# Power Laws in the Information Production Process: Lotkaian Informetrics

Leo Egghe

# Power Laws in the Information Production Process: Lotkaian Informetrics

**Leo Egghe**

# Working together to grow
# libraries in developing countries

www.elsevier.com | www.bookaid.org | www.sabre.org

**ELSEVIER**  **BOOK AID** International  Sabre Foundation

# Power Laws in the Information Production Process:
# Lotkaian Informetrics

# Library and Information Science

Series Editor: *Bert R. Boyce*
School of Library & Information Science
Louisiana State University, Baton Rouge

My wife is thanking me for the
many quiet evenings at the time
of the writing of this book.

To Ute, because she asked for it.

# PREFACE

Explaining and hence understanding is one of the key characteristics of human beings. Explaining is making logical, mathematical deductions based on a minimum of unexplained properties, called axioms. Indeed, without axioms, one is not able to make deductions. How these axioms are selected is the only non-explainable part of the theory. In fact, different choices are possible leading to different, in themselves consistent, theories which conflict when considered together. A typical example is the construction of the different types of geometries, Euclidean and non-Euclidean geometries which are in contradiction when considered together but which have their own applications.

In this book the object of study is a two-dimensional information production process, i.e. where one has sources (e.g. journals, authors, words, ...) which produce (or have) items (e.g. respectively articles, publications, words occurring in texts, ...) and in which one considers different functions describing quantitatively the production quantities of the different sources. All functions can be reduced to one type of function, namely the size-frequency function f: such a function gives, for every $n = 1$, 2, 3, ..., the number $f(n)$, being the number of sources with n items. This is the framework of study and in this framework we want to explain as many regularities that one encounters in the literature as possible, using the size-frequency function f.

As explained above, we need at least one axiom. The only axiom used in this book is that the size-frequency function f is Lotkaian, i.e. a power function of the form $f(n)=C/n^{\alpha}$, where $C>0$ and $\alpha^3\ 0$, hence also implying that the function f decreases. The name comes from its introduction into the literature by Alfred Lotka in 1926, see Lotka (1926). Based on this one assumption, one can be surprised about the enormous amounts of regularities that can be explained. In all cases the parameter $\alpha$ turns out to be crucial and is capable of, dependent on the different values that $\alpha$ can take, explaining different shapes of one phenomenon.

In this book we encounter explanations in the following directions: other informetric functions that are equivalent with Lotka's law (e.g. Zipf's law), concentration theory

(theory of inequality), fractal theory, modelling systems in which items can have multiple sources (as is the case in the system: articles written by several authors), modelling citation distributions. These models are developed from Chapter II on.

Although the law of Lotka is an axiom in this book, in Chapter I we investigate what other models are capable of "explaining" Lotka's law. The only purpose for these developments is to understand why Lotka's law is choosen in this book (and not another type of size-frequency function such as, e.g., an exponential function) and hence to make the choice acceptable (although this is not needed, strictly speaking). For the same reason we also give an overview of situations where Lotka's law is encountered, being the majority of the situations (including the "new" situations as networks, including the Internet).

The author is basing himself on the publications that he has written the past 20 years on the subject but also benifits from the work of many others. To them my sincerest thanks. The author is especially grateful to Professor Ronald Rousseau (co-winner of the 2001 Derek De Solla Price Award) with whom he co-authored several papers but with whom he also had numerous long discussions (often by phone) on the different topics described in this book.

We strongly hope this book will serve the informetrics community in the sense that it shows the logical links between many (at first sight unrelated) informetrics aspects. The informetrician, having read this book, can use it each time he/she encounters a new informetric phenomenon (often in the form of a data set) in the sense that one can investigate if the phenomenon shows regularities that are (or can be) explained using the arguments given in this book. The mathematical knowledge required is limited to elementary mathematics such as first-year calculus. Other, more advanced topics, are introduced in this book.

The author is indebted to the Limburgs Universitair Centrum (LUC) and the University of Antwerp (UA) for their support in doing informetric research: in LUC, the author is chief librarian and coordinator of the research project "bibliometrics" while in UA, he is professor in the School of Library and Information Science, where he teaches the courses on informetrics and on information retrieval. The author thanks

Mr. M. Pannekoeke (LUC) for the excellent typing and organization of this manuscript.

Leo Egghe

Diepenbeek, Belgium

Summer 2004

# TABLE OF CONTENTS

# INTRODUCTION

The most facinating aspect of informetrics is the study of what we could call two-dimensional informetrics. In this discipline one considers sources (e.g. journals, authors, ...) and items (being produced by a source - e.g. articles) and their interrelations. By this we mean the description of the link that exists between sources and items. Without the description of this link we would have two times a one dimensional informetrics study, one for the sources and one for the items. Essentially in two-dimensional informetrics the link between sources and items is described by two possible functions: a size-frequency function f and a rank-frequency function g. Although one function can be derived from the other, they are different descriptions of two-dimensionality. A size-frequency function f describes the number $f(n)$ of sources with $n = 1, 2, 3, ...$ items while a rank-frequency function g describes the number $g(r)$ of items in the source on rank $r = 1, 2, 3, ...$ (where the sources are ranked in decreasing order of the number of items they contain (or produce)). So, in essence, in f and g, the role of sources and items are interchanged. This is called duality, hence f and g can be considered as dual functions.

Rank-frequency functions are well-known in the literature, especially in the economics and linguistics literature, where one usually considers Pareto's law and Zipf's law, respectively, being power laws. Less encountered in the literature (except in information sciences) is the size-frequency function. If studied, one supposes in most cases also a power law for such a function, i.e. a function of the type $f(n) = C/n^\alpha$ with $\alpha \geq 0$. Such a function is then called the law of Lotka referring to its introduction in the informetrics literature in 1926, see Lotka (1926). The law of Lotka gives rise to a variety of derived results in informetrics, the description of them being the subject of this book. That we choose a size-frequency function as the main study-object is explained e.g. by its simplicity in formulation (in the discrete setting simpler than a rank-frequency function since the latter uses ranks which have been derived from the "sizes" n but also in the continuous setting, where sizes and ranks are taken in an interval, the formulation of the size-frequency function is more appealing and direct). A size-frequency function also allows for a study of fractional quantities (see Chapter VI), needed e.g. in the description of two-dimensional informetrics in which