

Relations between the shape of a size-frequency distribution and the shape of a rank-frequency distribution

Peer-reviewed author version

EGGHE, Leo & Waltman, Ludo (2011) Relations between the shape of a size-frequency distribution and the shape of a rank-frequency distribution. In: INFORMATION PROCESSING & MANAGEMENT, 47(2). p. 238-245.

DOI: 10.1016/j.ipm.2010.03.009

Handle: <http://hdl.handle.net/1942/11518>

# RELATIONS BETWEEN THE SHAPE OF A SIZE-FREQUENCY DISTRIBUTION AND THE SHAPE OF A RANK-FREQUENCY DISTRIBUTION

L. Egghe<sup>(\*)</sup> and L. Waltman<sup>(\*\*)</sup>

<sup>(\*)</sup>Universiteit Hasselt, Campus Diepenbeek, Agoralaan, B-3590 Diepenbeek, Belgium

<sup>(\*\*)</sup>Centre for Science and Technology Studies, Leiden University, P.O. Box 905, 2300 AX Leiden, The Netherlands

## ABSTRACT

We study the dependence of the shape of the rank-frequency distribution  $g$  on the shape of the size-frequency distribution  $f$  and vice versa. We show mathematically that  $g$  is convexly decreasing if and only if  $f$  is monotonically decreasing and that  $g$  has an S-shape (i.e.,  $g$  is first convexly decreasing and then concavely decreasing) if and only if  $f$  is first increasing and then decreasing.

To illustrate our mathematical results, we empirically analyze size- and rank-frequency distributions of the number of articles and of the impact factor of journals in various scientific fields. We find that most of the size-frequency distributions that we examine are first increasing and then decreasing. Most of the rank-frequency distributions that we examine have an S-shape. However, the concave part of the S-shape is sometimes very small.

## I. INTRODUCTION

An important topic in informetric research is the study of informetric distributions, such as distributions of authors, citations, or publications. In empirical work, there are two ways in which informetric distributions are commonly presented, namely as size-frequency distributions and as rank-frequency distributions. Both approaches to presenting informetric distributions convey the same information. As is well known, many informetric distributions approximately follow Lotka's law. For these distributions, the size- and rank-frequency presentations look similar, that is, they both show a decreasing power law. However, there are also informetric distributions that do not follow Lotka's law, and for these distributions the size- and rank-frequency presentations may look quite different. In this paper, we study this phenomenon. More specifically, we study, both mathematically and empirically, how size- and rank-frequency distributions

are related to each other. We also briefly touch upon the modeling of non-Lotkaian informetric distributions. We do so by presenting a mathematical analysis of a generalization of Zipf's law recently proposed by Mansilla et al. (2007).

The definitions of size- and rank-frequency distributions can be given in the context of information production processes (IPPs) (e.g., Egghe (2005a)). IPPs are systems consisting of sources that have, or produce, items. An example is given by journals that have (publish) articles. Another example is given by journals that have (receive) citations. Many more examples can be found in Chapter 1 in Egghe (2005a).

The size-frequency distribution  $f$  is defined as  $f(n)$  being the number ( $> 0$ ) of sources with  $n$  items ( $n = 1, 2, \dots$ ). If we rank the sources in decreasing order of their number of items and if we denote by  $r$  their ranks ( $r = 1, 2, \dots$ ), then the rank-frequency distribution  $g$  is defined as  $g(r)$  being the number of items in the source on rank  $r$ . So in the first example  $f(n)$  is the number of journals with  $n$  articles. If we rank the journals in decreasing order  $r$  of their number of articles, then  $g(r)$  is the number of articles in the journal on rank  $r$ . Replacing "articles" by "citations" yields the definitions of  $f(n)$  and  $g(r)$  in the second example.

It is clear that there is a general relation between the size-frequency distribution  $f$  and the rank-frequency distribution  $g$ . Denoting by  $g^{-1}$  the inverse function of  $g$ , we have by definition of  $f$  and  $g$

$$r = \sum_{n'=n}^{\infty} f(n') = g^{-1}(n) \quad (1)$$

where  $n = g(r)$ . Note that (1) defines a strictly decreasing function in  $n$ , which means that  $g$ , the inverse function of  $g^{-1}$ , indeed exists.

In the above examples (and in the examples in Chapter 1 in Egghe (2005a)),  $n$  is a positive whole number (a so-called natural number, i.e.,  $n \in \mathcal{N}$ ). However, we can generalize the IPP framework to cases where  $n$  need not be a whole number. This is needed for the following case, which we study in this paper. If we take the two examples of IPPs given above (i.e., journals and their number of articles and journals and their number of citations) and we divide the number of citations of a journal by the number of articles of a journal, then we obtain the impact factor (IF) of a journal. (Hence, journals and their IFs can be seen as an IPP derived from two other IPPs.)

In general IFs are not whole numbers. Hence, in the case of IFs, the definitions of the size-frequency distribution  $f$  and the rank-frequency distribution  $g$  cannot be given as above and (1) also cannot be used. Indeed, it does not make much sense to define  $f$  as

the number of journals with a certain IF. This is because IFs range in  $Q^+$ , the set of positive rational numbers. The solution to this problem is well known. We have to adopt the framework of continuous variables and treat  $f$  and  $g$  as density functions (in the same way as density functions of continuous variables are used in probability theory).

We now define  $f$  to be the size-frequency distribution where for every  $n \in R^+$ ,  $f(n)$  is the density ( $> 0$ ) of sources with  $n$  items, that is, for every  $m, n \in R^+$ ,  $m < n$ ,

$$\int_m^n f(n') dn' \quad (2)$$

denotes the number of sources with between  $m$  and  $n$  items (e.g., the number of journals with an IF between  $m$  and  $n$ ).

The corresponding rank-frequency distribution  $g$  is defined as

$$r = \int_n^\infty f(n') dn' = g^{-1}(n) \quad (3)$$

where  $n = g(r)$ . Equation (3) is a continuous version of (1). If  $n$  is a whole number, then the use of (3) rather than (1) can be convenient for calculatory reasons. In the case of “derived item values”, such as IFs, we have to use (3). Note that (3) implies that  $g^{-1}$  is strictly decreasing and hence that  $g$ , the inverse function of  $g^{-1}$ , indeed exists. Equation (3) defines  $g$  given  $f$ , but it also determines  $f$  given  $g$ , since (3) is equivalent with

$$f(n) = -\frac{1}{g'(g^{-1}(n))} \quad (4)$$

given that  $g(0) = \infty$ .

In earlier work by the first author (Egghe (2005a)), Lotkaian models for size-frequency distributions were studied as the basic functions in informetric research. In a Lotkaian framework, size- and rank-frequency distributions are both decreasing power laws. Although a Lotkaian framework is highly useful in many areas of informetric research, empirical data sometimes shows significant deviations from Lotkaian models. The empirical data studied in this paper illustrates this phenomenon. The data yields size-frequency distributions that in many cases do not approximate decreasing power laws. Instead, the distributions tend to be first increasing and then decreasing.<sup>1</sup> For such data, the use of Lotkaian models is not appropriate and a more general approach is needed. In

---

<sup>1</sup> Other examples of this phenomenon are books and their number of circulations and articles and their number of authors (oral communication by R. Rousseau).

this paper, we explore such an approach by studying the relation between size- and rank-frequency distributions without assuming a Lotkaian framework.

The paper is organized as follows. In the next section, we present a mathematical analysis of the relation between the shape of the size-frequency distribution  $f$  and the shape of the rank-frequency distribution  $g$ . We show that  $g$  is convexly decreasing if and only if  $f$  is monotonically decreasing and that  $g$  has an S-shape (i.e.,  $g$  is first convexly decreasing and then concavely decreasing) if and only if  $f$  is first increasing and then decreasing. In the third section, we empirically analyze size- and rank-frequency distributions of the number of articles and of the IF of journals in various scientific fields. We show examples of size-frequency distributions that are monotonically decreasing as well as of size-frequency distributions that are first increasing and then decreasing. We also show the corresponding rank-frequency distributions. Some rank-frequency distributions are convexly decreasing, while others have an S-shape. In the fourth section, we briefly consider the modeling of non-Lotkaian informetric distributions. We mathematically study a generalization of Zipf's law recently proposed by Mansilla et al. (2007), and we show how, depending on a parameter, this generalized Zipf's law yields either a convexly decreasing rank-frequency distribution or an S-shaped rank-frequency distribution.

## II. MATHEMATICAL ANALYSIS

We first need some lemmas on general injective functions  $g$  (i.e., for which  $g^{-1}$  exists).

### Lemma II.1

$g$  is strictly decreasing if and only if  $g^{-1}$  is strictly decreasing.

Proof :

$g$  is strictly decreasing if and only if, for all values  $r_1, r_2: r_1 < r_2 \Leftrightarrow g(r_1) > g(r_2)$ .

Denoting  $g(r_1) = n_1$  and  $g(r_2) = n_2$ , this is equivalent with  $g^{-1}(n_1) < g^{-1}(n_2) \Leftrightarrow n_1 > n_2$ .

Hence,  $g^{-1}$  is strictly decreasing.  $\square$

A similar proof can be given for strictly increasing functions  $g$  and with the word "strictly" omitted.

### Lemma II.2

Let  $g$  be decreasing. Then  $g$  is convex if and only if  $g^{-1}$  is convex. Also,  $g$  is concave if and only if  $g^{-1}$  is concave.

Proof :

$g$  is convex if and only if, for all values  $r_1, r_2$  and all values  $\lambda \in ]0, 1[$  :  
 $g(\lambda r_1 + (1-\lambda)r_2) \leq \lambda g(r_1) + (1-\lambda)g(r_2)$ . Since  $g$  is decreasing, we have by Lemma II.1 that  $g^{-1}$  is decreasing. Hence, the above inequality is equivalent with  
 $g^{-1}(g(\lambda r_1 + (1-\lambda)r_2)) \geq g^{-1}(\lambda g(r_1) + (1-\lambda)g(r_2))$ . Denoting  $g(r_1) = n_1$  and  $g(r_2) = n_2$ , we obtain  $\lambda g^{-1}(n_1) + (1-\lambda)g^{-1}(n_2) \geq g^{-1}(\lambda n_1 + (1-\lambda)n_2)$ , which proves that  $g^{-1}$  is convex. The proof of the second assertion is similar.  $\square$

The above lemma is not true if  $g$  is not decreasing. Indeed, a similar proof as the one above shows that if  $g$  is increasing and convex, then  $g^{-1}$  is concave and, similarly, that if  $g$  is increasing and concave, then  $g^{-1}$  is convex.

### Lemma II.3

Let  $g$  be decreasing. Then  $g$  has an S-shape, first convex and then concave, if and only if  $g^{-1}$  has an S-shape, first concave and then convex.

Proof :

Let  $g$  be defined on the interval  $[0, T]$  (we can even take  $[0, \infty[$  if  $T = \infty$ ). Suppose that  $g$  has an S-shape, first convex and then concave. Then there exists a number  $r_1 \in ]0, T[$  such that the restriction of  $g$  to the interval  $[0, r_1]$ , denoted  $g|_{[0, r_1]}$ , is convex and such that the restriction of  $g$  to the interval  $[r_1, T]$ , denoted  $g|_{[r_1, T]}$ , is concave. Since  $g|_{[0, r_1]}$  is decreasing and convex, we have by Lemma II.2 that  $(g|_{[0, r_1]})^{-1}$  is convex on the interval  $[g(r_1), g(0)]$ . Since  $g|_{[r_1, T]}$  is decreasing and concave, we have by Lemma II.2 that  $(g|_{[r_1, T]})^{-1}$  is concave on the interval  $[g(T), g(r_1)]$ . Hence we have that  $g^{-1}$  has an S-shape, first concave and then convex. The proof of the reverse assertion is similar.  $\square$

We now prove two theorems on shape relations between the size-frequency distribution  $f$  and the rank-frequency distribution  $g$ .

### Theorem II.4

$f$  is decreasing if and only if  $g$  is convex.

Proof :

From (3) we have

$$(g^{-1})'(n) = -f'(n) \quad (5)$$

and hence

$$(g^{-1})''(n) = -f''(n) \quad (6)$$

From (6) it follows that  $f$  is decreasing if and only if  $g^{-1}$  is convex. Since  $g^{-1}$  is decreasing (by (3)), we have by Lemma II.1 that  $g$  is decreasing. By Lemma II.2,  $g^{-1}$  is convex if and only if  $g$  is convex.  $\square$

### **Theorem II.5**

$f$  is first increasing and then decreasing if and only if  $g$  has an S-shape, first convex and then concave.

Proof :

By (6),  $f$  is first increasing and then decreasing if and only if  $g^{-1}$  has an S-shape, first concave and then convex. Since  $g$  is decreasing, we have by Lemma II.3 that  $g^{-1}$  has an S-shape, first concave and then convex, if and only if  $g$  has an S-shape, first convex and then concave.  $\square$

Without making additional assumptions, we cannot say more about the dependence of the shape of the rank-frequency distribution  $g$  on the shape of the size-frequency distribution  $f$  and vice versa (e.g., if  $g$  has an S-shape, then what is the location of the inflection point of  $g$ ?). We also cannot say more about the shape of  $\ln g$  based on the shape of  $g$  (e.g., a convex function  $g$  can lead to a convex function  $\ln g$  or to a function  $\ln g$  that has an S-shape).

## **III. EMPIRICAL ILLUSTRATION**

In this section, we provide an empirical illustration of our mathematical results on shape relations between size- and rank-frequency distributions. We use data from Thomson Reuters' Journal Citation Reports (JCR) for 2008. We focus on the number of articles that a journal has published and on the IF of a journal. This data allows us to examine different types of distributions. We also looked at the number of citations that a journal has received. However, the resulting size-frequency distributions all turned out to be

monotonically decreasing, which is not very interesting for the purpose of illustrating our mathematical results.

We analyze data for nine scientific fields. A field is defined by a JCR subject category or, in the case of chemistry, computer science, physics, and psychology, by a number of JCR subject categories taken together. Some summary statistics for the nine fields that we consider are reported in Table 1. As can be seen in the table, both the distribution of the number of articles that a journal has published and the distribution of the IF vary widely among fields. Of course, it is well known that on average IFs are much higher in, for example, biochemistry & molecular biology than in mathematics. However, even if we correct for such scale differences, different fields are still characterized by quite different distributions. This is indicated by the coefficient of variation and the skewness in Table 1. (The coefficient of variation, defined as the standard deviation divided by the mean, is a scale-invariant measure of the dispersion of a distribution. The coefficient of variation can also be interpreted as a measure of concentration (e.g., Chapter 4 in Egghe (2005a)). The skewness is a measure of the asymmetry of a distribution and is scale-invariant as well.)

Table 1. Summary statistics for nine scientific fields.  $N$  denotes the number of journals, CV denotes the coefficient of variation, and Skew. denotes the skewness.

Field	$N$	Number of articles				Impact factor			
		Mean	Median	CV	Skew.	Mean	Median	CV	Skew.
Biochemistry & molec. biol.	266	180.6	108	1.7	7.4	3.7	2.6	1.1	3.9
Chemistry	441	283.0	153	1.4	3.5	2.2	1.4	1.2	4.2
Computer science	391	75.8	48	1.1	2.9	1.4	1.1	0.8	2.5
Economics	207	51.8	36	1.0	3.4	1.0	0.8	0.8	2.2
Mathematics	206	83.4	53	1.5	5.7	0.7	0.6	0.7	3.1
Neurosciences	213	138.1	80	1.3	3.9	3.4	2.7	1.3	3.9
Pharmacology & pharmacy	213	137.5	94	0.9	2.3	2.9	2.3	1.1	4.9
Physics	311	356.7	144	1.9	4.8	2.2	1.3	1.5	5.4
Psychology	453	52.0	36	0.9	2.4	1.7	1.2	1.0	4.0

In the rest of this section, we focus on three fields in particular, namely chemistry, economics, and mathematics. The distributions characterizing these three fields are quite different. Together, the three fields can be regarded as representative for the nine fields listed in Table 1.

We first look at the way in which the number of articles that a journal has published is distributed in each of the three fields. The size-frequency distributions and the corresponding rank-frequency distributions are shown in Figure 1. As can be seen in the figure, the size-frequency distribution is monotonically decreasing in the case of chemistry, while it is first increasing and then decreasing in the case of economics and mathematics. Hence, based on Theorems II.4 and II.5, the rank-frequency distribution



should be convex in the case of chemistry and first convex and then concave in the case of economics and mathematics. The rank-frequency distributions shown in Figure 1 indeed have these shapes. However, in the case of mathematics, the concave part of the rank-frequency distribution is rather difficult to see. This is an example of a more general observation that we made by examining the distributions obtained for all nine fields listed in Table 1. It turns out that the concave part of a rank-frequency distribution is sometimes very small. Because of this, it can be difficult to distinguish between rank-frequency distributions that have an S-shape and rank-frequency distributions that do not have an S-shape.

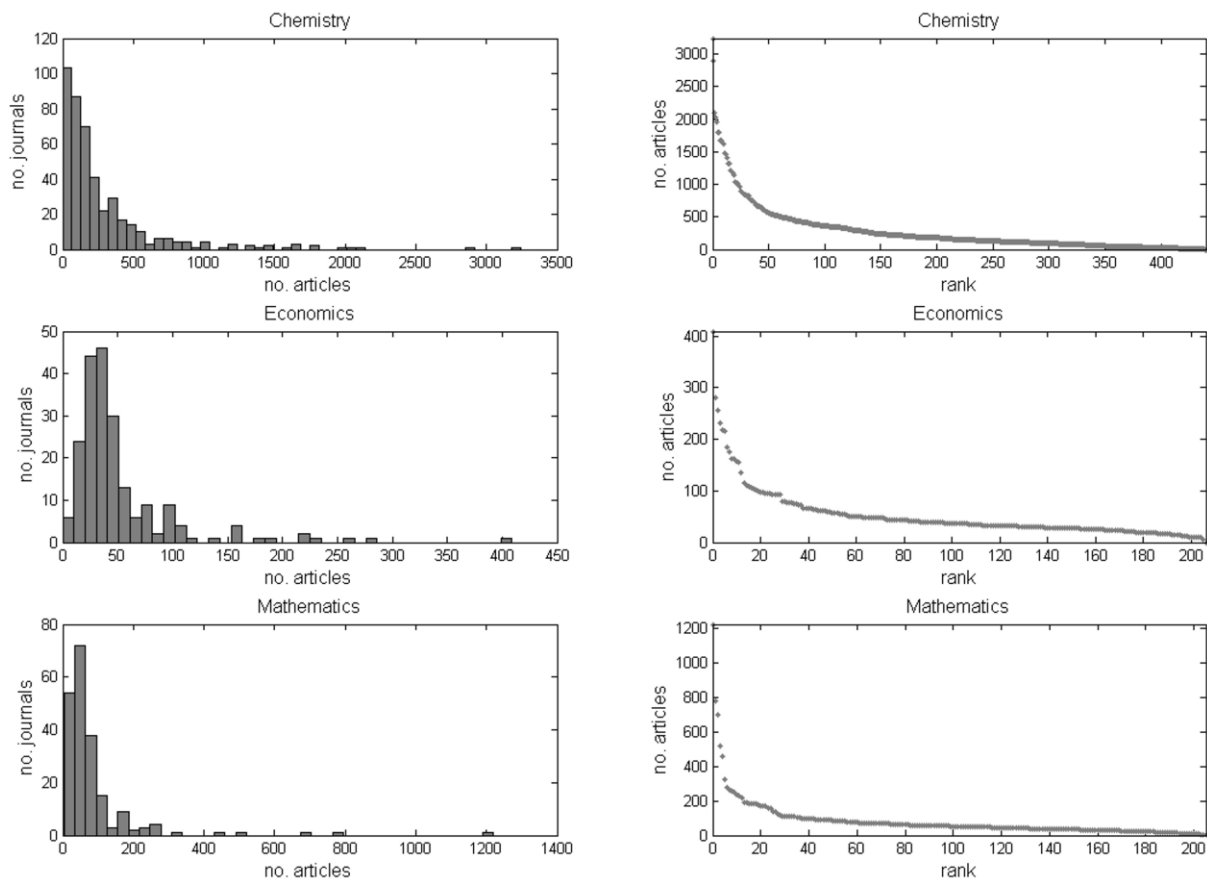


Figure 1. Size- and rank-frequency distributions of the number of articles that a journal has published.

We now turn to the distribution of the IF (see also Mansilla et al. (2007); Waltman and van Eck (2009)). The size-frequency distributions and the corresponding rank-frequency distributions for chemistry, economics, and mathematics are shown in Figure 2. IF size-frequency distributions generally seem to be first increasing and then decreasing. For some fields, however, the increasing part of the size-frequency distribution is almost negligible. This is for example the case for chemistry. (If we had used somewhat wider histogram bins, the increasing part of the size-frequency distribution for chemistry would

not even have been visible in Figure 2.) It follows from Theorem II.5 that the rank-frequency distributions for chemistry, economics, and mathematics should all have an S-shape. The rank-frequency distributions for economics and mathematics shown in Figure 2 clearly have an S-shape. In the case of chemistry, the S-shape of the rank-frequency distribution is much more difficult to see. However, this makes perfect sense. Since the increasing part of the size-frequency distribution for chemistry is very small, one would expect (based on Theorem II.4) that the rank-frequency distribution for chemistry is almost completely convex. This is indeed what we see in Figure 2.

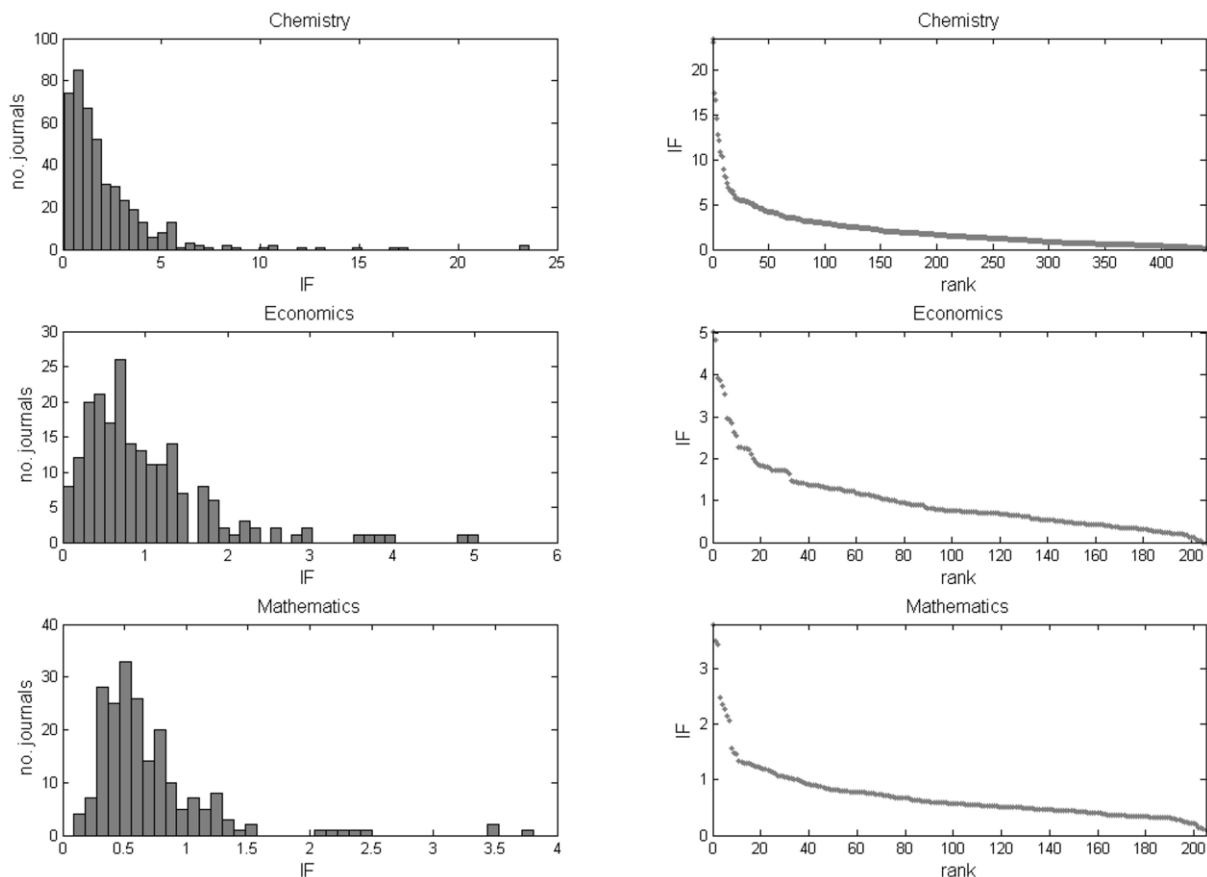


Figure 2. Size- and rank-frequency distributions of the IF of a journal.

#### IV. MODELING S-SHAPED RANK-FREQUENCY DISTRIBUTIONS

In the previous section, we have shown examples of size-frequency distributions that are first increasing and then decreasing. The corresponding rank-frequency distributions have an S-shape. Clearly, a size-frequency distribution that is first increasing and then decreasing does not follow Lotka's law. Similarly, an S-shaped rank-frequency distribution does not follow Zipf's law. Hence, to model such size- and rank-frequency

distributions in a satisfactory way, one needs a framework that is more flexible than the framework offered by the laws of Lotka and Zipf.

In this section, we study a recent proposal by Mansilla et al. (2007). Mansilla et al. are concerned with the modeling of rank-frequency distributions of IFs. They propose to use a generalization of Zipf's law given by

$$g(r) = K \frac{(N+1-r)^b}{r^a} \quad (7)$$

where  $a > 0$ ,  $b \geq 0$ , and  $K > 0$  are parameters,  $N$  is the total number of sources, and  $r = 1, \dots, N$ . If  $b = 0$ , (7) reduces to Zipf's law (and hence the corresponding size-frequency distribution is Lotkaian, see Egghe (2005a)). If  $a = b$ , (7) reduces to a function proposed by Lavalette (1996). Hence, (7) generalizes not only Zipf's law but also the function proposed by Lavalette. We note that (7) is also used by Campanario (in press-a, in press-b) and Martínez-Mekler et al. (2009).

We study (7) in a continuous setting. The following theorem states that, depending on  $b$ , (7) is either convex or S-shaped.

#### **Theorem IV.1**

Let  $g(r)$  denote the function in (7) with domain  $]0, N+1[$  and with  $a > 0$ ,  $b \geq 0$ , and  $K > 0$ . Then,

- (i)  $g(r)$  is strictly decreasing;
- (ii)  $g(r)$  has an S-shape, first convex and then concave, if  $0 < b < 1$ ;
- (iii)  $g(r)$  is convex if either  $b = 0$  or  $b \geq 1$ .

A proof of the theorem is provided in the appendix.

When fitting (7) to empirical IF data for various scientific fields, Mansilla et al. (2007) find for most fields that  $0 < b < 1$ . For a few fields they find that  $b \geq 1$ . Based on Theorem IV.1, this means that most of the fields studied by Mansilla et al. are characterized by an S-shaped IF rank-frequency distribution. This is in agreement with our empirical findings reported in the previous section.

A disadvantage of the rank-frequency distribution  $g$  in (7) is that there does not seem to exist a closed-form expression for the corresponding size-frequency distribution  $f$ . However, using Theorems II.4, II.5, and IV.1, we can at least derive some properties of  $f$ . It follows from Theorems II.4 and IV.1 that  $f$  is monotonically decreasing if either  $b = 0$  or  $b \geq 1$ , and it follows from Theorems II.5 and IV.1 that  $f$  is first increasing and then decreasing if  $0 < b < 1$ .

## V. CONCLUSION

We have mathematically analyzed the dependence of the shape of the rank-frequency distribution  $g$  on the shape of the size-frequency distribution  $f$  and vice versa. It turns out that  $g$  is convexly decreasing if and only if  $f$  is monotonically decreasing and that  $g$  has an S-shape (i.e.,  $g$  is first convexly decreasing and then concavely decreasing) if and only if  $f$  is first increasing and then decreasing.

Most size-frequency distributions in informetric research are monotonically decreasing. In this paper, however, we have empirically studied two exceptions to this rule, namely size-frequency distributions of the number of articles and of the IF of journals. For some fields these distributions are monotonically decreasing, but for most fields they are first increasing and then decreasing. (In the case of IFs, the increasing part of the distribution is sometimes quite small and may therefore not be visible in histograms with wide bins, such as in Beirlant et al. (2007) and Schwartz and Lopez Hellin (1996).) As one would expect based on our mathematical results, for most fields rank-frequency distributions of the number of articles and of the IF of journals have an S-shape. However, the concave part of the S-shape is sometimes very small.

We have also studied a generalization of Zipf's law recently proposed by Mansilla et al. (2007). It turns out that, depending on a parameter, this generalized Zipf's law yields either a convexly decreasing rank-frequency distribution or an S-shaped rank-frequency distribution. This flexibility explains why the proposal of Mansilla et al. is well suited for modeling rank-frequency distributions of IFs.

A question that remains is why some size-frequency distributions are monotonically decreasing while others are first increasing and then decreasing. Answering this question requires more insight into the underlying process that determines the shape of a size-frequency distribution. In the case of a Lotkaian (and hence monotonically decreasing) size-frequency distribution, it is sometimes suggested that a "success breeds success" mechanism or a mechanism based on exponential growth could be responsible for the shape of the distribution (e.g., Egghe (2005a, 2005b), Naranan (1970)). In a similar way, one could try to come up with a plausible mechanism that causes size-frequency distributions to be first increasing and then decreasing.<sup>2</sup> Related to this, one could try to build a model that explains functions such as the one proposed by Mansilla et al. (2007). We leave these issues for future research.

---

<sup>2</sup> For IF distributions, such a mechanism is studied by Egghe (2009). Egghe first points out that IFs are averages and then claims that, as a consequence of the central limit theorem, size-frequency distributions of IFs approximate normal distributions (for a similar reasoning, see van Raan (2006, p. 413)). Waltman and van Eck (2009) argue that Egghe's reasoning relies on unrealistic assumptions.

## REFERENCES

- J. Beirlant, W. Glänzel, A. Carbonez, and H. Leemans (2007). Scoring research output using statistical quantile plotting. *Journal of Informetrics*, 1(3), 185–192.
- J.M. Campanario (in press-a). Distribution of changes in impact factors over time. *Scientometrics*.
- J.M. Campanario (in press-b). Distribution of ranks of articles and citations in journals. *Journal of the American Society for Information Science and Technology*.
- L. Egghe (2005a). *Power Laws in the Information Production Process: Lotkaian Informetrics*. Elsevier, Oxford, UK.
- L. Egghe (2005b). The power of power laws and an interpretation of Lotkaian informetric systems as self-similar fractals. *Journal of the American Society for Information Science and Technology*, 56(7), 669–675.
- L. Egghe (2009). Mathematical derivation of the impact factor distribution. *Journal of Informetrics*, 3(4), 290–295.
- D. Lavalette (1996). Facteur d’impact: Impartialité ou impuissance? Internal Report, INSERM U350, Institut Curie, Paris.
- R. Mansilla, E. Köppen, G. Cocho and P. Miramontes (2007). On the behavior of journal impact factor rank-order distribution. *Journal of Informetrics*, 1(2), 155–160.
- G. Martínez-Mekler, R.A. Martínez, M.B. del Río, R. Mansilla, P. Miramontes and G. Cocho (2009). Universality of rank-ordering distributions in the arts and sciences. *PLoS ONE*, 4(3), e4791.
- S. Naranan (1970). Bradford’s law of bibliography of science: An interpretation. *Nature*, 227, 631–632.
- S. Schwartz and J. Lopez Hellin (1996). Measuring the impact of scientific publications. The case of the biomedical sciences. *Scientometrics*, 35(1), 119–132.
- A.F.J. van Raan (2006). Statistical properties of bibliometric indicators: Research group indicator distributions and correlations. *Journal of the American Society for Information Science and Technology*, 57(3), 408–430.
- L. Waltman and N.J. van Eck (2009). Some comments on Egghe’s derivation of the impact factor distribution. *Journal of Informetrics*, 3(4), 363–366.

## APPENDIX

In this appendix, we provide a proof of Theorem IV.1.

Let  $g(r)$  denote the function in (7) with domain  $]0, N + 1[$  and with  $a > 0$ ,  $b \geq 0$ , and  $K > 0$ . The first derivative of  $g(r)$  is given by

$$g'(r) = K \frac{(N+1-r)^{b-1} [(a-b)r - a(N+1)]}{r^{a+1}} \quad (\text{A1})$$

In the domain  $]0, N + 1[$ ,  $g'(r) < 0$  for all  $a > 0$ ,  $b \geq 0$ , and  $K > 0$ . Hence,  $g(r)$  is strictly decreasing for all  $a > 0$ ,  $b \geq 0$ , and  $K > 0$ . This proves part (i) of Theorem IV.1.

The second derivative of  $g(r)$  is given by

$$g''(r) = K \frac{(N+1-r)^{b-2}}{r^{a+2}} T(r) \quad (\text{A2})$$

where

$$T(r) = (a-b)(a-b+1)r^2 - 2a(a-b+1)(N+1)r + a(a+1)(N+1)^2 \quad (\text{A3})$$

In the domain  $]0, N + 1[$ ,  $g''(r)$  has the same sign as  $T(r)$ . Notice that  $T(r)$  is a quadratic equation (or a linear equation in case  $b = a$  or  $b = a + 1$ ).  $T(0)$  and  $T(N + 1)$  are given by

$$T(0) = a(a+1)(N+1)^2 \quad (\text{A4})$$

and

$$T(N+1) = b(b-1)(N+1)^2 \quad (\text{A5})$$

Hence,  $T(0) > 0$  for all  $a > 0$  and all  $b$ . The sign of  $T(N+1)$  depends on  $b$ .

We first consider the case in which  $a > 0$  and  $0 < b < 1$ . In this case, it follows from (A5) that  $T(N+1) < 0$ . Hence,  $T(r)$  is a quadratic (or linear) equation with  $T(0) > 0$  and  $T(N+1) < 0$ . It is clear that  $T(r)$  must have exactly one root in the interval  $]0, N + 1[$ . Let this root be denoted by  $r_1$ . For  $r \in ]0, r_1[$ ,  $T(r) > 0$  and consequently also  $g''(r) > 0$ . For  $r \in ]r_1, N + 1[$ ,  $T(r) < 0$  and consequently also  $g''(r) < 0$ . This means that  $g(r)$  has an S-shape, first convex and then concave. This proves part (ii) of Theorem IV.1.

We now consider the case in which  $a > 0$  and either  $b = 0$  or  $b \geq 1$ . In this case, it follows from (A5) that  $T(N+1) \geq 0$ . Hence,  $T(r)$  is a quadratic (or linear) equation with  $T(0) > 0$  and  $T(N+1) \geq 0$ . The discriminant of  $T(r)$  equals

$$D = 4ab(a-b+1)(N+1)^2 \quad (\text{A6})$$

$T(r)$  does not have a root in the interval  $]0, N+1[$ . To show this, we distinguish the following four cases:

- (i) If  $b = a$  or  $b = a+1$ ,  $T(r)$  is a linear equation. Since  $T(0) > 0$  and  $T(N+1) \geq 0$ ,  $T(r)$  does not have a root in the interval  $]0, N+1[$ .
- (ii) If  $b > a+1$ , (A6) yields  $D < 0$ . Hence,  $T(r)$  has no roots at all.
- (iii) If  $b = 0$ , (A6) yields  $D = 0$ . Hence,  $T(r)$  has one root. It follows from (A5) that this root is given by  $r_1 = N+1$ . This means that  $T(r)$  does not have a root in the interval  $]0, N+1[$ .
- (iv) If  $1 \leq b < a+1$  and  $b \neq a$ , (A6) yields  $D > 0$ . Hence,  $T(r)$  has two roots. One root of  $T(r)$  is given by

$$r_1 = \frac{a(a-b+1) + \sqrt{ab(a-b+1)}}{(a-b)(a-b+1)}(N+1) \quad (\text{A7})$$

Let the other root of  $T(r)$  be denoted by  $r_2$ . Based on (A7), it is not difficult to see that  $r_1 < 0$  or  $r_1 > N+1$ . Since  $T(0) > 0$  and  $T(N+1) \geq 0$ , it follows from this that  $r_2 < 0$  or  $r_2 \geq N+1$ . Hence,  $T(r)$  does not have a root in the interval  $]0, N+1[$ .

We have now shown that  $T(r)$  does not have a root in the interval  $]0, N+1[$ . Hence, for  $r \in ]0, N+1[$ ,  $T(r) > 0$  and consequently also  $g''(r) > 0$ . This means that  $g(r)$  is convex. This proves part (iii) of Theorem IV.1.