

On Hirsch's h , Egghe's g and Kosmulski's $h(2)$

QUENTIN L. BURRELL

Isle of Man International Business School, Douglas, Isle of Man

In recent issues of the ISSI Newsletter, EGGHE [2006A] proposed the g -index and KOSMULSKI [2006] the $h(2)$ -index, both claimed to be improvements on the original h -index proposed by HIRSCH [2005]. The aim of this paper is to investigate the inter-relationships between these measures and also their time dependence using the stochastic publication/citation model proposed by BURRELL [1992, 2007A]. We also make some tentative suggestions regarding the relative merits of these three proposed measures.

Introduction

The proposal by HIRSCH [2005] to introduce a single index to quantify a scientist's published research impact created an unprecedented response from the scientometric community. Within a year we saw responses from, among others, BANKS [2006], BORNMANN & DANIEL [2005], BRAUN & AL. [2005], BURRELL [2006], CRONIN & MEHO [2006], BATISTA & AL. [2006], EGGHE [2006 C], EGGHE & ROUSSEAU [2006], GLÄNZEL [2006A, 2006B], LIANG [2006], POPOV [2005], ROUSSEAU [2006A, 2006B] and VAN RAAN [2006]. Some of these have sought to demonstrate empirical applications of the index, some to extend its applicability and others to provide mathematical models, most notably EGGHE & ROUSSEAU [2006], GLÄNZEL [2006A] and BURRELL [2006, 2007A].

Others have proposed alternative measures, similar to or based upon the h -index. The "size" of the h -core [ROUSSEAU, 2006B] and the A -index [JIN, 2006; ROUSSEAU, 2006B] have been analysed by BURRELL [2007B] using a stochastic model for the publication/citation process originally proposed by BURRELL [1992, 2007A]. The idea that it might be more appropriate to use the h -rate, rather than the h -index, has been argued by BURRELL [2007C] based upon the work of LIANG [2006]. In this paper we use this same stochastic model to investigate aspects of the HIRSCH [2005] h -index, the EGGHE [2006A, 2006B] g -index and the KOSMULSKI [2006] $h(2)$ -index.

Received December 5, 2007

Address for correspondence:

QUENTIN L. BURRELL

The Nunnery, Old Castletown Road, Douglas, Isle of Man IM2 1QB, via United Kingdom

E-mail: q.burrell@ibs.ac.im

0138–9130/US \$ 20.00

Copyright © 2008 Akadémiai Kiadó, Budapest

All rights reserved

The various indexes

According to the preprint of HIRSCH [2005], the h-index for an author is that integer h such that h of his/her papers have at least h citations each, while the rest have fewer than h citations. Actually this is not quite well-defined, see the print version of HIRSCH [2005], GLÄNZEL [2006A] and ROUSSEAU [2006B], since there is ambiguity if there are several papers with the same number of citations at h . To get round this, let us introduce the following.

Notation. Write $f(n; T)$ for the number of an author's papers receiving exactly n citations by time T , and $N(n; T)$ for the number of an author's papers that have received at least n citations by time T so that

$$N(n; T) = \sum_{j=n}^{\infty} f(j; T)$$

and note that $N(n; T)$ decreases with increasing n .

Remarks. Here and later we include the time parameter T explicitly since one of our aims is to consider how the indexes develop in time as dynamic processes.

Definition 1. Hirsch's h-index at time T is, for any particular author, the integer $h(T)$ satisfying

$$h(T) = \max \{n : n \leq N(n, T)\} \quad (1)$$

For instance, if this maximal $n = 25$, say, then at least 25 of the author's papers have received 25 or more citations each while fewer than 26 have received 26 or more citations each. Note that this is an empirical measure, requiring observation of the actual values of $N(n; T)$.

Kosmulski's $h(2)$ index also concentrates on the number of most cited sources, but requires a much higher level of citation attraction.

Definition 2. Kosmulski's $h(2)$ -index at time T is, for any particular author, the integer $h_2(T)$ satisfying

$$h_2(T) = \max \{n : n \leq N(n^2, T)\} \quad (2)$$

For instance, if the maximal n is $n = 5$ then at least 5 papers have attracted 25 citations each, or more, while fewer than 6 have attracted at least 36 citations each. Clearly $h_2(T) \leq h(T)$ and we shall see that in most circumstances $h_2(T)$ is considerably less than $h(T)$.

Note. To incorporate the time dependence explicitly, we have modified Kosmulski's notation slightly, writing h_2 rather than $h(2)$ since this latter could be confused with our notation for the h-index at time $T = 2$.

Egghe's g -index is rather different from both h and h_2 in that it switches attention from the number of most productive papers to the actual number of citations attracted by these most productive papers. In this sense, it is related to the size of the Hirsch core,

see ROUSSEAU [2006B] and BURRELL [2007B], and Jin's A-index, see JIN [2006], ROUSSEAU [2006B] and BURRELL [2007B].

In our notation, the total number of citations received by those publications receiving n citations each is $nf(n;T)$ and the total number received by all those receiving at least n citations each is then given by

$$C(n;T) = \sum_{j \geq n} jf(j;T)$$

Let us refer to this as *the size of the n -core* and note that $C(n;T)$ decreases with increasing n . In particular, $C(h(T);T)$ gives the size of the Hirsch core at time T , denoted by $C(T)$ in BURRELL [2007B].

Another novelty in Egghe's approach is that a paper's rank, as determined by its number of received citations, is included explicitly. If we write $r = r(n)$ for the rank of a paper receiving n citations then in this notation Egghe's g -index can be defined as follows.

Definition 3. Egghe's g -index at time T is, for any particular author, the integer $g(T)$ satisfying

$$g(T) = \max \{r : r(n)^2 \leq C(n;T)\} \quad (3)$$

Thus if $g(T) = 25$ then the total number of citations received by those papers on rank ≤ 25 is at least $25^2 = 625$ while for those on rank ≤ 26 the total citation count is less than $26^2 = 676$.

Actually, it will be convenient for us to modify this slightly. In EGGHE [2006A, 2006B] an author's papers are ranked according to the number of citations received with the convention that when there are several papers with the same number of citations they are ranked serially, with no indication as to how they are to be ordered. (This last point is of no concern in determining the author's index, only when we wish to decide which papers are included and which excluded from any sort of core.) Our approach is slightly different in that all papers receiving the same number of citations are given the same rank. This conforms with standard practice in probability and economics, although it is not crucial in what follows, which is an illustrative analysis.

For any $n = 0, 1, 2, \dots$ we define the rank of n as the number of papers receiving at least n citations, and note that this is defined whether or not there are any papers receiving (exactly) n citations. In other words, Egghe's $r = r(n)$ is almost equivalent to what we have denoted by $N(n;T)$. We thus slightly modify the earlier definition to:

Definition 3*. Egghe's g -index at time T is, for any particular author, the integer $g(T)$ satisfying

$$g(T) = \max \{N(n;T) : N(n;T)^2 \leq C(n;T)\} \quad (4)$$

Remarks. (i) It is clear from the respective definitions that $h_2(T) \leq h(T) \leq g(T)$.

(ii) Note that when we talk of time T , this refers to the time that has elapsed since the start of a particular author's publication career so that it is important to realise that when we are considering several authors, "now" or "the current time" may correspond to different values of T .

(iii) Note that $f(0;T)$ gives the size of the zero class, the number of an author's papers that have received no citations by time T ; $N(0;T)$ gives the total number of publications by time T ; $C(0;T)$ gives the author's total number of received citations by time T .

The stochastic model

Here we just recap the essentials of the model and refer the reader to BURRELL [2007A] for full details. The basic idea is that an author publishes papers at certain times and that these papers subsequently attract citations following their publication, where both the publication and citation accumulation processes are random. We further assume that some papers are more citable than others so that the citation rate varies between different publications. The basic ideas behind this were originally put forward in BURRELL [1992]. The precise technical assumptions, without the mathematical details, are:

Assumptions

1. From the start of his/her publishing career (at time zero), an author publishes papers according to a Poisson process of rate θ which gives the mean number of publications per unit time, called the *publication rate*.

2. Any particular publication acquires citations according to a Poisson process of rate Λ , where Λ varies from paper to paper. Here Λ denotes the mean number of citations to the paper per unit time following publication, called the *citation rate*.

3. The citation rate Λ for this author varies over the set of his/her publications according to a gamma distribution of index $\nu \geq 1$ and scale parameter $\alpha > 0$.

See BURRELL [2007A] for the precise details.

In all of the following analysis, the basic result concerns the distribution of the number of citations garnered (up to the current time) for a typical paper by this author, whenever it was published during his/her career. This is given in the following:

Theorem [BURRELL, 2007A]

Under the assumptions of the model, the distribution of X_T , the number of citations to a randomly chosen paper by time T , is given by

$$P(X_T = r) = \frac{\alpha}{(\nu-1)T} B\left(\frac{T}{\alpha+T}; r+1, \nu-1\right) \text{ for } r = 0, 1, 2, \dots \quad (5)$$

where $B(x; a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_0^x y^{a-1} (1-y)^{b-1} dy$

is the cumulative distribution function of a beta distribution (of the first kind) with parameters a and b .

To apply the theoretical model, we need to re-interpret the empirical quantities $f(n;T)$, $N(n;T)$ and $C(n;T)$ introduced in the previous section as the expected values of the theoretical quantities. In terms of the stochastic model we have:

Proposition.

$$(i) f(n;T) = \frac{\alpha\theta}{(v-1)} B\left(\frac{T}{\alpha+T}; n+1, v-1\right) \text{ for } n = 0, 1, 2, \dots \quad (6)$$

$$(ii) N(n;T) = \theta T \left(1 - \frac{\alpha\theta}{(v-1)T} \sum_{j=0}^{n-1} B\left(\frac{T}{\alpha+T}; j+1, v-1\right) \right) \text{ for } n = 0, 1, 2, \dots \quad (7)$$

$$(iii) C(n;T) = \theta T \left(\frac{vT}{2\alpha} - \frac{\alpha}{(v-1)T} \sum_{j=0}^{n-1} n B\left(\frac{T}{\alpha+T}; j+1, v-1\right) \right) \text{ for } n = 0, 1, 2, \dots \quad (8)$$

Proof. See the Appendix.

Remark. Given that our model provides explicit formulae, it is worth noting that, according to the model,

$$\begin{aligned} N(0;T) &= \text{Expected total number of papers published by time } T \\ &= \theta T \end{aligned} \quad (9)$$

and

$$\begin{aligned} C(0;T) &= \text{Expected total number of citations received by time } T \\ &= \frac{1}{2}\theta T^2(v/\alpha) = \frac{1}{2}\theta\mu T^2 \end{aligned} \quad (10)$$

(For this last result, see the proof of the Proposition in BURRELL [2007B].)

From the above it follows that:

Proposition 2.

Given the model assumptions, $g(T)$ is defined if and only if $\theta > v/2\alpha = \mu/2$

This is an interesting result. It says that an author needs a publication rate at least half as big as his/her mean citation rate in order to have a defined g -index. For instance, if an author is publishing 2 papers a year, but these papers are attracting on average 5 citations a year, then the g -index is not defined, at least according to the theoretical model. That the (empirical) g -index is not always well-defined was acknowledged by EGGHE [2006B] in a "Note added in proof". His proposed solution was that in such cases "fictitious articles with 0 citations have to be added (until the g -index can be determined)". We do not pursue this in what follows, merely note that the (theoretical) g -index is not necessarily defined.

Numerical investigations

The theoretical model involves four parameters:

- (i) The author's publication rate, θ .
- (ii) The gamma parameters, ν and α .
- (iii) The length of the author's publishing career to the current time, T .

For this investigation we will mainly be concerned with the time dependence of the various indexes and how they are inter-related, illustrating these using particular examples. Specifically, we will look at publication rates of $\theta = 2, 5$ and 10 . Assuming a publication rate of two papers per annum might be thought of as typical in a field such as mathematics, but it would be viewed as low in fields such as biomedicine where there is much collaborative work. Taking the "moderate" publication rate of $\theta = 5$ papers per annum, this would certainly be high in mathematics while in scientometrics, it might be a reasonable value for several of the best-known contributors! Similarly, whether a publication rate of $\theta = 10$ is low, medium or high will depend very much on the discipline.

For the gamma parameters we will focus on the mean citation rate $\mu = \nu/\alpha$, taking as an illustrative value $\mu = 5$. This could be thought of as a low, medium or high citation rate, depending very much on the subject context. BURRELL [2007A] found that the value of μ was much more important in calculations using the model rather than the individual values of α and ν . For our illustrative examples we will use $\alpha = 1$ and $\nu = 5$.

- (i) Determination of the h and $h(2)$ -indexes

Note that, according to (1) and (2), determination of $h(T)$ and $h_2(T)$ requires the calculation of $N(n;T)$ and $N(n^2;T)$ respectively, for $n = 0, 1, 2, \dots$ and compare it with the value of n in each case, identifying the point when the inequality between the two is reversed. (Recall that $N(n;T)$ is the expected number of the author's papers receiving exactly n citations by time T .) Calculation of $N(n;T)$, and hence of $N(n^2;T)$, is straightforward from the formula (7) using any package that incorporates calculation of the cumulative distribution function of the beta distribution. As there is essentially only one quantity, namely $N(n;T)$ to be computed, the two indexes are most readily determined from spreadsheet output from calculation of $N(n;T)$ and is illustrated in Table 1 for the case $\theta = 2, \mu = 5$ and $T = 2, 4, 6, 8, 10$. Note first that, as expected, for fixed T the value of $N(n;T)$ decreases as n increases; for fixed n it increases as T increases. In particular note that $N(0;T) = 2T$ in each case, in accordance with (9). The values of the two indexes have been highlighted in Table 1, $h(T)$ in bold, $h_2(T)$ in bold italics. For instance, reading down the entries for $N(n;6)$ we find $N(8;6) = 8.08 > 8$ while $N(9;6) = 7.62 < 9$ so that by (1), $h(6) = 8$. Similarly, $N(3;6) = 10.50 > 3^2$ while $N(4;6) = 10.01 < 4^2$ so that from (2), $h_2(6) = 3$.

Table 1. Tabulation of $N(n;T)$ in the case $\theta = 2, \mu = 5$

n	T = 2	T = 4	T = 6	T = 8	T = 10
0	4.00	8.00	12.00	16.00	20.00
1	3.57	7.50	11.50	15.50	19.50
2	3.03	7.00	11.00	15.00	19.00
3	2.58	6.51	10.50	14.50	18.50
4	2.17	6.03	10.01	14.00	18.00
5	1.79	5.56	9.52	13.51	17.50
6	1.47	5.10	9.03	13.01	17.01
7	1.19	4.66	8.55	12.52	16.51
8	0.95	4.24	8.08	12.04	16.02
9	0.76	3.84	7.62	11.55	15.53
10	0.60	3.47	7.17	11.08	15.04
11	0.47	3.12	6.74	10.61	14.56
12	0.36	2.80	6.32	10.15	14.08
13	0.28	2.50	5.91	9.69	13.60
14	0.21	2.22	5.52	9.25	13.13

(ii) Determination of the g-index

As remarked before, g is rather different to the other two indexes as it involves not just the number of most productive papers but also the numbers of citations carried by them. For any given values of the parameters, we have to determine $g(T)$ according to (4) which involves calculation of $N(n;T)$, $N(n;T)^2$ and $C(n;T)$ for a range of values of n using the formulae (7) and (8). Evaluation of each of these is again straightforward in any statistical package allowing evaluation of the cumulative distribution function of the beta distribution. Rather than give this numerical determination, note that for the approximate determination of the g-index by graphical means, it is easy to plot $N(n;T)^2$ and $C(n;T)$ against n. Where these intersect, we can read off the value of $N(n;T)^2$ and hence find the appropriate $N(n;T) = g(T)$. We illustrate this as in Figure 1 (a)-(c) for $T = 5$ and various values of θ .

In Figure 1(a) with $\theta = 10$, the point of intersection appears to be at about $n = 13$, and we can read off the approximate value of the ordinate. In fact, calculations show that $N(13;5)^2 = 425.78 < C(13;5) = 436.82$ while $N(12; 5)^2 = 504.33 > C(12; 5) = 459.21$ so that $g(5) = N(13;5) = 20.63$.

In the case of Figure 1(b) the point of intersection is at about $n = 7$ and again we can read off the approximate value of the ordinate. Calculations now show that $N(7;5)^2 = 271.38 < C(7;5) = 279.40$ while $N(6;5)^2 = 311.04 > C(6;5) = 287.32$. Hence we have $g(5) = N(7;5) = 16.47$.

Note that in Figure 1(c) there is no point of intersection because here we have an example where $\theta < \mu/2$, in which case, from Proposition 2 we have that the g-index does not exist.

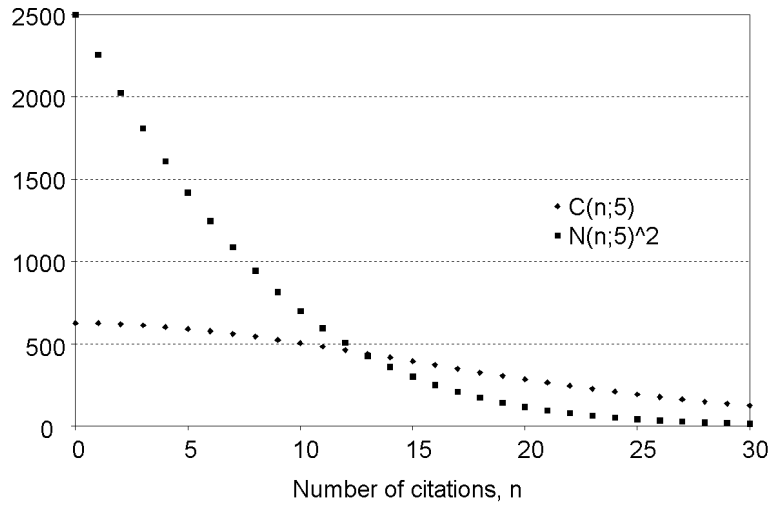


Figure 1a. Graphical determination of $g(5)$, with $\theta = 10$, $\mu = 5$

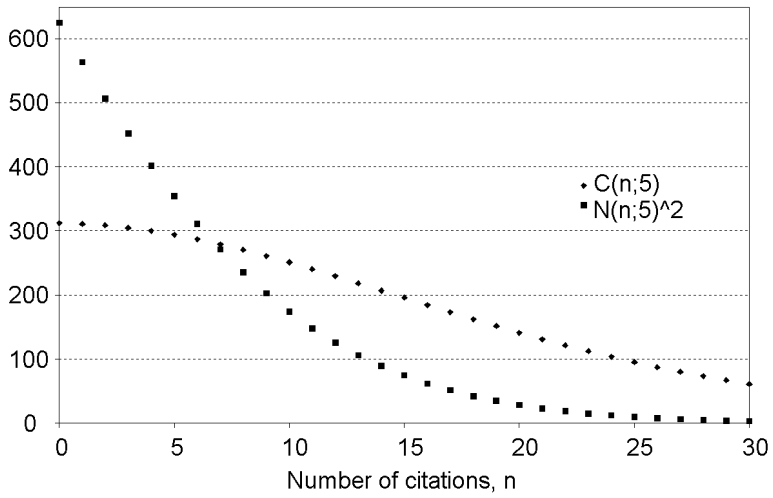


Figure 1b. Graphical determination of $g(5)$, with $\theta = 5$, $\mu = 5$

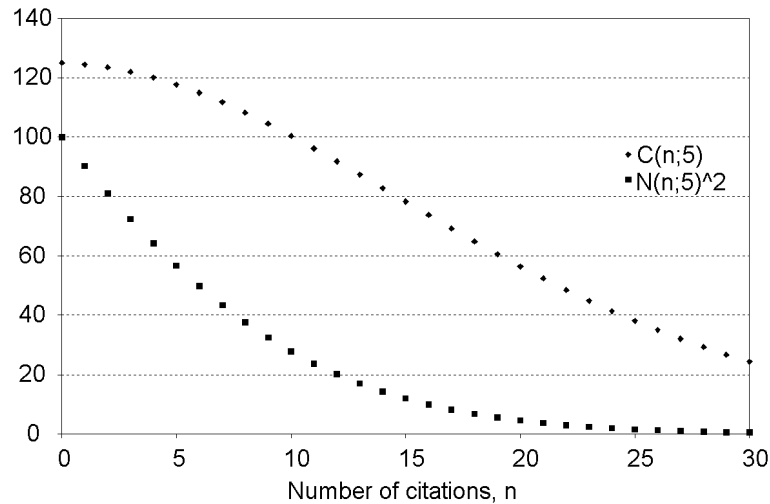


Figure 1c. Graphical determination of $g(5)$, with $\theta = 2$, $\mu = 5$

Remark. It is worth noting that, unlike $h(T)$ and $h_2(T)$, the *theoretical* g -index as determined by our model need not be an integer since it arises as the expected value of $N(n;T)$, a random variable. (See the above examples.)

(iii) Time dependence of the indexes

In Figure 2(a, b) we give plots of $h(T)$, $g(T)$ and $h_2(T)$ against time for the two examples where the g -index exists together with the fitted regression lines. Note that the fitted lines have been constrained to pass through the origin, a reasonable requirement in all cases. In BURRELL [2007A] we showed that with our theoretical model the h -index is almost directly proportional to time, thus supporting the original conjecture of HIRSCH [2005]. Here we find that the same is true also for the g -index. For the $h(2)$ -index, the situation is rather less clear cut. A certain linearity is again evident, but the R^2 values are much smaller. One reason for this is that the increase of $h_2(T)$ with T is slow and, by its construction, $h_2(T)$ is integer-valued so that it increases in a sequence of discrete steps.

(iv) Inter-relationships between the indexes

The results of the previous section suggest that each index is approximately a simple multiple of time. From this it follows that each index should be (approximately) directly proportional to each of the others. Simple graphical and statistical analysis confirms this. (Details can be obtained from the author.)

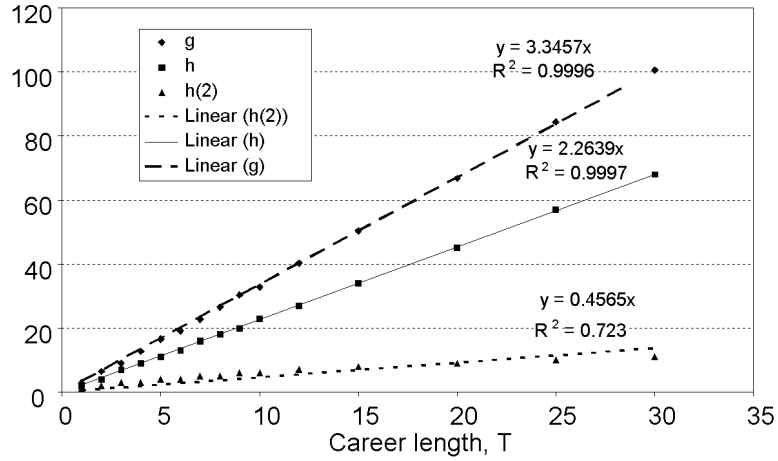


Figure 2a. The case $\theta = 5, \mu = 5$

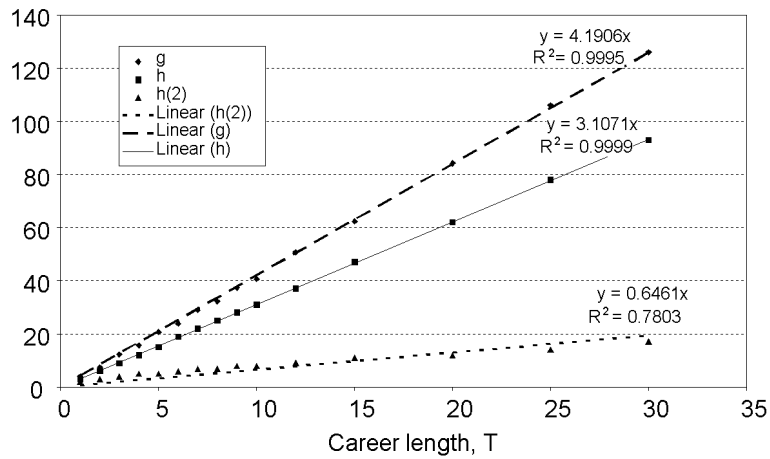


Figure 2b. The case $\theta = 10, \mu = 5$

The direct proportionality of g and h is predicted in EGGHE's [2006B] paper, using a non-stochastic model based upon an assumed Lotka form for the distribution of citations. EGGHE [2006B] further suggests using the constant of proportionality $g(T)/h(T)$ as a possibly interesting measure to use when comparing the outputs of different scientists. We agree that this would seem to be a possibly fruitful line of inquiry for empirical scientometric studies.

Concluding remarks

Within the confines of the model, we find that all three of an author's indexes are (approximately) directly proportional to the current career length and hence that they are proportional to each other. However, this is based on selected numerical calculations based on a particular model rather than a proper mathematical analysis and hence further theoretical as well as more extensive numerical work is required. On the other hand, if these preliminary findings are found to be more generally true, then each index is essentially measuring the same thing – only the constant of proportionality is different. From a practical point of view, which index is to be preferred? As we have that for any author, $g \geq h \geq h(2)$ this means that, in principle, accurate determination of the g-index requires more work than does the h-index, which in turn requires more work than the h(2)-index. On the other hand, EGGHE [2006B] shows that in comparative empirical studies the g and h indexes can lead to different perspectives. It would seem that, at this stage at least, all three measures are worth exploring further, particularly in empirical studies.

*

Some of the results in this paper were originally presented at the 9th International Science and Technology Indicators Conference [BURRELL, 2006] and the 11th International Conference of the International Society for Scientometrics and Informetrics [BURRELL, 2007D].

References

- BANKS, M. G. (2006), An extension of the Hirsch index: Indexing scientific topics and compounds. *Scientometrics*, 69 (1) : 161–168.
- BATISTA, P. D., CAMPITELI, M. G., KINOCHI, O. (2006), Is it possible to compare researchers with different scientific interests? *Scientometrics*, 68 (1) : 179–189.
- BORNHANN, L., DANIEL, H.-D. (2005), Does the h-index for ranking of scientists really work? *Scientometrics*, 65 (3) : 391–392.
- BURRELL, Q. L. (1992), A simple model for linked informetric processes. *Information Processing and Management*, 28 : 637–645.
- BURRELL, Q. L. (2006), Hirsch's h-index: a preliminary stochastic model. *Book of Abstracts: 9th International Science and Technology Indicators Conference*, 7-9 September 2006, Leuven, Belgium, 26–28. Katholieke Universiteit, Leuven.
- BURRELL, Q. L. (2007A), Hirsch's h-index: a stochastic model. *Journal of Informetrics*, 1 (1) : 16–25.
- BURRELL, Q. L. (2007B), On the h-index, the size of the Hirsch core and Jin's A- index. *Journal of Informetrics*, 1 (2) : 170–177.
- BURRELL, Q. L. (2007C), Hirsch index or Hirsch rate? Some thoughts arising from Liang's data. *Scientometrics*, 73 (1) : 19–28.
- BURRELL, Q. L. (2007D), Hirsch's h-index and Egghe's g-index. Paper presented at *ISSI 2007 – 11th International Conference of the International Society for Scientometrics and Informetrics*. Serrano Central Campus, Spanish Research Council, Madrid, Spain, 24-28 June, 2007.
- CRONIN, B., MEHO, L. (2006), Using the h-index to rank influential information scientists. *Journal of the American Society for Information Science and Technology*, 57 (9) : 1275–1278.

BURRELL: On Hirsch's h, Egghe's g and Kosmulski's h(2)

- EGGHE, L. (2006A), An improvement of the h-index: the g-index. *ISSI Newsletter*, 2 (1) : 8–9.
- EGGHE, L. (2006B), Theory and practice of the g-index. *Scientometrics*, 69 (1) : 131–152.
- EGGHE, L. (2006C), Dynamic h-index: the Hirsch index in function of time. *Journal of the American Society for Information Science and Technology* (To appear.)
- EGGHE, L., ROUSSEAU, R. (2006), An informetric model for the h-index. *Scientometrics*, 69 (1) : 121–129.
- GLÄNZEL, W. (2006A), On the H-index – a mathematical approach to a new measure of publication activity and citation impact. *Scientometrics*, 67 (2) : 315–321.
- GLÄNZEL, W. (2006B), On the opportunities and limitations of the H-index. *Science Focus*, 1 (1) : 10–11. (In Chinese.)
- HIRSCH, J. E. (2005), 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. (Also available in preprint form as arXiv:physics/0508113, accessible at <http://xxx.arxiv.org/abs/physics/0508025>.)
- JIN, BH. (2006), H-index: an evaluation indicator proposed by scientist. *Science Focus*, 1 (1) : 8–9. (In Chinese.)
- KOSMULSKI, M. (2006), A new Hirsch-type index saves time and works equally well as the original h-index. *ISSI Newsletter*, 2 (3) : 4–6.
- LIANG, L. (2006), h-index sequence and h-index matrix: Constructions and applications. *Scientometrics*, 69 (1) : 153–159.
- ROSS, S. (1996), *Stochastic Processes*. 2nd edition. New York: John Wiley.
- ROUSSEAU, R. (2006A), A case study: evolution of JASIS' Hirsch index. *Science Focus*, 1 (1) : 16–17, (in Chinese). (English version available at E-LIS, code 5430.)
- ROUSSEAU, R. (2006B), New developments related to the Hirsch index. *Science Focus*, 1 (4) : 23–25, (in Chinese). (English version available at E-LIS, code 6736.)
- VAN RAAN, A. J. (2006), Comparison of the Hirsch-index with standard bibliometric indicators and with peer judgment for 147 chemistry research groups. *Scientometrics*, 67 (3) : 491–502.

Appendix

Proof of Proposition

$$\begin{aligned}
 \text{(i)} \quad f(n;T) &= E[\# \text{ papers receiving } n \text{ citations by time } T] \\
 &= E[\# \text{ papers by time } T]P(\text{paper receives } n \text{ citations}) \\
 &= \theta T P(X_T = n) \\
 &= \frac{\alpha\theta}{(v-1)} B\left(\frac{T}{\alpha+T}; n+1, v-1\right) \text{ for } n = 0, 1, 2, \dots
 \end{aligned}$$

In the above we have used standard results for the mean of a binomial distribution, the mean of a Poisson process and (5)

$$\text{(ii)} \quad N(n;T) = E[\# \text{ papers receiving at least } n \text{ citations by time } T] = \sum_{j=n}^{\infty} f(j;T)$$

$$= \frac{\alpha\theta}{(v-1)} \sum_{j \geq n} B\left(\frac{T}{\alpha+T}; j+1, v-1\right)$$

$$= E[\text{total } \# \text{ papers by time } T]$$

$$- E[\# \text{ papers receiving less than } n \text{ citations by time } T]$$

$$= \theta T - \frac{\alpha\theta}{(v-1)} \sum_{j=0}^{n-1} B\left(\frac{T}{\alpha+T}; j+1, v-1\right)$$

$$= \theta T \left(1 - \frac{\alpha}{(v-1)T} \sum_{j=0}^{n-1} B\left(\frac{T}{\alpha+T}; j+1, v-1\right) \right)$$

$$\text{(iii)} \quad C(n;T) = E[\text{total } \# \text{ citations for those papers receiving at least } n \text{ citations each}]$$

$$= \sum_{j \geq n} j f(j;T)$$

$$\theta T \left(\frac{vT}{2\alpha} - \frac{\alpha}{(v-1)T} \sum_{j=0}^{n-1} j B\left(\frac{T}{\alpha+T}; j+1, v-1\right) \right)$$