{(0)}, \quad (2)$$

$$\mathbf{y}^{(1)} = (I + C^T) \mathbf{y}^{(1)} = (I + C^T) M^T W \mathbf{y}^{(0)}. \quad (3)$$

Similarly, if we start with the equation  $\mathbf{y} \leftarrow M^T \mathbf{x}$ , the successive iterations proceed as below.

$$\mathbf{y}^{(1)} = M^T \mathbf{x}^{(0)}, \quad (4)$$

$$\mathbf{y}^{(1)} = (I + C^T) \mathbf{y}^{(1)} = (I + C^T) M^T \mathbf{x}^{(0)}, \quad (5)$$

$$\mathbf{x}^{(1)} = W \mathbf{y}^{(1)} = W(I + C^T) M^T \mathbf{x}^{(0)}. \quad (6)$$

Proceeding similarly, at the  $k$ -th iteration the author and paper vectors are given by,

$$\mathbf{x}^{(k)} = [W(I + C^T) M^T]^k \mathbf{x}^{(0)}, \quad (7)$$

$$\mathbf{y}^{(k)} = [(I + C^T) M^T W]^k \mathbf{y}^{(0)}. \quad (8)$$

### B. Proof of convergence of author scores and paper scores

In this section, we will prove the following theorem.

**Theorem 1.** *The sequences  $\mathbf{x}^{(k)}$  and  $\mathbf{y}^{(k)}$ ,  $k = 0, 1, 2, \dots$  converge to the limits  $\mathbf{x}^*$  and  $\mathbf{y}^*$  respectively. Moreover,  $\mathbf{x}^*$  is the principal eigenvector of the matrix  $W(I + C^T) M^T$  and  $\mathbf{y}^*$  is the principal eigenvector of the matrix  $(I + C^T) M^T W$ .*

Further, both  $\mathbf{x}^*$  and  $\mathbf{y}^*$  are non-negative and non-zero vectors.

*Proof:* From the above discussion we have,  $\mathbf{x}^{(k)} = P^k \mathbf{x}^{(0)}$  and  $\mathbf{y}^{(k)} = Q^k \mathbf{y}^{(0)}$ , where  $P = W(I + C^T) M^T$  and  $Q = (I + C^T) M^T W$ . Note that  $P$  is an  $m \times m$  square matrix, whereas  $Q$  is an  $n \times n$  square matrix. Moreover,  $\mathbf{x}^{(k+1)} = P^{k+1} \mathbf{x}^{(0)} = P \cdot P^k \mathbf{x}^{(0)} = P \mathbf{x}^{(k)}$ . If the author score  $\mathbf{x}^{(k)}$  converges to the vector  $\mathbf{x}^*$  in the limit when  $k \rightarrow \infty$ , then this vector should satisfy  $P \mathbf{x}^* = \mathbf{x}^*$ . This means that  $\mathbf{x}^*$  is an eigenvector of  $P$ , with the corresponding eigenvalue being 1. Similarly, if the paper score  $\mathbf{y}^{(k)}$  converges to the vector  $\mathbf{y}^*$  in the limit when  $k \rightarrow \infty$ , then  $\mathbf{y}^*$  must be an eigenvector of  $Q$ , with the corresponding eigenvalue being 1.

To prove that a non-negative eigenvalue and a non-negative eigenvector exists, we use the following theorem from linear algebra.

**Theorem 2** (PERRON-FROBENIUS THEOREM). [17], [11], [9] *Let  $A = (a_{ij})$  be an  $n \times n$  non-negative matrix, meaning that  $a_{ij} \geq 0, \forall i, j : 1 \leq i, j \leq n$ . Then the following statements hold.*

- 1) *A has a real eigenvalue  $c \geq 0$  such that  $c > |c'|$  for all other eigenvalues  $c'$ .*
- 2) *There is an eigenvector  $v$  with non-negative real components corresponding to the largest eigenvalue  $c : Av = cv, v_i \geq 0, 1 \leq i \leq n$ , and  $v$  is unique up to multiplication by a constant.*
- 3) *If the largest eigenvalue  $c$  is equal to 1, then for any starting vector  $\mathbf{x}^{(0)} \neq 0$  with non-negative components, the sequence of vectors  $A^k \mathbf{x}^{(0)}$  converge to a vector in the direction of  $v$  as  $k \rightarrow \infty$ .*

Thus, by the Perron-Frobenius theorem, the author and paper scores both converge to unique non-negative vectors  $\mathbf{x}^*$  and  $\mathbf{y}^*$ , after repeated applications of the update rules. These two vectors are the limiting values of the author and paper scores. Moreover, none of the vectors  $\mathbf{x}^*$  and  $\mathbf{y}^*$  can be the zero vector. At least one of their components must be non-zero, because the initial vectors  $\mathbf{x}^{(0)}$  and  $\mathbf{y}^{(0)}$  are the all-1 vectors, and at each iteration the vectors are normalized. So the sum of their components add up to 1. ■

### C. Time-complexity of each iteration

We have to multiply some matrices as can be seen from equations (7) and (8). Computing the product of a  $m \times n$  matrix and a  $n \times p$  matrix requires  $O(mnp)$  time. Computing the matrices  $W(I + C^T) M^T$  and  $(I + C^T) M^T W$  takes  $O(mn(m+n))$  time each. Hence, each iteration can be done in  $O(mn(m+n))$  time.

## VI. EXPERIMENTAL ANALYSIS

We compute the author and paper scores for some graphs and show that they match with our intuition.

### A. Example 1

For the graph in Figure 1, the author and paper vectors, publication matrix and citation matrix are given below. Note that for this example,  $m = 4, n = 5$ .

$$\mathbf{x} = [1, 1, 1, 1], \mathbf{y} = [1, 1, 1, 1, 1],$$

$$M = \begin{bmatrix} 1 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 & 1 \end{bmatrix}, C = \begin{bmatrix} 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}.$$

The author scores after 10 iterations are given below. These scores converge (up to 3 decimal places) as there is no change between two successive iterations.

$$\mathbf{x} = [0.259, 0.132, 0.289, 0.320],$$

$$\mathbf{y} = [0.082, 0.141, 0.264, 0.123, 0.390].$$

Intuitively it is clear that author  $a_4$  and paper  $p_5$  should get the highest scores. Author  $a_4$  has written 4 papers  $p_1, p_3, p_4, p_5$  and some of them have high paper scores. Similarly,  $p_5$  has been written by 3 authors  $a_1, a_3, a_4$  and cited by 3 papers  $p_1, p_3, p_4$ , some of which have high scores.  $a_2$  gets the lowest score as it has written only 2 papers  $p_2, p_4$ , none of which have high paper scores. Similarly,  $p_1$  gets the lowest score since it has no citations, although it has two authors  $a_1, a_4$ .

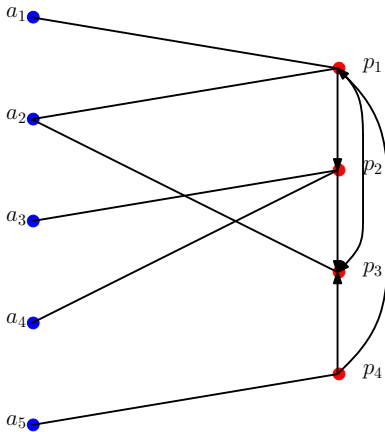


Fig. 2. The publication and citation graphs for Example 2 showing authors, papers and citations.

### B. Example 2

For the graph in Figure 2, the author and paper vectors, publication matrix and citation matrix are given below. Note that for this example,  $m = 5, n = 4$ .

$$\mathbf{x} = [1, 1, 1, 1, 1], \mathbf{y} = [1, 1, 1, 1],$$

$$M = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, C = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 \end{bmatrix}.$$

The author scores after 10 iterations are given below. These scores converge (up to 3 decimal places) as there is no change between two successive iterations.

$$\mathbf{x} = [0.106, 0.590, 0.152, 0.152, 0.000],$$

$$\mathbf{y} = [0.212, 0.304, 0.484, 0.000].$$

Intuitively it is clear that author  $a_2$  and paper  $p_3$  should get the highest scores. Author  $a_2$  has written two papers  $p_1$  and  $p_3$ . In turn,  $p_3$  has been cited by papers  $p_1, p_2$  and  $p_4$ . On the other hand,  $p_4$  gets the lowest paper score 0, since it is not cited by any paper.  $a_5$  gets the lowest author score since it has only written paper  $p_4$ , which has a score 0. These scores matches with our intuition.

## VII. EXTENSIONS

### A. Weights on edges of the publication graph

In Section V we assumed that all authors contribute equally to a paper. However, this is not true in practice. Different co-authors have different contribution to a paper. We can easily incorporate this feature in our CITE index, by slightly modifying the publication matrix  $M$ . In Section V,  $M$  was a 0-1 matrix. For the weighted version, we construct a matrix  $N$ , where edge weight  $n_{ij}$  denotes the contribution of author  $i$  in writing paper  $j$ . An author having higher contribution has higher weight, compared to an author having lesser contribution. The new matrix  $N$  might not be a 0-1 matrix.  $W'$  is the matrix corresponding to  $N$ . Hence,  $x_i = \sum_{j=1}^n \left( \frac{n_{ij}}{\sum_{i=1}^m n_{ij}} \right) y_j = \sum_{j=1}^n w'_{ij} y_j$ , where  $w'_{ij} = \frac{n_{ij}}{\sum_{i=1}^m n_{ij}}$  is the weight associated with the paper  $p_j$ . Now,  $\mathbf{x}$  can be written as  $\mathbf{x} \leftarrow W' \mathbf{y}$ , where  $W'$  is the weight matrix whose  $(i, j)^{th}$  entry is  $w'_{ij}$ .

Now the author and paper vectors at the  $k$ -th iteration can be written as,

$$\mathbf{x}^{(k)} = [W'(I + C^T)N^T]^k \mathbf{x}^{(0)},$$

$$\mathbf{y}^{(k)} = [(I + C^T)N^T W']^k \mathbf{y}^{(0)}.$$

### B. Reputation of authors

Our CITE index can be modified to include the reputation of authors. Each author can rank other authors who he/she believes has done original, ground breaking work. We define the *author reputation graph*  $G_R = (V_R, E_R)$ , similar to the citation graph. The vertices are the set of authors, i.e.,  $V_R = \mathcal{A}$ . Thus, the number of vertices in this graph is  $m$ . A directed edge from node  $i$  to node  $j$  of weight  $r_{ij}$  exists, if author  $i$  has rated author  $j$  with a score  $r_{ij}$ . The *author reputation matrix*  $R$  is defined similarly. For an author  $a$ ,  $REP(a)$  is defined as the set of authors who have ranked author  $a$ . In other words,  $REP(a) = \{b \in \mathcal{A} : (b, a) \in E_R\}$ .

On incorporating the reputation matrix, another rule is added to the list. For each author  $a_i$ , add to its  $a$ -score  $x_i$ ,

the sum of the  $a$ -scores of all the authors who have rated author  $a_i$ . In other words,  $x_i = x_i + \sum_{k \in REP(i)} x_k = x_i + \sum_{k=1}^m r_{ki} x_k$ , for  $i = 1, \dots, m$ .

Let the initial author vector and paper vector be  $\mathbf{x}^{(0)}$  and  $\mathbf{y}^{(0)}$  respectively. If we start with the equation  $\mathbf{x} \leftarrow W\mathbf{y}$ , the successive iterations proceed as below.

$$\begin{aligned}\mathbf{x}^{(1)} &= W\mathbf{y}^{(0)}, \\ \mathbf{x}^{(1)} &= (I + R^T)\mathbf{x}^{(1)} = (I + R^T)W\mathbf{y}^{(0)} \\ \mathbf{y}^{(1)} &= M^T\mathbf{x}^{(1)} = M^T(I + R^T)W\mathbf{y}^{(0)}, \\ \mathbf{y}^{(1)} &= (I + C^T)\mathbf{y}^{(1)} = (I + C^T)M^T(I + R^T)W\mathbf{y}^{(0)}.\end{aligned}$$

Proceeding similarly, at the  $k$ -th iteration the author and paper vectors are given by,

$$\begin{aligned}\mathbf{x}^{(k)} &= [(I + R^T)W(I + C^T)M^T]^k \mathbf{x}^{(0)}, \\ \mathbf{y}^{(k)} &= [(I + C^T)M^T(I + R^T)W]^k \mathbf{y}^{(0)}.\end{aligned}$$

### C. Evaluation of journals and conferences

One important question is how to measure the quality of scientific journals and conferences. If we have the author scores and paper scores of all authors and papers published in a journal/conference, we can use the average author score and the average paper score as a metric for determining the quality of the journal/conference.

## VIII. CONCLUSION AND FUTURE DIRECTIONS

In this paper, we have proposed a new citation metric CITE<sub>X</sub> to judge the quality of authors and papers in scientific publications. CITE<sub>X</sub> assigns scores to authors and papers, with higher scores indicating more importance, by analyzing the link structures in the publication graph and the citation graph. We also considered some extensions to the basic scheme. Here are some future directions to work on.

- In a real-world scenario, authors and papers will be added over time. Dynamically modifying the scores from the current scores in an incremental fashion is a challenging problem.
- Applying this framework to the setting of customer recommendation of products will be an interesting idea. Here the nodes of the graph are customers and products. A customer can give a score to a product, which is like the weighted version of the publication graph.
- This is an example of an *interdependent network*, where there are two graphs – the collaboration/reputation graph and the citation graph. In addition, there is a publication graph which records the cross-edges between the nodes in the two graphs. Extending CITE<sub>X</sub> to other interdependent networks will be an interesting direction to think about.
- Using further parameters such as time of publication and age of authors, in addition to the link structure will require further thoughts.
- Can this technique be extended to *directly* assign scores to journals and conferences, rather than doing it *indirectly*, as in section VII-C?

- The analytical power of eigenvector-based methods is not yet fully understood. It would be interesting to pursue this question in the context of the algorithm presented here. Considering random graph models that contain enough structure to capture certain global properties of the model is a promising direction.

## REFERENCES

- [1] Pablo D. Batista, Mônica G. Campiteli, and Osame Kinouchi. Is it possible to compare researchers with different scientific interests? *Scientometrics*, 68(1):179–189, 2006.
- [2] Carl Bergstrom. Measuring the value and prestige of scholarly journals. *College & Research Libraries News*, 68(5):3146, 2007.
- [3] Carl T Bergstrom, Jevin D West, and Marc A Wiseman. The Eigenfactor<sup>TM</sup> metrics. *The Journal of Neuroscience*, 28(45):11433–11434, 2008.
- [4] Johan Bollen, Herbert Van de Sompel, Joan A. Smith, and Richard Luce. Toward alternative metrics of journal impact: A comparison of download and citation data. *Inf. Process. Manage.*, 41(6):1419–1440, 2005.
- [5] Johan Bollen, Marko A Rodriguez, and Herbert Van de Sompel. Journal status. *Scientometrics*, 69(3):669–687, 2006.
- [6] Sergey Brin and Lawrence Page. The anatomy of a large-scale hypertextual web search engine. *Computer networks and ISDN systems*, 30(1):107–117, 1998.
- [7] Peng Chen, Huafeng Xie, Sergei Maslov, and Sidney Redner. Finding scientific gems with google’s pagerank algorithm. *Journal of Informetrics*, 1(1):8–15, 2007.
- [8] Philip M Davis. Eigenfactor: Does the principle of repeated improvement result in better estimates than raw citation counts? *Journal of the American Society for Information Science and Technology*, 59(13):2186–2188, 2008.
- [9] David Easley and Jon Kleinberg. *Networks, crowds, and markets: Reasoning about a highly connected world*. Cambridge University Press, 2010.
- [10] Leo Egghe. Theory and practise of the  $g$ -index. *Scientometrics*, 69(1):131–152, 2006.
- [11] Georg Ferdinand Frobenius. *Über Matrizen aus nicht negativen Elementen*. Königliche Akademie der Wissenschaften, 1912.
- [12] Eugene Garfield. Citation indexes for science. *Science*, 122(3159):108–111, 1955.
- [13] Jorge E Hirsch. An index to quantify an individual’s scientific research output. *Proceedings of the National academy of Sciences of the United States of America*, 102(46):16569–16572, 2005.
- [14] Glen Jeh and Jennifer Widom. Simrank: a measure of structural-context similarity. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 538–543. ACM, 2002.
- [15] Jon M. Kleinberg. Authoritative sources in a hyperlinked environment. *J. ACM*, 46(5):604–632, 1999.
- [16] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The PageRank citation ranking: Bringing order to the web. Technical report, Stanford InfoLab, 1999.
- [17] Oskar Perron. Zur theorie der matrices. *Mathematische Annalen*, 64(2):248–263, 1907.
- [18] Gabriel Pinski and Francis Narin. Citation influence for journal aggregates of scientific publications: Theory, with application to the literature of physics. *Inf. Process. Manage.*, 12(5):297–312, 1976.
- [19] Yang Sun and C Lee Giles. Popularity weighted ranking for academic digital libraries. In *Proceedings of the 29th European conference on IR research*, pages 605–612. Springer-Verlag, 2007.
- [20] Dylan Walker, Huafeng Xie, Koon-Kiu Yan, and Sergei Maslov. Ranking scientific publications using a model of network traffic. *Journal of Statistical Mechanics: Theory and Experiment*, 2007(06):P06010, 2007.
- [21] Chun-Ting Zhang. The e-index, complementing the h-index for excess citations. *PLoS ONE*, 5(5):1–22, 2009.
- [22] Ding Zhou, Sergey A Orshanskiy, Hongyuan Zha, and C Lee Giles. Co-ranking authors and documents in a heterogeneous network. In *Data Mining, 2007. ICDM 2007. Seventh IEEE International Conference on*, pages 739–744. IEEE, 2007.