

A time-aware citation-based model for evaluating scientific products

[Extended Abstract]

Gianna M. Del Corso^{*}
Dipartimento di Informatica
Università di Pisa
Largo Pontecorvo, 3, 56127 Pisa, Italy
delcorso@di.unipi.it

F. Romani
Dipartimento di Informatica
Università di Pisa
Largo Pontecorvo, 3, 56127 Pisa, Italy
romani@di.unipi.it

ABSTRACT

Recently a citation-based model for ranking scientific journals and papers together with authors has been proposed [4, 5]. In that model, papers, authors and journals mutually contribute to the attribution of ranking score to each other. The rank of each subject is computed as the stationary distribution of a suitable Markov chain. In this paper, we add into the model a factor accounting for the aging of papers. In particular, the importance of each paper slowly decreases as time elapses unless a fresh citation arrives conferring new importance to the cited papers. The experimental part shows the effectiveness of the introduction of the time into the model. In fact, new papers gain importance with respect to older ones, sustaining in this way new trends of research with respect to subjects popular years ago.

Categories and Subject Descriptors

G.1.3 [Numerical Linear Algebra]: Eigenvalues and Eigenvectors

1. INTRODUCTION

Evaluation of scientific research has always been a very important problem. Recently, the number of scientific journals and papers has increased at an almost exponential rate [14] making the task of using and evaluating scientific literature much harder than in the past. For example, researchers now rely on search engines such as Google Scholar to choose what to read or what to cite. This problem does not affect only researchers but also funding agencies, university administrators, and reviewers called to evaluate productivity of researchers and institutions. Most of the time it is impossible to give an in-depth evaluation of the research

^{*}Corresponding author.

performed by a scholar or institution, and it is becoming common to use indirect indicators of quality.

Among such indicators, the most popular are currently those based on citation analysis which allow a quick, simple and objective evaluation of a large amount of data when peer review is not practicable.

In the literature one can find many different metrics for evaluating papers, journals, or researchers. The reason is that there are many possible different purposes for ranking. For instance, the ranking of journals is interesting for librarians to decide on subscriptions and for authors to decide where to publish. The ranking of papers is becoming useful for untangling the maze of papers published everyday, and decide what to read or what to cite. Likewise, it is becoming common to evaluate scholars on the basis of their scientific productivity for distributing funds, or even for hiring people.

Among the different methods proposed in the literature for ranking scientific research we can distinguish between methods based on citation statistics — such as Impact Factor (IF) (see [9] and references therein for an historical review), simple Citation Count, the MCQ by the American Mathematical Society [2] — and methods based on approaches similar to Google PageRank, such as Eigenfactor [3], SCImago [13] and others [12, 11].

Metrics based on citation statistics are easy to compute but not all the scientific community agrees on the effectiveness of these metrics to capture concepts such as reputation or influence. Metrics based on PageRank-like techniques seem more appealing since the effectiveness of PageRank for ranking web pages is proved by everyday use, and citations in a paper have a similar role as links in web pages. The main idea is that not all citations are equal and that, rather than the number, one should consider the “quality” of citations.

In [4, 5], jointly with D.A. Bini, we propose an integrated three-class model for the ranking of papers, authors, and journals loosely inspired by the PageRank algorithm. In our model papers, authors, and journals represent three distinct classes that mutually contribute to the attribution of a ranking score to each element of each class. The idea is that to evaluate an author we consider not only the quality of the journals where his/her papers have been published, but also the quality of every single paper he/she authored. In addition, we take into account also the “quality” of the co-authors. In fact, an important author who writes a joint paper with a less important one, expresses a sort of trust-

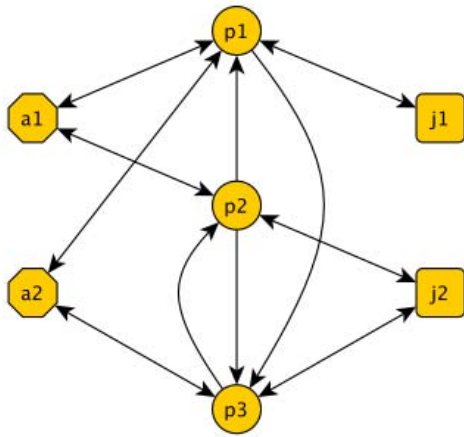


Figure 1: A graph where we have different nodes for each category. We have three papers, two authors and two journals.

ing vote by conferring to that author more visibility in the scientific community. Similarly, to evaluate the quality of a paper we consider the quality of the journal where the paper is published, the citations received, and at the reputation of its authors. Also, when evaluating a journal we take into account not only the cross-citations among journals — as done by many methods such as Impact Factor [8], Eigenfactor [3], and others [6, 11] — but also the quality of every single paper published there, and the authoritativeness of the authors who published on that journal.

In this paper we first review the basic model introduced in [4, 5] and then we modify the model introducing a freshness decay factor for citations. In particular, we take into account also the time of publication of the papers, and we account for the decrease of the importance of a paper over time. In our framework, papers, naturally start to lose importance right after publication, unless a fresh citation is received underlying a renewed interest in the paper. As a side effect we have that even recent papers that have not yet gathered enough citations have a chance to rank in higher position respect to old papers that received a greater number of citations in the past but have not been cited recently.

The paper is organized as follows. Section 2 describes the basic model already presented in [4, 5]. In Section 2.1 we show how a time-aware mechanism can be incorporated into the model. Section 3 contains the results of an extensive experimentation carried on on a synthetic dataset of one-million papers, half a million authors and 5,000 journals.

2. THE BASIC MODEL

Assume we are given n_P papers together with their bibliographic data. More precisely, of each paper we know the authors, the journal where the paper is published, and the list of citations contained in the paper. With this information we construct a graph with three different kinds of nodes (see Figure 1). We associate with this graph three matrices, one for each kind of nodes: the matrix F which

records which journal has published each paper, the matrix K which stores information about authorship, and the matrix H which records the citation structure among papers. In particular, let n_J be the total number of distinct journals where the n_P papers are published, and let n_A denote the number of distinct authors who authored the n_P papers. We define $F = (f_{i,j})$ as the $n_J \times n_P$ binary matrix such that

$$f_{i,j} = \begin{cases} 1 & \text{if paper } j \text{ is published in journal } i \\ 0 & \text{otherwise,} \end{cases}$$

$K = (k_{i,j})$ as the $n_A \times n_P$ binary matrix such that

$$k_{i,j} = \begin{cases} 1 & \text{if author } i \text{ has written paper } j \\ 0 & \text{otherwise,} \end{cases}$$

and $H = (h_{i,j})$ as the $n_P \times n_P$ matrix such that

$$h_{i,j} = \begin{cases} 1 & \text{if paper } i \text{ has paper } j \text{ in its reference list} \\ 0 & \text{otherwise.} \end{cases}$$

In the example of Figure 1 we have

$$F = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \end{bmatrix} \quad K = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \quad H = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}.$$

We can combine these three matrices to obtain the following 3×3 block matrix

$$A = \begin{bmatrix} FHF^T & FK^T & F \\ KF^T & KK^T & K \\ F^T & K^T & H \end{bmatrix} \quad (1)$$

of size $N = n_J + n_A + n_P$. For the example in Figure 1 we have

$$A = \begin{bmatrix} 0 & 1 & 1 & 1 & 1 & 0 & 0 \\ 1 & 2 & 1 & 1 & 0 & 1 & 1 \\ \hline 1 & 1 & 2 & 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 2 & 1 & 0 & 1 \\ \hline 1 & 0 & 1 & 1 & 0 & 0 & 1 \\ 0 & 1 & 1 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 \end{bmatrix}.$$

Each block of this matrix expresses the relationship between the subjects belonging to the three classes of *Journals*, *Authors* and *Papers*. More precisely, the entry (i, j) of the block FHF^T contains the number of citations that the papers published in journal i received from the papers published in journal j ; the entry (i, j) of the block FK^T contains the number of papers that author j has published in journal i ; the entry (i, j) of the block KK^T contains the number of papers co-authored by authors i and j .

We can scale the rows of A to obtain a row-stochastic matrix P , that is $Pe = e$, where $e = (1, \dots, 1)^T$. Then, we compute the ranking score of the subjects as the left eigenvector corresponding to the eigenvalue 1,

$$\pi^T = \pi^T P.$$

More precisely, numbering the subjects from 1 to N , the rank value (or importance) π_j of subject j is the weighted sum of the importances π_i of all the other subjects i which are in relation with j , where the weights are $p_{i,j}$, that is

$$\pi_j = \sum_{i=1}^N \pi_i p_{i,j}.$$

The row stochasticity of P implies that the overall amount of importance that a subject i transfers to the other subjects coincides with the importance of i . In other words, the amount of importance in the system is neither created nor destroyed.

To guarantee the existence and uniqueness of a solution we need A to be irreducible. Under this condition, it is always possible to find a scaling technique such that the matrix P can be constructed. The Perron Frobenius theorem [10] guarantees the existence of a unique vector π , such that $\pi_i > 0$ and $\sum_i \pi_i = 1$. We refer to π as the *Perron vector* of P . Moreover, in order to have nice convergence properties of iterative algorithms for the computation of π we need A to be aperiodic.

Note that working with the stochastic matrix P rather than computing the dominant eigenvector of A has advantages also from a numerical point of view. In fact, the approximation of the dominant eigenvector is done using an iterative procedure, and we do not need to perform a normalization at each step to limit the growth of the entries of the intermediate vectors.

We observe that P can be viewed as the transition matrix of a Markov chain so that π turns out to be the stationary distribution of the chain. The irreducibility of P makes the chain ergodic and ensures the existence and uniqueness of a stationary distribution.

2.1 Introducing time into the model

In the model just described, the amount of importance a paper confers to a cited paper does not depend on the time of publication. In most of the methods in citation analysis [8, 3] the dependence on the time is enforced just considering citations to papers only in a restricted time window. These models tend to favor papers that gather citations immediately after their publication, even if the declared intent was to make the rank current [8]. However, in many fields such as mathematics or economics, most of the citations occur up to ten years after publication [1], so that the rank computed will be based only on a small percentage of the citation activity, missing most of the citations.

Our idea is different because we introduce the concept that the value of the citations to papers change over the time. This means that papers that do not receive citations lose importance as time elapses. Conversely, old papers that are continuously cited over the years do not lose importance and as a side effect they confer authority also to the papers in their reference list. We can observe that papers with many citations in the past but no longer cited in the present time are penalized by this model, while recent papers highly cited have a chance to rank in higher positions even if their absolute citation count is lower.

Assume that for all the papers we know the year of publication, and let t_i be the time paper p_i was published. Assume that the papers are reordered in such a way $t_{i+1} \geq t_i$, $i = 1, 2, \dots, n_p$. Then in equation 1 we replace matrix H with the matrix $H_T = TH$, where T is the diagonal matrix with $T_{ii} = f(t_i)$, and $f: [t_0, t_{\max}] \rightarrow [0, 1]$ is a non increasing function. In Section 2.2 we will discuss the choice of the decay function $f(t)$.

To force irreducibility and then guarantee the existence of a unique stationary distribution, we introduce an additional node (see [4, 5, 7]) for each class of subjects. In particular,

we add a dummy paper which is cited by all the papers and cites back all the papers except itself. We also assume that the dummy paper is written by the dummy author and is published in the dummy journal. Mathematically, this corresponds to considering the matrices \widehat{H}_T , \widehat{K} and \widehat{F} obtained from H_T , K and F as follows,

$$\widehat{H}_T = \left[\begin{array}{c|c} H_T & \mathbf{e} \\ \hline \mathbf{e}^T & 0 \end{array} \right], \quad \widehat{K} = \left[\begin{array}{c|c} K & \mathbf{0} \\ \hline \mathbf{0}^T & 1 \end{array} \right], \quad \widehat{F} = \left[\begin{array}{c|c} F & \mathbf{0} \\ \hline \mathbf{0}^T & 1 \end{array} \right],$$

and replacing H , K and F in (1) with \widehat{H}_T , \widehat{K} and \widehat{F} , respectively. It is easy to prove that the Markov chain described by

$$\widehat{A} = \widehat{A}(t) = \left[\begin{array}{ccc} \widehat{F}\widehat{H}_T\widehat{F}^T & \widehat{F}\widehat{K}^T & \widehat{F} \\ \widehat{K}\widehat{F}^T & \widehat{K}\widehat{K}^T & \widehat{K} \\ \widehat{F}^T & \widehat{K}^T & \widehat{H}_T \end{array} \right]$$

is irreducible and aperiodic.

We already discussed the importance of scaling the rows of A to obtain a row-stochastic matrix. The simplest strategy is dividing each row of A by the sum of the entries in the row. A more flexible strategy, introduced in [4, 5], consists in performing a separate normalization of each block of A . That is, each block of A is normalized to yield nine row-stochastic matrices; then these matrices are compounded with weights $\Gamma = (\gamma_{i,j})_{i,j=1,3}$, where Γ is row stochastic, into a new stochastic matrix. The entries of Γ can be used to weight the amount of importance that each class (*Journal*, *Authors*, and *Papers*) transfers to the other classes. An in-depth discussion about the different possible normalization strategies of the single blocks is presented in [5] where a proposal for the normalization of each block is discussed. [7] discusses the role of the weight parameters Γ .

Denote by

$$Q = \left[\begin{array}{ccc} J_J & J_A & J_P \\ A_J & A_A & A_P \\ P_J & P_A & P_P \end{array} \right], \quad (2)$$

the matrix obtained from the corresponding block in matrix \widehat{A} . Each block is row-stochastic, and for example J_J is the stochastic matrix obtained by the normalization of $\widehat{F}\widehat{H}_T\widehat{F}^T$.

The notation used in (2) shows the role of each block with respect to the classes *Journals*, *Authors* and *Papers*. For instance, the entries of block J_A weight the amount of importance that *Journals* transfer to *Authors*.

Let $\Gamma = (\gamma_{i,j})$ be a 3×3 row-stochastic matrix, then the matrix

$$P = \left[\begin{array}{ccc} \gamma_{1,1} J_J & \gamma_{1,2} J_A & \gamma_{1,3} J_P \\ \gamma_{2,1} A_J & \gamma_{2,2} A_A & \gamma_{2,3} A_P \\ \gamma_{3,1} P_J & \gamma_{3,2} P_A & \gamma_{3,3} P_P \end{array} \right]. \quad (3)$$

is row-stochastic and its entries $p_{i,j} \geq 0$ express the amount of importance that subject i transfers to subject j . The parameters $\gamma_{i,j}$ can be used to tune the role that each class has with respect to the other classes. For instance, choosing $\gamma_{3,3}$ greater than $\gamma_{2,3}$ and $\gamma_{1,3}$ means that the importance of papers comes more from the citations they receive rather than from the importance of their authors or of the journals where they are published.

In [7] different choices of the weight matrix Γ are discussed in detail showing that the choice of a weighting criterion rather than another can change the behavior of our ranking algorithm. In the same paper a discussion of the probabilistic interpretation of the model and of the role of the dummy

players is given. From that discussion it turns out that a good choice for the weighting matrix is

$$\Gamma = \frac{1}{N} \begin{bmatrix} n_J & n_A & n_P \\ n_J & n_A & n_P \\ n_J & n_A & n_P \end{bmatrix},$$

where $N = n_J + n_A + n_P$ is the size of matrix P in (3). In fact, as discussed in [7], with this weighting strategy, the average value of a paper a journal or an author is the same. We used such a Γ in the experiments reported in Section 3.

2.2 The decay function

In the previous section we introduced a time-aware mechanism for weighting the citations to papers. The idea is that of multiplying the citation matrix H by a diagonal matrix T such that $T_{ii} = f(t_i)$ where t_i is the time paper p_i was issued. In this section we motivate the choice for the decay function f used in the experimental part.

The aging of citations follows an exponential decay rule, where

$$f(t) = \exp(-\alpha(t_c - t)),$$

where t_c is the current time and α is a constant obtained from the half-life decay time ρ , that is the time required to halve the value of a citation with the relation $\exp(-\alpha\rho) = 1/2$. This means that each citation contributes to the rank of cited papers fully only at the time of publication, and as time elapses the importance of citations decreases.

The exponential function is well suited for describing the decay of importance of citations released by a paper p , because the value at time t does not depend on the actual value of t but on the time elapsed since the publication of p . This characteristic makes it very easy to update the matrix H_T at the new current time $t_{c_1} = t_c + \delta$. In fact, if no papers are published in the time range $[t_c, t_c + \delta]$, we have that each citation has decreased its importance by a factor $\exp(-\alpha\delta)$. If H_{T_1} denotes the citation matrix at time t_{c_1} we have $H_{T_1} = \exp(-\alpha\delta)H_T$.

The difference between the use of the exponential decay function rather than a fixed time window is that our model is able to keep track of the past. In particular, we have differences in those situations where the importance of a paper is recognized only years after its publication. In methods such as Impact Factor or Eigenfactor the citations to those papers do not contribute to the ranking of the journal every time a paper is published outside the time window. Despite the delayed recognition of some papers might seem a marginal fact because we are in general interested in evaluating recent papers, this can produce differences in the rank of journals and authors.

3. EXPERIMENTAL RESULTS

[4, 5, 7] report the results of several tests on real and synthetic data for the model without time. In this section we present some results obtained scaling the citation matrix with the diagonal time matrix T as described in the previous sections. The experiments were performed on synthetic data since real datasets are either not publicly available and usable, or so incomplete that the characteristics of the bibliographic items do not correspond to those recognized in real

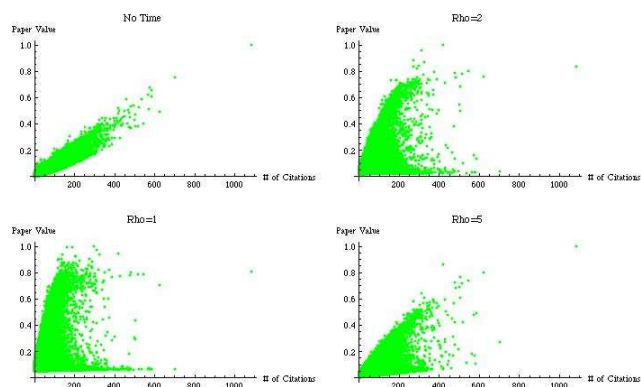


Figure 2: Dependence of the rank of a paper from the number of received citations. In the first plot without the decay in time of the importance of citations, the dependence of the rank is almost linear. The remaining three plots report the cases where $\rho = 1, 2, 5$. In these cases papers which receive a large number of citations do not have necessarily a large rank.

cases. In this respect, a generative model for building up synthetic matrices describing the subjects journals, authors and papers was proposed in [7]. The synthetic data produced agree with the properties observed on real datasets, allowing us to test the algorithm on a larger set of data, where we can evaluate the robustness of our ideas on special critical situations.

The purpose of this experimentation is to see how the rank of the subjects, either being papers, authors or journals is affected by the introduction of exponential decay in the importance of citations.

The generative algorithm presented in [7] was used to produce a dataset with one million of papers, half a million of authors and 5,000 journals, which respects the proportion of the cardinality of the classes in real databases [2].

Figure 2 represents the dependence of the rank of papers from the number of received citations. The first plot represents the situation where we have not introduced the time-decay factor. In that case the rank of a paper is strictly related to the number of citations received and we observe an almost linear dependence. When we introduce the decay in time of the importance of citations, we see that it is no more true that best papers are those receiving a greater number of citations because it is more important to receive fresh citations rather than to have received many of them in the past. The last three plots in Figure 2 represent the situation for $\rho = 1, 2, 5$, where ρ is the half-life parameter expressed in years. In particular $\rho = 1$ means that the importance of a citing paper halves each year, while $\rho = 5$ denotes that it is only after 5 years that a citation halves its importance. We see that as ρ increases, the dependence on the number of citations starts to appear. Repeating the same tests using as function $f(t)$ a step function which removes citations to papers older than two or five years, will produce plots similar to those obtained without any decay function.

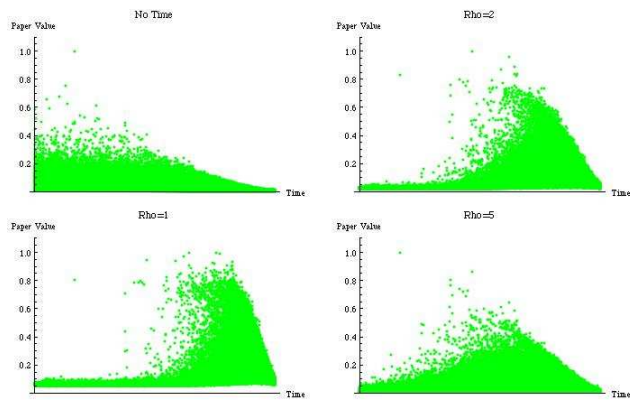


Figure 3: Rank of papers versus time of publications. We assume a total ordering of papers respect to the time of publication. The top-left plot, depicts the rank of papers with respect to their time of publication. We note that the papers published recently all have a small rank since they have not received many citations yet. The situation changes when we introduce a decay in time of the importance of citations. The extreme situation is when the importance of citations halves every year (corresponding to $\rho = 1$). Old papers lose importance while recent papers become more important. As the half-life decay factor increases citation count gains over freshness.

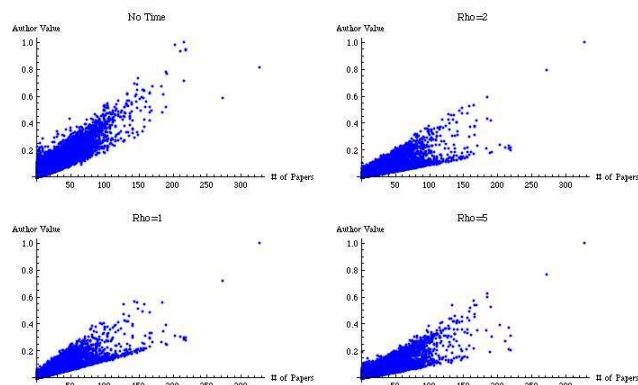


Figure 4: Dependence of the rank of authors on the number of papers written. We see that with the introduction of the decay in time of the importance of citations we have a more spread plot, even if more productive authors obtain greater rank values.

Figure 3 represents the rank of papers plotted versus their time of publication, publications ranging from 1980 to 2008. We see in the top-left plot, where the rank does not depend on time, that new papers have a small rank since they have received only few citations. This situation is corrected by the introduction of the time decay function which boosts the rank of recent papers. This mechanism is adequate in all the situations where we are called on to evaluate the productivity of young researchers or in general when we want to value the timing of research. In fact, the aging of citations allows new papers and young researchers to emerge from the maze of papers published.

The introduction of time into the model affects also the rank of authors (see Figure 4). For example, the dependence of the importance of an author on the number of papers written is modified. It is however still true that the more an author publishes the more important he/she becomes.

The introduction into the model of the decay in time leaves almost unchanged the plots for the category journals. Figure 5 depicts the dependence of the rank of journals on the mean number of citations received by each paper, but we have similar situations when we consider the dependence of the rank of journals on the mean value of papers or authors, or on the overall number of papers published on that journal. Of course, even if the cloud-like shape does not change, this does not mean that the rank of single journals is the same for the four models. To analyze how the rank of journals changes we should follow the evolution in time of the importance of top journals.

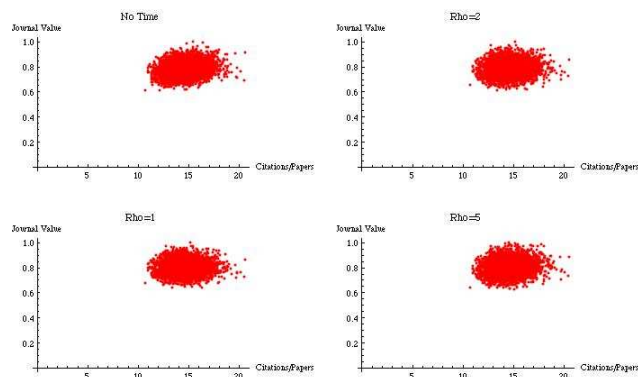


Figure 5: Dependence of the rank of a journal from the average number of citations received by each paper published on that journal. The situation does not change much with the introduction of the decay in time.

4. CONCLUSIONS

In this paper we have introduced a decay in time of the importance of citations and we have tested the new model on synthetic data simulating a collection of 1 Million papers, half a million authors and 5,000 journals. The modification introduced is effective in boosting recent papers with respect to old papers that received many citations in the past but no more citations in recent years. To the contrary, seminal papers published many years ago but still cited do not lose importance because of the fresh citations.

There remain interesting theoretical questions that need to be further investigated as done for the model where citations do not age. A question that we plan to tackle is to establish a relation between the rank of two chains one obtained from the other by a rank-one perturbation of matrix H_T .

5. ACKNOWLEDGMENTS

The authors would like to thank Dario A. Bini for the many useful discussions about the time-aware model.

We also thank the anonymous referees whose detailed comments helped improving the quality of the presentation.

6. REFERENCES

- [1] R. Adler, J. Ewing, and P. Taylor. Citation statistics. Technical report, IMU-ICIAM-IMS, June 2008. <http://www.mathunion.org/fileadmin/IMU/Report/CitationStatistics.pdf>.
- [2] AMS. MathSciNet, Mathematical Reviews on the Web. <http://www.ams.org/mathscinet/>.
- [3] C. T. Bergstrom. Eigenfactor: Measuring the value of prestige of scholarly journals. *C&RL News*, 68(5), 2007.
- [4] D. A. Bini, G. M. D. Corso, and F. Romani. Evaluating scientific products by means of citation-based models: a first analysis and validation. *Electron. Trans. Numer. Anal.*, 22:1–16, 2008.
- [5] D. A. Bini, G. M. D. Corso, and F. Romani. A combined approach for evaluating papers, authors and scientific journals. *Journal of Computational and Applied Mathematics*, 2009. to appear.
- [6] J. Bollen, M. A. Rodriguez, and H. V. de Sompel. Journal status. *Scientometrics*, 69(3):669–687, 2006.
- [7] G. M. D. Corso and F. Romani. Versatile weighting strategies for a citation-based research evaluation model. Technical Report TR-09-04, Dipartimento di Informatica, Università di Pisa, March 2009.
- [8] E. Garfield. Citation analysis as a tool in journal evaluation. *Science*, 178:471, 1972.
- [9] E. Garfield. The history and meaning of journal impact factor. *Journal of the American Medical Association*, 293:90–93, 2006.
- [10] L. Hogben, editor. *Handbook of Linear Algebra*. Chapman and Hall, 2007.
- [11] I. Palacios-Huerta and O. Volij. The measurement of intellectual influence. *Econometrica*, 72(3):963–977, 2004.
- [12] G. Pinski and F. Narin. Citation influence for journal aggregates of scientific publications- theory, with applications to literature of physics. *Information Processing & Management*, 12:297–312, 1976.
- [13] SCImago. Sjr- scimago journal & country rank, 2007. <http://www.scimagojr.com>.
- [14] M. Stringer, M. Sales-Pardo, and L. A. N. Amaral. Effectiveness of journal ranking schemes as a tool for locating information. *PLoS ONE*, 3(2):e1683, 2008. <http://dx.doi.org/10.1371/journal.pone.0001683>.