

1 **Extreme inequalities of citation counts in environmental sciences**

2 Deepthi Chimalakonda¹, Alex R. Cook^{2,3}, L. Roman Carrasco^{1,*}

3 ¹ Department of Biological Sciences, National University of Singapore, Singapore

4 ² Saw Swee Hock School of Public Health, National University of Singapore, Singapore, Singapore.

5 ³ Yale-NUS College, National University of Singapore, Singapore.

6 Author for correspondence: dbctrl@nus.edu.sg

7

8 **Abstract**

9 Well-established scientists are expected to be more likely to have their work recognised than early-
10 career individuals and thus receive more citations. Estimating the degree of inequality in citation
11 counts in environmental sciences can help identify the dynamics behind citation inequalities.

12 Using the scientific profiles of researchers in the Google Scholar database, we estimated the
13 inequality in the distribution of citations in the disciplines of evolutionary biology, conservation
14 biology and ecology. The data were modelled using short-tailed (exponential) and long-tailed power-
15 law (Pareto) distributions. The inequality in performance in each distribution was assessed using Gini
16 coefficients.

17 Citations counts per researcher presented Gini coefficients of 0.82–0.89, indicating extreme
18 inequality. The results suggest that the reinforcement in citation counts due to seniority and
19 previous success might be very strong. To produce meaningful comparisons of actual research
20 impact using citation counts, factors such as lab size, collaborations or role in articles should ideally
21 be controlled for.

22

23 **Keywords:** h-index; Matthew effect; preferential attachment; research merit; rich get richer; role-
24 based h-index.

25

26 Introduction

27 Citation analysis is a bibliometric method increasingly used to assess the research output of
28 universities, scientists, journals and even countries. Citations are used as basis for evaluating
29 researchers for either positions or tenure, awarding grants and in determining the rank of
30 universities [1-4]. Citation counts and the h-index are popular indicators of scientific merit [5,6]—
31 where the h-index is the measure of the number of the researcher's articles that have at least h
32 citations [1]. Although it has been found that the h-index is a better predictor of future achievement
33 than citation numbers [7], citation counts still remain as one of the most commonly used measures
34 of the research impact of individuals [8].

35 Citation counts may however not fully demonstrate the scientific merit of the researcher [9]: they
36 often include self-citations and negative citations [10,11], might benefit from multiple-author
37 publication practices [12,13] and they are not time-standardised to control for the fact that
38 researchers who have been publishing for a long time tend to have more citations due to the
39 additive effect resulting from the increase in the absolute number of published articles and the
40 number of citations per article.

41 A second group of more intangible factors influencing citations are related to the previous success
42 and seniority of the scientist. In other words, scientists who achieve fame are more likely to have
43 their work recognised, for instance receiving more citations or collaborations, than early-career
44 individuals producing similar quality of work [14,15]. This phenomenon has been termed
45 “preferential attachment” or “the rich get richer effect”. This means that resources are distributed
46 proportionally to what an individual already has [16]. Other factors can also contribute to widen the
47 gap between a few successful scientists and the rest: opportunities to collaborate with other top
48 scientists in landmark articles, attainment of larger grants, having larger teams or attracting better
49 quality students. These factors are especially prominent in environmental related sciences where
50 principal investigators construct large groups and where multiple-author papers are common.

51 Preferential attachment has been widely studied in the field of economics, especially in connection
52 to web links, wealth and population distributions [17]. Power laws tend to explain the pattern in
53 these “rich get richer” models [18], where inequality is quite high. For instance the wealth
54 distribution of the Forbes 400 list of US richest people follows a power law [19]. Similarly the income
55 and wealth distributions in US and UK has been shown to follow a distribution with the high-end tail
56 following a power law [20]. Power laws have also been used to show that the citations of academic
57 articles are proportional to the number of citations that the article already has [18]. Such a process
58 creates wider inequalities in the frequency distribution of the variable of interest. This inequality can
59 be measured by using the Gini co-efficient, initially conceived to measure inequalities in incomes
60 distributions. A Gini co-efficient of zero suggests perfect equality and one indicates maximum
61 inequality [21].

62 Although the combined influence of reinforcement factors due to previous success is expected, it is
63 very hard to estimate their influence on the inequality in citation counts. Here we aim to estimate
64 the degree of inequality in citation counts for several environmental science disciplines using the
65 recently available Google Scholar research profiles data.

66

67 **Methods**

68 Using the scientific profiles of researchers in the Google Scholar database, we extracted the citation
69 counts for all scholars in various disciplines using the labels “ecology”, “conservation biology” and
70 “evolutionary biology”.

71 The data were modelled using short-tailed (exponential) and long-tailed power-law (Pareto)
72 distributions in the statistical environment R using the package VGAM [22]. A preliminary analysis
73 indicated that Pareto type IV distributions obtained the best fit of all the Pareto types and was
74 subsequently used for the analysis. The exponential (1) and Pareto type IV distributions (2) are
75 expressed by the following formulas respectively:

$$76 \quad F(y) = \lambda e^{-\lambda y} \quad (1)$$

$$77 \quad F(y) = 1 - [1 + ((y - \mu) / \sigma)^{1/\gamma}]^{-\alpha} \quad (2)$$

78 Where λ is the rate, μ the location, σ the scale, α the shape and γ the inequality parameters
79 respectively.

80 The models were fit to the data using Markov Chain Monte Carlo methods with a Gibbs sampler. The
81 fit of the distributions was compared using the Akaike Information Criterion.

82 Pareto distributions can be used to characterize income distribution through their association with
83 the Gini coefficient [21]. We employed the following formula to estimate the Gini coefficient of the
84 Pareto IV functions [23]:

$$85 \quad G = 1 - \left(\frac{\mu + 2\sigma\alpha B(2\alpha - \gamma, \gamma + 1)}{\mu + \sigma\alpha B(\alpha - \gamma, \gamma + 1)} \right),$$

86 where B denotes the Beta function.

87

88 **Results**

89 Pareto type IV provided the best fit for all the datasets (Table 1, Figure 1). Citations per researcher in
90 Ecology, Evolutionary Biology and Conservation biology presented Gini coefficients of 0.89, 0.83 and
91 0.82 respectively, indicating a very high inequality. The inequality in citations in different fields of
92 study was also obvious from our results where twenty percent of the scientists had more than eighty
93 percent of all the citations (82% for ecology and conservation biology and 83% for evolutionary
94 biology).

95 The citations within top cited environmental scientists present themselves high inequality. For
96 instance, the most cited scientist in conservation biology has been cited 56,618 times and out of
97 these 22,107 times correspond to a statistical book highly used for model analysis in the field. In a
98 similar fashion, the top most cited scientist in evolutionary biology with 192,602 citations accrues a
99 large share of his citations through the construction of widely used software for evolutionary
100 genetics analysis.

101 Discussion

102 Our results show that there is extreme inequality in the citation counts in environmental sciences.
103 For illustration, the greatest income inequality in the world occurs in South Africa, with a Gini
104 coefficient of 0.7 [24] which is still below the 0.82–0.89 values obtained for citation counts
105 distributions.

106 Extreme inequalities in environmental sciences could be increased due to factors related to the
107 facility to collaborate with other successful scientists and to be invited as co-author in landmark
108 articles or the attraction of the best students and postdoctoral researchers. Other institutional
109 factors could be the promotion process or grant allocations. In addition to these, search engines
110 such as Google Scholar and Web of Science tend to give higher order of appearance to results with
111 high citations, creating a bias for researchers to read and cite these papers. At any rate, teasing out
112 the actual research impact from the reinforcement in inequalities due to exogenous factors would
113 need further research.

114 With views to produce more meaningful comparisons of scientific output, measures to alleviate the
115 extreme inequalities in citation numbers would be needed. For instance to account for the facility to
116 obtain collaborations and co-authorship of established scientists, variants of the h-index have been
117 proposed: h-index methods that take into account the number of co-authors and relative position of
118 the authors while assigning the rank could be employed [25]. Other role-based indices have also
119 been suggested, where single author publications get higher ranking than multiple author ones or
120 where only major contributions such as first and corresponding authors are considered for the
121 estimation of the h-index [26,27]. These methods, however, would still not be able to account for
122 institutional factors such as promotion, lab size, grant allocation and capacity to attract exceptional
123 students.

124

125 Acknowledgements

126 D.C. and L.R.C. are thankful to the Singaporean MOE Tier 1 grant WBS R154000556112.

127

128 References

- 129 1. Hirsch JE (2005) An index to quantify an individual's scientific research output. Proceedings of the
130 National academy of Sciences of the United States of America 102: 16569.
- 131 2. Moed HF, Burger WJM, Frankfort JG, Van Raan AFJ (1985) The use of bibliometric data for the
132 measurement of university research performance. Research Policy 14: 131-149.
- 133 3. Johnes J, Johnes G (1995) Research funding and performance in UK university departments of
134 economics: a frontier analysis. Economics of Education Review 14: 301-314.
- 135 4. King DA (2004) The scientific impact of nations. Nature 430: 311-316.
- 136 5. King J (1987) A review of bibliometric and other science indicators and their role in research
137 evaluation. Journal of information science 13: 261-276.
- 138 6. Radicchi F, Fortunato S, Castellano C (2008) Universality of citation distributions: Toward an
139 objective measure of scientific impact. Proceedings of the National Academy of Sciences
140 105: 17268-17272.

- 141 7. Bornmann L, Daniel H-D (2005) Does the h-index for ranking of scientists really work?
142 *Scientometrics* 65: 391-392.
- 143 8. Adam D (2002) Citation analysis: The counting house. *Nature* 415: 726-729.
- 144 9. Garfield E (1979) Is citation analysis a legitimate evaluation tool? *Scientometrics* 1: 359-375.
- 145 10. Fassoulaki A, Paraskeva A, Papilas K, Karabinis G (2000) Self-citations in six anaesthesia journals
146 and their significance in determining the impact factor. *British Journal of Anaesthesia* 84:
147 266-269.
- 148 11. Fowler J, Aksnes D (2007) Does self-citation pay? *Scientometrics* 72: 427-437.
- 149 12. Persson O, Glänzel W, Danell R (2004) Inflationary bibliometric values: The role of scientific
150 collaboration and the need for relative indicators in evaluative studies. *Scientometrics* 60:
151 421-432.
- 152 13. Hsu J-w, Huang D-w (2011) Correlation between impact and collaboration. *Scientometrics* 86:
153 317-324.
- 154 14. Merton RK (1968) The Matthew effect in science. *Science* 159: 56-63.
- 155 15. Eom Y-H, Fortunato S (2011) Characterizing and modeling citation dynamics. *PLoS one* 6: e24926.
- 156 16. Barabási A-L, Albert R (1999) Emergence of scaling in random networks. *Science* 286: 509-512.
- 157 17. Jeong H, Néda Z, Barabási A-L (2003) Measuring preferential attachment in evolving networks.
158 *EPL (Europhysics Letters)* 61: 567.
- 159 18. Newman ME (2005) Power laws, Pareto distributions and Zipf's law. *Contemporary physics* 46:
160 323-351.
- 161 19. Levy M, Solomon S (1997) New evidence for the power-law distribution of wealth. *Physica A:
162 Statistical Mechanics and its Applications* 242: 90-94.
- 163 20. Drăgulescu A, Yakovenko VM (2001) Exponential and power-law probability distributions of
164 wealth and income in the United Kingdom and the United States. *Physica A: Statistical
165 Mechanics and its Applications* 299: 213-221.
- 166 21. Gini C (1921) Measurement of inequality of incomes. *The Economic Journal* 31: 124-126.
- 167 22. R Core Development Team (2010) R: A Language and Environment for Statistical Computing.
168 Vienna, Austria: R Foundation for Statistical Computing.
- 169 23. Chotikapanich D (2008) Modeling income distributions and Lorenz curves [electronic resource]:
170 Springer.
- 171 24. The World Bank (2011) South Africa Overview. The World Bank. Accessed on 12/9/2013 at:
172 <http://www.worldbank.org/en/country/southafrica/overview>.
- 173 25. Wan J-k, Hua P-h, Rousseau R (2007) The pure h-index: calculating an author's h-index by taking
174 co-authors into account. *COLLNET Journal of Scientometrics and Information Management*
175 1: 1-5.
- 176 26. Hu X, Rousseau R, Chen J (2010) In those fields where multiple authorship is the rule, the h-index
177 should be supplemented by role-based h-indices. *Journal of Information Science* 36: 73-85.
- 178 27. Schreiber M (2009) A case study of the modified Hirsch index hm accounting for multiple
179 coauthors. *Journal of the American Society for Information Science and Technology* 60:
180 1274-1282.

181

182

183

184 **Tables and figures captions**

185 Table 1. AIC values for exponential and Pareto IV models fitted to the citation counts in
186 environmental sciences.

Field	Model	AIC
Evolutionary Biology	Pareto IV	20075
	Exponential	21871
Ecology	Pareto IV	43531
	Exponential	47112
Conservation Biology	Pareto IV	33143
	Exponential	35879

187

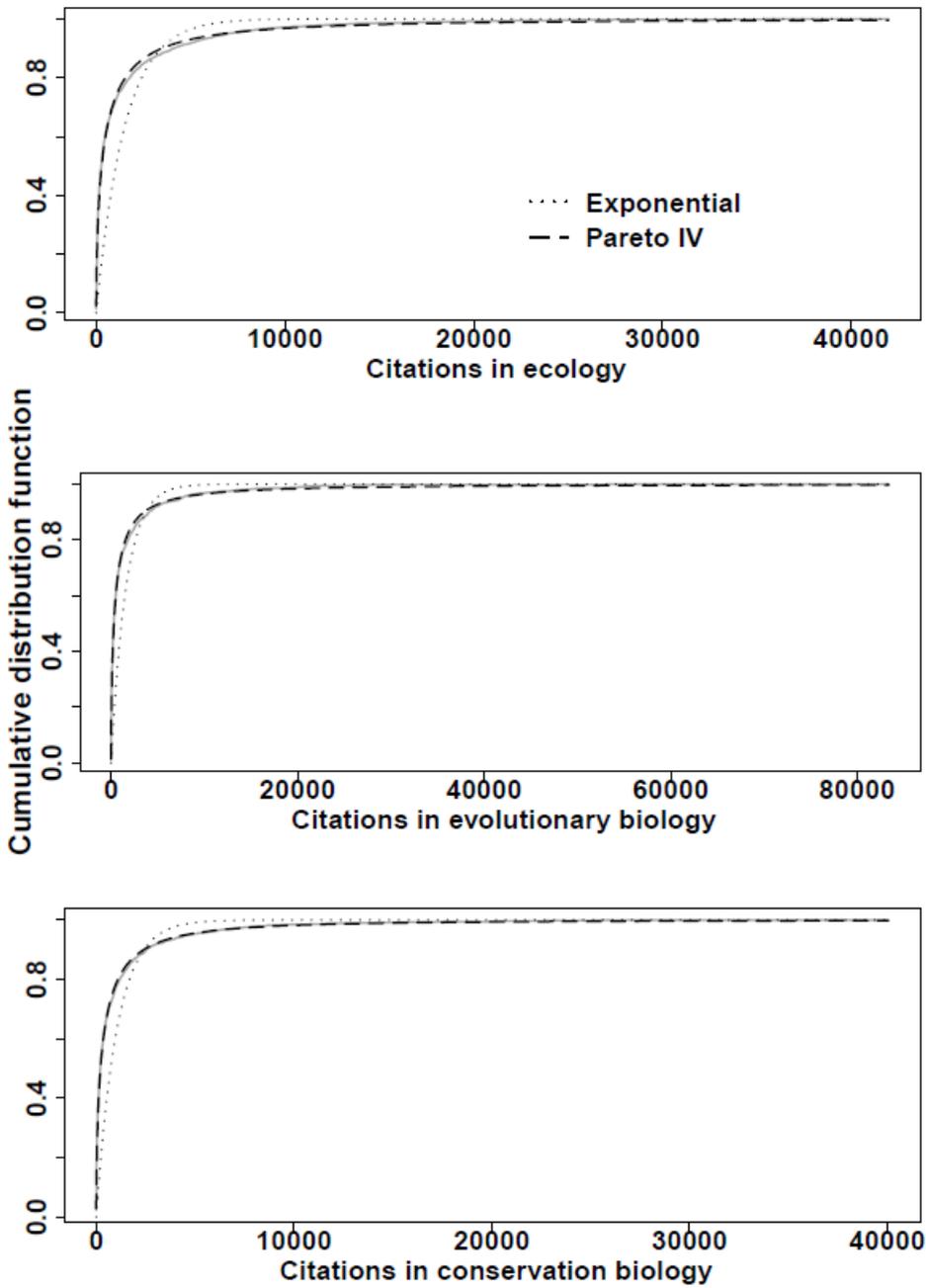
188 Table 2. Parameter values and Gini coefficients of the Pareto IV models fitted to the data.

Field	Location	Scale	Inequality	Shape	Gini coefficient
Evolutionary Biology	0	6.79	1.43	1.78	0.89
Ecology	0	7.80	1.65	2.86	0.83
Conservation Biology	0	7.52	1.62	2.96	0.82

189

190

191 Figure 1. Comparison of Pareto IV and exponential fits to citations in ecology, conservation biology
192 and evolutionary ecology.



193

194