

J. R. Statist. Soc. A (2016)
179, Part 1, pp. 1–63

Statistical modelling of citation exchange between statistics journals

Cristiano Varin,

Università Ca' Foscari, Venezia, Italy

Manuela Cattelan

Università degli Studi di Padova, Italy

and David Firth

University of Warwick, Coventry, UK

[*Read before The Royal Statistical Society at a meeting organized by the General Applications Section on Wednesday, May 13th, 2015, Professor P. Clarke in the Chair*]

Summary. Rankings of scholarly journals based on citation data are often met with scepticism by the scientific community. Part of the scepticism is due to disparity between the common perception of journals' prestige and their ranking based on citation counts. A more serious concern is the inappropriate use of journal rankings to evaluate the scientific influence of researchers. The paper focuses on analysis of the table of cross-citations among a selection of statistics journals. Data are collected from the *Web of Science* database published by Thomson Reuters. Our results suggest that modelling the exchange of citations between journals is useful to highlight the most prestigious journals, but also that journal citation data are characterized by considerable heterogeneity, which needs to be properly summarized. Inferential conclusions require care to avoid potential overinterpretation of insignificant differences between journal ratings. Comparison with published ratings of institutions from the UK's research assessment exercise shows strong correlation at aggregate level between assessed research quality and journal citation 'export scores' within the discipline of statistics.

Keywords: Bradley–Terry model; Citation data; Export score; Impact factor; Journal ranking; Research evaluation; Stigler model

1. Introduction

The problem of ranking scholarly journals has arisen partly as an economic matter. When the number of scientific journals started to increase, librarians were faced with decisions about which journal subscriptions should consume their limited economic resources; a natural response was to be guided by the relative importance of different journals according to a published or otherwise agreed ranking. Gross and Gross (1927) proposed the counting of citations received by journals as a direct measure of their importance. Garfield (1955) suggested that the number of citations received should be normalized by the number of citable items published by a journal. This idea is at the origin of the *impact factor*, which is the best-known index for ranking journals. Published since the 1960s, the impact factor is 'an average citation rate per published article' (Garfield, 1972).

Address for correspondence: Cristiano Varin, Department of Environmental Sciences, Informatics and Statistics, Ca' Foscari University, Zeta Building, Via Torino 155, 30170 Mestre, Italy.
E-mail: cristiano.varin@unive.it

© 2015 The Authors Journal of the Royal Statistical Society: Series A (Statistics in Society) 0964–1998/16/179001
Published by John Wiley & Sons Ltd on behalf of the Royal Statistical Society.
This is an open access article under the terms of the Creative Commons Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

The impact factor of the journals where scholars publish has also been employed—improperly, many might argue—in appointing to academic positions, in awarding research grants and in ranking universities and their departments. The ‘San Francisco declaration on research assessment’ (<http://am.ascb.org/dora>, 2013) and the Institute of Electrical and Electronics Engineers position statement on ‘Appropriate use of bibliometric indicators for the assessment of journals, research proposals, and individuals’ (Institute of Electrical and Electronics Engineers Board of Directors, 2013) are just two of the most recent authoritative standpoints regarding the risks of automatic, metric-based evaluations of scholars. Typically, only a small fraction of all published articles accounts for most of the citations that are received by a journal (Seglen, 1997). Single authors should ideally be evaluated on the basis of their own outputs and not through citations of other papers that have appeared in the journals where their papers have been published (Seglen, 1997; Adler *et al.*, 2009; Silverman, 2009). As stated in a recent *Science* editorial about impact factor distortions (Alberts, 2013),

‘... the leaders of the scientific enterprise must accept full responsibility for thoughtfully analyzing the scientific contributions of other researchers. To do so in a meaningful way requires the actual reading of a small selected set of each researcher’s publications, a task that must not be passed by default to journal editors’.

Indicators derived from citations received by papers written by a particular author (e.g. Bornmann and Marx (2014)) can be a useful complement for evaluation of trends and patterns of that author’s impact, but not a substitute for the reading of papers.

Journal rankings based on the impact factor often differ substantially from common perceptions of journal prestige (Theoharakis and Skordia, 2003; Arnold and Fowler, 2011). Various causes of such discrepancy have been pointed out. First, there is the phenomenon that more ‘applied’ journals tend to receive citations from other scientific fields more often than do journals that publish theoretical work. This may be related to uncounted ‘indirect citations’ arising when methodology that is developed in a theoretical journal is then popularized by papers published in applied journals accessible to a wider audience and thus receiving more citations than the original source (Journal-Ranking.com, 2007; Putirka *et al.*, 2013). Second is the short time period that is used for computation of the impact factor, which can be completely inappropriate for some fields, in particular for mathematics and statistics (van Nierop, 2009; Arnold and Fowler, 2011). Finally, there is the risk of manipulation, whereby authors might be asked by journal editors to add irrelevant citations to other papers published in their journal (Sevinc, 2004; Frandsen, 2007; Archambault and Larivière, 2009; Arnold and Fowler, 2011). According to a large survey published in *Science* (Wilhite and Fong, 2012), about 20% of academics in social science and business fields have been asked to ‘pad their papers with superfluous references to get published’ (van Noorden, 2012). The survey data also suggest that junior faculty members are more likely to be pressured to cite superfluous papers. Recently, Thomson Reuters has started to publish the impact factor both with and without journal self-citations, thereby allowing evaluation of the contribution of self-citations to the impact factor calculation. Moreover, Thomson Reuters has occasionally excluded journals with an excessive self-citation rate from the ‘Journal citation reports’ (JCRs).

Given these criticisms, it is not surprising that the impact factor and other ‘quantitative’ journal rankings have given rise to substantial scepticism about the value of citation data. Several proposals have been developed in the bibliometric literature to overcome the weaknesses of the impact factor; examples include the *article influence score* (Bergstrom, 2007; West, 2010), the *H-index* for journals (Braun *et al.*, 2006; Pratelli *et al.*, 2012), the *source-normalized impact per paper* index (Waltman *et al.*, 2013) and methods based on percentile rank classes (Marx and Bornmann, 2013).

The aforementioned *Science* editorial (Alberts, 2013) reports that

‘... in some nations, publication in a journal with an impact factor below 5.0 is officially of zero value’.

In the latest edition (2013) of the JCR, the only journal with an impact factor larger than 5 in the category ‘Statistics and probability’ was the *Journal of the Royal Statistical Society, Series B*, with impact factor 5.721. The category ‘Mathematics’ achieved still lower impact factors, with the highest value there in 2013 being 3.08 for *Communications on Pure and Applied Mathematics*. Several bibliometric indicators have been developed, or adjusted, to allow for cross-field comparisons, e.g. Leydesdorff *et al.* (2013) and Waltman and Van Eck (2013), and could be considered to alleviate unfair comparisons. However, our opinion is that comparisons between different research fields will rarely make sense, and that such comparisons should be avoided. Research fields differ very widely, e.g. in terms of the frequency of publication, the typical number of authors per paper and the typical number of citations made in a paper, as well as in the sizes of their research communities. Journal homogeneity is a minimal prerequisite for a meaningful statistical analysis of citation data (Lehmann *et al.*, 2009).

Journal citation data are unavoidably characterized by substantial variability (e.g. Amin and Mabe (2000)). A clear illustration of this variability, suggested by the Associate Editor for this paper, comes from an early editorial of *Briefings in Bioinformatics* (Bishop and Bird, 2007) announcing that this new journal had received an impact factor of 24.37. However, the editors noted that a very large fraction of the journal’s citations came from a single paper; if that paper were to be dropped, then the journal’s impact factor would decrease to about 4. The variability of the impact factor is inherently related to the heavy-tailed distribution of citation counts. Averaged indicators such as the impact factor are clearly unsuitable for summarizing highly skew distributions. Nevertheless, quantification of uncertainty is typically lacking in published rankings of journals. A recent exception is Chen *et al.* (2014) who employed a bootstrap estimator for the variability of journal impact factors. Also the source-normalized impact per paper indicator that was published by Leiden University’s Centre for Science and Technology Studies based on the Elsevier Scopus database, and available on line at www.journalindicators.com, is accompanied by a ‘stability interval’ computed via a bootstrap method. See also Hall and Miller (2009, 2010) and references therein for more details on statistical assessment of the authority of rankings.

The impact factor was developed to identify which journals have the greatest influence on subsequent research. The other metrics that are mentioned in this paper originated as possible improvements on the impact factor, with the same aim. Palacios-Huerta and Volij (2004) listed a set of properties that a ranking method which measures the intellectual influence of journals, by using citation counts, should satisfy. However, the list of all desirable features of a ranking method might reasonably be extended to include features other than citations, depending on the purpose of the ranking. For example, when librarians decide which journals to take, they should consider the influence of a journal in one or more research fields, but they may also take into account its cost effectiveness. The Web site www.journalprices.com, which is maintained by Professor Ted Bergstrom and Professor Preston McAfee, ranks journals according to their price per article, price per citation and a composite index.

A researcher when deciding where to submit a paper most probably considers each candidate journal’s record of publishing papers on similar topics, and the importance of the journal in the research field; but he or she may also consider the speed of the reviewing process, the typical time between acceptance and publication of the paper, possible page charges, and the likely effect on his or her own career. Certain institutions and national evaluation agencies publish rankings of journals which are used to evaluate researcher performance and to inform the hiring of new faculty members. For various economics and management-related disciplines,

the ‘*Journal quality list*’, which is compiled by Professor Anne-Wil Harzing and is available at www.harzing.com/jql.htm, combines more than 20 different rankings made by universities or evaluation agencies in various countries. Such rankings typically are based on bibliometric indices, expert surveys or a mix of both.

Modern technologies have fostered the rise of alternative metrics such as ‘webometrics’ based on citations on the Internet or numbers of downloads of articles. Recently, interest has moved from Web citation analysis to social media usage analysis. In some disciplines the focus is now towards broader measurement of research impact through the use of Web-based quantities such as citations in social media sites, newspapers, government policy documents and blogs. This is mainly implemented at the level of individual articles (see for example the Altmetric service (Adie and Roe, 2013) which is available at www.altmetric.com), but the analysis may also be made at journal level. Along with the advantages of timeliness, availability of data and consideration of different sources, such measures also have certain drawbacks related to data quality, possible bias and data manipulation (Bornmann, 2014).

A primary purpose of the present paper is to illustrate the risks of overinterpretation of insignificant differences between journal ratings. In particular, we focus on the analysis of the exchange of citations between a relatively homogeneous list of journals. Following Stigler (1994), we model the table of cross-citations between journals in the same field by using a Bradley–Terry model (Bradley and Terry, 1952) and thereby derive a ranking of the journals’ ability to ‘export intellectual influence’ (Stigler, 1994). Although the Stigler approach has desirable properties and is sufficiently simple to be promoted also outside the statistics community, there have been rather few published examples of application of this model since its first appearance; Stigler *et al.* (1995) and Liner and Amin (2004) are two notable examples of its application to the journals of economics.

We pay particular attention to methods that summarize the uncertainty in a ranking produced through the Stigler model-based approach. Our focus on properly accounting for ‘model-based uncertainty in making comparisons’ is close in spirit to Goldstein and Spiegelhalter (1996). We propose to fit the Stigler model with the quasi-likelihood method (Wedderburn, 1974) to account for interdependence between the citations exchanged between pairs of journals, and to summarize estimation uncertainty by using quasi-variances (Firth and de Menezes, 2005). We also suggest the use of the ranking lasso penalty (Masarotto and Varin, 2012) when fitting the Stigler model, to combine the benefits of shrinkage with an enhanced interpretation arising from automatic presentational grouping of journals with similar merits.

The paper is organized as follows. Section 2 describes the data collected from the *Web of Science* database compiled by Thomson Reuters; then, as preliminary background to the paper’s main content on journal rankings, Section 3 illustrates the use of cluster analysis to identify groups of statistics journals sharing similar aims and types of content. Section 4 provides a brief summary of journal rankings published by Thomson Reuters in the JCRs. Section 5 discusses the Stigler method and applies it to the table of cross-citations between statistics journals. Section 6 compares journal ratings based on citation data with results from the UK research assessment exercise, and Section 7 collects some concluding remarks.

The citation data set and the computer code used for the analyses written in the R language (R Core Team, 2015) are available from

<http://wileyonlinelibrary.com/journal/rss-datasets>

2. The *Web of Science* database

The database that was used for our analyses is the 2010 edition of the *Web of Science* that was produced by Thomson Reuters. The citation data contained in the database are employed to

compile the JCRs, whose science edition summarizes citation exchange between more than 8000 journals in science and technology. Within the JCR, scholarly journals are grouped into 171 overlapping subject categories. In particular, in 2010 the ‘Statistics and probability’ category comprised 110 journals. The choice of the journals that are encompassed in this category is to some extent arbitrary. The Scopus database, which is the main commercial competitor of the *Web of Science*, in 2010 included in its statistics and probability category 105 journals, but only about two-thirds of them were classified in the same category within the *Web of Science*. The statistics and probability category contains also journals related to fields such as econometrics, chemistry, computational biology, engineering and psychometrics.

A severe criticism of the impact factor relates to the time period that is used for its calculation. The standard version of the impact factor considers citations received to articles published in the previous 2 years. This period is too short to reach the peak of citations of an article, especially in mathematical disciplines (Hall, 2009). van Nierop (2009) found that articles published in statistics journals typically reach the peak of their citations more than 3 years after publication; as reported by the JCR, the median age of the articles cited in this category is more than 10 years. Thomson Reuters acknowledges this issue and computes a second version of the impact factor using citations to papers published in the previous 5 years. Recent published alternatives to the impact factor, to be discussed in Section 4, also count citations to articles that appeared in the previous 5 years. The present paper considers citations of articles published in the previous 10 years, to capture the influence, over a more substantial period, of work published in statistical journals.

A key requirement for the methods that are described here, as well as in our view for any sensible analysis of citation data, is that the journals jointly analysed should be as homogeneous as possible. Accordingly, analyses are conducted on a subset of the journals from the statistics and probability category, among which there is a relatively high level of citation exchange. The selection is obtained by discarding journals in probability, econometrics, computational biology, chemometrics and engineering, and other journals that are not sufficiently related to the majority of the journals in the selection. Furthermore, journals recently established, and thus lacking a record of 10 years of citable items, also are dropped. The final selection consists of the 47 journals that are listed in Table 1. Obviously, the methods that are discussed in this paper can be similarly applied to other selections motivated by different purposes. For example, a statistician who is interested in applications to economics might consider a different selection with journals of econometrics and statistical methodology, discarding instead journals oriented towards biomedical applications.

The JCR database supplies detailed information about the citations that are exchanged between pairs of journals through the *cited journal table* and the *citing journal table*. The cited journal table for journal i contains the number of times that articles published in journal j during 2010 cite articles published in journal i in previous years. Similarly, the citing journal table for journal i contains the number of times that articles published in journal j in previous years were cited in journal i during 2010. Both of the tables contain some very modest loss of information. In fact, all journals that cite journal i are listed in the cited journal table for journal i only if the number of citing journals is less than 25. Otherwise, the cited journal table reports only those journals that cite journal i at least twice in *all past years*, thus counting also citations to papers that were published earlier than the decade 2001–2010 considered here. Remaining journals that cite journal i only once in all past years are collected in the category ‘all others’. Information on journals cited only once is similarly treated in the citing journal table.

Cited and citing journal tables allow construction of the cross-citation matrix $\mathbf{C} = (c_{ij})$, where c_{ij} is the number of citations from articles published in journal j in 2010 to papers published

Table 1. List of selected statistics journals, with abbreviations used in the paper

<i>Journal name</i>	<i>Abbreviation</i>
<i>American Statistician</i>	AmS
<i>Annals of the Institute of Statistical Mathematics</i>	AISM
<i>Annals of Statistics</i>	AoS
<i>Australian and New Zealand Journal of Statistics</i>	ANZS
<i>Bernoulli</i>	Bern
<i>Biometrical Journal</i>	BioJ
<i>Biometrics</i>	Bcs
<i>Biometrika</i>	Bka
<i>Biostatistics</i>	Biost
<i>Canadian Journal of Statistics</i>	CJS
<i>Communications in Statistics—Simulation and Computation</i>	CSSC
<i>Communications in Statistics—Theory and Methods</i>	CSTM
<i>Computational Statistics</i>	CmpSt
<i>Computational Statistics and Data Analysis</i>	CSDA
<i>Environmental and Ecological Statistics</i>	EES
<i>Environmetrics</i>	Envr
<i>International Statistical Review</i>	ISR
<i>Journal of Agricultural, Biological and Environmental Statistics</i>	JABES
<i>Journal of the American Statistical Association</i>	JASA
<i>Journal of Applied Statistics</i>	JAS
<i>Journal of Biopharmaceutical Statistics</i>	JBS
<i>Journal of Computational and Graphical Statistics</i>	JCGS
<i>Journal of Multivariate Analysis</i>	JMA
<i>Journal of Nonparametric Statistics</i>	JNS
<i>Journal of the Royal Statistical Society, Series A</i>	JRSS-A
<i>Journal of the Royal Statistical Society, Series B</i>	JRSS-B
<i>Journal of the Royal Statistical Society, Series C</i>	JRSS-C
<i>Journal of Statistical Computation and Simulation</i>	JSCS
<i>Journal of Statistical Planning and Inference</i>	JSPI
<i>Journal of Statistical Software</i>	JSS
<i>Journal of Time Series Analysis</i>	JTSA
<i>Lifetime Data Analysis</i>	LDA
<i>Metrika</i>	Mtka
<i>Scandinavian Journal of Statistics</i>	SJS
<i>Stata Journal</i>	StataJ
<i>Statistical Methods in Medical Research</i>	SMMR
<i>Statistical Modelling</i>	StMod
<i>Statistica Neerlandica</i>	StNee
<i>Statistical Papers</i>	StPap
<i>Statistical Science</i>	StSci
<i>Statistica Sinica</i>	StSin
<i>Statistics</i>	Stats
<i>Statistics and Computing</i>	StCmp
<i>Statistics in Medicine</i>	StMed
<i>Statistics and Probability Letters</i>	SPL
<i>Technometrics</i>	Tech
<i>Test</i>	Test

in journal i in the chosen time window ($i = 1, \dots, n$). In our analyses, $n = 47$, the number of selected statistics journals, and the time window is the previous 10 years. In the rest of this section we provide summary information about citations made and received by each statistics journal at aggregate level, whereas Sections 3 and 5 discuss statistical analyses derived from citations exchanged by pairs of journals.

Table 2 shows the citations made by papers published in each statistics journal in 2010 to papers published in other journals in the decade 2001–2010, as well as the citations that the papers published in each statistics journal in 2001–2010 received from papers published in other journals in 2010. The same information is visualized in the bar plots of Fig. 1. Citations made and received are classified into three categories, namely journal self-citations from a paper published in a journal to another paper in the same journal, citations to or from journals in the list of selected statistics journals and citations to or from journals not in the selection.

The total numbers of citations reported in the second and fifth columns of Table 2 include citations given or received by all journals included in the *Web of Science* database, not only those in the field of statistics. The totals are influenced by journals' sizes and by the citation patterns of other categories to which journals are related. The number of references to articles published in 2001–2010 ranges from 275 for citations made in *Statistical Modelling*, which has a small size publishing around 350–400 pages per year, to 4022 for *Statistics in Medicine*, which is a large journal with size ranging from 3500 to 6000 pages annually in the period examined. The number of citations from a journal to articles in the same journal is quite variable and ranges from 0.8% of all citations for *Computational Statistics* to 24% for *Stata Journal*. On average, 6% of the references in a journal are to articles appearing in the same journal and 40% of references are addressed to journals in the list, including journal self-citations. The *Journal of the Royal Statistical Society, Series A*, has the lowest percentage of citations to other journals in the list, at only 10%. Had we kept the whole 'Statistics and probability' category of the JCR, that percentage would have risen, by just 2 points to 12%; most of the references appearing in the *Journal of the Royal Statistical Society, Series A*, are to journals outside the statistics and probability category.

The number of citations received ranges from 168 for *Computational Statistics* to 6602 for *Statistics in Medicine*. Clearly, the numbers are influenced by the size of the journal. For example, the small number of citations received by *Computational Statistics* relates to only around 700 pages published per year by that journal. The citations received are influenced also by the citation patterns of other subject categories. In particular, the number of citations that are received by a journal oriented towards medical applications benefits from communication with a large field including many high impact journals. For example, around 75% of the citations received by *Statistics in Medicine* came from journals outside the list of statistics journals, mostly from medical journals. On average, 7% of the citations received by journals in the list came from the same journal and 40% were from journals in the list.

As stated already, the statistics journals on which we focus have been selected from the statistics and probability category of the JCR, with the aim of retaining those which communicate more. The median fraction of citations from journals discarded from our selection to journals in the list is only 4%, whereas the median fraction of citations received by non-selected journals from journals in the list is 7%. An important example of an excluded journal is *Econometrica*, which was ranked in leading positions by all the published citation indices. *Econometrica* had only about 2% of its references addressed to other journals in our list, and received only 5% of its citations from journals within our list.

3. Clustering journals

Statistics journals have different stated objectives, and different types of content. Some journals emphasize applications and modelling, whereas others focus on theoretical and mathematical developments, or deal with computational and algorithmic aspects of statistical analysis. Applied journals are often targeted to particular areas, such as statistics for medical applications, or

Table 2. Citations made, Citing, and received, Cited, in 2010 to or from articles published in 2001–2010†

<i>Journal</i>	<i>Citing</i>			<i>Cited</i>		
	<i>Total</i>	<i>Self</i>	<i>Stat</i>	<i>Total</i>	<i>Self</i>	<i>Stat</i>
AmS	380	0.11	0.43	648	0.07	0.29
AIMS	459	0.04	0.36	350	0.05	0.57
AoS	1663	0.17	0.48	3335	0.09	0.47
ANZS	284	0.02	0.35	270	0.02	0.34
Bern	692	0.03	0.29	615	0.04	0.39
BioJ	845	0.07	0.50	664	0.08	0.42
Bcs	1606	0.12	0.49	2669	0.07	0.45
Bka	872	0.09	0.57	1713	0.04	0.60
Biost	874	0.06	0.41	1948	0.03	0.22
CJS	419	0.04	0.51	362	0.04	0.60
CSSC	966	0.03	0.43	344	0.08	0.48
CSTM	1580	0.06	0.41	718	0.13	0.59
CmpSt	371	0.01	0.33	168	0.02	0.38
CSDA	3820	0.13	0.45	2891	0.17	0.40
EES	399	0.10	0.34	382	0.10	0.23
Envr	657	0.05	0.27	505	0.06	0.27
ISR	377	0.05	0.21	295	0.07	0.32
JABES	456	0.04	0.26	300	0.05	0.27
JASA	2434	0.10	0.41	4389	0.05	0.44
JAS	1248	0.03	0.31	436	0.08	0.33
JBS	1132	0.09	0.33	605	0.16	0.33
JCGS	697	0.06	0.44	870	0.05	0.43
JMA	2167	0.09	0.49	1225	0.15	0.52
JNS	562	0.03	0.52	237	0.07	0.65
JRSS-A	852	0.05	0.15	716	0.05	0.24
JRSS-B	506	0.11	0.51	2554	0.02	0.42
JRSS-C	731	0.02	0.30	479	0.03	0.34
JSCS	736	0.04	0.43	374	0.09	0.45
JSPI	3019	0.08	0.44	1756	0.13	0.54
JSS	1361	0.07	0.21	1001	0.09	0.17
JTSA	327	0.08	0.32	356	0.07	0.41
LDA	334	0.06	0.57	247	0.09	0.59
Mtka	297	0.07	0.56	264	0.08	0.59
SJS	493	0.02	0.50	562	0.02	0.60
StataJ	316	0.24	0.36	977	0.08	0.11
SMMR	746	0.04	0.33	813	0.03	0.18
StMod	275	0.03	0.41	237	0.03	0.35
StNee	325	0.01	0.24	191	0.02	0.31
StPap	518	0.03	0.35	193	0.08	0.42
StSci	1454	0.03	0.29	924	0.05	0.35
StSin	1070	0.04	0.57	935	0.05	0.54
Stats	311	0.02	0.47	254	0.02	0.43
StCmp	575	0.04	0.46	710	0.03	0.24
StMed	4022	0.16	0.42	6602	0.10	0.24
SPL	1828	0.08	0.36	1348	0.11	0.46
Tech	494	0.09	0.37	688	0.06	0.38
Test	498	0.01	0.61	243	0.03	0.54

†Columns are total citations, Total, proportion of citations that are journal self-citations, Self, and proportion of citations that are to or from statistics journals, Stat, including journal self-citations. Journal abbreviations are as in Table 1.

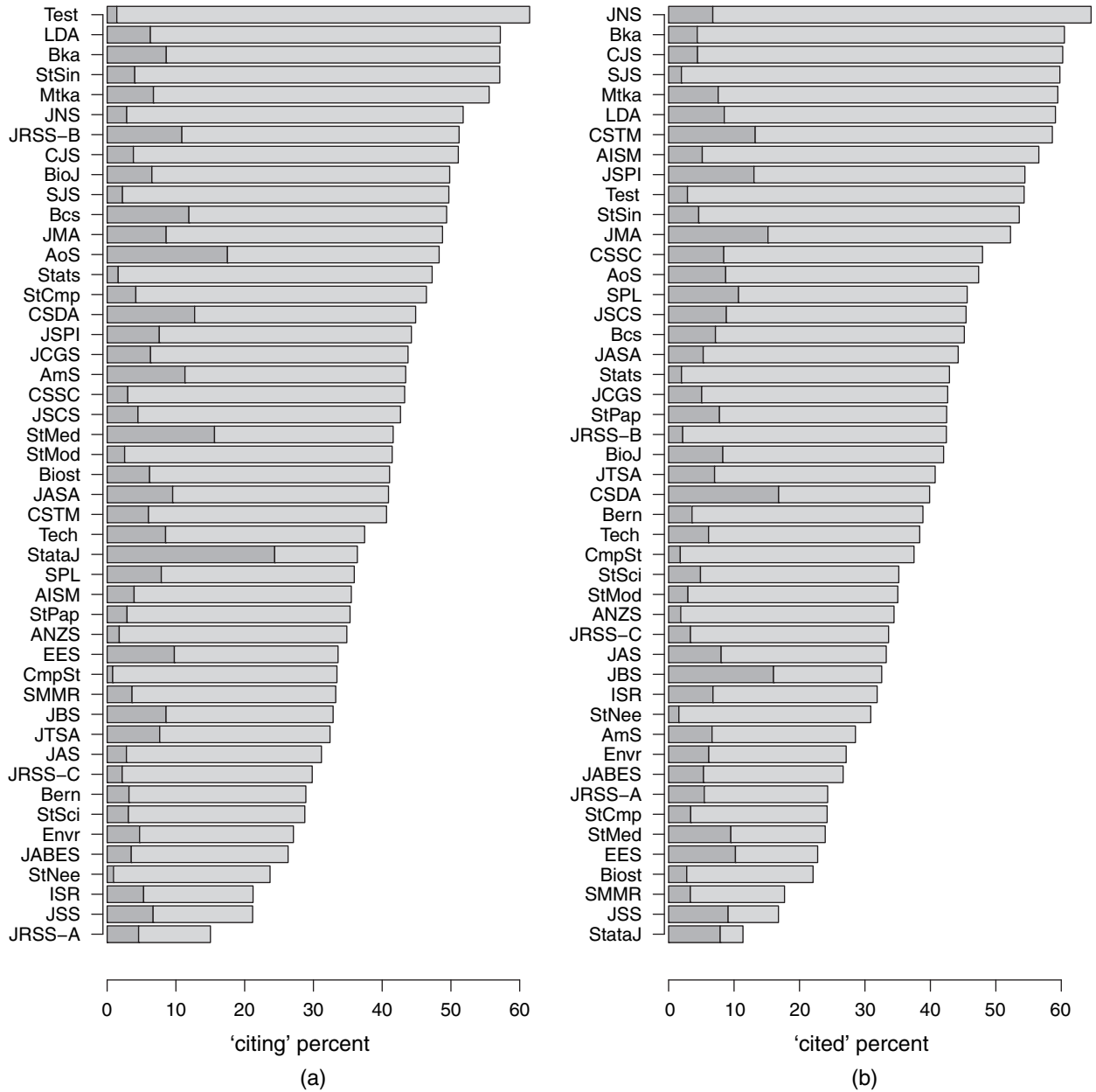


Fig. 1. Bar plots of (a) citations made and (b) citations received for the statistics journals selected, as listed in Table 2 based on the 2010 JCR: for each journal, the bar displays the percentage of self-citations (■) and the percentage of citations made or received to or from other statistics journals in the list (□)

for environmental sciences. Therefore, it is quite natural to consider whether the cross-citation matrix \mathbf{C} allows the identification of groups of journals with similar aims and types of content. Clustering of scholarly journals has been extensively discussed in the bibliometric literature and a variety of clustering methods have been considered. Examples include the hill climbing method (Carpenter and Narin, 1973), k -means (Boyack *et al.*, 2005) and methods based on graph theory (Leydesdorff, 2004; Liu *et al.*, 2012).

Consider the total number t_{ij} of citations exchanged between journals i and j ,

$$t_{ij} = \begin{cases} c_{ij} + c_{ji}, & \text{for } i \neq j, \\ c_{ii}, & \text{for } i = j. \end{cases} \quad (1)$$

Among various possibilities—see, for example, Boyack *et al.* (2005)—the distance between two journals can be measured by quantity $d_{ij} = 1 - \rho_{ij}$, where ρ_{ij} is the Pearson correlation coefficient of variables t_{ik} and t_{jk} ($k = 1, \dots, n$), i.e.

$$\rho_{ij} = \frac{\sum_{k=1}^n (t_{ik} - \bar{t}_i)(t_{jk} - \bar{t}_j)}{\sqrt{\left\{ \sum_{k=1}^n (t_{ik} - \bar{t}_i)^2 \sum_{k=1}^n (t_{jk} - \bar{t}_j)^2 \right\}}},$$

with $\bar{t}_i = \sum_{k=1}^n t_{ik}/n$. Among the many available clustering algorithms, we consider a hierarchical agglomerative cluster analysis with complete linkage (Kaufman and Rousseeuw, 1990). The clustering process is visualized through the dendrogram in Fig. 2. Visual inspection of the dendrogram suggests cutting it at distance 0.6, thereby obtaining eight clusters, two of which are singletons. The clusters identified are grouped in brackets in Fig. 2.

We comment first on the groups and later on the singletons, following the order of the journals in Fig. 2. The first group, (1), includes a large number of general journals concerned with theory and methods of statistics, but also with applications. Among others, the group includes the *Journal of Time Series Analysis*, the *Journal of Statistical Planning and Inference* and *Annals of the Institute of Statistical Mathematics*.

The second group, (2), contains the leading journals in the development of statistical theory and methods: *Annals of Statistics*, *Biometrika*, the *Journal of the American Statistical Association* and the *Journal of the Royal Statistical Society*, Series B. The group includes also other methodological journals such as *Bernoulli*, the *Scandinavian Journal of Statistics* and *Statistica Sinica*. It is possible to identify some natural subgroups: the *Journal of Computational and Graphical Statistics* and *Statistics and Computing*; *Biometrika*, the *Journal of the Royal Statistical Society*, Series B, and the *Journal of the American Statistical Association*; *Annals of Statistics* and *Statistica Sinica*.

The third group, (3), comprises journals mostly dealing with computational aspects of statistics, such as *Computational Statistics and Data Analysis*, *Communications in Statistics—Simulation and Computation*, *Computational Statistics* and the *Journal of Statistical Computation and Simulation*. Other members of the group with a less direct orientation towards computational methods are *Technometrics* and the *Journal of Applied Statistics*.

The fourth group, (4), includes just two journals both of which publish mainly review articles, namely the *American Statistician* and the *International Statistical Review*.

The fifth group, (5), comprises the three journals specializing in ecological and environmental applications: the *Journal of Agricultural, Biological and Environmental Statistics*, *Environmental and Ecological Statistics* and *Environmetrics*.

The last group, (6), includes various journals emphasizing applications, especially to health sciences and similar areas. It encompasses journals oriented towards biological and medical applications such as *Biometrics* and *Statistics in Medicine*, and also journals publishing papers about more general statistical applications, such as the *Journal of the Royal Statistical Society*, Series A and C. The review journal *Statistical Science* also falls into this group; it is not grouped together with the other two review journals already mentioned. Within the group there are some natural subgroupings: *Statistics in Medicine* with *Statistical Methods in Medical Research*; and *Biometrics* with *Biostatistics*.

Finally, and perhaps not surprisingly, the two singletons are the software-oriented *Journal of Statistical Software* and *Stata Journal*. The latter is, by some distance, the most remote journal in the list according to the measure of distance that is used here.

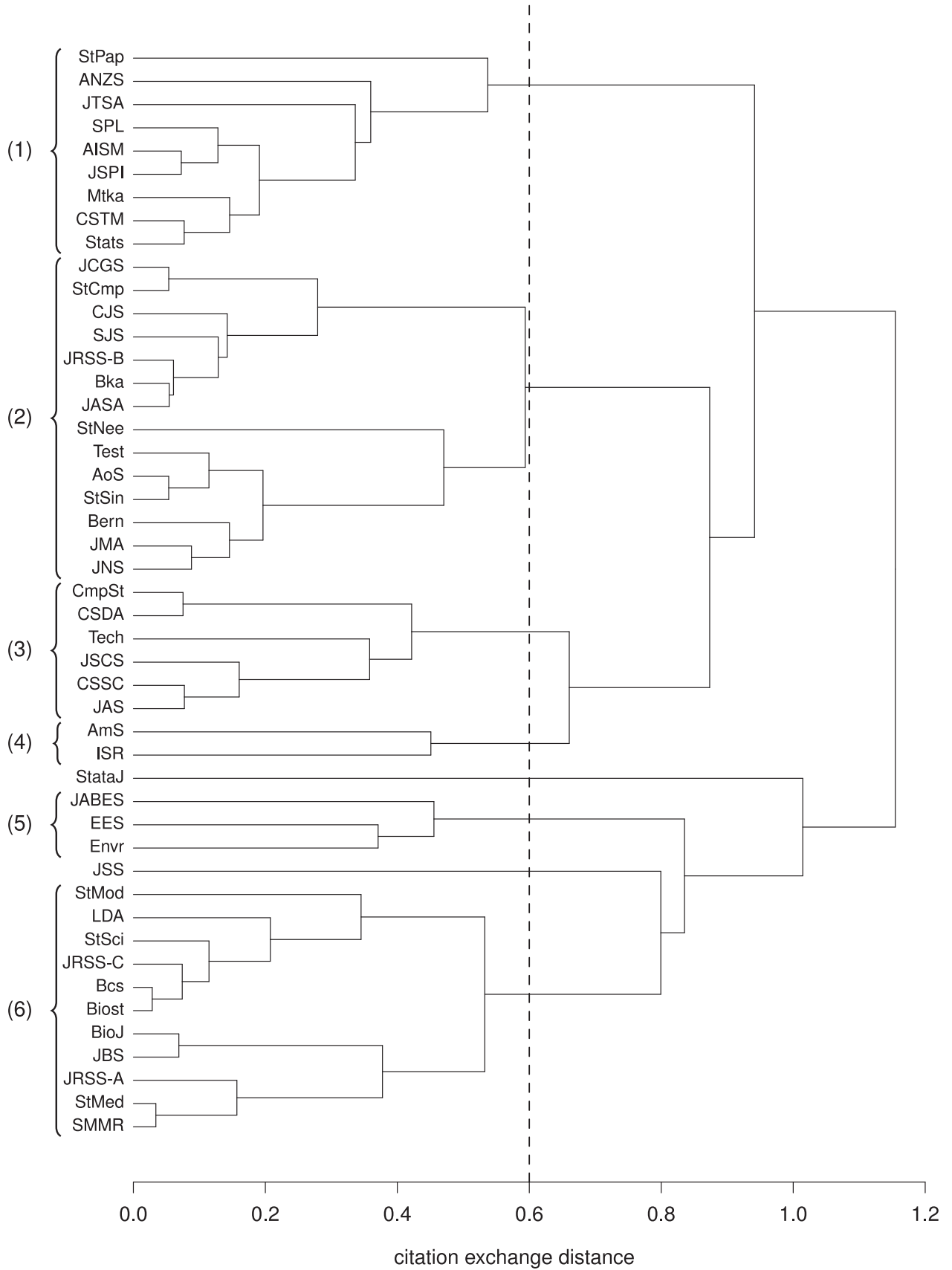


Fig. 2. Dendrogram of complete-linkage hierarchical cluster analysis: clusters obtained by cutting the dendrogram at distance 0.6

4. Ranking journals

The Thomson Reuters JCR Web site annually publishes various rating indices, the best-known being the already mentioned impact factor. Thomson Reuters also publishes the *immediacy index*, which describes the average number of times that an article is cited in the year of its publication. The immediacy index is unsuitable for evaluating statistics journals, but it could be worthy of attention in fields where citations occur very quickly, e.g. some areas of neuroscience and other life sciences.

It is well known in the bibliometric literature that the calculation of the impact factor contains some important inconsistencies (Glänzel and Moed, 2002). The numerator of the impact factor includes citations to all items, whereas the number of citable items in the denominator excludes letters to the editor and editorials; such letters are an important element of some journals, notably medical journals. The inclusion of self-citations, defined as citations from a journal to articles in the same journal, exposes the impact factor to possible manipulation by editors. Indeed, Sevinc (2004), Frandsen (2007) and Wilhite and Fong (2012) have reported instances where authors were asked to add irrelevant references to their articles, presumably with the aim of increasing the impact factor of the journal. As previously mentioned, recently Thomson Reuters has made available also the impact factor without journal self-citations. Journal self-citations can also be a consequence of authors' preferring to cite papers that are published in the same journal instead of equally relevant papers published elsewhere, particularly if they perceive such self-citation as likely to be welcomed by the journal's editors. Nevertheless, the potential for such behaviour should not lead to the conclusion that self-citations are always unfair. Many self-citations are likely to be genuine, especially since scholars often select a journal for submission of their work according to the presence of previously published papers on related topics.

The *eigenfactor score* and the derived *article influence score* (Bergstrom, 2007; West, 2010) have been proposed to overcome the limitations of the impact factor. Both the eigenfactor and the article influence score are computed over a 5-year time period, with journal self-citations removed to eliminate possible sources of manipulation. The idea underlying the eigenfactor score is that the importance of a journal relates to the time that is spent by scholars in reading that journal. As stated by Bergstrom (2007), it is possible to imagine that a scholar starts reading an article selected at random. Then, the scholar randomly selects another article from the references of the first paper and reads it. Afterwards, a further article is selected at random from the references that were included in the previous one and the process may go on *ad infinitum*. In such a process, the time that is spent in reading a journal might reasonably be regarded as an indicator of that journal's importance.

Apart from modifications that are needed to account for special cases such as journals that do not cite any other journal, the eigenfactor algorithm is summarized as follows. The eigenfactor is computed from the normalized citation matrix $\tilde{\mathbf{C}} = (\tilde{c}_{ij})$, whose elements are the citations c_{ij} from journal j to articles published in the previous 5 years in journal i divided by the total number of references in j in those years, $\tilde{c}_{ij} = c_{ij} / \sum_{i=1}^n c_{ij}$. The diagonal elements of $\tilde{\mathbf{C}}$ are set to 0, to discard self-citations. A further ingredient of the eigenfactor is the vector of normalized numbers of articles $\mathbf{a} = (a_1, \dots, a_n)^T$, with a_i being the number of articles published by journal i during the 5-year period divided by the number of articles published by all journals considered. Let \mathbf{e}^T be the row vector of 1s, so that $\mathbf{a}\mathbf{e}^T$ is a matrix with all identical columns \mathbf{a} . Then

$$\mathbf{P} = \lambda \tilde{\mathbf{C}} + (1 - \lambda) \mathbf{a}\mathbf{e}^T$$

is the transition matrix of a Markov process that assigns probability λ to a random movement in

the journal citation network, and probability $1 - \lambda$ to a random jump to any journal; for jumps of the latter kind, destination journal attractiveness is simply proportional to size.

The damping parameter λ is set to 0.85, just as in the PageRank algorithm at the basis of the Google search engine; see Brin and Page (1998). The leading eigenvector ψ of \mathbf{P} corresponds to the steady state fraction of time spent reading each journal. The eigenfactor score EF_i for journal i is defined as ‘the percentage of the total weighted citations that journal i receives’, i.e.

$$EF_i = 100 \frac{[\tilde{\mathbf{C}}\psi]_i}{\sum_{i=1}^n [\tilde{\mathbf{C}}\psi]_i}, \quad i = 1, \dots, n,$$

where $[\mathbf{x}]_i$ denotes the i th element of vector \mathbf{x} . See www.eigenfactor.org/methods.pdf for more details of the methodology behind the eigenfactor algorithm.

The eigenfactor ‘measures the total influence of a journal on the scholarly literature’ (Bergstrom, 2007) and thus it depends on the number of articles that are published by a journal. The article influence score AI_i of journal i is instead a measure of the per-article citation influence of the journal, obtained by normalizing the eigenfactor as follows:

$$AI_i = 0.01 \frac{EF_i}{a_i}, \quad i = 1, \dots, n.$$

Distinctive aspects of the article influence score with respect to the impact factor are

- (a) the use of a formal stochastic model to derive the journal ranking and
- (b) the use of bivariate data—the cross-citations c_{ij} —in contrast with the univariate citation counts that are used by the impact factor.

An appealing feature of the article influence score is that citations are weighted according to the importance of the source, whereas the impact factor counts all citations equally (Franceschet, 2010). Accordingly, the bibliometric literature classifies the article influence score as a measure of journal ‘prestige’ and the impact factor as a measure of journal ‘popularity’ (Bollen *et al.*, 2006). Table 3 summarizes some of the main features of the ranking methods that are discussed in this section and also of the Stigler model that will be discussed in Section 5 below.

The rankings of the selected statistics journals according to impact factor, impact factor without journal self-citations, 5-year impact factor, immediacy index and article influence score

Table 3. Characteristics of the journal rankings derived from the JCR†

<i>Ranking</i>	<i>Citation period (years)</i>	<i>Stochastic model</i>	<i>Data</i>	<i>Excludes self-citation</i>	<i>Global or local</i>
II	1	None	Univariate	No	Global
IF	2	None	Univariate	No	Global
IFno	2	None	Univariate	Yes	Global
IF5	5	None	Univariate	No	Global
AI	5	Markov process	Bivariate	Yes	Global
SM	10	Bradley–Terry	Bivariate	Yes	Local

†Rankings are the immediacy index II, impact factor IF, impact factor without self-citations, IFno, 5-year impact factor, IF5, article influence score AI and the Stigler model studied in this paper, SM. The ‘Data’ column indicates whether the data used are bivariate cross-citation counts or only univariate citation counts. ‘Global or local’ relates to whether a ranking is ‘local’ to the main journals of statistics, or ‘global’ in that it is applied across disciplines.

Table 4. Rankings of selected statistics journals based on the JCR, 2010 edition†

Rank	Results according to the following scores:					
	II	IF	IFno	IF5	AI	SM
1	JSS	JRSS-B	JRSS-B	JRSS-B	JRSS-B	JRSS-B
2	Biost	AoS	Biost	JSS	StSci	AoS
3	SMMR	Biost	AoS	StSci	JASA	Bka
4	StCmp	JSS	JRSS-A	JASA	AoS	JASA
5	AoS	JRSS-A	JSS	Biost	Bka	Bcs
6	EES	StSci	StSci	AoS	Biost	JRSS-A
7	JRSS-B	StMed	StMed	StataJ	StataJ	Bern
8	JCGS	JASA	JASA	SMMR	StCmp	SJS
9	StMed	StataJ	StataJ	JRSS-A	JRSS-A	Biost
10	BioJ	StCmp	StCmp	Bka	JSS	JCGS
11	CSDA	Bka	SMMR	StCmp	Bcs	Tech
12	StSci	SMMR	Bka	StMed	Bern	AmS
13	JRSS-A	Bcs	EES	Bcs	JCGS	JTSA
14	StSin	EES	Bcs	Tech	SMMR	ISR
15	JBS	Tech	Tech	JCGS	Tech	AIMS
16	StataJ	BioJ	BioJ	EES	SJS	CJS
17	Bcs	JCGS	JCGS	CSDA	StMed	StSin
18	Envr	CSDA	Test	SJS	Test	StSci
19	Bka	JBS	AIMS	AmS	CJS	LDA
20	JMA	Test	Bern	JBS	StSin	JRSS-C
21	Tech	JMA	StSin	Bern	JRSS-C	StMed
22	JASA	Bern	LDA	JRSS-C	AmS	ANZS
23	JRSS-C	AmS	JMA	BioJ	JMA	StCmp
24	ISR	AIMS	CSDA	JABES	EES	StataJ
25	JNS	StSin	SJS	JMA	JTSA	SPL
26	Test	LDA	ISR	CJS	LDA	StNee
27	Bern	ISR	JBS	Test	BioJ	Envr
28	JABES	SJS	AmS	StMod	StMod	JABES
29	JSPI	Envr	Envr	StSin	CSDA	Mtka
30	SJS	JABES	StMod	LDA	JABES	StMod
31	AmS	StMod	CJS	Envr	AIMS	JSPI
32	AIMS	JSPI	JABES	JTSA	ANZS	SMMR
33	StMod	CJS	JTSA	ISR	ISR	BioJ
34	Mtka	JTSA	JSPI	ANZS	JSPI	JMA
35	StNee	JRSS-C	ANZS	JSPI	Envr	EES
36	StPap	ANZS	StPap	AIMS	JBS	CSDA
37	SPL	StPap	Mtka	Stats	StNee	JNS
38	ANZS	Mtka	JRSS-C	Mtka	CmpSt	CmpSt
39	LDA	Stats	Stats	CmpSt	JNS	Stats
40	JTSA	CmpSt	CmpSt	StNee	Stats	Test
41	JSCS	JSCS	JSCS	JSCS	Mtka	CSTM
42	CJS	JNS	JNS	StPap	JSCS	JSS
43	CmpSt	SPL	SPL	SPL	StPap	JBS
44	CSTM	CSTM	CSTM	JNS	SPL	JSCS
45	Stats	CSSC	StNee	JAS	CSTM	CSSC
46	JAS	StNee	CSSC	CSTM	CSSC	StPap
47	CSSC	JAS	JAS	CSSC	JAS	JAS

†Columns correspond to the immediacy index II, impact factor IF, impact factor without self-citations IFno, 5-year impact factor IF5, article influence score AI and the Stigler model SM. Braces indicate groups identified by the ranking lasso.

are reported in the second to sixth columns of Table 4. The substantial variation between those five rankings is the first aspect that leaps to the eye; these different published measures clearly do not yield a common, unambiguous picture of the journals' relative standings.

A diffuse opinion within the statistical community is that the four most prestigious statistics journals are (in alphabetic order) *Annals of Statistics*, *Biometrika*, the *Journal of the American Statistical Association* and the *Journal of the Royal Statistical Society, Series B*. See, for example, the survey about how statisticians perceive statistics journals that is described in Theoharakis and Skordia (2003). Accordingly, a minimal requirement for a ranking of acceptable quality is that the four most prestigious journals should occupy prominent positions. Following this criterion, the least satisfactory ranking is, as expected, that based on the immediacy index, which ranks the *Journal of the American Statistical Association* only 22nd and *Biometrika* just a few positions ahead at 19th.

In the three versions of impact factor ranking, the *Journal of the Royal Statistical Society, Series B*, always occupies first position, the *Annals of Statistics* ranges between second and sixth, the *Journal of the American Statistical Association* between fourth and eighth, and *Biometrika* between 10th and 12th. The two software journals have quite high impact factors: the *Journal of Statistical Software* is ranked between second and fifth by the three different impact factor versions, whereas *Stata Journal* is between seventh and ninth. Other journals ranked highly according to the impact factor measures are *Biostatistics* and *Statistical Science*.

Among the indices that are published by Thomson Reuters, the article influence score yields the most satisfactory ranking with respect to the four leading journals mentioned above, all of which stand within the first five positions.

All the indices discussed in this section are constructed by using the complete *Web of Science* database, thus counting citations from journals in other fields as well as citations between statistics and probability journals.

5. The Stigler model

Stigler (1994) considered the export of intellectual influence from a journal to determine its importance. The export of influence is measured through the citations that are received by the journal. Stigler assumed that the log-odds that journal i exports to journal j rather than vice versa are equal to the difference of the journals' *export scores*,

$$\text{log-odds}(\text{journal } i \text{ is cited by journal } j) = \mu_i - \mu_j, \tag{2}$$

where μ_i is the export score of journal i . In Stephen Stigler's words 'the larger the export score, the greater the propensity to export intellectual influence'. The Stigler model is an example of the Bradley–Terry model (Bradley and Terry, 1952; David, 1963; Agresti, 2013) for paired comparison data. According to equation (2), the citation counts c_{ij} are realizations of binomial variables C_{ij} with expected value

$$E(C_{ij}) = t_{ij}\pi_{ij}, \tag{3}$$

where $\pi_{ij} = \exp(\mu_i - \mu_j) / \{1 + \exp(\mu_i - \mu_j)\}$ and t_{ij} is the total number of citations exchanged between journals i and j , as defined in equation (1).

The Stigler model has some attractive features.

- (a) *Statistical modelling*: similarly to the eigenfactor and the derived article influence score, the Stigler method is based on stochastic modelling of a matrix of cross-citation counts. The

methods differ regarding the modelling perspective—a Markov process for the eigenfactor *versus* a Bradley–Terry model in the Stigler method—and, perhaps most importantly, the use of formal statistical methods. The Stigler model is calibrated through well-established statistical fitting methods, such as maximum likelihood or quasi-likelihood (see Section 5.1), with estimation uncertainty summarized accordingly (Section 5.3). Moreover, Stigler model assumptions are readily checked by the analysis of suitably defined residuals, as described in Section 5.2.

- (b) *The size of the journals is not important.* Rankings based on the Stigler model are not affected by the numbers of papers published. As shown by Stigler (1994), page 102, if two journals are merged into a single journal then the odds in favour of that ‘super’ journal against any third journal is a weighted average of the odds for the two separate journals against the third. Normalization for journal size, which is explicit in the definitions of various impact factor and article influence measures, is thus implicit for the Stigler model.
- (c) *Journal self-citations are not counted.* In contrast with the standard impact factor, rankings based on journal export scores μ_i are not affected by the risk of manipulation through journal self-citations.
- (d) *Only citations between journals under comparison are counted.* If the Stigler model is applied to the list of 47 statistics journals, then only citations between these journals are counted. Such an application of the Stigler model thus aims unambiguously to measure influence within the research field of statistics, rather than combining that with potential influence on other research fields. As noted in Table 3, this property differentiates the Stigler model from the other ranking indices published by Thomson Reuters, which use citations from all journals in potentially any fields to create a ‘global’ ranking of all scholarly journals. Obviously it would be possible also to recompute more ‘locally’ the various impact factor measures and/or eigenfactor-based indices, by using only citations exchanged between the journals in a restricted set to be compared.
- (e) *The citing journal is taken into account.* Like the article influence score, the Stigler model measures journals’ relative prestige, because it is derived from bivariate citation counts and thus takes into account the source of each citation. The Stigler model decomposes the cross-citation matrix \mathbf{C} differently, though; it can be re-expressed in log-linear form as the ‘quasi-symmetry’ model,

$$E(C_{ij}) = t_{ij} \exp(\alpha_i + \beta_j), \quad (4)$$

in which the export score for journal i is $\mu_i = \alpha_i - \beta_i$.

- (f) *Lack-of-fit assessment:* Stigler *et al.* (1995) and Liner and Amin (2004) observed increasing lack of fit of the Stigler model when additional journals that trade little with those already under comparison are included in the analysis. Ritzberger (2008) stated bluntly that the Stigler model ‘suffers from a lack of fit’ and dismissed it—incorrectly, in our view—for that reason. We agree instead with Liner and Amin (2004) who suggested that statistical lack-of-fit assessment is another positive feature of the Stigler model that can be used, for example, to identify groups of journals belonging to different research fields, journals which should perhaps not be ranked together. Certainly the existence of principled lack-of-fit assessment for the Stigler model should not be a reason to prefer other methods for which no such assessment is available.

See also Table 3 for a comparison of properties of the ranking methods that are considered in this paper.

5.1. Model fitting

Maximum likelihood estimation of the vector of journal export scores $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)^\top$ can be obtained through standard software for fitting generalized linear models. Alternatively, specialized software such as the R package `BradleyTerry2` (Turner and Firth, 2012) is available through the Comprehensive R Archive Network repository. Since the Stigler model is specified through pairwise differences of export scores $\mu_i - \mu_j$, model identification requires a constraint, such as a ‘reference journal’ constraint $\mu_1 = 0$ or the sum constraint $\sum_{i=1}^n \mu_i = 0$. Without loss of generality we use the latter constraint in what follows.

Standard maximum likelihood estimation of the Stigler model would assume that citation counts c_{ij} are realizations of independent binomial variables C_{ij} . Such an assumption is likely to be inappropriate, since research citations are not independent of one another in practice; see Cattelan (2012) for a general discussion on handling dependence in paired comparison modelling. The presence of dependence between citations can be expected to lead to the well-known phenomenon of overdispersion. A simple way to deal with overdispersion is provided by the method of quasi-likelihood (Wedderburn, 1974). Accordingly, we consider a ‘quasi-Stigler’ model,

$$\begin{aligned} E(C_{ij}) &= t_{ij}\pi_{ij}, \\ \text{var}(C_{ij}) &= \phi t_{ij}\pi_{ij}(1 - \pi_{ij}), \end{aligned} \quad (5)$$

where $\phi > 0$ is the dispersion parameter. Let \mathbf{c} be the vector that is obtained by stacking all citation counts c_{ij} in some arbitrary order, and let \mathbf{t} and $\boldsymbol{\pi}$ be the corresponding vectors of totals t_{ij} and expected values π_{ij} respectively. Then estimates of the export scores are obtained by solving the quasi-likelihood estimating equations

$$\mathbf{D}^\top \mathbf{V}^{-1}(\mathbf{c} - \mathbf{t}\boldsymbol{\pi}) = 0, \quad (6)$$

where \mathbf{D} is the Jacobian of $\boldsymbol{\pi}$ with respect to the export scores $\boldsymbol{\mu}$, and $\mathbf{V} = \mathbf{V}(\boldsymbol{\mu})$ is the diagonal matrix with elements $\text{var}(C_{ij})/\phi$. Under the assumed model (5), quasi-likelihood estimators are consistent and asymptotically normally distributed with variance–covariance matrix $\phi(\mathbf{D}^\top \mathbf{V}^{-1} \mathbf{D})^{-1}$. The dispersion parameter is usually estimated via the squared Pearson residuals as

$$\hat{\phi} = \frac{1}{m - n + 1} \sum_{i < j}^n \frac{(c_{ij} - t_{ij}\hat{\pi}_{ij})^2}{t_{ij}\hat{\pi}_{ij}(1 - \hat{\pi}_{ij})}, \quad (7)$$

where $\hat{\boldsymbol{\pi}}$ is the vector of estimates $\hat{\pi}_{ij} = \exp(\hat{\mu}_i - \hat{\mu}_j) / \{1 + \exp(\hat{\mu}_i - \hat{\mu}_j)\}$, with $\hat{\mu}_i$ being the quasi-likelihood estimate of the export score μ_i , and $m = \sum_{i < j} \mathbf{1}(t_{ij} > 0)$ the number of pairs of journals that exchange citations. Well-known properties of quasi-likelihood estimation are robustness against misspecification of the variance matrix \mathbf{V} and optimality within the class of linear unbiased estimating equations.

The estimate of the dispersion parameter that is obtained here, for the model applied to statistics journal cross-citations between 2001 and 2010, is $\hat{\phi} = 1.76$, indicative of overdispersion. The quasi-likelihood estimated export scores of the statistics journals are reported in Table 5 and will be discussed later in Section 5.4.

5.2. Model validation

An essential feature of the Stigler model is that the export score of any journal is a constant.

In particular, in model (2) the export score of journal i is not affected by the identity of the citing journal j . Citations that are exchanged between journals can be seen as results of contests between opposing journals and the residuals for contests involving journal i should not exhibit any relationship with the corresponding estimated export scores of the ‘opponent’ journals j . With this in mind, we define the *journal residual* r_i for journal i as the standardized regression coefficient derived from the linear regression of Pearson residuals involving journal i on the estimated export scores of the corresponding opponent journals. More precisely, the i th journal residual is defined here as

$$r_i = \frac{\sum_{j=1}^n \hat{\mu}_j r_{ij}}{\sqrt{\left(\hat{\phi} \sum_{j=1}^n \hat{\mu}_j^2\right)}},$$

where r_{ij} is the Pearson residual for citations of i by j ,

$$r_{ij} = \frac{c_{ij} - t_{ij}\hat{\pi}_{ij}}{\sqrt{\{t_{ij}\hat{\pi}_{ij}(1 - \hat{\pi}_{ij})\}}}.$$

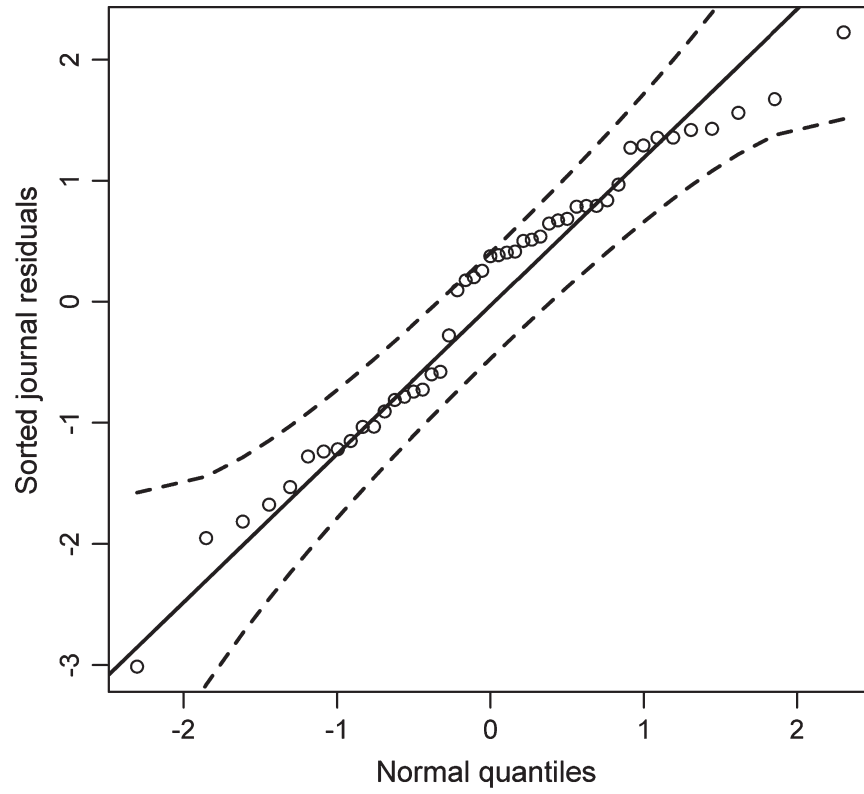
The journal residual r_i indicates the extent to which i performs systematically better than predicted by the model either when the opponent j is strong, as indicated by a positive-valued journal residual for i , or when the opponent j is weak, as indicated by a negative-valued journal residual for i . The journal residuals thus provide a basis for useful diagnostics, targeted specifically at readily interpretable departures from the model assumed.

Under the assumed quasi-Stigler model, journal residuals are approximately realizations of standard normal variables and are unrelated to the export scores. The normal probability plot of the journal residuals displayed in Fig. 3(a) indicates that the normality assumption is indeed approximately satisfied. The scatter plot of the journal residuals r_i against estimated export scores $\hat{\mu}_i$ in Fig. 3(b) shows no clear pattern; there is no evidence of correlation between journal residuals and export scores. As expected on the basis of approximate normality of the residuals, only two journals—i.e. 4.3% of journals—have residuals that are larger in absolute value than 1.96. These journals are *Communications in Statistics—Theory and Methods* ($r_{\text{CSTM}} = 2.23$) and *Test* ($r_{\text{Test}} = -3.01$). The overall conclusion from this graphical inspection of journal residuals is that the assumptions of the quasi-Stigler model appear to be essentially satisfied for the data that are used here.

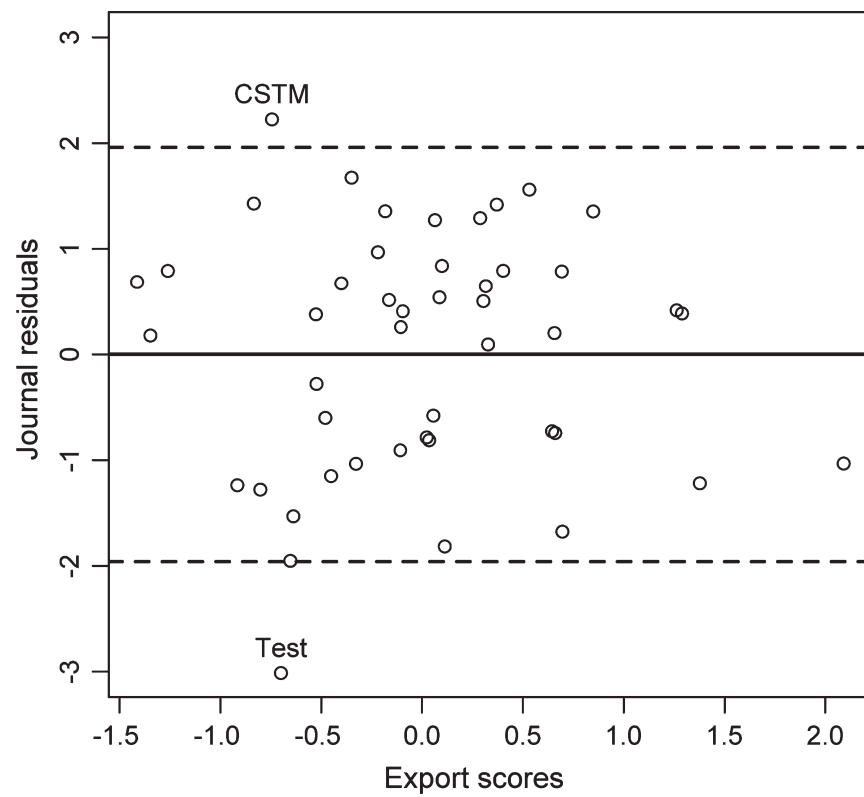
5.3. Estimation uncertainty

Estimation uncertainty is commonly unexplored, and is rarely reported, in relation to the various published journal rankings. Despite this lacuna, many academics have produced vibrant critiques of ‘statistical citation analyses’, although such analyses are actually rather non-statistical. Recent research in the bibliometric field has suggested that uncertainty in estimated journal ratings might be estimated via bootstrap simulation; see the already mentioned Chen *et al.* (2014) and the ‘stability intervals’ for the source-normalized impact per paper index. A key advantage of the Stigler model over other ranking methods is straightforward quantification of the uncertainty in journal export scores.

Since the Stigler model is identified through pairwise differences, uncertainty quantification requires the complete variance matrix of $\hat{\mu}$. Routine reporting of such a large variance matrix is



(a)



(b)

Fig. 3. (a) Normal probability plot of journal residuals with 95% simulation envelope and (b) scatter plot of journal residuals *versus* estimated journal export scores

Table 5. Journal ranking based on the Stigler model using data from the JCR 2010 edition†

Rank	Journal	SM	QSE	SMgrouped	Rank	Journal	SM	QSE	SMgrouped
1	JRSS-B	2.09	0.11	1.87	25	SPL	-0.09	0.09	-0.04
2	AoS	1.38	0.07	1.17	26	StNee	-0.10	0.25	-0.04
3	Bka	1.29	0.08	1.11	27	Envr	-0.11	0.18	-0.04
4	JASA	1.26	0.06	1.11	28	JABES	-0.16	0.23	-0.04
5	Bcs	0.85	0.07	0.65	29	Mtka	-0.18	0.17	-0.04
6	JRSS-A	0.70	0.19	0.31	30	StMod	-0.22	0.21	-0.04
7	Bern	0.69	0.15	0.31	31	JSPI	-0.33	0.07	-0.31
8	SJS	0.66	0.12	0.31	32	SMMR	-0.35	0.16	-0.31
9	Biost	0.66	0.11	0.31	33	BioJ	-0.40	0.12	-0.31
10	JCGS	0.64	0.12	0.31	34	JMA	-0.45	0.08	-0.36
11	Tech	0.53	0.15	0.31	35	EES	-0.48	0.25	-0.36
12	AmS	0.40	0.18	0.04	36	CSDA	-0.52	0.07	-0.36
13	JTSA	0.37	0.20	0.04	37	JNS	-0.53	0.15	-0.36
14	ISR	0.33	0.25	0.04	38	CmpSt	-0.64	0.22	-0.36
15	AISM	0.32	0.16	0.04	39	Stats	-0.65	0.18	-0.36
16	CJS	0.30	0.14	0.04	40	Test	-0.70	0.15	-0.36
17	StSin	0.29	0.09	0.04	41	CSTM	-0.74	0.10	-0.36
18	StSci	0.11	0.11	-0.04	42	JSS	-0.80	0.19	-0.36
19	LDA	0.10	0.17	-0.04	43	JBS	-0.83	0.16	-0.36
20	JRSS-C	0.09	0.15	-0.04	44	JSCS	-0.92	0.15	-0.36
21	StMed	0.06	0.07	-0.04	45	CSSC	-1.26	0.14	-0.88
22	ANZS	0.06	0.21	-0.04	46	StPap	-1.35	0.20	-0.88
23	StCmp	0.04	0.15	-0.04	47	JAS	-1.41	0.15	-0.88
24	StataJ	0.02	0.33	-0.04					

†Columns are the quasi-likelihood estimated Stigler model export scores SM with associated quasi-standard errors QSE, and estimated export scores after grouping by lasso, SMgrouped.

impracticable for brevity. A neat solution is provided through the presentational device of quasi-variances (Firth and de Menezes, 2005), constructed in such a way as to allow approximate calculation of any variance of a difference, $\text{var}(\hat{\mu}_i - \hat{\mu}_j)$, as if $\hat{\mu}_i$ and $\hat{\mu}_j$ were independent:

$$\text{var}(\hat{\mu}_i - \hat{\mu}_j) \simeq \text{qvar}_i + \text{qvar}_j, \quad \text{for all choices of } i \text{ and } j.$$

Reporting the estimated export scores with their quasi-variances, then, is an economical way to allow approximate inference on the significance of the difference between any two journals' export scores. The quasi-variances are computed by minimizing a suitable penalty function of the differences between the true variances, $\text{var}(\hat{\mu}_i - \hat{\mu}_j)$, and their quasi-variance representations $\text{qvar}_i + \text{qvar}_j$. See Firth and de Menezes (2005) for details.

Table 5 reports the estimated journal export scores computed under the sum constraint $\sum_{i=1}^n \mu_i = 0$ and the corresponding quasi-standard errors, defined as the square root of the quasi-variances. Quasi-variances are calculated by using the R package `qvcalc` (Firth, 2012). For illustration, consider testing whether the export score of *Biometrika* is significantly different from that of the *Journal of the American Statistical Association*. The z -test statistic as approximated through the quasi-variances is

$$z \simeq \frac{\hat{\mu}_{\text{Bka}} - \hat{\mu}_{\text{JASA}}}{\sqrt{(\text{qvar}_{\text{Bka}} + \text{qvar}_{\text{JASA}})}} = \frac{1.29 - 1.26}{\sqrt{(0.08^2 + 0.06^2)}} = 0.30.$$

The 'usual' variances for those two export scores in the sum-constrained parameterization are respectively 0.0376 and 0.0344, and the covariance is 0.0312; thus the 'exact' value of the z -statistic in this example is

$$z = \frac{1.29 - 1.26}{\sqrt{\{0.0376 - 2(0.0312) + 0.0344\}}} = 0.31,$$

so the approximation based on quasi-variances is quite accurate. In this case the z -statistic suggests that there is insufficient evidence to rule out the possibility that *Biometrika* and the *Journal of the American Statistical Association* have the same ability to ‘export intellectual influence’ within the 47 statistics journals in the list.

5.4. Results

We proceed now with interpretation of the ranking based on the Stigler model. It is reassuring that the four leading statistics journals that were mentioned previously are ranked in the first four positions. The *Journal of the Royal Statistical Society, Series B*, is ranked first with a remarkably larger export score than the second-ranked journal, the *Annals of Statistics*: the approximate z -statistic for the significance of the difference of their export scores is 5.44. The third position is occupied by *Biometrika*, closely followed by the *Journal of the American Statistical Association*.

The fifth-ranked journal is *Biometrics*, followed by the *Journal of the Royal Statistical Society, Series A*, *Bernoulli*, the *Scandinavian Journal of Statistics*, *Biostatistics*, the *Journal of Computational and Graphical Statistics* and *Technometrics*.

The ‘centipede’ plot in Fig. 4 visualizes the estimated export scores along with the 95% comparison intervals with limits $\hat{\mu}_i \pm 1.96 \text{QSE}(\hat{\mu}_i)$, where ‘QSE’ denotes the quasi-standard error. The centipede plot highlights the outstanding position of the *Journal of the Royal Statistical Society, Series B*, and indeed of the four top journals whose comparison intervals are well separated from those of the remaining journals. However, the most striking general feature is the substantial uncertainty in most of the estimated journal scores. Many of the small differences that appear between the estimated export scores are not statistically significant.

5.5. Ranking in groups with lasso

Shrinkage estimation offers notable improvement over standard maximum likelihood estimation when the target is simultaneous estimation of a vector of mean parameters; see, for example, Morris (1983). It seems natural to consider shrinkage estimation also for the Stigler model. Masarotto and Varin (2012) fitted Bradley–Terry models with a lasso-type penalty (Tibshirani, 1996) which, in our application here, forces journals with close export scores to be estimated at the same level. The method, which is termed the ranking lasso, has the twofold advantages of shrinkage and enhanced interpretation, because it avoids overinterpretation of small differences between estimated journal export scores.

For a given value of a bound parameter $s \geq 0$, the ranking lasso method fits the Stigler model by solving the quasi-likelihood equations (6) with an L_1 -penalty on all the pairwise differences of export scores, i.e

$$\mathbf{D}^T \mathbf{V}^{-1} (\mathbf{c} - \mathbf{t}\boldsymbol{\pi}) = \mathbf{0}, \quad \text{subject to } \sum_{i < j}^n w_{ij} |\mu_i - \mu_j| \leq s \text{ and } \sum_{i=1}^n \mu_i = 0, \quad (8)$$

where the w_{ij} are data-dependent weights discussed below.

Quasi-likelihood estimation is obtained for a sufficiently large value of the bound s . As s decreases to 0, the L_1 -penalty causes journal export scores that differ little to be estimated at the same value, thus producing a ranking in groups. The ranking lasso method can be interpreted as a generalized version of the fused lasso (Tibshirani *et al.*, 2005).

Since quasi-likelihood estimates coincide with maximum likelihood estimates for the corres-

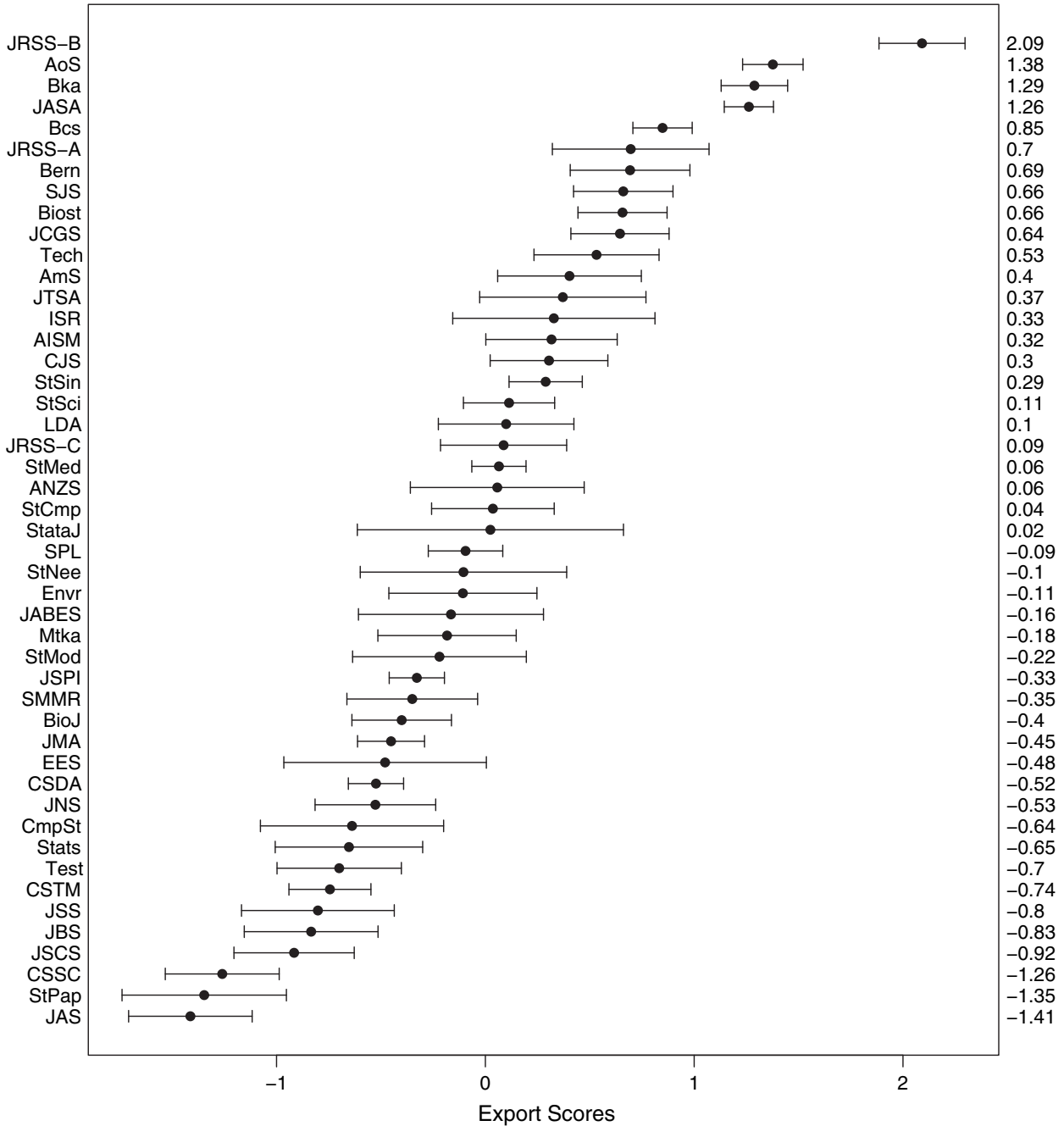


Fig. 4. Centipede plot of estimated journal export scores and 95% comparison intervals based on the JCR 2010 edition: the error bar limits are $\hat{\mu}_i \pm 1.96 \text{QSE}(\hat{\mu}_i)$, with the estimated export scores $\hat{\mu}_i$ marked (●)

ponding exponential dispersion model, ranking lasso solutions can be computed as penalized likelihood estimates. Masarotto and Varin (2012) obtained estimates of the adaptive ranking lasso by using an augmented Lagrangian algorithm (Nocedal and Wright, 2006) for a sequence of bounds s ranging from complete shrinkage ($s = 0$)—i.e. all journals have the same estimated export score—to the quasi-likelihood solution ($s = \infty$).

Many researchers (e.g. Fan and Li (2001) and Zou (2006)) have observed that lasso-type penalties may be too severe, thus yielding inconsistent estimates of the non-zero effects. In the ranking lasso context, this means that, if the weights w_{ij} in problem (8) are all identical, then the pairwise differences $\mu_i - \mu_j$ whose ‘true’ value is non-zero might not be consistently

estimated. Among various possibilities, an effective way to overcome the drawback is to resort to the adaptive lasso method (Zou, 2006), which imposes a heavier penalty on small effects. Accordingly, the adaptive ranking lasso employs weights that are equal to the reciprocal of a consistent estimate of $\mu_i - \mu_j$, such as $w_{ij} = |\hat{\mu}_i^{(\text{QLE})} - \hat{\mu}_j^{(\text{QLE})}|^{-1}$, with $\hat{\mu}_i^{(\text{QLE})}$ being the quasi-likelihood estimate of the export score for journal i .

Lasso tuning parameters are often determined by cross-validation. Unfortunately, the interjournal ‘tournament’ structure of the data does not allow the identification of internal replication; hence it is not clear how cross-validation can be applied to citation data. Alternatively, tuning parameters can be determined by minimization of suitable information criteria. The usual Akaike information criterion is not valid with quasi-likelihood estimation because the likelihood function is formally unspecified. A valid alternative is based on the Takeuchi information criterion TIC (Takeuchi, 1976) which extends the Akaike information criterion when the likelihood function is misspecified. Let $\hat{\boldsymbol{\mu}}(s) = (\hat{\mu}_1(s), \dots, \hat{\mu}_n(s))^T$ denote the solution of problem (8) for a given value of the bound s . Then the optimal value for s is chosen by minimization of

$$\text{TIC}(s) = -2\hat{l}(s) + 2 \text{tr}\{\mathbf{J}(s)\mathbf{I}(s)^{-1}\},$$

where $\hat{l}(s) = l\{\hat{\boldsymbol{\mu}}(s)\}$ is the misspecified log-likelihood of the Stigler model

$$l(\boldsymbol{\mu}) = \sum_{i < j}^n c_{ij}(\mu_i - \mu_j) - t_{ij} \ln\{1 + \exp(\mu_i - \mu_j)\}$$

computed at $\hat{\boldsymbol{\mu}}(s)$, $\mathbf{J}(s) = \text{var}\{\nabla l(\boldsymbol{\mu})\}|_{\boldsymbol{\mu}=\hat{\boldsymbol{\mu}}(s)}$ and $\mathbf{I}(s) = -E\{\nabla^2 l(\boldsymbol{\mu})\}|_{\boldsymbol{\mu}=\hat{\boldsymbol{\mu}}(s)}$. Under the assumed quasi-Stigler model, $\mathbf{J}(s) = \phi\mathbf{I}(s)$ and the TIC-statistic reduces to

$$\text{TIC}(s) = -2\hat{l}(s) + 2\phi p,$$

where p is the number of distinct groups formed with bound s . The dispersion parameter ϕ can be estimated as in equation (7). The effect of overdispersion is inflation of the Akaike information criterion model dimension penalty.

Fig. 5 displays the path plot of the ranking lasso, and Table 5 reports estimated export scores corresponding to the solution identified by TIC. See also Table 4 for a comparison with the Thomson Reuters published rankings. The path plot of Fig. 5 visualizes how the estimates of the export scores vary as the degree of shrinkage decreases, i.e. as the bound s increases. The plot confirms the outstanding position of the *Journal of the Royal Statistical Society, Series B*, the leader in the ranking at any level of shrinkage. Also *Annals of Statistics* keeps the second position for about three-quarters of the path before joining the paths of *Biometrika* and the *Journal of the American Statistical Association*. *Biometrics* is solitary in fifth position for almost the whole of its path. The TIC-statistic identifies a sparse solution with only 10 groups. According to TIC, the five top journals are followed by a group of six further journals, namely the *Journal of the Royal Statistical Society, Series A*, *Bernoulli*, the *Scandinavian Journal of Statistics*, *Biostatistics*, the *Journal of Computational and Graphical Statistics* and *Technometrics*. However, the main conclusion from this ranking lasso analysis is that many of the estimated journal export scores are not clearly distinguishable from one another.

6. Comparison with results from the UK research assessment exercise

6.1. Background

In the UK, the quality of the research that is carried out in universities is assessed periodically by

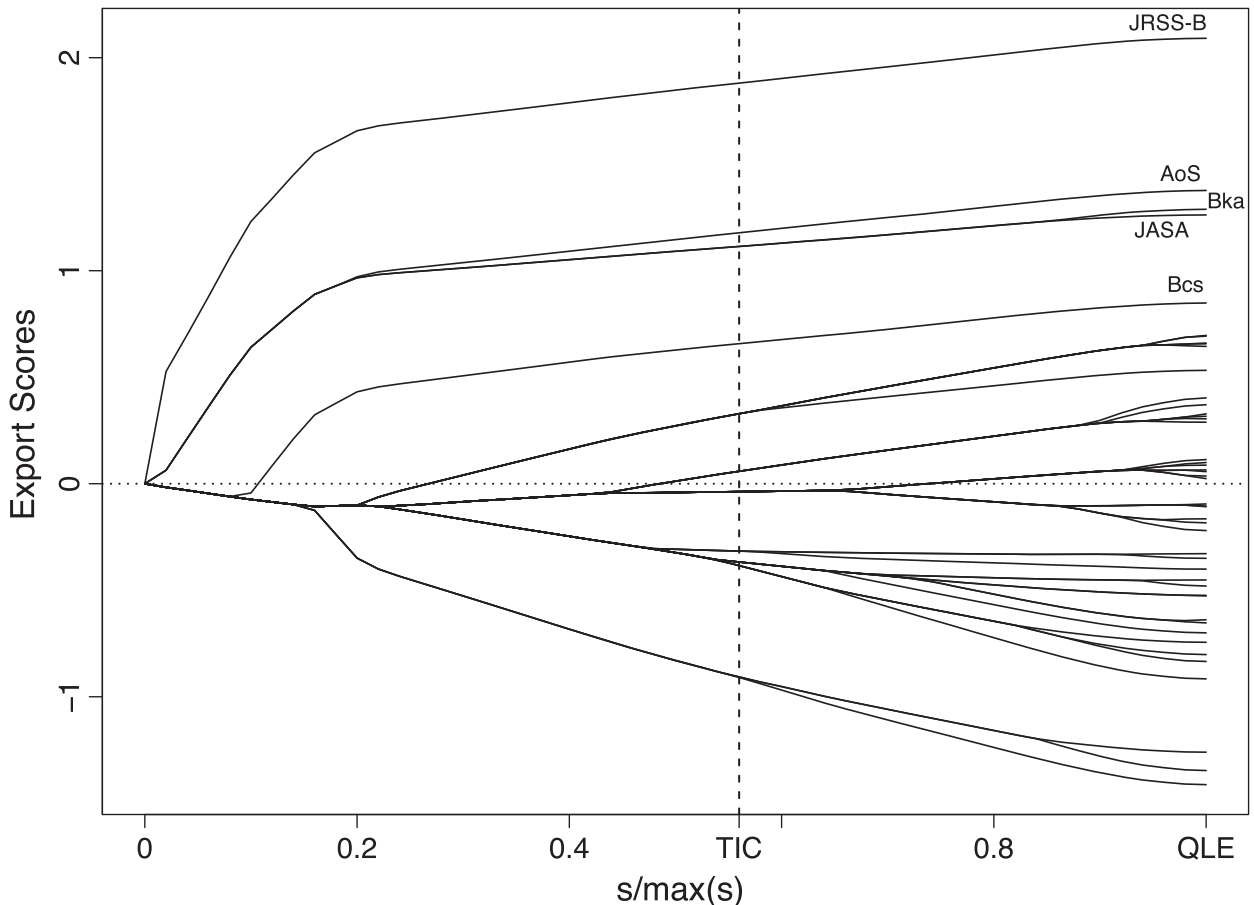


Fig. 5. Path plot of adaptive ranking lasso analysis based on the JCR 2010 edition: QLE, quasi-likelihood estimate; TIC, Takeuchi information criterion

the government-supported funding councils, as a primary basis for future funding allocations. At the time of writing, the most recent such assessment to be completed was the 2008 RAE, full details of which are on line at www.rae.ac.uk. The next such assessment to report, at the end of 2014, will be the similar ‘research excellence framework’. Each unit of assessment is an academic ‘department’, corresponding to a specified research discipline. In the 2008 RAE, ‘Statistics and operational research’ was one of 67 such research disciplines; in contrast the 2014 research excellence framework has only 36 separate discipline areas identified for assessment, and research in statistics will be part of a new and much larger ‘Mathematical sciences’ unit of assessment. The results from the 2008 RAE are therefore likely to provide the last opportunity to make a directly statistics-focused comparison with journal rankings.

The word ‘department’ in the 2008 RAE refers to a discipline-specific group of researchers submitted for assessment by a university, or sometimes by two universities together: a department in the 2008 RAE need not be an established academic unit within a university, and indeed many of the 2008 RAE statistics and operational research departments were actually groups of researchers working in university departments of mathematics or other disciplines.

It is often argued that the substantial cost of assessing research outputs through review by a panel of experts, as was done in the 2008 RAE, might be reduced by employing suitable metrics based on citation data. See, for example, Jump (2014). Here we briefly explore this in quite a specific way, through data on journals rather than on the citations that are attracted by individual research papers submitted for assessment.

The comparisons to be made here can also be viewed as exploring an aspect of ‘criterion validity’ of the various journal ranking methods: if highly ranked journals tend to contain high quality research, then there should be evidence through strong correlations, even at the ‘department’ level of aggregation, between expert panel assessments of research quality and journal ranking scores.

6.2. Data and methods

We examine only Sub-panel 22, ‘Statistics and operational research’ of the 2008 RAE. The specific data used here are

- (a) the detailed ‘RA2’ (research outputs) submissions made by departments to the 2008 RAE (these list up to four research outputs per submitted researcher) and
- (b) the published 2008 RAE results on the assessed quality of research outputs, namely the ‘outputs subprofile’ for each department.

From the RA2 data, only research outputs categorized in the 2008 RAE as ‘journal article’ are considered here. For each such article, the journal’s name is found in the ‘publisher’ field of the data. A complication is that the name of any given journal can appear in many different ways in the RA2 data, e.g. ‘*Journal of the Royal Statistical Society B*’ and ‘*Journal of the Royal Statistical Society Series B: Statistical Methodology*’, and the International Standard Serial Number codes as entered in the RA2 data are similarly unreliable. Unambiguously resolving all of the many different representations of journal names proved to be the most time-consuming part of the comparison exercise that is reported here.

The 2008 RAE outputs subprofile for each department gives the assessed percentage of research outputs at each of five quality levels, these being ‘world leading’ (shorthand code ‘4*’), ‘internationally excellent’ (shorthand ‘3*’), then ‘2*’, ‘1*’ and ‘U’ (unclassified). For example, the outputs subprofile for University of Oxford, the highest-rated statistics and operational research submission in the 2008 RAE, is

4*	3*	2*	1*	U
37.0	49.5	11.4	2.1	0.

Our focus will be on the fractions at the 4* and 3* quality levels, since those are used as the basis for research funding. Specifically, in the comparisons that are made here the RAE ‘score’ used will be the percentage at 4* plus a third of the percentage at 3*, computed from each department’s 2008 RAE outputs subprofile. Thus, for example, Oxford’s 2008 RAE score is calculated as $37.0 + 49.5/3 = 53.5$. This scoring formula is essentially that used since 2010 to determine funding council allocations; we have considered also various other possibilities, such as simply the percentage at 4*, or the percentage at 3* or higher, and found that the results below are not sensitive to this choice.

For each of the journal ranking methods listed in Table 3, a bibliometrics-based comparator score per department is then constructed in a natural way as follows. Each RAE-submitted journal article is scored individually, by for example the impact factor of the journal in which it appeared; and those individual article scores are then averaged across all of a department’s RAE-submitted journal articles. For the averaging, we use the simple arithmetic mean of scores; an exception is that Stigler model export scores are exponentiated before averaging, so that they are positive valued like the scores for the other methods considered. Use of the median was considered as an alternative to the mean; it was found to produce very similar results, which accordingly will not be reported here.

A complicating factor for the simple scoring scheme just described is that journal scores were not readily available for all the journals named in the RAE submissions. For the various ‘global’ ranking measures (see Table 3), scores were available for the 110 journals in the JCR ‘Statistics and probability’ category, which covers approximately 70% of the RAE-submitted journal articles to be scored. For the Stigler model as used in this paper, though, only the subset of 47 statistics journals that are listed in Table 1 are scored; and this subset accounts for just under half of the RAE-submitted journal articles. In what follows we have ignored all articles that appeared in unscored journals, and used the rest. To enable a more direct comparison with the use of Stigler model scores, for each of the global indices we computed also a restricted version of its mean score for each department, i.e. restricted to using scores for only the 47 statistics journals from Table 1.

Of the 30 departments submitting work in ‘Statistics and operational research’ to the 2008 RAE, four turned out to have substantially less than 50% of their submitted journal articles in the JCR ‘Statistics and probability’ category of journals. The data from those four departments, which were relatively small groups and whose RAE-submitted work was mainly in operational research, have been omitted from the following analysis.

The statistical methods that are used below to examine department level relationships between the RAE scores and journal-based scores are simply correlation coefficients and scatter plots. Given the arbitrary nature of data availability for this particular exercise, anything more sophisticated would seem inappropriate.

6.3. Results

Table 6 shows, for bibliometrics-based mean scores based on each of the various journal ranking measures discussed in this paper, the computed correlation with departmental RAE score. The main features of Table 6 are as follows.

- (a) The article influence and Stigler model scores correlate more strongly with RAE results than do scores based on the other journal ranking measures.
- (b) The various global measures show stronger correlation with the RAE results when they are used only to score articles from the 47 statistics journals of Table 1, rather than to score everything from the larger set of journals in the JCR ‘Statistics and probability’ category.

The first of these findings unsurprisingly gives clear support to the notion that the use of bivariate citation counts, which take account of the source of each citation and hence lead to measures of journal ‘prestige’ rather than ‘popularity’, is important if a resultant ranking of journals

Table 6. 2008 RAE score for research outputs in 26 UK ‘Statistics and operational research’ departments: Pearson correlation with departmental mean scores derived from the various journal rating indices based on the 2010 JCR

<i>Journals scored</i>	<i>Results for the following journal scoring methods:</i>						
	<i>II</i>	<i>IF</i>	<i>IFno</i>	<i>IF5</i>	<i>AI</i>	<i>SM</i>	<i>SMgrouped</i>
All of the JCR ‘Statistics and probability’ category	0.34	0.47	0.49	0.50	0.73	—	—
Only the 47 statistics journals listed in Table 1	0.34	0.69	0.70	0.73	0.79	0.81	0.82

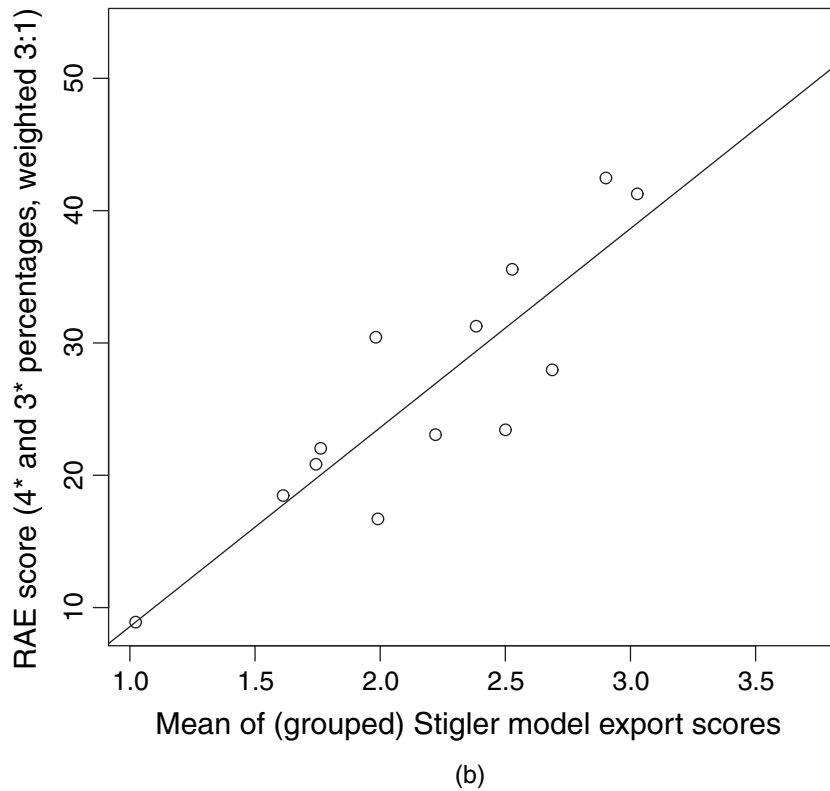
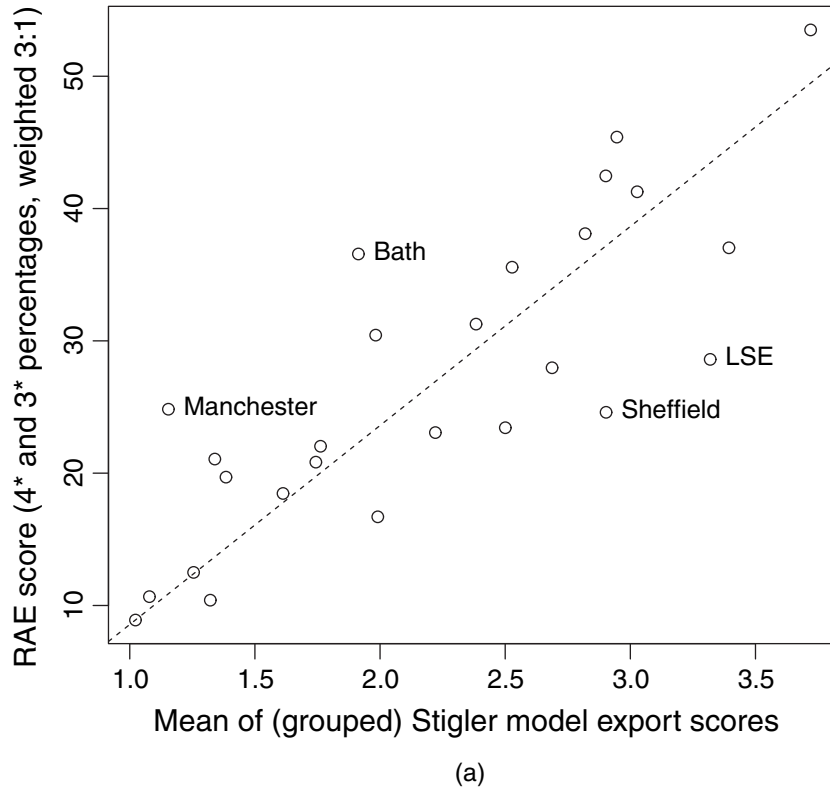


Fig. 6. (a) Scatter plot of the 2008 RAE outcome (scores derived from the published RAE ‘outputs’ sub-profiles) against averaged Stigler model journal export scores for RAE-submitted papers (the 26 plotted points are the main ‘Statistics and operational research’ groups at UK universities; four outliers from a straight line fit are highlighted) and (b) a subset of the same scatter plot: just the 13 research groups for which papers published in the 47 journals in Table 1 formed the majority of their RAE-submitted research outputs; the straight line shown in both panels is the least squares fit to these 13 points

should relate strongly to the perceived quality of published articles. The second finding is more interesting: for good agreement with departmental RAE ratings; it can be substantially better to score only those journals that are in a relatively homogeneous subset than to use all the scores that might be available for a larger set of journals. In the present context, for example, citation patterns for research in probability are known to differ appreciably from those in statistics, and global scoring of journals across these disciplines would tend not to rate highly even the very best work in probability.

The strongest correlations found in Table 6 are those based on journal export scores from the Stigler model, from columns ‘SM’ and ‘SM grouped’ of Table 5. The departmental means of grouped export scores from the ranking lasso method correlate most strongly with RAE scores, which is a finding that supports the notion that small estimated differences between journals are likely to be spurious. Fig. 6(a) shows the relationship between RAE score and the mean of ‘SM-grouped’ exponentiated journal export scores, for the 26 departments whose RAE-submitted journal articles were predominantly in the JCR ‘Statistics and probability’ category; the correlation as reported in Table 6 is 0.82. The four largest outliers from a straight line relationship are identified in the plot, and it is notable that all of those four departments are such that the ratio

$$\frac{\text{number of RAE outputs in the 47 statistics journals of Table 1}}{\text{total number of RAE-submitted journal articles}} \quad (9)$$

is less than $\frac{1}{2}$. Thus the largest outliers are all departments for which the majority of RAE-submitted journal articles are not actually scored by our application of the Stigler model, and this seems entirely to be expected. Fig. 6(b) plots the same scores but now omitting all the 13 departments whose ratio (9) is less than $\frac{1}{2}$. The result is, as expected, much closer to a straight line relationship; the correlation in this restricted set of the most ‘statistical’ departments increases to 0.88.

Some brief remarks on interpretation of these findings appear in Section 7. 5 below. The data and R language code for this comparison are included in this paper’s supplementary Web materials.

7. Concluding remarks

7.1. *The role of statistical modelling in citation analysis*

In his Presidential address at the 2011 Institute of Mathematical Statistics Annual Meeting about controversial aspects of measuring research performance through bibliometrics, Professor Peter Hall concluded that

‘As statisticians we should become more involved in these matters than we are. We are often the subject of the analyses discussed above, and almost alone we have the skills to respond to them, for example by developing new methodologies or by pointing out that existing approaches are challenged. To illustrate the fact that issues that are obvious to statisticians are often ignored in bibliometric analysis, I mention that many proponents of impact factors, and other aspects of citation analysis, have little concept of the problems caused by averaging very heavy tailed data. (Citation data are typically of this type.) We should definitely take a greater interest in this area’ (Hall, 2011).

The model-based approach to journal ranking that is discussed in this paper is a contribution in the direction that Professor Hall recommended. Explicit statistical modelling of citation data has two important merits: first, transparency, since model assumptions need to be clearly stated and can be assessed through standard diagnostic tools; secondly, the evaluation and reporting of uncertainty in statistical models can be based on well-established methods.

7.2. The importance of reporting uncertainty in journal rankings

Many journals' Web sites report the latest journal impact factor and the journal's corresponding rank in its category. Very small differences in the reported impact factor often imply large differences in the corresponding rankings of statistics journals. Statisticians should naturally be concerned about whether such differences are significant. Our analyses conclude that many of the apparent differences between estimated export scores are insignificant, and thus differences in journal ranks are often not reliable. The clear difficulty of discriminating between journals on the basis of citation data is further evidence that the use of journal rankings for evaluation of individual researchers will often—and perhaps always—be inappropriate.

In view of the uncertainty in rankings, it makes sense to ask whether the use of 'grouped' ranks such as those that emerge from the lasso method of Section 5.5 should be universally advocated. If the rankings or associated scores are to be used for prediction, then the usual arguments for shrinkage methods apply and such grouping, to help to eliminate apparent but spurious differences between journals, is likely to be beneficial; predictions based on grouped ranks or scores are likely to be at least as good as those made without the grouping, as indeed we found in Section 6.3 in connection with the 2008 RAE outcomes. For presentational purposes, though, the key requirement is at least some indication of the amount of uncertainty, and ungrouped estimates coupled with realistically wide intervals, as in the centipede plot of Fig. 4, will often suffice.

7.3. A 'read papers' effect?

Discussion papers read to the Society at meetings organized by the Research Section of the Royal Statistical Society are a distinctive aspect of the *Journal of the Royal Statistical Society, Series B*. It is natural to ask whether there is a 'read papers effect' which might explain the prominence of that journal under the metric used in this paper. During the study period 2001–2010, the *Journal of the Royal Statistical Society, Series B*, published in total 446 articles, 36 of which were papers read to the Society. Half of these papers were published during the three years 2002–2004. The *Journal of the Royal Statistical Society, Series B*, received in total 2554 citations from papers published in 2010, with 1029 of those citations coming from other statistics journals in the list. Despite the fact that papers read to the Society were only 8.1% of all published *Journal of the Royal Statistical Society, Series B*, papers, they accounted for 25.4% (649/2554) of all citations received by the *Journal of the Royal Statistical Society, Series B*, in 2010, and 23.1% (238/1029) of the citations from the other statistics journals in the list.

Papers read to the Society are certainly an important aspect of the success of the *Journal of the Royal Statistical Society, Series B*. However, not all such papers contribute strongly to the citations received by the journal. In fact, a closer look at citation counts reveals that the distribution of the citations received by papers read to the Society is very skew, not differently from what happens for 'standard' papers. The most cited read paper published in 2001–2010 was Spiegelhalter *et al.* (2002), which alone received 11.9% of all *Journal of the Royal Statistical Society, Series B*, citations in 2010, and 7.4% of those received from other statistics journals in the list. About 75% of the remaining discussion papers published in the study period each received less than 0.5% of the 2010 *Journal of the Royal Statistical Society, Series B*, citations.

A precise quantification of the 'read paper' effect is difficult. Refitting the Stigler model dropping the citations that were received by these papers seems an unfair exercise. Proper evaluation of the effect would require removal also of the citations received by other papers derived from papers read to the Society and published either in the *Journal of the Royal Statistical Society, Series B*, or elsewhere.

7.4. Possible extensions

7.4.1. Fractioned citations

The analyses that are discussed in this paper are based on the total numbers c_{ij} of citations exchanged by pairs of journals in a given period and available through the JCRs. One potential drawback of this approach is that citations are all counted equally, irrespective of the number of references contained in the citing paper. Some recent papers in the bibliometric literature (e.g. Zitt and Small (2008), Moed (2010), Leydesdorff and Opthof (2010) and Leydesdorff and Bornmann (2011)) suggest that the impact factor and other citation indices should be recomputed by using fractional counting, in which each citation is counted as $1/n$ with n being the number of references in the citing paper. Fractional counting is a natural expedient to take account of varying lengths of reference lists in papers; for example, a typical review article contains many more references than does a short, technical research paper. The Stigler model extends easily to handle such fractional counting, e.g. through the quasi-symmetry formulation (4); and the rest of the methodology described here would apply with straightforward modifications.

7.4.2. Evolution of export scores

This paper discusses a ‘static’ Stigler model fitted to data extracted from a single JCR edition. A natural extension would be to study the evolution of citation exchange between pairs of journals over several years, through a dynamic version of the Stigler model. A general form for such a model is

$$\log\text{-odds}(\text{journal } i \text{ is cited by journal } j \text{ in year } t) = \mu_i(t) - \mu_j(t),$$

where each journal’s time-dependent export score $\mu_i(t)$ is assumed to be a separate smooth function of t . Such a model would not only facilitate the systematic study of time trends in the relative intellectual influence of journals; it would also ‘borrow strength’ across years to help to smooth out spurious variation, whether it be ‘random’ variation arising from the allocation of citing papers to a specific year’s JCR edition, or variation caused by transient, idiosyncratic patterns of citation. A variety of such dynamic extensions of the Bradley–Terry model have been developed in other contexts, especially the modelling of sports data; see, for example, Fahrmeir and Tutz (1994), Glickman (1999), Knorr-Held (2000) and Cattelan *et al.* (2013).

7.5. Citation-based metrics and research assessment

From the strong correlations found in Section 6 between the 2008 RAE outcomes and journal ranking scores, it is tempting to conclude that the expert review element of such a research assessment might reasonably be replaced, mainly or entirely, by automated scoring of journal articles based on the journals in which they have appeared. Certainly Fig. 6 indicates that such scoring, when applied to the main journals of statistics, can perform quite well as a predictor of RAE outcomes for research groups whose publications have appeared mostly in those journals.

The following points should be noted, however.

- (a) Even with correlation as high as 0.88, as in Fig. 6(b), there can be substantial differences between departments’ positions based on RAE outcomes and on journal scores. For example, in Fig. 6(b) there are two departments whose mean scores based on our application of the Stigler model are between 1.9 and 2.0 and thus essentially equal, but their computed RAE scores, at 16.7 and 30.4, differ very substantially indeed.
- (b) High correlation was achieved by scoring only a relatively homogeneous subset of all the journals in which the RAE-submitted work appeared. Scoring a wider set of journals,

to cover most or all of the journal articles appearing in the 2008 RAE ‘Statistics and operational research’ submissions, leads to much lower levels of agreement with RAE results.

In relation to point (a) it could of course be argued that, in cases such as the two departments mentioned, the 2008 RAE panel of experts were wrong, or it could be that the difference that was seen between those two departments in the RAE results is largely attributable to the 40% or so of journal articles for each department that were not scored because they were outside the list in Table 1. Point (b), in contrast, seems more clearly to be a severe limitation on the potential use of journal scores in place of expert review. The use of cluster analysis as in Section 3, in conjunction with expert judgements about which journals are ‘core’ to disciplines and subdisciplines, can help to establish relatively homogeneous subsets of journals that might reasonably be ranked together; but comparison across the boundaries of such subsets is much more problematic.

The analysis that is described in this paper concerns journals. It says nothing directly about the possible use of citation data on individual research outputs, as were made available to several of the review panels in the 2014 research excellence framework for example. For research in mathematics or statistics it seems clear that such data on recent publications carry little information, mainly because of long and widely varying times taken for good research to achieve ‘impact’ through citations; indeed, the mathematical sciences subpanel in the 2014 research excellence framework chose not to use such data at all. Our analysis does, however, indicate that any counting of citations to inform assessment of research quality should at least take account of the source of each citation.

Acknowledgements

The authors are grateful to Alan Agresti, Mike Titterington, the referees, the Series A Joint Editor and Associate Editor, and the Editor for discussion papers, for helpful comments on earlier versions of this work. The kind permission of Thomson Reuters to distribute the 2010 JCR cross-citation counts is also gratefully acknowledged.

This work was supported by the UK Engineering and Physical Sciences Research Council through Centre for Research in Statistical Methodology grant EP/D002060/1, by University of Padua grant CDPA131553 and by an *Iride* grant from the Dipartimento di Scienze Ambientali, Informatica e Statistica, Università Ca’ Foscari.

References

- Adie, E. and Roe, W. (2013) Altmetric: enriching scholarly content with article-level discussion and metrics. *Learnd Publish.*, **26**, 11–17.
- Adler, R., Ewing, J. and Taylor, P. (2009) Citation statistics (with discussion). *Statist. Sci.*, **24**, 1–14.
- Agresti, A. (2013) *Categorical Data Analysis*, 3rd edn. New York: Wiley.
- Alberts, B. (2013) Impact factor distortions. *Science*, **340**, 787.
- Amin, M. and Mabe, M. (2000) Impact factors: use and abuse. *Perspect. Publish.*, **1**, 1–6.
- Archambault, E. and Larivière, V. (2009) History of the journal impact factor: contingencies and consequences. *Scientometrics*, **79**, 635–649.
- Arnold, D. N. and Fowler, K. K. (2011) Nefarious numbers. *Not. Am. Math. Soc.*, **58**, 434–437.
- Bergstrom, C. (2007) Eigenfactor: measuring the value and the prestige of scholarly journals. *Coll. Res. Lib. News*, **68**, 314–316.
- Bishop, M. and Bird, C. (2007) BIB’s first impact factor is 24.37. *Brief. Bioinform.*, **8**, 207.
- Bollen, J., Rodriguez, M. A. and de Sompel, H. V. (2006) Journal status. *Scientometrics*, **69**, 669–687.
- Bornmann, L. (2014) Do altmetrics point to the broader impact of research?: an overview of benefits and disadvantages of altmetrics. *J. Informetr.*, **8**, 895–903.
- Bornmann, L. and Marx, W. (2014) How to evaluate individual researchers working in the natural and life sciences meaningfully?: a proposal of methods based on percentiles of citations. *Scientometrics*, **98**, 487–509.
- Boyack, K. W., Klavans, R. and Börner, K. (2005) Mapping the backbone of science. *Scientometrics*, **64**, 351–374.

- Bradley, R. A. and Terry, M. E. (1952) The rank analysis of incomplete block designs: I, The method of paired comparisons. *Biometrika*, **39**, 324–345.
- Braun, T., Glänzel, W. and Schubert, A. (2006) A Hirsch-type index for journals. *Scientometrics*, **69**, 169–173.
- Brin, S. and Page, L. (1998) The anatomy of a large-scale hypertextual web search engine. *Comput. Netwks ISDN Syst.*, **30**, 107–117.
- Carpenter, M. P. and Narin, F. (1973) Clustering of scientific journals. *J. Am. Soc. Inform. Sci.*, **24**, 425–436.
- Cattelan, M. (2012) Models for paired comparison data: a review with emphasis on dependent data. *Statist. Sci.*, **27**, 412–433.
- Cattelan, M., Varin, C. and Firth, D. (2013) Dynamic Bradley–Terry modelling of sports tournaments. *Appl. Statist.*, **62**, 135–150.
- Chen, K.-M., Jen, T.-H. and Wu, M. (2014) Estimating the accuracies of journal impact factor through bootstrap. *J. Inform.*, **8**, 181–196.
- David, H. A. (1963) *The Method of Paired Comparisons*. New York: Hafner.
- Fahrmeir, L. and Tutz, G. (1994) Dynamic stochastic models for time-dependent ordered paired comparison systems. *J. Am. Statist. Ass.*, **89**, 1438–1449.
- Fan, J. and Li, R. (2001) Variable selection via nonconcave penalized likelihood and its oracle properties. *J. Am. Statist. Ass.*, **96**, 1348–1360.
- Firth, D. (2012) qvcalc: quasi variances for factor effects in statistical models. *R Package Version 0.8-8*. (Available from CRAN.R-project.org/package=qvcalc.)
- Firth, D. and de Menezes, R. X. (2005) Quasi-variances. *Biometrika*, **91**, 65–80.
- Franceschet, M. (2010) Ten good reasons to use the Eigenfactor metrics. *Inform. Process. Mangmnt*, **46**, 555–558.
- Frandsen, T. F. (2007) Journal self-citations—analysing the JIF mechanism. *J. Informetr.*, **1**, 47–58.
- Garfield, E. (1955) Citation indices for Science. *Science*, **122**, 108–111.
- Garfield, E. (1972) Citation analysis as a tool in journal evaluation. *Science*, **178**, 471–479.
- Glänzel, W. and Moed, H. F. (2002) Journal impact measures in bibliometric research. *Scientometrics*, **53**, 171–193.
- Glickman, M. E. (1999) Parameter estimation in large dynamic paired comparison experiments. *Appl. Statist.*, **48**, 377–394.
- Goldstein, H. and Spiegelhalter, D. J. (1996) League tables and their limitations: statistical issues in comparisons of institutional performance (with discussion). *J. R. Statist. Soc. A*, **159**, 385–443.
- Gross, P. L. K. and Gross, E. M. (1927) College libraries and chemical education. *Science*, **66**, 385–389.
- Hall, P. G. (2009) Comment: Citation statistics. *Statist. Sci.*, **24**, 25–26.
- Hall, P. G. (2011) ‘Ranking our excellence,’ or ‘assessing our quality,’ or whatever.... *Inst. Math. Statist. Bull.*, **40**, 12–14.
- Hall, P. and Miller, H. (2009) Using the bootstrap to quantify the authority of an empirical ranking. *Ann. Statist.*, **37**, 3929–3959.
- Hall, P. and Miller, H. (2010) Modeling the variability of rankings. *Ann. Statist.*, **38**, 2652–2677.
- Institute of Electrical and Electronics Engineers Board of Directors (2013) IEEE position statement on ‘Appropriate use of bibliometric indicators for the assessment of journals, research proposals, and individuals’. Institute of Electrical and Electronics Engineers.
- Journal-Ranking.com (2007) *Present Ranking Endeavors*. Red Jasper. (Available from www.journal-ranking.com/ranking/web/content/intro.html.)
- Jump, P. (2014) Light dose of metrics could ease REF pain. *Times Higher Educ.*, no. 2178, Nov. 13th, 11. (Available from www.timeshighereducation.co.uk/news/regular-diet-of-metrics-lite-may-make-full-ref-more-palatable/2016912.article.)
- Kaufman, L. and Rousseeuw, P. J. (1990) *Finding Groups in Data: an Introduction to Cluster Analysis*. Hoboken: Wiley.
- Knorr-Held, L. (2000) Dynamic rating of sports teams. *Statistician*, **49**, 261–276.
- Lehmann, S., Lautrup, B. E. and Jackson, A. D. (2009) Comment: Citation statistics. *Statist. Sci.*, **24**, 17–20.
- Leydesdorff, L. (2004) Clusters and maps of science based on bi-connected graphs in Journal Citation Reports. *J. Documentn*, **60**, 371–427.
- Leydesdorff, L. and Bornmann, L. (2011) How fractional counting of citations affects the impact factor: normalization in terms of differences in citation potentials among fields of science. *J. Am. Soc. Inform. Sci. Technol.*, **62**, 217–229.
- Leydesdorff, L. and Opthof, T. (2010) Scopus’ Source Normalized Impact per Paper (SNIP) versus the Journal Impact Factor based on fractional counting of citations. *J. Am. Soc. Inform. Sci. Technol.*, **61**, 2365–2369.
- Leydesdorff, L., Radicchi, F., Bornmann, L., Castellano, C. and de Nooy, W. (2013) Field-normalized impact factors (IFs): a comparison of rescaling and fractionally counted IFs. *J. Am. Soc. Inform. Sci. Technol.*, **64**, 2299–2309.
- Liner, G. H. and Amin, M. (2004) Methods of ranking economics journals. *Atl. Econ. J.*, **32**, 140–149.
- Liu, X., Glänzel, W. and De Moor, B. (2012) Optimal and hierarchical clustering of large-scale hybrid networks for scientific mapping. *Scientometrics*, **91**, 473–493.
- Marx, W. and Bornmann, L. (2013) Journal impact factor: ‘the poor man’s citation analysis’ and alternative approaches. *Eur. Sci. Editing*, **39**, 62–63.

- Masarotto, G. and Varin, C. (2012) The ranking lasso and its application to sport tournaments. *Ann. Appl. Statist.*, **6**, 1949–1970.
- Moed, H. F. (2010) Measuring contextual citation impact of scientific journals. *J. Informetr.*, **4**, 265–277.
- Morris, C. N. (1983) Parametric empirical Bayes inference: theory and applications. *J. Am. Statist. Ass.*, **78**, 47–65.
- van Nierop, E. (2009) Why do statistics journals have low impact factors? *Statist. Neerland.*, **63**, 52–62.
- Nocedal, J. and Wright, S. J. (2006) *Numerical Optimization*, 2nd edn. New York: Springer.
- van Noorden, R. (2012) Researchers feel pressure to cite superfluous papers. *Nat. News*, Feb. 12th.
- Palacios-Huerta, I. and Volij, O. (2004) The measurement of intellectual influence. *Econometrica*, **72**, 963–977.
- Pratelli, L., Baccini, A., Barbaresi, L. and Marcheselli, M. (2012) Statistical analysis of the Hirsch Index. *Scand. J. Statist.*, **39**, 681–694.
- Putirka, K., Kunz, M., Swainson, I. and Thomson, J. (2013) Journal Impact Factors: their relevance and their influence on society-published scientific journals. *Am. Mineral.*, **98**, 1055–1065.
- R Core Team (2015) *R: a Language and Environment for Statistical Computing*. Vienna: R Foundation for Statistical Computing.
- Ritzberger, K. (2008) A ranking of journals in economics and related fields. *Germ. Econ. Rev.*, **9**, 402–430.
- Seglen, P. O. (1997) Why the impact factor of journals should not be used for evaluating research. *Br. Med. J.*, **314**, 498–502.
- Sevinc, A. (2004) Manipulating impact factor: an unethical issue or an editor's choice? *Swiss Med. Wkly*, **134**, 410.
- Silverman, B. W. (2009) Comment: Citation statistics. *Statist. Sci.*, **24**, 21–24.
- Spiegelhalter, D. J., Best, N. G., Carlin, B. P. and van der Linde, A. (2002) Bayesian measures of model complexity and fit (with discussion). *J. R. Statist. Soc. B*, **64**, 583–639.
- Stigler, S. M. (1994) Citation patterns in the journals of statistics and probability. *Statist. Sci.*, **9**, 94–108.
- Stigler, G. J., Stigler, S. M. and Friedland, C. (1995) The journals of economics. *J. Polit. Econ.*, **103**, 331–359.
- Takeuchi, K. (1976) Distribution of informational statistics and a criterion of model fitting (in Japanese). *Suri-Kagaku*, **153**, 12–18.
- Theoharakis, V. and Skordia, M. (2003) How do statisticians perceive statistics journals? *Am. Statist.*, **57**, 115–123.
- Tibshirani, R. (1996) Regression shrinkage and selection via the lasso. *J. R. Statist. Soc. B*, **58**, 267–288.
- Tibshirani, R., Saunders, M., Rosset, S., Zhu, J. and Knight, K. (2005) Sparsity and smoothness via the fused lasso. *J. R. Statist. Soc. B*, **67**, 91–108.
- Turner, H. and Firth, D. (2012) Bradley-Terry models in R: the `BradleyTerry2` package. *J. Statist. Softw.*, **48**, 1–21.
- Waltman, L. and Van Eck, N. J. (2013) Source normalized indicators of citation impact: an overview of different approaches and an empirical comparison. *Scientometrics*, **96**, 699–716.
- Waltman, L., van Eck, J. N., van Leeuwen, T. N. and Visser, M. S. (2013) Some modifications to the SNIP journal impact indicator. *J. Informetr.*, **7**, 272–285.
- Wedderburn, R. W. M. (1974) Quasi-likelihood, generalized linear models, and the Gauss-Newton method. *Biometrika*, **61**, 439–447.
- West, J. D. (2010) Eigenfactor: ranking and mapping scientific knowledge. *PhD Dissertation*. University of Washington, Seattle.
- Wilhite, A. W. and Fong, E. A. (2012) Coercive citation in academic publishing. *Science*, **335**, 542–543.
- Zitt, M. and Small, H. (2008) Modifying the journal impact factor by fractional citation weighting: the audience factor. *J. Am. Soc. Inform. Sci. Technol.*, **59**, 1856–1860.
- Zou, H. (2006) The adaptive lasso and its oracle properties. *J. Am. Statist. Ass.*, **101**, 1418–1429.

Discussion on the paper by Varin, Cattelan and Firth

David Colquhoun (University College London)

It is a pleasure to propose the vote of thanks for a paper that puts yet another nail in the coffin of the journal impact factor (JIF).

There are two classes of reasons to deplore JIFs. One is that they are statistically dubious, and that is what Varin and his colleagues develop. It has been obvious for a long time that it is statistically illiterate to characterize very skew distributions by their mean. And it is statistically illiterate to present point estimates with no indication of their uncertainty. The existence of so many different methods for ranking journals, each of which gives different answers, renders them useless.

Supporting information

Additional 'supporting' may be found in the on-line version of this article:

'Supplement to "Statistical modelling of citation exchange between statistics journals"'.