# f-value: Measuring an article's scientific impact

## E. Fragkiadaki[1], G. Evangelidis[1], N. Samaras[1], D. A. Dervos[2]

[1]Department of Applied Informatics, University of Macedonia Economic and Social Sciences, 54006 Thessaloniki, Greece

[2]Department of Information Technology, Alexander Technology Educational Institute (ATEI) of Thessaloniki, 57400 Sindos, Greece

**Abstract**

The f-value is a new method that measures the importance of a research article by taking into account all citations received, directly and indirectly, up to depth n. This method considers all information present in a Citation Graph in order to produce a ranking of the articles. Apart from the mathematical equation that calculates the f-value, we also present the corresponding algorithm with its implementation, plus an experimental comparison of f-value with some known indicators of an article's scientific importance, namely, the number of citations and the Page Rank for citation analysis. Finally, we discuss the similarities and differences among the indicators.

## 1   Introduction

The importance of Citation Analysis has become more obvious during the past few years. The vast increase of scientific production made it very difficult for scientists to keep track of publications they might be interested in. Many methods have been developed to rank scientific journals/conferences, authors and scientific publications by measuring their importance.

The most widely used ranking method for journals/conferences is the Impact Factor proposed by Eugene Garfield (**?** ). The ranking is based on the number of citations received by the articles published in the journal/conference in question.

In order to measure the importance of a researcher's work, other metrics have been proposed that use the collection of all articles a researcher has (co-) authored, plus the sum of all direct citations received. Such indexes are the h-Index (16), g-Index (9), and their variations.

For example, there have been variations of the h-index that take into account (a) the total number of citations included in the Hirsch-core (a-index, r-index) (17), (b) the age of the publications included in the Hirsch-core (ar-index) (17) (c) the age of the publications of an author (contemporary h-index) (21), (d) the age of the citations (trend h-index) (21), (e) the combination of the above two (age-decaying h-index) (18), (f) not only the citations inside the Hirsch-core but also the ones received by publications currently not included in the Hirsch-core (tapered h-index) (2).

There have been some variations of the g-index as well, like the gr-index and the grat-index (15). There has also been a proposal for an index that combines the h-index and the g-index called hg-index (1) that treats some of the disadvantages of both indices.

The importance of a scientific publication is most commonly measured based on the number of citations it has received. A different approach was proposed by Rousseau (20), which claims that publications mentioned in the Reference List have an impact on the publication in question, and also, recently, there has been a proposal for applying the philosophy of Page Rank (5) on a Citation Graph (19). Finally, the Cascading Citations Indexing Framework approach (6; 8; 7) suggests that citations should be addressed at the (article, author) level in order to rank the contribution of each author's scientific work.

We suggest a new method for measuring the importance of a research article, the f-value. We produce a ranking of the publications included in the CiteSeer bibliographic database(14) and compare our results with the ones obtained by competitive methods.

In Section 2, the Number of Citations, the Cascading Citations Indexing Framework, and the Page Rank for citation graphs approaches are presented. Section 3 describes the basic concept of the f-value and in Section 4, we justify the selection of the specific reducing factor used in the calculation of the f-value. The paper continuous by presenting the f-value algorithm in Section 5 and the different rankings produced by three different methods in Section 6. Section 7 describes the similarities and differences of the f-value with the competitive methods, and, finally, the last section concludes the paper.

## 2   Related Work

A citation graph is a representation of the relationships that exist between research articles based on the references that each article provides. In Figure 1, articles are shown as nodes of a directed graph. In this example there are 7 articles labeled A to G.

The edges of the graph represent references among articles. For example, the edge leaving node B can be interpreted as "article B references article D". The incoming edges are the direct citations received by a specific article. For article D we can state that "article D receives one direct citation from article B".

## 2.1   Number of Citations

This method produces a ranking of scientific publications based on the number of citations they receive. It is by far the most simplistic approach, but, it is widely used. For
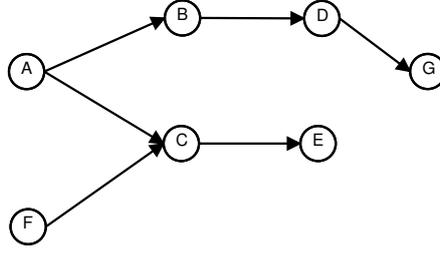
Fig. 1: Citation Graph A

example, in the citation graph of Figure 1, articles A and F receive 0 citations, articles B, D, E and G receive 1 citation each, and article C receives 2 citations.

## 2.2   The Cascading Citations Indexing Framework ($c^2$-IF)

The fundamental concept in the $c^2$-IF approach is the n-gen citation. According to $c^2$-IF, direct citations like the ones discussed in the previous section are called 1-gen citations. If we carefully examine the graph we observe that article D also receives an indirect citation from article A through article B. This is considered to be a 2-gen citation. In general, an n-gen citation exists between a source article S and a target article T, if there is a directed path in the citation graph from node S to node T. In the example of Figure 1, the greatest gen citation present is of rank 3, from article A to article G, through the citation path A → B → D → G.

According to $c^2$-IF, the citations that a (article, author) pair receives can be calculated up to depth n, thus, producing a number of distinct values. So, if we choose to consider the citations up-to depth 3, the following values will be calculated: 1-gen citations, 2-gen citations, and 3-gen citations. These values are stored in a table called Medal Standings Output (MSO).

Unlike the $c^2$-IF approach, self-citations are not excluded from the citation counts to be considered next. Equivalently, the present work considers the citations to be at the <article>, rather than at the <article, author>, level. We also stress that the $c^2$-IF approach is not to be considered as a ranking method but merely a framework that extends the citation indexing paradigm to include 2-, 3-, ..., k-gen citations.

## 2.3   Page Rank

The original Page Rank (5) produces a ranking of web pages by taking into account the number and importance of pages linking to each web page. The formula used by the Page Rank algorithm is

$$PR(A) = (1-d) + d * \sum_i \frac{PR(T_i)}{C(T_i)} \qquad (1)$$

where $PR(T_i)$ is the Page Rank value of page $T_i$ linking to page A whose Page Rank value we wish to calculate, and $C(T_i)$ is the number of outbound links of page $T_i$. Finally, $d$ is the damping factor. In order to better explain the damping factor, we should firstly give a general description of the concept of Page Rank.

The Page Rank algorithm is based on the Random Surfer model which states that a person, the "random surfer", navigates through the web randomly, by clicking on links

present on a web page. So, how high a web page ranks has to do with the probability that this "random surfer" eventually visits the web page in question. The probability increases as the number of incoming links increases and the effect is even more intense if these links come from web pages which score high, thus having themselves high probability to be visited. But, there is always a chance that our "random surfer" gets bored and chooses to simply leave, a reaction indicated by the damping factor, which on the original paper was chosen to be 0.85. In most discussions about Page Rank, 0.85 is the value used for the damping factor, but, there is at least one paper that we know of that examines the behavior of the original Page Rank algorithm when different values are chosen (4). So, for the most common value of the damping factor, Equation 1 actually becomes

$$PR(A) = 0.15 + 0.85 * \sum_i \frac{PR(T_i)}{C(T_i)} \qquad (2)$$

(19) suggest a variation of the original Page Rank algorithm applied to citation graphs. In that paper, the authors apply equation 1 by choosing 0.5 instead of 0.15 as the stopping probability. They choose the specific value based on an empirical study that states that researchers will probably not follow 6 articles and stop but only two.

## 3   f-value description

The Cascading Citations Indexing Framework defines the n-gen citations as a means of acknowledging the importance of a research article based not only on its direct influence (number of 1-gen citations) but also on the influence that articles citing that article have on their scientific field. In this paper, we introduce the f-value, a new way for quantifying the importance of a research article, which considersthe accumulated importance of all articles that have based their scientific contribution on the article in question, directly or indirectly.

Let us consider the following example. We have 6 articles, labeled A to F related as shown in Figure 2, thus producing the MSO table shown in Table 1.
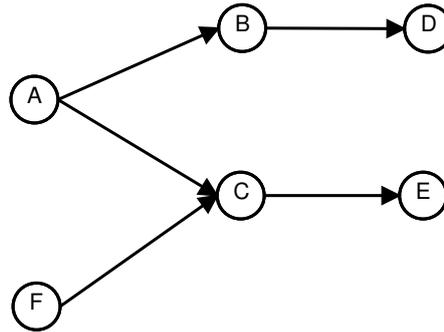


Fig. 2: Citation Graph B

A possible way to calculate the f-value of an article A by taking into account the indirect citations could be

$$f(A) = 1 + (f(A_1) + f(A_2) + ... + f(A_n)) \qquad (3)$$

| Article | 1-gen citations | 2-gen citations |
|:-------:|:---------------:|:---------------:|
| C | 2 | 0 |
| E | 1 | 2 |
| D | 1 | 1 |
| B | 1 | 0 |
| A | 0 | 0 |
| F | 0 | 0 |

Tab. 1: MSO Table for Citation Graph B

where $f(A)$ is the f-value of article A, and $A_i$, $i = 1..n$ are the articles citing article A. According to the equation, the minimum f-value for a published article is 1. Thus, the f-value of article A is 1 plus the sum of the f-values of all articles citing article A.

By performing the calculations for the citation graph of 2, we produce the graph shown in Figure 3, with the number on top of the nodes representing the f-values for the corresponding articles.
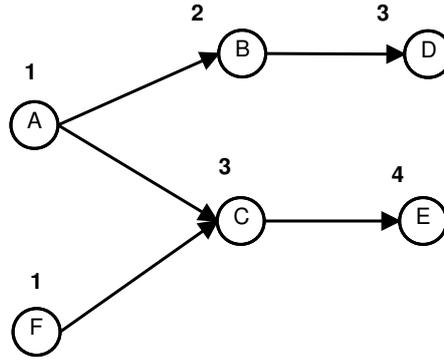


Fig. 3: f-values for Citation Graph B

Such an approach results to each article eventually receiving thus much credit as the sum of the credit received by all articles that cite it, making no distinction between direct or indirect citations. This is also obvious by examining the results shown in Figure 3. The f-value of each article is 1 plus the f-values of all direct citations. Of special interest are the f-values of articles C and D which are both 3. This means that based on Equation (3) these two articles are equally important even though article C has received 2 1-gen citations and article D has received one 1-gen citation and one 2-gen citation.

So there must be some factor that will assist us in differentiating direct and indirect citations. This is going to be a value that will reduce the cascaded f-value passed to an article's direct citations. We have concluded that, for the database used in this paper, this factor should be the number $\frac{1}{2.2}$ for reasons that we are going to explain in Section 4. Here is the new equation that calculates the f-value of an article:

$$f(A) = 1 + \frac{1}{2.2} * (f(A_1) + f(A_2) + ... + f(A_n)) \tag{4}$$

Applying this equation to the citation graph of Figure 2 we get the graph of Figure 4, in which the calculated f-values clearly distinguish among the different combinations
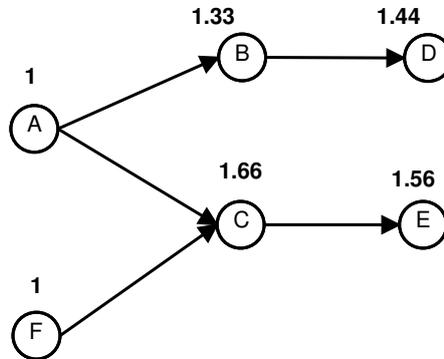
of n-gen citations received by the articles.

Fig. 4: f-values for Citation Graph B

## 4    Determining the reducing factor

For the calculations performed in the previous section we have used a reducing factor equal to $\frac{1}{2.2}$. We consider this value to be a good representative of the relation between 1-gen and 2-gen citations based on previous work that we have performed (10) with the specific database. More specifically in that paper we presented an algorithm that calculates and stores all the citations present in a bibliographic database, and produces the MSO table for all the distinct (article, author) pairs up-to depth 3. We tested the algorithm with data provided by the CiteSeer bibliographic database (14**?**). A relational database was used to store part of the original data from CiteSeer as well as the calculated values, the 1-gen, 2-gen, 3-gen citations (plus all individual citation paths) and the MSO table.

### 4.1    Data used

We chose the CiteSeer database because:

- It indexes a sufficient number of research articles and is not limited to certain journals

- It mostly covers the scientific area of Computer and Information Science

- it uses the Open Access Initiative (OAI) format, which is XML based.

A sample record is shown in Figure 5. For simplicity, only the identifiers that are used by the algorithm are listed.

Each article is defined by a unique *<identifier>* tag generated by CiteSeer, as shown in Figure 5. Other fields required by the algorithm are the title (*<dc:title>* tag) and the list of references included in each article (*<oai_citeseer:relation>* tag).

### 4.2    Preprocessing

The original data consisted of the entire CiteSeer database; a total of 72 files, each holding 10,000 articles with their corresponding bibliographic details. Articles appearing in the list of references of a particular article are also part of the CiteSeer database.

```
<record>
<header>
<identifier>oai:CiteSeerPSU:number#</identifier>
</header>
<metadata>
<dc:title>The Title</dc:title>
<oai_citeseer:pubyear>Publication Year</oai_citeseer:pubyear>
<oai_citeseer:relation type="References">
<oai_citeseer:uri>oai:CiteSeerPSU:number#</oai_citeseer:uri>
</oai_citeseer:relation>
<oai_citeseer:relation type="References">
<oai_citeseer:uri>oai:CiteSeerPSU:number#</oai_citeseer:uri>
</oai_citeseer:relation>
</oai_citeseer:oai_citeseer>
</metadata>
</record>
```

Fig. 5: CiteSeer Record

In order to retrieve the necessary information and to store it in the relational database we developed a parsing algorithm.

During the parsing process certain errors occurred, mainly concerning articles with insufficient information. For the algorithms presented here, articles lacking information about their authors (26,040 in total) or their publication year (280,098 in total) where excluded from the procedure.

## 4.3 $c^2$-IF algorithm Results and Statistical Analysis

The $c^2$-IF algorithm presented in (10) calculates the numbers of direct and indirect citations present in a Citation Graph, up to a pre-specified depth (in this case up to depth 3). Moreover, it stores in the relational database all the paths in the citation graph that produce these citations thus giving us complete knowledge of the graph. We note that the database stores information about 410,205 articles, with 265,563 identified authors and 1,245,171 direct references among the articles.

Before executing the algorithm we made sure that no cycles existed in the directed citation graph. A cycle indicates that there are articles inside our database that cite chronologically younger articles. This can occur either because the bibliographic information supplied for some articles is erroneous or because some citations refer to articles not yet published officially at the time when the citing articles are published. In order to avoid such anomalies that can lead to inaccurate results, we excluded from the procedure citations received from an article earlier than its publication year. Thus, the direct references among articles in the database were reduced from 1,245,171 to 1,000,077.

After the execution of the algorithm, 1,000,077 1-gen citations, 4,095,493 2-gen citations and 14,924,150 3-gen citations were detected among the articles and that many paths were stored in the database. An interesting fact is that from the 410,025 articles originally included in the database only 133,658 receive at least one citation. To gain a better understanding of our data we calculated the summary statistics for each n-gen (n=1, 2, 3) citation type ( see Table 2).

If we compare the mean to the median we observe that in all 3 cases the median is lower than the mean. This means that even though the means are high they are mostly affected by a small number of observations with high values. This hypothesis is proven true if we examine the quartile information. For example, for 1-gen citations we find

that at least 75% of the observations in our database are smaller than the corresponding mean value, whereas, the maximum value is 1,280 which is much larger than the usual values calculated for articles. Even greater are the differences for 2-gen citations and 3-gen citations.

|        | 1-gen | 2-gen  | 3-gen  |
|--------|-------|--------|--------|
| mean   | 7.48  | 30.64  | 111.7  |
| SD     | 18.98 | 139.36 | 774.38 |
| min    | 1     | 0      | 0      |
| 25%    | 1     | 0      | 0      |
| median | 3     | 2      | 0      |
| 75%    | 7     | 15     | 18     |
| max    | 1,280 | 12,186 | 82,182 |

Tab. 2: Summary Statistics for 1-gen, 2-gen and 3-gen citations

Finally we identified the ratios

$$\frac{\text{number of 2-gen citations}}{\text{number of 1-gen citations}} \tag{5}$$

and

$$\frac{\text{number of 3-gen citations}}{\text{number of 2-gen citations}} \tag{6}$$

for all articles in our database and we calculated the corresponding summary statistics shown in Table 3.

|        | 2-gen / 1-gen | 3-gen / 2-gen |
|--------|---------------|---------------|
| mean   | 2.22          | 1.54          |
| SD     | 4.92          | 2.48          |
| min    | 0             | 0             |
| 25%    | 0             | 0             |
| median | 1.00          | 0.91          |
| 75%    | 2.643         | 2.10          |
| max    | 454           | 227           |

Tab. 3: Summary Statistics for ratios 5 and 6

We observe that on average for each 1-gen citation an article receives from within our database, it also receives 2.22 2-gen citations and for each 2-gen citation it receives 1.54 3-gen citations. This is an expected result since according to the definition of n-gen citations, the (n+1)-gen citations an article receives is the sum of all 1-gen citations received by the n-gen citations of the article. For example the 2-gen citations received by an article are the sum of all 1-gen citations received by the articles directly citing the article in question (1-gen citations). We also mention that there are 44,280 articles for which we can not calculate ratio 6 because the number of 2-gen citations they have received so far is 0.

Based on these statistical data we chose to use $1/2.2$ as a reducing factor for the calculation of f-value. We expect this value to differ among scientific areas or bibliographic databases.

# 5 f-value algorithm

In this section we present an algorithm that calculates the f-values of all articles in our bibliographic database. This algorithm requires a finite number of iterations to calculate the f-values.

The algorithm receives as input the list of articles to be processed (***I***), the ***Article Direct Citations (ADC)*** data structure which includes for each article the list of articles that cite it, and, the ***Article F-Values (AFV)*** data structure which includes the articles that need to be processed plus their current f-value and a flag that denotes whether this value has changed since the last iteration. In other words, if we denote an article by $R_x$, then for a database with m articles, the list of all articles that need to be processed is ***I***=[$R_1$, $R_2$, $R_3$, ..., R]. Let $CR_x$ denote the list of articles that reference $R_x$. Thus, $CR_x$ is a subset of ***I*** and the Article Direct Citations (ADC) data structure is ***ADC***=[$CR_1$, $CR_2$ ,$CR_3$ , ... , $CR_m$]. Additionally, for each article $R_x$, let $VR_x$ denote the information required for this article during the execution of the algorithm. This information consists of the f-value calculated so far for this article and of a flag indicating whether the f-value changed since the last iteration of the algorithm. Thus, $VR_x$= [fval=1, changed=0] for every article $R_x$in the beginning of the algorithm. Finally, the Article F Values structure is ***AFV***=[$VR_1$, $VR_2$, ..., $VR_m$]. The algorithm returns the AFV structure with the calculated f-values for all articles in the database.

During the first iteration of the algorithm, all articles have an f-value equal to 1. At each iteration, the algorithm calculates the f-values of all articles in the database based on the f-values calculated during the previous iteration and records whether any f-value has changed between the two iterations. If there is at least one changed value the algorithm requires one more iteration because that change could propagate to more articles in the following iteration. If there is no f-value change then all f-values have been calculated and the the algorithm terminates.

[H]

```
1 Input:
2    I list of articles to be processed
3    ADC data structure with direct citations of each
4 article
5    AFV data structure with initial f-values and
6 flags
7 Output:
8    AFV data structure with calculated f-values and
9 flags
10
11 ADC = remove_cycles(ADC)
12 NChanged = 0
13 first = true
14 while (first || NChanged > 0) do
15   first = false
16   NChanged = 0
17   PREV_AFV = AFV
18   foreach R in I do
19    prev_fval = AFV[R][fval]
20    AFV[R][fval] = 1
21    RCIT = ADC[R]
22    for T in RCIT do
23     AFV[R][fval] = AFV[R][fval] +
24 /2.2*PREV_AFV[T][fval]
25    if AFV[R][fval] != prev_fval then
26     AFV[R][changed] = 1
27     NChanged = NChanged + 1
```

```
28    else
29      AFV[R][changed] = 0
```
    f-value algorithm

In order to avoid possible errors in the execution of the algorithm we must ensure that no cycles exist in the collection of articles stored in our database. Since the algorithm calculates the f-value of an article based on the f-values of the articles that cite it if there is a cycle the algorithm will enter an infinite loop.

## 6    Experimental Results

In order to compare the three different methods for measuring an article's scientific impact, we tested them against our database and report the obtained rankings per method. Recall that only 133,658 out of 410,025 articles listed in our database actually receive at least one 1-gen citation. In addition, there are 203,607 articles that do not give any citation, 38,100 of which receive citations from other articles while the rest do not give or receive any citations. Apart from presenting the rankings, the tables are complemented with the $c^2$-IF Information about the n-gen citations received by the articles up to depth 3. This information derives from the $c^2$-IF algorithm originally introduced at (10). The algorithm was modified for the needs of the present paper. Table 4, shows the top 10 articles according to the received number of citations.

| | | | | $c^2$-IF Information | | |
|---|---|---|---|---|---|---|
| Rank | Article Title | Pub. Year | Num. of Citations | 1-gen | 2-gen | 3-gen |
| 1 | Graph-Based Algorithms for Boolean Function Manipulation | 1986 | 1,280 | 1,280 | 7,057 | 31,724 |
| 2 | Optimization by Simulated Annealing | 1983 | 1,027 | 1,027 | 4,508 | 17,090 |
| 3 | Congestion Avoidance and Control | 1988 | 879 | 879 | 12,186 | 92,182 |
| 4 | A Method for Obtaining Digital Signatures and Public-Key Cryptosystems | 1978 | 867 | 867 | 7,678 | 42,807 |
| 5 | Statecharts: A Visual Formalism For Complex Systems | 1987 | 803 | 803 | 3,590 | 12,045 |
| 6 | Random Early Detection Gateways for Congestion Avoidance | 1993 | 762 | 762 | 6,244 | 31,185 |
| 7 | Fast Algorithms for Mining Association Rules | 1994 | 735 | 735 | 3,681 | 12,688 |
| 8 | Tcl and the Tk Toolkit | 1994 | 700 | 700 | 4,726 | 22,976 |
| 9 | Chord: A Scalable Peer-to-Peer Lookup Service for Internet Applications | 2001 | 610 | 610 | 1,672 | 1,351 |
| 10 | Mining Association Rules between Sets of Items in Large Databases | 1993 | 594 | 594 | 5,178 | 22,961 |

Tab. 4: Number of Citations: Top 10 ranked articles

In order to test the Page Rank algorithm for citation graphs against our bibliographic database, we used an implementation written by Vincent Krï¿œeutler in Python, which is based on an essay by David Austin published at the American Mathematical Society portal (3). The implementation of the Page Rank algorithm as a package was imported to a Python script created for handling the reading/writing from/to the database and transforming the data into an acceptable by the package format. The results are shown in Table 5.

| Rank | Article Title | Pub. Year | PR value | $c^2$-IF Information | | |
|---|---|---|---|---|---|---|
| | | | | 1-gen | 2-gen | 3-gen |
| 1 | Optimization by Simulated Annealing | 1983 | $686,054 * 10^{-9}$ | 1,027 | 4,508 | 17,090 |
| 2 | Graph-Based Algorithms for Boolean Function Manipulation | 1986 | $662,149 * 10^{-9}$ | 1,280 | 7,057 | 31,724 |
| 3 | New Directions in Cryptography | 1976 | $576,792 * 10^{-9}$ | 422 | 5,224 | 34,203 |
| 4 | A Method for Obtaining Digital Signatures and Public-Key Cryptosystems | 1978 | $526,387 * 10^{-9}$ | 867 | 7,678 | 42,807 |
| 5 | Congestion Avoidance and Control | 1988 | $461,410 * 10^{-9}$ | 879 | 12,186 | 92,182 |
| 6 | Applications Of Circumscription To Formalizing Common Sense Knowledge | 1986 | $323,209 * 10^{-9}$ | 226 | 2,611 | 16,881 |
| 7 | Tcl and the Tk Toolkit | 1994 | $315,861 * 10^{-9}$ | 700 | 4,726 | 22,976 |
| 8 | Implementing Mathematics with The Nuprl Proof Development System | 1986 | $309,963 * 10^{-9}$ | 398 | 3,858 | 17,598 |
| 9 | Statecharts: A Visual Formalism For Complex Systems | 1987 | $309,718 * 10^{-9}$ | 803 | 3,590 | 12,045 |
| 10 | A Timeout-Based Congestion Control Scheme for Window Flow-Controlled Networks | 1986 | $304,607 * 10^{-9}$ | 26 | 1,486 | 21,789 |

Tab. 5: Page Rank: Top 10 ranked articles

Algorithm 5 was implemented and executed against our database. Table 6 shows information about the top 10 ranked articles.

| Rank | Article Title | Pub. Year | f-value | $c^2$-IF Information | | |
|---|---|---|---|---|---|---|
| | | | | 1-gen | 2-gen | 3-gen |
| 1 | Congestion Avoidance and Control | 1988 | 258,534 | 879 | 12,186 | 92,182 |
| 2 | Design and Implementation of the Sun Network Filesystem | 1985 | 234,037 | 296 | 4,299 | 39,239 |
| 3 | The UNIX Time-Sharing System | 1974 | 224,167 | 127 | 1,405 | 14,236 |
| 4 | A Scheme for Real-Time Channel Establishment in Wide-Area Networks | 1990 | 192,736 | 421 | 5,172 | 46,302 |
| 5 | A Timeout-Based Congestion Control Scheme for Window Flow-Controlled Networks | 1986 | 181,751 | 26 | 1,486 | 21,789 |
| 6 | A Fast File System for UNIX | 1984 | 148,843 | 83 | 1,610 | 13,429 |
| 7 | New Directions in Cryptography | 1976 | 138,137 | 422 | 5,224 | 34,203 |
| 8 | An Open Operating System for a Single-User Machine | 1979 | 114,979 | 12 | 878 | 9,894 |
| 9 | Using Sparse Capabilities in a Distributed Operating System | 1986 | 109,455 | 51 | 523 | 5,418 |
| 10 | Why Aren't Operating Systems Getting Faster As Fast As Hardware? | 1989 | 103,480 | 149 | 2,451 | 19,929 |

Tab. 6: f-value: Top 10 ranked articles

Finally, table 7 shows the summary statistics for all three methods.

| | Number of Citations | Page Rank | f-value |
|---|---|---|---|
| mean | 7.48 | $2,451 * 10^{-9}$ | 43.06 |
| SD | 18.98 | $4,258 * 10^{-9}$ | 1,221 |
| min | 1 | $1,788 * 10^{-9}$ | 1 |
| 25% | 1 | $1,788 * 10^{-9}$ | 1 |
| median | 3 | $1,788 * 10^{-9}$ | 1 |
| 75% | 7 | $2,011 * 10^{-9}$ | 1.66 |
| max | 1,280 | $686,954 * 10^{-9}$ | 258,534 |

Tab. 7: Summary Statistics

## 7   Discussion

In this section, we present the similarities and differences of the three methods. In addition, we attempt to interpret the experimental results we obtained.

The Number of Citations, a measure used traditionally in citation analysis, plays an important role in all methods. In the Page Rank method, the direct citations a publication receives are referred to as inbound links to its node in the citation graph and they are similarly used in the f-value method.

In general, the latter two methods are based on the assumption that the use of the Number of Citations as a measurement of the importance of a scientific publication is insufficient. The resulting ranking is solely based on the direct impact the article has without taking into account its present state (whether it remains in the researchers' preferences) or its derived contribution (the impact it has on the research in the specific scientific field). The f-value method and the Page Rank method appear to be very similar in nature, thus, before elaborating on their experimental results, we discuss their main differences and similarities. These are summarized in the following:

1. *The logic behind the equation:* Page Rank focuses on a person (the "random scientist") moving from article to article randomly by choosing to read next an article that appears as a citation in the List of References of the article she reads. All cited articleshave the same probability to be selected. The f-value is not based on such a probability, but on the cumulative value of the n-gen citations that an article has received.

2. *How are citations treated:* Page Rank for Citation graphs divides equally the value of an article among its cited articles. Such a division implies that among two articles with equal values, A and B, if A cites 10 articles and B cites 20 articles, then articles cited by A will receive twice as much recognition than articles cited by B, just because A has cited fewer articles. Since we cannot assume that cited articles have less impact when they are encountered in longer reference lists, we claim that this division of value does not correspond to a real world behavior, thus, it is not included in the calculations of an article's f-value.

3. *The damping factor:* In the f-value calculation there is no damping factor. Instead, there is a reducing factor used to decrease the accumulated value of the n-gen citations. This factor has been chosen to be $\frac{1}{2.2}$ (see Section 4). In addition, the f-value also has a minimum value of 1 for all articles. The f-value of an articlealways increases as more articles cite directly and/or indirectly the article in question.

Even though the equations used in the calculation of the Page Rank for Citation Analysis and the f-value appear similar, the logic behind each implementation is different.

We now proceed with a discussion of the experimental results in an effort to better understand the differences and similarities among the three methods. Examining the top 10 ranked articles based on the Number of Citations (Table 4), it is very interesting to notice the $c^2$-IF information provided, especially for the top four ranked articles. We observe that according to this method, the "Congestion Avoidance and Control" article is ranked 3rd, because it has received fewer direct citations than the two articles above it. On the other hand, if we examine the $c^2$-IF information, we can clearly see that it has received considerably more 2-gen citations and 3-gen citations than the first and second ranked articles. The same is true to a lesser extent for the fourth ranked article. But, this information is not taken under consideration for this ranking.

Table 5, shows the top 10 articles based on the PageRank method along with the corresponding $c^2$-IF Information. The ranking is different here, and, by inspecting the $c^2$-IF information of the the top two articles, we observe that the first ranked article has less 1-gen, 2-gen and even 3-gen citations than the second ranked article. This ordering can only be explained if we consider the way Page Rank values are calculated. Apparently, the "Optimization by Simulated Annealing" article has received less 1-gen, 2-gen and 3-gen citations than the second article as an absolute number, but, the prestige (Page Rank value) of the articles that cite it played an important role in the calculations. In addition, the number of citations provided by the citing articles also affected the result. So, we have to assume that although the up to 3-gen citations of the first article are fewer than the ones received by the "Graph-Based Algorithms for Boolean Function Manipulation" article, they are either of higher value and/or have a smaller number of outbound links.

The f-value results are presented in Table 6 along with the corresponding $c^2$-IF information. Let us examine the first ranked article. This article was ranked third according to the Number of Citations. This is explained by the fact that the calculation of the f-value is exchaustive in nature and takes into consideration all the knowledge present in the citation graph. In other words, an article's f-value increases as it receives more citations at each depth, all the way to the longestcitation path.

Finally, Table 8 shows all articles listed in tables 4, 5 and 6 along with their $c^2$-IF information. The articles are ordered by their f-value rank. Again, we observe that the rankings vary significantly depending on the method used.

The first method, Number of Citations, only takes into account the direct impact an articles has based on the number of citations it receives. On the other hand, the Page Rank method does not take into account the direct impact alone but it also considers, to some extent, the added value provided by the citing articles of the article in question. We should point out though that the Page Rank method is not an exchaustive method, that is, for the calculation of the importance of a research article one does not traverse the entire citation graph. Finally, in the f-value method the indirect impact an article has is fully accumulated in the calculations. The whole citation graph is traversed and the value of each article is partially propagated to all articles that it cites, thus producing an exchaustive method that uses all the information present in the citation graph.

The f-value method is based on historical data, that is, it is dependent on the dataset. It is very likely that the reducing factor will be different for different datasets. A different reducing factor is expected to alter the resulting ranking, but the extend at which the ranking is affected requires more research.

| Article Title | Pub. Year | Ranks | | | c²-IF Information | | |
|---|---|---|---|---|---|---|---|
| | | f-value | Number of Citations | Page Rank | 1-gen | 2-gen | 3-gen |
| Congestion Avoidance and Control | 1988 | 1 | 3 | 5 | 879 | 12,186 | 92,182 |
| Design and Implementation of the Sun Network Filesystem | 1985 | 2 | 75 | 20 | 296 | 4,299 | 39,239 |
| The UNIX Time-Sharing System | 1974 | 3 | 498 | 39 | 127 | 1,405 | 14,236 |
| A Scheme for Real-Time Channel Establishment in Wide-Area Networks | 1990 | 4 | 26 | 11 | 421 | 5,172 | 46,302 |
| A Timeout-Based Congestion Control Scheme for Window Flow-Controlled Networks | 1986 | 5 | 7,365 | 10 | 26 | 1,486 | 21,789 |
| A Fast File System for UNIX | 1984 | 6 | 1,126 | 139 | 83 | 1,610 | 13,429 |
| New Directions in Cryptography | 1976 | 7 | 23 | 3 | 422 | 5,224 | 34,203 |
| An Open Operating System for a Single-User Machine | 1979 | 8 | 18,272 | 268 | 12 | 878 | 9,894 |
| Using Sparse Capabilities in a Distributed Operating System | 1986 | 9 | 2,608 | 323 | 51 | 523 | 5,418 |
| Why Aren't Operating Systems Getting Faster As Fast As Hardware? | 1989 | 10 | 365 | 182 | 149 | 2,451 | 19,929 |
| A Method for Obtaining Digital Signatures and Public-Key Cryptosystems | 1978 | 19 | 4 | 4 | 867 | 7,678 | 42,807 |
| Applications Of Circumscription To Formalizing Common Sense Knowledge | 1986 | 71 | 143 | 6 | 226 | 2,611 | 16,881 |
| Graph-Based Algorithms for Boolean Function Manipulation | 1986 | 76 | 1 | 2 | 1,280 | 7,057 | 31,724 |
| Random Early Detection Gateways for Congestion Avoidance | 1993 | 129 | 6 | 34 | 762 | 6,244 | 31,185 |
| Tcl and the Tk Toolkit | 1994 | 131 | 8 | 7 | 700 | 4,726 | 22,976 |
| Implementing Mathematics with The Nuprl Proof Development System | 1986 | 156 | 35 | 8 | 398 | 3,858 | 17,598 |
| Optimization by Simulated Annealing | 1983 | 168 | 2 | 1 | 1,027 | 4,508 | 17,090 |
| Mining Association Rules between Sets of Items in Large Databases | 1993 | 249 | 10 | 23 | 594 | 5,178 | 22,961 |
| Statecharts: A Visual Formalism For Complex Systems | 1987 | 326 | 5 | 9 | 803 | 3,590 | 12,045 |
| Fast Algorithms for Mining Association Rules | 1994 | 531 | 7 | 25 | 735 | 3,681 | 12,688 |
| Chord: A Scalable Peer-to-Peer Lookup Service for Internet Applications | 2001 | 2,621 | 9 | 150 | 610 | 1,672 | 1,351 |

Tab. 8: Summarized results of Top article rankings based on all three approaches

# 8   Conclusions

Based on the Cascading Citations Indexing Framework, we proposed a new method for measuring the importance of a research article. The f-value method calculates a unique value for each article that takes into consideration the n-gen citations received by the specific article. We developed an algorithm that calculates the f-value for all articles in a bibliographic database, and we experimentaly compared our method to two other popular methods.

Future work on this field will: (a) try to incorporate other aspects of the $c^2$-IF to the calculation of f-value, like the self-citations and the chords, (b) examine the impact the different values of the reducing factor have on the final ranking of the articles, and, (c) examine whether there can be a unified f-value for interdisciplinary articles.

# References

[1] Alonso, S., Cabrerizo, F. J., Herrera-Viedma, E., and Herrera, F.: hg-index: a new index to characterize the scientific output of researchers based on the h- and g-indices, Scientometrics, 2009

[2] Anderson, Thomas R., Hankin, Robin K. S., and Killworth, Peter D.: Beyond the Durfee square $Enhancing the h-index to score total publication output, Scientometrics 76(3), volume 76, 577588, 2008$

[3] Austin, David: How Google Finds Your Needle in the Web's Haystack , 2006

[4] Boldi, Paolo, Santini, Massimo, and Vigna, Sebastiano: PageRank: Functional dependencies, ACM Transactions on Information Systems 27(4), volume 27, ACM, 1â23, 2009

[5] Brin, S., and Page, L.: The Anatomy of a Large-Scale Hypertextual Web Search Engine, Computer Networks and ISDN Systems, 107â117, 1998

[6] Dervos, D., and Kalkanis, T.: cc-IFF: A Cascading Citations Impact Factor Framework for the Automatic Ranking of Research Publications, 3rd IEEE International Workshop on Intelligent Data Acquisition and Advanced Computer Systems: Technology and Applications (IDAACS 2005), Sofia, Bulgaria, September 2005

[7] Dervos, D., and Klimis, L.: Exploiting Cascading Citations for Retrieval, Proc. of the ASSIST 2008 Annual Meeting, October 2008

[8] Dervos, D., Samaras, N., Evangelidis, G., and Folias, T.: A New Framework for the Citation Indexing Paradigm, Proc. of the ASSIST 2006 Annual Meeting, Austin, Texas, USA, November 2006

[9] Egghe, L.: Theory and practise of the g-index, Scientometrics 69(1), volume 69, 131â152, 2006

[10] Fragkiadaki, E., Evangelidis, G., Samaras, N., and Dervos, D.: Cascading Citations Indexing Framework Algorithm Implementation and Testing, Informatics, Panhellenic Conference on 0, volume 0, IEEE Computer Society, 70â74, 2009

[11] Garfield, E.: Citation indexes for science. A new dimension in documentation through association of ideas, Science 122, volume 122, 1123â1127, 1955

[12] Garfield, E.: Journal impact factor: a brief review, CMAJ 161(8), volume 161, 979â980, October 1999

[13] Garfield, E.: The Agony and the Ecstasy - The History and Meaning of the Journal Impact Factor., sep 2005

[14] Giles, C. Lee, Bollacker, Kurt D., and Lawrence, Steve: CiteSeer: An Automatic Citation Indexing System, , ACM Press, 89â98, 1998

[15] Guns, R., and Rousseau, R.: Real and rational variants of the h-index and the g-index, Journal of Informetrics 3(1), volume 3, 64â71, January 2009

[16] Hirsch, J.E.: An index to quantify an individual's scientific research output, Proceedings of the National Academy of Sciences, volume 102, National Acad Sciences, 16569â16572, 2005

[17] Jin, BiHui, Liang, LiMing, Rousseau, Ronald, and Egghe, Leo: The R - and AR -indices: Complementing the h -index, Chinese Science Bulletin 52(6), volume 52, 855â863, 03 2007

[18] Katsaros, D., Sidiropoulos, A., and Manopoulos, Y.: Age Decaying H-Index for Social Network of Citations, SAW Proceedings of the BIS 2007 Workshop on Social Aspects of the Web, Poznan, Poland, April 27, 2007, volume 245, CEUR-WS.org, 2007

[19] Ma, Nan, Guan, Jiancheng, and Zhao, Yi: Bringing PageRank to the citation analysis, Information Processing and Management 44(2), volume 44, Pergamon Press, Inc., 800â810, 2008

[20] Rousseau, R.: The Gozinto Theorem: Using Citations to Determine Influences on a Scientific Publication, Scientometrics 11(3-4), volume 11, 217â229, 1987

[21] Sidiropoulos, A., Katsaros, D., and Manolopoulos, Y.: Generalized Hirsch h-index for disclosing latent facts in citation networks, Scientometrics 72(2), volume 72, 253â280, 2007