

Probabilities for encountering genius, basic, ordinary or insignificant papers based on the cumulative nth citation distribution

LEO EGGHE^{a,b}

^a Universiteit Hasselt (UHasselt), Diepenbeek (Belgium)

^b Universiteit Antwerpen (UA), Campus Drie Eiken, Wilrijk (Belgium)

This article calculates probabilities for the occurrence of different types of papers such as genius papers, basic papers, ordinary papers or insignificant papers. The basis of these calculations are the formulae for the cumulative nth citation distribution, being the cumulative distribution of times at which articles receive their nth ($n = 1,2,3,\dots$) citation.

These formulae (proved in previous papers) are extended to allow for different aging rates of the papers. These new results are then used to define different importance classes of papers according to the different values of n, in function of time t. Examples are given in case of a classification into four parts: genius papers, basic papers, ordinary papers and (almost) insignificant papers.

The fact that, in these examples, the size of each class is inversely related to the importance of the journals in this class is proved in a general mathematical context in which we have an arbitrary number of classes and where the threshold values of n in each class are defined according to the natural law of Weber-Fechner.

Introduction

Classifying papers (journal articles, articles in conference proceedings,...) according to their importance for a scientific discipline is an extremely difficult task. Of course, visibility can be measured according to the number of citations a paper receives in the

Received June 21, 2006

Address for correspondence:

LEO EGGHE

Universiteit Hasselt (UHasselt), Agoralaan, B-3590 Diepenbeek, Belgium

E-mail: leo.egghe@uhasselt.be

0138-9130/US \$ 20.00

Copyright © 2007 Akadémiai Kiadó, Budapest

All rights reserved

time t (years) after its publication. Often, visibility, as measured by citations received, is used as a substitute for the measurement of its importance.

In this paper we will use the above mentioned measurement technique for developing a mathematical model for importance classes of papers. More specifically we will use the models of the cumulative n^{th} citation distribution, developed earlier in EGGHE (2000) and EGGHE & RAO (2001) ($n = 1, 2, 3, \dots$). This is the distribution of times at which an article receives its n^{th} citation (diachronous model). It is clear that n in function of time t , used as variables, are important values in determining the degree of visibility of a paper, hence of its importance for the scientific community.

Using the model (an assumption, being the basic aging distribution)

$$c(t) = ba^t \quad (1)$$

for the density function of citations to an article, t time (e.g. years) after its publication ($t > 0$, $0 < a < 1$, $b > 0$, a, b constants) and the Lotka distribution (again an assumption, being a basic size-frequency function – see EGGHE (2005))

$$\varphi(M) = \frac{E}{M^\alpha} \quad (2)$$

for the distribution of articles that (eventually) receive M citations ($M \geq 1$), hence amongst the evercited articles), we arrived in EGGHE & RAO (2001), based on the theory developed in EGGHE (2000), at the formula (for $n = 1, 2, 3, \dots$, $t > 0$):

$$\Phi_n(t) = \left(\frac{1-a^t}{n} \right)^{\alpha-1} \quad (3)$$

for the cumulative fraction of papers with n or more citations at time t , fraction with respect to the set of ever cited papers. For the sake of completeness, we will shortly repeat the proof of the above result in the Appendix. For more on the n^{th} (sometimes $n = 1$) citation distribution, we refer to GLÄNZEL (1992), ROUSSEAU (1994) and BURRELL (2001, 2002).

The proof of the result (3) supposes that α and a are constants. This is reasonable for α since the size-frequency function describing the fraction of papers that (eventually) receive M citations (amongst the set of ever cited papers) can be supposed to be fixed for a fixed discipline (which we assume here). Of course, even within each fixed discipline, one has all kinds of papers (of the different importance or visibility degree as we want to describe in this paper), yielding different aging curves (1), hence different ‘ a ’s. In the next section we will extend (3), allowing for different a -values, hereby giving a formula for the cumulative n^{th} citation dsitribution, incorporating the different types of articles (as e.g. expressed by different a values).

The third section then generally presents the way to subdivide the total set of papers into different visibility or importance classes. An example can be given of 4 classes describing the classes of genius papers, basic papers, ordinary papers and insignificant papers (or at least little significant papers). The size of these classes can be calculated according to the probability for a paper to belong to one of such (disjoint) classes. Examples are given.

The last section then proves some mathematical properties of these classes. Using the above mentioned models we can show that for $\alpha \geq 2$, in all cases of i classes, the class of least significant papers is always the largest. We furthermore prove the following result: If the importance classes are constructed using exponential thresholds for the value of n (i.e. dividing the set of ever cited articles according to citation levels that increase exponentially – a natural choice as we will explain in the sequel), then we show that the size of each class is inversely related to its importance level, i.e. the less important a paper is, this paper then belongs to the larger importance class and vice-versa: the most important papers belong to the smallest classes. This, at first sight, evident fact can be proved based on the aging law (1), the law of Lotka (2) and the law of Weber-Fechner and explains the elitary character of visibility or importance and gives exact probabilities for a paper to belong to one of such classes, solving the problem posed in this paper.

The n^{th} citation distribution for variable values of a

Formula (3) denotes the probability (amongst the ever cited papers) that a paper has n or more citations at time t . Here a is given and fixed. Hence, in the extended environment where a is variable in the interval $]0,1[$ we could also denote (3) as

$$P(n, t | a) = \Phi_n(t) \quad (4)$$

being the conditional probability (amongst the ever cited papers) that a paper has n or more citations at time t , *given* the fixed value a of the aging rate.

By definition of conditional probability we have

$$P(n,t,a) = P(n,t | a)P(a) \quad (5)$$

where $P(n,t,a)$ is the joint probability in (n,t,a) and $P(a)$ is the probability that a paper has aging rate a . We will assume that $P(a) = 1$ for $a \in]0,1[$ expressing that all values of a are equally possible (and note that $\int_0^1 P(a)da = 1$). We use this as a first approximation

noting that in each scientific discipline low values of a are common for ordinary papers which receive citations for a relatively short period of time and that high values of a are common for rather insignificant papers (with a relative low constant number of

citations, e.g. one every year) but also for basic or genius papers which keep on attracting citations over a very long time period.

Allowing in this way different values of a we have that

$$P(n,t) = \int_0^1 P(n,t,a)da \quad (6)$$

is the marginal probability (amongst the ever cited papers) that a paper has n or more citations at time t (averaged over the aging rate a).

Combining (3), (4), (5) and the supposed $P(a) = 1, \forall a \in]0,1[$, we obtain

$$\begin{aligned} P(n,t) &= \int_0^1 \left(\frac{1-a^t}{n} \right)^{\alpha-1} da \\ P(n,t) &= \frac{1}{n^{\alpha-1}} \int_0^1 (1-a^t)^{\alpha-1} da \end{aligned} \quad (7)$$

As remarked by one of the referees, the expression in (7) can be evaluated as follows:

$$\begin{aligned} I &= \int_0^1 (1-a^t)^{\alpha-1} da \\ &= \frac{1}{t} \int_0^1 (1-y)^{\alpha-1} y^{\frac{1}{t}-1} dy \end{aligned}$$

upon substituting $y = a^t$. This is the classical beta function and hence we can write

$$I = \frac{1}{t} \frac{\Gamma(\alpha)\Gamma\left(\frac{1}{t}\right)}{\Gamma\left(\alpha + \frac{1}{t}\right)} = \frac{\Gamma(\alpha)\Gamma\left(\frac{1}{t}+1\right)}{\Gamma\left(\alpha + \frac{1}{t}\right)} \quad (8)$$

For $\alpha = 2$, the most common value for Lotka's exponent (and a "turning" point in Lotkaian informetrics – see EGGHE (2005)) we have for (8)

$$P(n,t) = \frac{1}{n} \left(1 - \frac{1}{t+1} \right) \quad (9)$$

Formula (9) is very simple and gives a good idea of the (n,t) -dependencies: $P(n,t)$ is naturally decreasing in n : it is linearly dependent on $\frac{1}{n}$ while $P(n,t)$ is naturally increasing in t : the functional relation being linear in $1 - \frac{1}{t+1}$. Note also that

$$\lim_{t \rightarrow \infty} P(n, t) = \lim_{t \rightarrow \infty} P(n, t, a) = \frac{1}{n} \quad (10)$$

since for $t \rightarrow \infty$, the joint probability distribution (5) becomes independent of a (by (3), (4) and (5) and since $0 < a < 1$). For $n = 1$, we have

$$P(1, t) = 1 - \frac{1}{t+1} \quad (11)$$

and hence $\lim_{t \rightarrow \infty} P(1, t) = 1$ as it should be a cumulative distribution function.

For t very large we can, of course, use (7) and (8) to obtain the more general formula (than (10)):

$$\lim_{t \rightarrow \infty} P(n, t) = \frac{1}{n^{\alpha-1}} \quad (12)$$

This is in agreement with (2) since (11) expresses the fraction (amongst the ever cited papers) of papers with, eventually (i.e. $t = \infty$), n or more citations, i.e. (use (2) and (A5)):

$$\begin{aligned} \int_n^\infty \phi(M) dM &= \int_n^\infty \frac{\alpha-1}{M^\alpha} dM \\ &= \frac{1}{n^{\alpha-1}} \\ &= \lim_{t \rightarrow \infty} P(n, t) \end{aligned}$$

It is now clear how to use (7) or (8) in the study of the different classes (of importance or visibility) of papers: the variables t and n allow for expressing how large or how small n is (the minimum number of received citations) and this for every time t selected. This clearly describes importance of papers at every measuring point t . We will make these intuitive ideas more clear in the next section.

Division of a set of papers into visibility or importance classes

We have now a general device to define visibility or importance classes for a set of papers (e.g. papers in a certain (vast) scientific field). Using (8) (or, more generally (6) or (7)) and fixing a time t we can introduce a decreasing sequence of number of citations as delimiters (thresholds) for the different classes:

$$n_1 > n_2 > \dots > n_i = 1$$

yielding i classes:

(1) The class of papers with n_1 or more citations at time t (i.e. the papers with a number of citations in the highest level: n_1 or more at t). The probability for a paper to belong to this class, hence its relative size equals, according to (7)

$$P_1 := P(n_1, t) = \frac{1}{n_1^{\alpha-1}} I(\alpha, t) \quad (13)$$

where we denoted

$$I(\alpha, t) = \int_0^1 (1 - a^t)^{\alpha-1} da \quad (14)$$

Hence, in the special case (8) we have

$$P_1 = \frac{1}{n_1} \left(1 - \frac{1}{t+1} \right) \quad (15)$$

(2) The class of papers with n_2 or more citations, but less than n_1 citations at time t (i.e. the papers in the second largest citation category). The probability for a paper to belong to this class, hence its relative size equals, according to (7)

$$P_2 := P(n_2, t) - P(n_1, t) = \left(\frac{1}{n_2^{\alpha-1}} - \frac{1}{n_1^{\alpha-1}} \right) I(\alpha, t) \quad (16)$$

which reduces, in the case (8), to

$$P_2 = \left(\frac{1}{n_2} - \frac{1}{n_1} \right) \left(1 - \frac{1}{t+1} \right) \quad (17)$$

(3) The next class has probability (hence relative size)

$$P_3 := P(n_3, t) - P(n_2, t)$$

$$= \left(\frac{1}{n_3^{\alpha-1}} - \frac{1}{n_2^{\alpha-1}} \right) I(\alpha, t) \quad (18)$$

or, in the special case (8):

$$P_3 = \left(\frac{1}{n_3} - \frac{1}{n_2} \right) \left(1 - \frac{1}{t+1} \right) \quad (19)$$

...

(i) The class of papers with 1 or more citations, but less than n_{i-1} citations at time t (i.e. the class of papers with the least citations at t , amongst the ever cited papers). The probability to belong to this class, hence its relative size equals

$$P_i := P(n_i, t) - P(n_{i-1}, t)$$

$$P_i = P(1, t) - P(n_{i-1}, t)$$

$$P_i = \left(1 - \frac{1}{n_{i-1}^{\alpha-1}} \right) I(\alpha, t) \quad (20)$$

or, in the special case (8):

$$P_i = \left(1 - \frac{1}{n_{i-1}} \right) \left(1 - \frac{1}{t+1} \right) \quad (21)$$

Let us give an example of $i = 4$. If $n_1 > n_2 > n_3 > n_4 = 1$ are well-chosen (for this, see the next section), we could classify the set of papers as

- (1) The set of genius papers whose relative size is given by (13) or (15),
- (2) The set of basic papers whose relative size is given by (16) or (14),
- (3) The set of ordinary papers whose relative size is given by (18) or (19),
- (4) The set of (almost) insignificant papers whose relative size is given by (20) or (21).

As an example, take $t = 9$, $n_1 = 200 > n_2 = 50 > n_3 = 10 > n_4 = 1$. We will also use $\alpha = 2$. Of course these cut-off values are discutable and at least they are dependent of the subject. Here it is just an abstract example.

We have

$$P_1 = 0.0045$$

$$P_2 = 0.0135$$

$$P_3 = 0.072$$

$$P_4 = 0.81$$

So, clearly, $P_1 < P_2 < P_3 << P_4$.

Hence, if we look at papers that are 9 years old, we define genius papers as papers receiving 200 or more citations in this period. They occur in 0.45% of the cases in this example (i.e. 0.45% of the ever cited papers is a genius paper). Basic papers, having between 50 and 199 citations in this 9-year period, comprise 1.35% of the total number of ever cited papers. Ordinary papers (number of citations between 10 and 49 in this 9-year period) comprise 7.2% of the ever cited papers and finally 81% of the ever cited papers have between 1 and 9 citations in this period (the class of almost insignificant papers). This last class should be completed with the papers that had not yet a citation at $t = 9$ (fraction $1 - (0.0045 + 0.0135 + 0.072 + 0.81) = 0.1$ of ever cited papers) and even with the papers without any citation ever, increasing the relative size of this last class (of rather invisible papers) to even more than 81%. Alternatively, the class of never cited papers could be handled as an extra class in the division of papers. In the next section, however, besides other results, we will show that, if $\alpha \geq 2$, even the class (i) itself is always the largest amongst the i defined classes showing that the used model of citedness inequality (e.g. Lotka's law (2)) mathematically implies that most papers are rather invisible (unimportant).

The choice of the values n_1, \dots, n_{i-1} is of course not determined nor is the choice of the time t . This is the power of the model: For any choice of values $n_1 > \dots > n_{i-1}, t$ we can use (7), (8) and (13)-(21), where we are able to compare the relative sizes of the determined visibility (importance) classes of papers.

Even without a fixed choice of the citation numbers n_j we are able to compare the class sizes. This will be done in the next section.

Mathematical properties of the defined article classes

It is clear that an interesting property to investigate is

$$P_1 < P_2 < \dots < P_{i-1} < P_i \quad (22)$$

Indeed, if (22) is true then we have the “logical” situation that importance classes increase in size if and only if the visibility (or importance) of the papers decreases.

We have the following trivial result.

Proposition 1:

If (22) is satisfied for a certain $t_0 > 0$ then (22) is valid for all $t > 0$.

Proof:

This follows readily from (13)-(21) since $I(\alpha, t)$ $\left(\text{or } 1 - \frac{1}{t+1} \right)$ is independent of the choices of n_1, \dots, n_{i-1} . \square

From the above it also follows that (22) (for any t) is valid if and only if we have the following inequalities

$$1 - \frac{1}{n_{i-1}^{\alpha-1}} > \frac{1}{n_{i-1}^{\alpha-1}} - \frac{1}{n_{i-2}^{\alpha-2}} \quad (23)$$

$$\frac{1}{n_{i-1}^{\alpha-1}} - \frac{1}{n_{i-2}^{\alpha-1}} > \frac{1}{n_{i-2}^{\alpha-1}} - \frac{1}{n_{i-3}^{\alpha-1}} \quad (24)$$

...

$$\frac{1}{n_3^{\alpha-1}} - \frac{1}{n_2^{\alpha-1}} > \frac{1}{n_2^{\alpha-1}} - \frac{1}{n_1^{\alpha-1}} \quad (25)$$

$$\frac{1}{n_2^{\alpha-1}} - \frac{1}{n_1^{\alpha-1}} > \frac{1}{n_1^{\alpha-1}} \quad (26)$$

We have the following results.

Proposition 2:

Independent of the choices of the values n_1, \dots, n_{i-1} , we always have, if $\alpha \geq 2$, that

$$P_{i-1} < P_i \quad (27)$$

In other words, the class of least cited papers is always larger than the class containing the second-to-least cited papers.

Proof:

We have to show, by (23), that

$$n_{i-1}^{\alpha-1} > \frac{2}{1 + \frac{1}{n_{i-2}^{\alpha-1}}} \quad (28)$$

But the right hand side of (28) is strictly smaller than 2 while $n_{i-1} > n_i = 1$ hence $n_{i-1} \geq 2$. If $\alpha \geq 2$ we have that $n_{i-1}^{\alpha-1} \geq 2$, whence the result. \square

The following theorem treats all inequalities and gives a sufficient condition for (22) to be true.

Theorem 3:

If

$$n_j > 2^{\alpha-1} n_{j+1} \quad (29)$$

for all $j = 1, \dots, i-1$, then (22) is valid.

Proof:

(26) requires

$$\frac{1}{n_2^{\alpha-1}} > \frac{2}{n_1^{\alpha-1}}$$

hence

$$n_1 > 2^{\alpha-1} n_2 \quad (30)$$

(in fact (30) is necessary and sufficient in order to have that $P_1 < P_2$).

To have the validity of one of the inequalities (24)-(25) we must have that (for $j \in \{2, \dots, i-2\}$)

$$\frac{1}{n_{j+1}^{\alpha-1}} - \frac{1}{n_j^{\alpha-1}} > \frac{1}{n_j^{\alpha-1}} - \frac{1}{n_{j-1}^{\alpha-1}}$$

hence

$$\frac{1}{n_{j+1}^{\alpha-1}} > \frac{2}{n_j^{\alpha-1}} - \frac{1}{n_{j-1}^{\alpha-1}}$$

It suffices to require

$$\frac{1}{n_{j+1}^{\alpha-1}} > \frac{2}{n_j^{\alpha-1}}$$

hence (29). Finally (23) is valid if and only if

$$1 > \frac{2}{n_{i-1}^{\alpha-1}} - \frac{1}{n_{i-2}^{\alpha-1}}$$

It suffices to require

$$1 > \frac{2}{n_{i-1}^{\alpha-1}}$$

hence

$$n_{i-1} > 2^{\frac{1}{\alpha-1}} = 2^{\frac{1}{\alpha-1}} n_i . \quad \square$$

This leads to the following important corollary.

Corollary 4:

If, for all $j = 1, \dots, i$

$$n_j = q^{i-j} \quad (31)$$

with

$$q > 2^{\frac{1}{\alpha-1}} \quad (32)$$

then (22) is valid.

Proof:

It suffices to prove (29). We have, for $j = 1, \dots, i-1$:

$$n_j = q^{i-j}$$

$$= q \cdot q^{i-j-1}$$

$$= q \cdot n_{j+1}$$

$$n_j > 2^{\frac{1}{\alpha-1}} n_{j+1}$$

proving (29). \square

Requirement (31) is very natural: it states that the sequence $1 = n_i, n_{i-1}, \dots, n_1$ must be exponentially increasing, i.e. that $\log(n_j)$ ($j = i, \dots, 1$) (any log can be used) is linearly increasing. This is related with the well-known law of Weber-Fechner (see e.g. EGGHE (2005)) stating that the sensation is proportional to the logarithm of the stimulus, i.e. that the sensation is linearly related to the logarithm of the stimulus. Here we need linearity in $\log(n_j)$ in function of the class numbers $j = i, \dots, 1$. It is indeed logical, in the construction of the importance classes, to allow for an exponential increase of number of citations when going from the classes of low number of citations to the ones with higher number of citations, so that the class index $i-j+1$ ($j = i, \dots, 1$) is linearly related with the logarithm of the (lower bound) of number of citations in this class (comparable with Weber-Fechner's law).

A similar construction is done when expressing distances from a certain point, e.g. in calculating the distance that library users live from e.g. a public library: one might start with classes with small distances from the library (say in the order of 500 meters) but for users living further away from the library it is best to construct classes, grouping users over several kilometers (it does not make sense to group users linearly in function of their distance to the library: constructing groups in a range [0,500m], [500m,1km] seems reasonable but not in a range [30km,30.5km] and so on). In the same way it is reasonable to define classes of not-so-important papers by using a small number of citations ($n = 1,2,3\dots$) but for basic or genius papers we need more citations than there are classes defined (see e.g. the example in the previous section).

Note that, if $\alpha = 2$ (the most common value of Lotka's exponent) we have the requirement (in Corollary 4) $q > 2$ in order to have the validity of (22). In this case one could also use $q = e > 2$, yielding the classical exponential function $x \rightarrow e^x$ in (31). This function can be used whenever

$$e > 2^{\frac{1}{\alpha-1}}$$

or

$$\alpha > 1.693$$

which is true in most cases.

Note also that, in order to have natural numbers for the threshold values n_j we need to round off the numbers $x = q^{j-j}$ (in Corollary 4) to the next higher entire number $\lceil x \rceil$ ($\lceil x \rceil$ denotes the smallest entire number larger than or equal to x , also called the floor function) and to check for (29).

Conclusions and open problems

In this paper we extended the model for the cumulative n^{th} citation distribution, developed in EGGHE (2000) and EGGHE & RAO (2001) to the universe of ever cited papers and allowing for variable aging rate a .

This new model is then used to define different classes of importance or visibility of papers going from genius papers to (almost) insignificant papers. The role played by Lotka's distribution in the model for the cumulative n^{th} citation distribution is extended here so that we are able to prove that, if class limits of n are chosen exponentially (in function of the class number), i.e. adopting a form of the law of Weber-Fechner, we always have that the class sizes increase with decreasing citedness of the papers. In this way we give a mathematical explanation of this "natural" and well-known phenomenon.

These results are also a mathematical support for the experimental result in GLÄNZEL et al. (2003) on genius papers (described there as papers with a high number of citations for t high, where the number of citations for t low are low).

As remarked by one of the referees such papers could also be called “Sleeping Beauties”, cf. VAN RAAN (2004), BURRELL (2005).

We leave open to find similar (or analogous) mathematical results based on the characteristics of highly cited papers (here the distinction between t low or high as in GLÄNZEL et al. (2003) is not made) as described in AKSNES (2003).

References

- AKSNES, D. W. (2003), Characteristics of highly cited papers. *Research Evaluation*, 12 (3) : 159–170.
- BURRELL, Q. L. (2001), Stochastic modelling of the first-citation distribution. *Scientometrics*, 52 (1) : 3–12.
- BURRELL, Q. L. (2002), The nth-citation distribution and obsolescence. *Scientometrics*, 53 (3) : 309–323.
- BURRELL, Q. L. (2005), Are “Sleeping Beauties” to be expected? *Scientometrics*, 65 (3) : 381–389.
- EGGHE, L. (2000), A heuristic study of the first-citation distribution. *Scientometrics*, 48 (3) : 345–359.
- EGGHE, L. (2005), *Power Laws in the Information Production Process: Lotkaian Informetrics*. Elsevier, Oxford (UK), 2005.
- EGGHE, L., I. K. RAVICHANDRA RAO (2001), Theory of first-citation distributions and applications. *Mathematical and Computer Modelling*, 34 : 81–90.
- GLÄNZEL, W. (1992), On some stopping times of citation processes. From theory to indicators. *Information Processing and Management*, 28 (1) : 53–60.
- GLÄNZEL, W., B. SCHLEMMER, B. THIJS (2003), Better late than never? On the chance to become highly cited only beyond the standard bibliometric time horizon. *Scientometrics*, 58 (3) : 571–586.
- ROUSSEAU, R. (1994), Double exponential models for first-citation processes. *Scientometrics*, 30 (1) : 213–227.
- VAN RAAN, A. F. J. (2004), Sleeping Beauties in science. *Scientometrics*, 59 (3) : 467–472.

Appendix

Proof of formula (3)

The cumulative distribution C of citations at time $t > 0$ is given by (using (1)):

$$C(t) = \int_0^t c(s)ds$$

$$C(t) = \int_0^t ba^s ds$$

$$C(t) = \frac{b}{\ln a} (a^t - 1) \quad (\text{A1})$$

Since $\lim_{t \rightarrow \infty} C(t)$ must be 1 we have, since $0 < a < 1$,

$$C(t) = 1 - a^t \quad (\text{A2})$$

A paper with M citations in total has $n = 1, 2, 3, \dots$ citations at time t if

$$MC(t) = n \quad (\text{A3})$$

Here we make the (simple) assumption that citations are spread out over papers proportional to their number of citations at $t = \infty$ (deterministic argument).

For all values $M' > M$ we evidently have

$$M'C(t) > n$$

hence these documents belong to the ones that receive their n^{th} citation before t . Their cumulative fraction is, according to (2),

$$\begin{aligned} \int_M^\infty \varphi(M') dM' &= \int_M^\infty \frac{E}{M'M'^\alpha} dM' \\ &= M^{1-\alpha} \end{aligned} \quad (\text{A4})$$

noting the fact that, since φ is a distribution, we have

$$\int_1^\infty \varphi(M') dM' = 1$$

and hence (supposing $\alpha > 1$ which can always be supposed for Lotka's law – see EGGHE (2005))

$$E = \alpha - 1 \quad (\text{A5})$$

Formulae (A2), (A3) and (A5) combined yield

$$\begin{aligned} \int_M^\infty \varphi(M') dM' &= \left(\frac{C(t)}{n} \right)^{\alpha-1} \\ &= \left(\frac{1-a^t}{n} \right)^{\alpha-1} \end{aligned} \quad (\text{A6})$$

But this is nothing but the cumulative n^{th} citation distribution (amongst the ever cited articles, since $M \geq 1$ in (2)), denoted Φ_n . Hence

$$\Phi_n(t) = \left(\frac{1-a^t}{n} \right)^{\alpha-1} \quad (\text{A7})$$

for $n = 1, 2, 3, \dots$.