



Policy Seminar 2018

Dipartimento di Economia e Management

“Dati segreti e classifiche pseudoscientifiche: la via italiana alla valutazione della ricerca”

Alberto Baccini – Università degli Studi di Siena



ROARS
Return On Academic ReSearch



Sommario

1. Valutazione della ricerca: assetto istituzionale
2. Valutazione della ricerca: lo stato dell'arte nel 2011
3. La via italiana alla valutazione della ricerca
4. Cronaca di un esperimento annunciato
5. Bibliometrics vs peer review: do they agree?
6. Concordanza o fallacia statistica?
7. Dati chiusi, concordanza non replicabile
8. Conclusioni

2. Valutazione della ricerca: assetto istituzionale

Le istituzioni dell'economia pianificata della ricerca?

Riforma Gelmini e successivi governi di centrosinistra hanno introdotto strumenti di governo e di controllo a distanza per le università ed i docenti

ANVUR non è un'agenzia autonoma/authority (come AGCM o AGCOM)

ANVUR non è un *quango* (quasi-autonomous non-governmental organization) che agisce a distanza dal governo (come HEFCE/ora Research England in UK)

ANVUR è un'agenzia governativa:

- (i) il suo consiglio è costituito da sette membri nominati dal ministro;
- (ii) agisce realizzando attività direttamente definite con decreti ministeriali.

Tra le istituzioni europee simili, come AERES in Francia o ANECA in Spagna, nessuna concentra così tanti poteri e funzioni.

In nessun altro paese occidentale è stato sviluppato un analogo controllo governativo delle scienze e delle università.

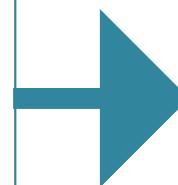
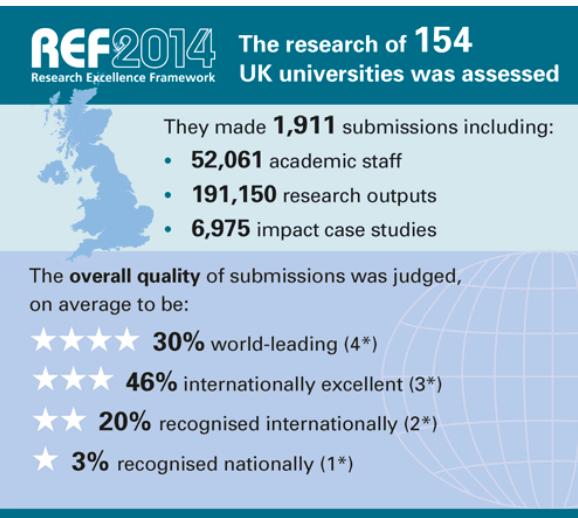
Modello organizzativo simile a quello della scienza nelle economie pianificate.

2. Valutazione della ricerca: lo stato dell'arte nel 2011

Research Excellence Framework

The Research Excellence Framework (REF) is the new system for assessing the quality of research in UK higher education institutions.

The [results](#) of the 2014 REF were published on 18 December 2014.



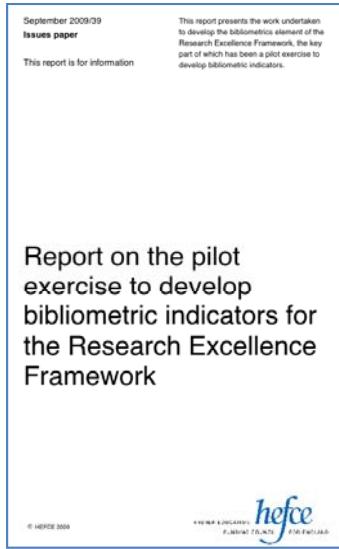
September 2009/39 Issues paper

This report is for information

This report presents the work undertaken to develop the bibliometrics element of the Research Excellence Framework, the key part of which has been a pilot exercise to develop bibliometric indicators.

September 2009/39

Report on the pilot exercise to develop bibliometric indicators for the Research Excellence Framework



Key points

8. Bibliometrics are not sufficiently robust at this stage to be used formulaically or to replace expert review in the REF. However there is considerable scope for citation information to be used to inform expert review.



The Australian Way

the-australian-way.de

ERA 2010: Ranking of Journals

The screenshot shows the official website for the Australian Research Council's ERA 2010 ranking of journals. The header features the Australian Government Australian Research Council logo, a large stylized 'ERA' logo, and links for ERA Login, RMS/GAMS Login, Site Map, and Contacts. A search bar is also present. The main navigation menu includes Home, About ARC, National Competitive Grants Program, Information for Applicants, Media, General Information, and Research Excellence. Below the menu is a banner with the word 'RESEARCH EXCELLENCE' and a collage of scientific images. The left sidebar contains links for Contacts, Current consultations, ERA 2010, ERA Newsletter, Frequently Asked Questions, Liaison Officers, Media Releases, Presentations, and Trial 2009. The main content area discusses the 'Tiers for the Australian Ranking of Journals' based on the 'Overall criterion: Quality of the papers'. It defines Tier A* as the best in its field, Tier A as having very high quality, Tier B as having solid but not outstanding reputation, and Tier C as peer-reviewed journals that do not meet higher tiers.

You are here: Home > Research Excellence > ERA 2010

Tiers for the Australian Ranking of Journals

Overall criterion: Quality of the papers

A*

Typically an A* journal would be one of the best in its field or subfield in which to publish and would typically cover the entire field/subfield. Virtually all papers they publish will be of a very high quality. These are journals where most of the work is important (it will really shape the field) and where researchers boast about getting accepted. Acceptance rates would typically be low and the editorial board would be dominated by field leaders, including many from top institutions.

A

The majority of papers in a Tier A journal will be of very high quality. Publishing in an A journal would enhance the author's standing, showing they have real engagement with the global research community and that they have something to say about problems of some significance. Typical signs of an A journal are lowish acceptance rates and an editorial board which includes a reasonable fraction of well known researchers from top institutions.

B

Tier B covers journals with a solid, though not outstanding, reputation. Generally, in a Tier B journal, one would expect only a few papers of very high quality. They are often important outlets for the work of PhD students and early career researchers. Typical examples would be regional journals with high acceptance rates, and editorial boards that have few leading researchers from top international institutions.

C

Tier C includes quality, peer reviewed, journals that do not meet the criteria of the higher tiers.

Content Last Modified: 25/09/09

30 maggio 2011

Kim Carr: «*There is clear and consistent evidence that the rankings were being deployed inappropriately within some quarters of the sector, in ways that could produce harmful outcomes [...]. [...] the removal of the ranks and the provision of the publication profile will ensure they will be used descriptively rather than prescriptively.*»



Kim Carr, the Australian Minister for Innovation, Industry, Science and Research



House of Commons
Science and Technology
Committee

Peer review in scientific publications

Eighth Report of Session 2010–12

Volume I: Report, together with formal minutes, oral and written evidence

Additional written evidence is contained in Volume II, available on the Committee website at www.parliament.uk/science

*Ordered by the House of Commons
to be printed 18 July 2011*

David Sweeney [Director HEFCE]: «*it is an underpinning element in the exercise that journal impact factors will not be used. I think we were very interested to see that in Australia, where they conceived an exercise that was heavily dependent on journal rankings, after carrying out the first exercise, they decided that alternative ways of assessing quality»*



International
Mathematical
Union
(IMU)



Joint Committee on Quantitative Assessment of Research

Citation Statistics

A report from the International Mathematical Union (IMU) in cooperation with the International Council of Industrial and Applied Mathematics (ICIAM) and the Institute of Mathematical Statistics (IMS)

Corrected version,
6/12/08

*“The idea that research assessment must be done using “simple and objective” methods is increasingly prevalent today. The “simple and objective” methods are broadly interpreted as bibliometrics, that is, citation data and the statistics derived from them. There is a belief that citation statistics are inherently more accurate because they substitute simple numbers for complex judgments, and hence overcome the possible subjectivity of peer review. But **this belief is unfounded.**”*



17 gennaio 2011

Du bon usage de la bibliométrie
pour l'évaluation individuelle des chercheurs

“Any bibliometric evaluation should be tightly associated to a close examination of a researcher’s work, in particular to evaluate its originality, an element that cannot be assessed through a bibliometric study.”

2. VQR, la via italiana alla valutazione della ricerca



Agenzia Nazionale di Valutazione del
sistema Universitario e della Ricerca

National Agency for the Evaluation of
Universities and Research Institutes

Bibliometria o peer review?

L'ANVUR ha indirizzato gli esperti verso una tecnica mista di strumenti bibliometrici e di *peer review*. In certe aree, quelle per le quali esistono dati citazionali certificati, è stata usata prevalentemente la bibliometria. Una bibliometria che ha fatto uso dell'impact factor della rivista sul quale un articolo di ricerca era stato pubblicato e del numero di citazioni ricevuto dall'articolo stesso, secondo una combinazione ben definita, riservando la *peer review* ai soli casi in cui i due indicatori fornivano risultati inconciliabili tra loro. Così è stato fatto per le aree di matematica e informatica, fisica, chimica, scienze della terra, agraria e veterinaria, biologia, medicina, ingegneria civile ed architettura, ingegneria industriale e dell'informazione, con la parziale eccezione di architettura. Mentre per le scienze umane, quali le scienze dell'antichità, filologiche e letterarie e storico artistiche, le scienze storico, filosofiche e psicologiche , la giurisprudenza, le scienze politiche e sociologiche e, in parte, le scienze economiche e statistiche, per le quali non sono disponibili dati citazionali, gli esperti hanno fatto ricorso a referee esterni.

Stefano Fantoni et Alessio Ancaiani, « La valutazione dei prodotti della ricerca ed il suo impatto sul sistema universitario e degli enti di ricerca: il caso italiano », *Mélanges de la Casa de Velázquez [En ligne]*, 44-2 | 2014, mis en ligne le 01 janvier 2018, consulté le 09 avril 2018. URL : <http://journals.openedition.org/mcv/5864>

Il “mix valutativo” della VQR 2004-2010

- Inedito metodo bibliometrico:

		Bibliometric Indicator			
		1	2	3	4
n. of citations	1	A	IR	IR	IR
	2	A	B	C	D
	3	A	B	C	D
	4	IR	IR	IR	D

Figure 2. The Bibliometric matrix.
Source: ANVUR.

- Informed peer review

Ma è lecito mescolare peer review e bibliometria?

National Agency for the Evaluation of
Universities and Research Institutes



Agenzia Nazionale di Valutazione del
sistema Universitario e della Ricerca

Evaluation of Research Quality



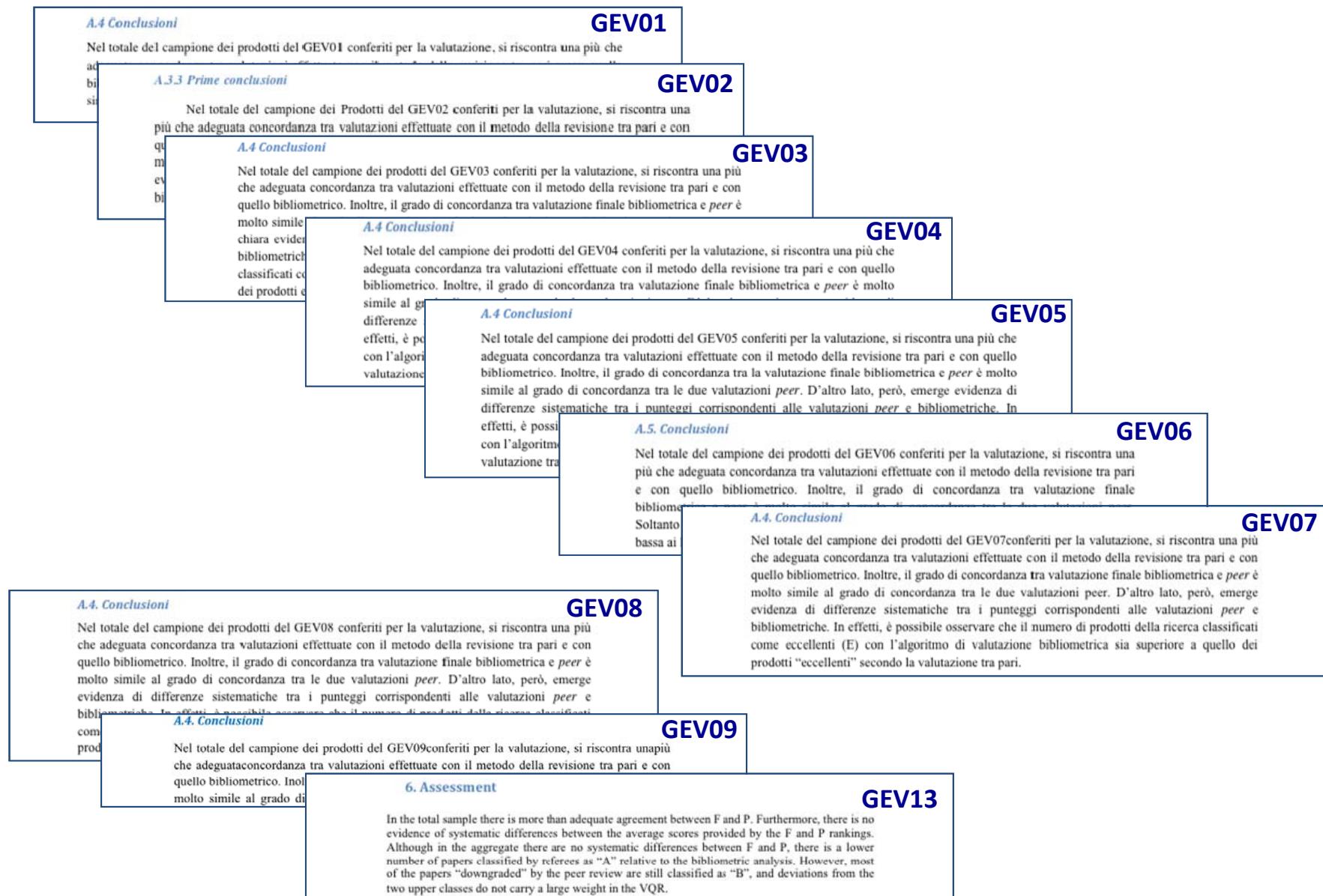
Valutazione Qualità della Ricerca

Appendice B. Il confronto tra valutazione *peer* e valutazione bibliometrica

I GEV che hanno utilizzato gli indicatori bibliometrici per la valutazione degli articoli indicizzati in ISI WoS e Scopus hanno selezionato, con un algoritmo di estrazione casuale in grado di garantire una buona copertura statistica di tutti i sub-GEV, un numero pari a circa il 10% degli articoli valutati bibliometricamente e li hanno sottoposti alla valutazione *peer*. L'obiettivo era un confronto tra le due metodologie di valutazione applicate allo stesso campione di articoli, per valutare il grado di corrispondenza dei risultati. Nel seguito, saranno presentati i risultati in forma sintetica e aggregata. Per confronti più puntuali si rimanda alla lettura dell'appendice apposita dei rapporti di area.

3. Cronaca di un esperimento annunciato

Conclusioni tutte uguali



Conclusioni tutte uguali

“Nel totale del campione dei prodotti del GEV_X conferiti per la valutazione, si riscontra una più che adeguata concordanza tra valutazioni effettuate con il metodo della revisione tra pari e con quello bibliometrico.”

Conclusioni tutte uguali ... o quasi

National Agency for the Evaluation of
Universities and Research Institutes



Agenzia Nazionale di Valutazione del
sistema Universitario e della Ricerca

Evaluation of Research Quality



Valutazione Qualità della Ricerca

A.4. Conclusioni

Nel totale del campione dei prodotti del GEV09 conferiti per la valutazione, si riscontra unapiù che adeguataconcordanza tra valutazioni effettuate con il metodo della revisione tra pari e con quello bibliometrico. Inoltre, il grado di concordanza tra valutazione finale bibliometrica e *peer* è molto simile al grado di concordanza tra le due valutazioni *peer*. D’altro lato, però, emerge evidenza di differenze sistematiche tra i punteggi corrispondenti alle valutazioni *peere* bibliometriche. In effetti, è possibile osservare che il numero di prodotti della ricerca classificati come eccellenti (E) con l’algoritmo di valutazione bibliometrica sia superiore a quello dei prodotti “eccellenti” secondo la valutazione tra pari.

Il grado di concordanza tra valutazioni *peer* e valutazioni bibliometriche è moderato in quasi tutti i macro-settori, risulta invece piuttosto elevato in Ingegneria informatica. Le differenze sistematiche tra i punteggi medi sono statisticamente significative e sempre di segno positivo (ossia, la valutazione bibliometrica è significativamente più favorevole in media rispetto a quella *peer*).

Facciamo uno zoom sul Rapporto di Area 09

National Agency for the Evaluation of
Universities and Research Institutes

anvur
Agenzia Nazionale di Valutazione del
sistema Universitario e della Ricerca

Evaluation of Research Quality

vQr
Valutazione Qualità della Ricerca

A.4. Conclusioni

Nel totale del campione dei prodotti del GEV09conferiti per la valutazione, si riscontra unapiù che adeguataconcordanza tra valutazioni effettuate con il metodo della revisione tra pari e con quello bibliometrico. Inoltre, il grado di concordanza tra valutazione finale bibliometrica e *peer* è molto simile al grado di concordanza tra le due valutazioni *peer*. D'altro lato, però, emerge evidenza di differenze sistematiche tra i punteggi corrispondenti alle valutazioni *peer* bibliometriche. In effetti, è possibile osservare che il numero di prodotti della ricerca classificati come eccellenti (E) con l'algoritmo di valutazione bibliometrica sia superiore a quello dei prodotti "eccellenti" secondo la valutazione tra pari.

Il grado di concordanza tra valutazioni *peer* e valutazioni bibliometriche è moderato in quasi tutti i macro-settori, risulta invece piuttosto elevato in Ingegneria informatica. Le differenze sistematiche tra i punteggi medi sono statisticamente significative e sempre di segno positivo (ossia, la valutazione bibliometrica è significativamente più favorevole in media rispetto a quella *peer*).

Nel totale del campione dei prodotti del GEV09conferiti per la valutazione, si riscontra unapiù che adeguataconcordanza tra valutazioni effettuate con il metodo della revisione tra pari e con quello bibliometrico. Inoltre, il grado di concordanza tra valutazione finale bibliometrica e *peer* è

Rapporto di Area 09

National Agency for the Evaluation of Universities and Research Institutes
anvur
Agenzia Nazionale di Valutazione del sistema Universitario e della Ricerca

Evaluation of Research Quality
vQr
Valutazione Qualità della Ricerca

A.4. Conclusioni

Nel totale del campione dei prodotti del GEV09conferiti per la valutazione, si riscontra unapiù che adeguataconcordanza tra valutazioni effettuate con il metodo della revisione tra pari e con quello bibliometrico. Inoltre, il grado di concordanza tra valutazione finale bibliometrica e *peer* è molto simile al grado di concordanza tra le due valutazioni *peer*. D'altro lato, però, emerge evidenza di differenze sistematiche tra i punteggi corrispondenti alle valutazioni *peer* bibliometriche. In effetti, è possibile osservare che il numero di prodotti della ricerca classificati come eccellenti (E) con l'algoritmo di valutazione bibliometrica sia superiore a quello dei prodotti "eccellenti" secondo la valutazione tra pari.

Il grado di concordanza tra valutazioni *peer* e valutazioni bibliometriche è moderato in quasi tutti i macro-settori, risulta invece piuttosto elevato in Ingegneria informatica. Le differenze sistematiche tra i punteggi medi sono statisticamente significative e sempre di segno positivo (ossia, la valutazione bibliometrica è significativamente più favorevole in media rispetto a quella *peer*).

Nel totale del campione dei prodotti del GEV09conferiti per la valutazione, si riscontra unapiù che adeguataconcordanza tra valutazioni effettuate con il metodo della revisione tra pari e con quello bibliometrico. Inoltre, il grado di concordanza tra valutazione finale bibliometrica e *peer* è

ma la concordanza è **più che adeguata o moderata?**

Il grado di concordanza tra valutazioni *peer* e valutazioni bibliometriche è moderato in quasi tutti i macro-settori, risulta invece piuttosto elevato in Ingegneria informatica. Le differenze

Facciamo uno zoom sul Rapporto di Area 09

National Agency for the Evaluation of Universities and Research Institutes

Agenzia Nazionale di Valutazione del sistema Universitario e della Ricerca

Evaluation of Research Quality

Valutazione Qualità della Ricerca

A.4. Conclusioni

Nel totale del campione dei prodotti del GEV09conferiti per la valutazione, si riscontra unapiù che adeguataconcordanza tra valutazioni effettuate con il metodo della revisione tra pari e con quello bibliometrico. Inoltre, il grado di concordanza tra valutazione finale bibliometrica e *peer* è molto simile al grado di concordanza tra le due valutazioni *peer*. D'altro lato, però, emerge evidenza di differenze sistematiche tra i punteggi corrispondenti alle valutazioni *peer* bibliometriche. In effetti, è possibile osservare che il numero di prodotti della ricerca classificati come eccellenti (E) con l'algoritmo di valutazione bibliometrica sia superiore a quello dei prodotti "eccellenti" secondo la valutazione tra pari.

Il grado di concordanza tra valutazioni *peer* e valutazioni bibliometriche è moderato in quasi tutti i macro-settori, risulta invece piuttosto elevato in Ingegneria informatica. Le differenze sistematiche tra i punteggi medi sono statisticamente significative e sempre di segno positivo (ossia, la valutazione bibliometrica è significativamente più favorevole in media rispetto a quella *peer*).

Nel totale del campione dei prodotti del GEV09conferiti per la valutazione, si riscontra unapiù che adeguataconcordanza tra valutazioni effettuate con il metodo della revisione tra pari e con quello bibliometrico. Inoltre, il grado di concordanza tra valutazione finale bibliometrica e *peer* è

Mancano degli spazi.

Non è che il rapporto dell'area 09 (**quella con la concordanza peggiore**), ha subito una correzione "last minute" per uniformarlo agli altri rapporti, con una sostituzione che richiedeva più caratteri?

Un rapporto, molti working papers e
anche un articolo scientifico



Appendice B. Il confronto tra valutazione *peer* e valutazione bibliometrica

I GEV che hanno utilizzato gli indicatori bibliometrici per la valutazione degli articoli indicizzati in ISI WoS e Scopus hanno selezionato, con un algoritmo di estrazione casuale in grado di garantire una buona copertura statistica di tutti i sub-GEV, un numero pari a circa il 10% degli articoli valutati bibliometricamente e li hanno sottoposti alla valutazione *peer*. L'obiettivo era un confronto tra le due metodologie di valutazione applicate allo stesso campione di articoli, per valutare il grado di corrispondenza dei risultati. Nel seguito, saranno presentati i risultati in forma sintetica e aggregata. Per confronti più puntuali si rimanda alla lettura dell'appendice apposita dei rapporti di area.

B.1 Il campionamento statistico

Un campione casuale di 9199 articoli su rivista passibili di valutazione bibliometrica è stato estratto dalla popolazione di 99.005 articoli, valutabili bibliometricamente e sottomessi alla valutazione nei GEV che hanno utilizzato indicatori bibliometrici. La popolazione è stata stratificata in base alla distribuzione dei prodotti all'interno dei sub-GEV individuati nelle varie Aree. Ai fini della stratificazione, gli articoli sono stati attribuiti ai sub-GEV sulla base del settore scientifico-disciplinare (SSD) nel quale sono stati valutati, escludendo i casi di articoli duplicati presentati da diversi autori all'interno di uno stesso strato campionario. Complessivamente, il campione include il 9,3% degli articoli sottoposti a valutazione bibliometrica nelle Aree "bibliometriche". L'estrazione è stata effettuata ai primi di settembre 2012, prima dell'inizio del processo di revisione *peer*, mediante una procedura casuale con il vincolo di selezionare una proporzione significativa di prodotti in ciascun sub-GEV. La Tabella B.1 riporta l'elenco dei GEV bibliometrici e, per ciascuno di essi, la dimensione della popolazione e del campione estratto in valori assoluti e in percentuale sulla popolazione.

Bibliometric and peer review methods for research evaluation: a methodological appraisement

Tindaro Cicero and Marco Malgarini and Carmela Anna Nappi and Franco Peracchi

ANVUR, ANVUR, ANVUR, Department of Economic and Finance, University of Rome Tor Vergata and EIEF

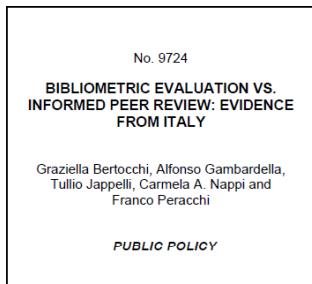
Online at <https://mpra.ub.uni-muenchen.de/50470/>
MPRA Paper No. 50470, posted 8 October 2013 19:30 UTC

2. Il campione statistico

Un campione casuale di 9.199 articoli su rivista passibili di valutazione bibliometrica è stato estratto dalla popolazione di 99.005 articoli valutabili bibliometricamente e sottomessi alla valutazione nelle cosiddette "arie bibliometriche", cioè nelle aree scientifiche che hanno utilizzato indicatori bibliometrici (scienze matematiche e informatiche, scienze fisiche, scienze chimiche, scienze della terra, scienze biologiche, scienze mediche, scienze agrarie e veterinarie, ingegneria civile e architettura, ingegneria industriale e dell'informazione, e scienze economiche e stistiche). La popolazione è stata stratificata in base alla distribuzione dei prodotti all'interno dei settori individuati nelle varie aree. Ai fini della stratificazione, gli articoli sono stati attribuiti ai settori sulla base del settore scientifico-disciplinare (SSD) nel quale sono stati valutati, eliminando le duplicazioni dovute alla presentazione di uno stesso articolo da parte di autori diversi all'interno di uno stesso strato campionario. Complessivamente, il campione include il 9,3% degli articoli sottoposti a valutazione bibliometrica nelle aree bibliometriche. L'estrazione è stata effettuata nel settembre 2012, prima dell'inizio del processo di revisione *peer*, mediante una procedura casuale

¹ Una precedente versione del lavoro è stata pubblicata come Appendice del Rapporto Finale ANVUR sulla Valutazione della Qualità della Ricerca 2004-2010, disponibile all'indirizzo <http://www.anvur.org/rapporto/>. Gli autori ringraziano il Professor Sergio Benedetto, coordinatore della VQR, per le numerose utili discussioni avute nel corso del lavoro. Un sentito ringraziamento va anche ai tecnici del CINECA che hanno messo a disposizione i dati. Ogni eventuale errore rimane ovviamente di esclusiva responsabilità degli autori.

DISCUSSION PAPER SERIES



WORKING PAPER NO. 344

Bibliometric Evaluation vs. Informed Peer Review: Evidence from Italy

Graziella Bertocchi, Alfonso Gambardella, Tullio Jappelli,
Carmela A. Nappi and Franco Peracchi

October 2013



DEMB Working Paper Series

N. 20

Bibliometric Evaluation vs. Informed Peer Review: Evidence from Italy

Graziella Bertocchi¹,
Alfonso Gambardella²,
Tullio Jappelli³,
Carmela A. Nappi⁴,
Franco Peracchi⁵

October 2013



DISCUSSION PAPER SERIES

IZA DP No. 7739

Bibliometric Evaluation vs. Informed Peer Review:
Evidence from Italy

Graziella Bertocchi
Alfonso Gambardella
Tullio Jappelli

Carmela A. Nappi
Franco Peracchi

November 2013



WORKING PAPER SERIES

Bibliometric Evaluation vs. Informed Peer Review: Evidence from Italy

Graziella Bertocchi, Alfonso Gambardella, Tullio
Jappelli, Carmela A. Nappi, Franco Peracchi

Working Paper 93

October 2013



Bibliometric evaluation vs. informed peer review: Evidence from Italy[†]

Graziella Bertocchi^a, Alfonso Gambardella^b, Tullio Jappelli^{c,d*}, Carmela A. Nappi^d,
Franco Peracchi^e

^aDepartment of Economics "Marco Biagi", University of Modena and Reggio Emilia, Viale Berengario, 51, 41121 Modena, Italy

^bDepartment of Management & Technology and CROS, Bocconi University, Via Rovigo, 1, 20136 Milan, Italy

^cDepartment of Economics and Statistics and CSEF, University of Naples Federico II, Via Cinthia, 21, 80126 Napoli, Italy

^dANVUR, Piazza Kennedy, 20, 00144 Rome, Italy

^eDepartment of Economics and Finance, University of Rome Tor Vergata, Via Columbia, 2, 00133 Rome, Italy

ARTICLE INFO
Article history:
Received 29 October 2013
Received in revised form 28 July 2014
Accepted 18 August 2014
Available online 18 September 2014

Keywords:
Research assessment
Informed peer review
Bibliometric evaluation
VQR

ABSTRACT

A relevant question for the organization of large-scale research assessments is whether bibliometric evaluation and informed peer review yield similar results. In this paper, we draw on the experience of the panel that evaluated Italian research in Economics, Management and Statistics during the national assessment exercise (VQR) relative to the period 2004–2010. We exploit the unique opportunity of studying a sample of 590 journal articles randomly drawn from a pool of nearly 5681 journal articles (out of nearly 12,000 journal titles) and measure the agreement between panel evaluations by bibliometric analysis and by informed peer review. In the total sample we find fair to good agreement between informed peer review and bibliometric analysis and absence of statistical bias between the two. We then discuss the nature, implications, and limitations of this correlation.

© 2014 Elsevier B.V. All rights reserved.

1. Introduction

Measuring research quality is a topic of growing interest to universities and research institutions. It has become a central issue in relation to the efficient allocation of public resources which, in many countries and especially in Europe, represent the main component of university funding. In the recent past, a number of countries – Austria, France, Italy, Netherlands, Scandinavia, UK – have introduced national assessment exercises to gauge the quality of academic research. We have also seen

a new trend in the way funds are being allocated to higher education in Europe, on the basis not only of actual costs but also, to promote excellence, academic performance. Examples of performance-based university research funding systems (OECD, 2010; Hicks, 2012; Reborn and Turrini, 2013) include the British Research Excellence Framework (REF) and the Italian Evaluation of Research Quality. Performance-based funding, however, comes with significant costs in terms of time and resources, and such costs may differ considerably across evaluation methods (Geuna and Martin, 2003; Martin, 2011).

The main criteria for evaluating research performance combine, in various ways, bibliometric indicators (Modell, 2005; Nicolaisen, 2007) and peer review (Bornmann, 2011). Bibliometric indicators

* The authors have been, respectively, president of the panel evaluating Italian

Assessing Italian research quality: A comparison between bibliometric evaluation and informed peer review

Graziella Bertocchi, Alfonso Gambardella, Tullio Jappelli, Carmela Nappi, Franco Peracchi 28 July 2014

Assessing the quality of academic research is important – particularly in countries where universities receive most of their funding from the government. This column presents evidence from an Italian research assessment exercise. Bibliometric analysis – based on the journal in which a paper was published and its number of citations – produced very similar evaluations of research quality to informed peer review. Since bibliometric analysis is less costly, it can be used to monitor research on a more continuous basis and to predict the outcome of future peer-reviewed assessments.



lavoce.info

DOMENICA 5 NOVEMBRE 2017

HOME

ARGOMENTI

DOSSIER

RUBRICHE

[Home](#) > [Argomenti](#) > [Scuola e università](#) > Bibliometria o “peer review” per valutare la ricerca?

SCUOLA E UNIVERSITÀ

Bibliometria o “peer review” per valutare la ricerca?

07.11.13

Graziella Bertocchi, Alfonso Gambardella, Tullio Jappelli, Carmela A. Nappi e Franco Peracchi



4. Bibliometrics vs peer review: do they agree?



Bibliometric evaluation vs. informed peer review: Evidence from Italy[☆]



Graziella Bertocchi^a, Alfonso Gambardella^b, Tullio Jappelli^{c,*}, Carmela A. Nappi^d, Franco Peracchi^e

^a Department of Economics "Marco Biagi", University of Modena and Reggio Emilia, Viale Berengario, 51, 41121 Modena, Italy

^b Department of Management & Technology and CRIOS, Bocconi University, Via Roentgen, 1, 20136 Milan, Italy

^c Department of Economics and Statistics and CSEF, University of Naples Federico II, Via Cinthia, 21, 80126 Napoli, Italy

^d ANVUR, Piazza Kennedy, 20, 00144 Rome, Italy

^e Department of Economics and Finance, University of Rome Tor Vergata, Via Columbia, 2, 00133 Rome, Italy

ARTICLE INFO

Article history:

Received 29 October 2013

Received in revised form 28 July 2014

Accepted 18 August 2014

Available online 18 September 2014

ABSTRACT

A relevant question for the organization of large-scale research assessments is whether bibliometric evaluation and informed peer review yield similar results. In this paper, we draw on the experience of the panel that evaluated Italian research in Economics, Management and Statistics during the national assessment exercise (VQR) relative to the period 2004–2010. We exploit the unique opportunity of studying a sample of 590 journal articles randomly drawn from a population of 5681 journal articles (out of nearly 12,000 journal and non-journal publications), which the panel evaluated both by bibliometric analysis and by informed peer review. In the total sample we find fair to good agreement between informed peer review and bibliometric analysis and absence of statistical bias between the two. We then discuss the nature, implications, and limitations of this correlation.

© 2014 Elsevier B.V. All rights reserved.

Keywords:

Research assessment

Informed peer review

Bibliometric evaluation

VQR

Table 11
Comparison between F and P .

Bibliometric (F)	Peer (P)				Total
	A	B	C	D	
A	98 49.49	72 36.36	19 9.60	9 4.55	198 100.00
B	11 10.78	56 54.90	26 25.49	9 8.82	102 100.00
C	4 3.88	25 24.27	39 37.86	35 33.98	103 100.00
D	3 1.60	21 11.23	45 24.06	118 63.10	187 100.00
Total	116 19.66	174 29.49	129 21.86	171 28.98	590 100.00

Note: The table tabulates the distribution of the journal articles in the sample by informed peer review and bibliometric evaluations, expressed through the merit classes. The elements on the main diagonal correspond to cases for which informed peer review and bibliometric evaluation coincide. The off-diagonal elements correspond to cases of disagreement between informed peer review and bibliometric evaluation.

Cohen's kappa

Cohen's kappa measures the agreement between two raters who each classify N items into C mutually exclusive categories. The first mention of a kappa-like statistic is attributed to Galton (1892);^[2] see Smeeton (1985).^[3]

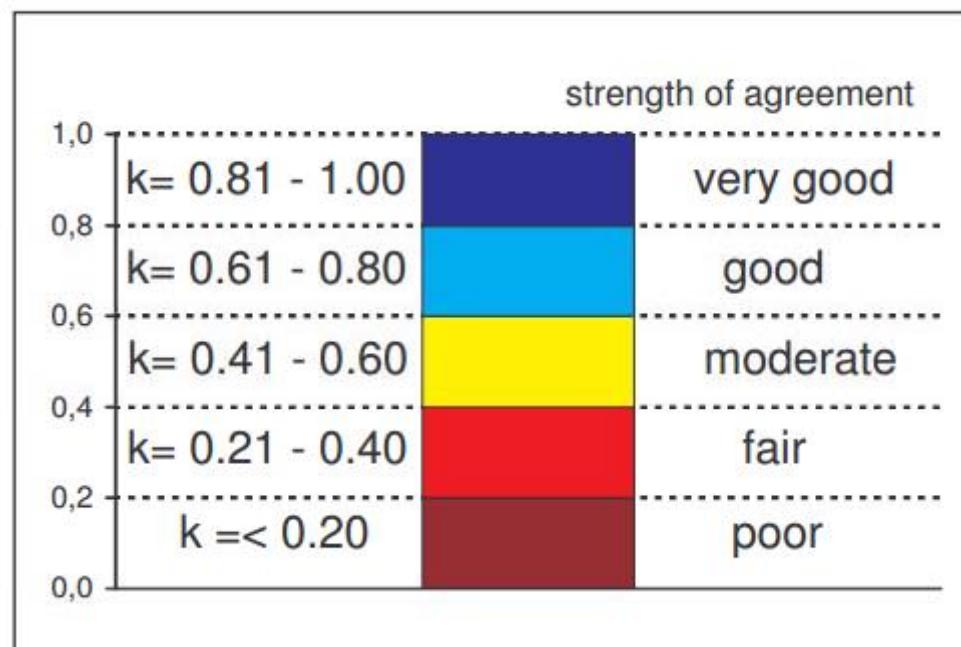
The definition of κ is:

$$\kappa \equiv \frac{p_o - p_e}{1 - p_e} = 1 - \frac{1 - p_o}{1 - p_e},$$

where p_o is the relative observed agreement among raters (identical to [accuracy](#)), and p_e is the hypothetical probability of chance agreement, using the observed data to calculate the probabilities of each observer randomly seeing each category. If the raters are in complete agreement then $\kappa = 1$. If there is no agreement among the raters other than what would be expected by chance (as given by p_e), $\kappa \leq 0$.



<i>K</i> values	Description
Landis and Koch (1977)	
<0.00	Poor
0.00–0.20	Slight
0.21–0.40	Fair
0.41–0.60	Moderate
0.61–0.80	Substantial
0.81–1.00	Almost perfect
Altman (1991)	
<0.20	Poor
0.21–0.40	Fair
0.41–0.60	Moderate
0.61–0.80	Good
0.81–1.00	Very good
Fleiss et al. (2003)	
<0.40	Poor
0.40–0.75	Fair to good
>0.75	Excellent
George and Mallery (2003)	
<0.51	Unacceptable
0.51–0.60	Poor
0.61–0.70	Questionable
0.71–0.80	Acceptable
0.81–0.90	Good
0.91–1.00	Excellent
Stemler and Tsai (2008)	
<0.50	Unacceptable
>0.50	Acceptable



Weighted Cohen's kappa

Weighted kappa [\[edit\]](#)

Weighted kappa lets you count disagreements differently^[15] and is especially useful when codes are ordered.^{[7]:66} Three matrices are involved, the matrix of observed scores, the matrix of expected scores based on chance agreement, and the weight matrix. Weight matrix cells located on the diagonal (upper-left to bottom-right) represent agreement and thus contain zeros. Off-diagonal cells contain weights indicating the seriousness of that disagreement. Often, cells one off the diagonal are weighted 1, those two off 2, etc.

The equation for weighted κ is:

$$\kappa = 1 - \frac{\sum_{i=1}^k \sum_{j=1}^k w_{ij} x_{ij}}{\sum_{i=1}^k \sum_{j=1}^k w_{ij} m_{ij}}$$

where k =number of codes and w_{ij} , x_{ij} , and m_{ij} are elements in the weight, observed, and expected matrices, respectively. When diagonal cells contain weights of 0 and all off-diagonal cells weights of 1, this formula produces the same value of kappa as the calculation given above.



Table 13Kappa statistic for the amount of agreement between *F* and *P* scores.

	Total sample (1)	Economics (2)	History (3)	Management (4)	Statistics (5)
<i>F</i> and <i>P</i> , linear weight kappa	0.54 (18.11)**	0.56 (11.94)**	0.32 (2.95)**	0.49 (8.91)**	0.55 (9.41)**
<i>F</i> and <i>P</i> , VQR weighted kappa	0.54 (17.29)**	0.56 (11.53)**	0.29 (2.56)**	0.50 (8.37)**	0.55 (9.18)**

Note: The table reports the kappa statistic and the associated z-value in parenthesis for the total sample and by research sub-area.

* Indicates significance at the 5% level.

** Indicates significance at the 1% level.

National Agency for the Evaluation of
Universities and Research InstitutesAgenzia Nazionale di Valutazione del
sistema Universitario e della Ricerca

Evaluation of Research Quality



Valutazione Qualità della Ricerca

For three of the four research areas, there is more than adequate agreement, while the agreement is somewhat lower for History.



ROARS

Return On Academic ReSearch



Lo strano caso delle concordanze della VQR

Alberto Baccini - 10 febbraio 2014

5

Nella VQR (aree 1-9 e 13) è stato usato un mix di strumenti per la valutazione del singolo prodotto di ricerca: la bibliometria e...

rendere in italiano il significato di *fair*: si potrebbe forse tradurre con discreto o modesto. O forse Accettabile, se si vuole usare un gergo di gradimento ad ANVUR. Certo è che è difficile sostenere che *fair* o *poor* possano essere tradotti in italiano con l'espressione usata da ANVUR nel rapporto: "più che adeguato".

Concordanza: “fair to good”. Ma quanto “good”?

Table 13

Kappa statistic for the amount of agreement between *F* and *P* scores.

	Total sample
	(1)
<i>F</i> and <i>P</i> , linear weight kappa	0.54 (18.11)**
<i>F</i> and <i>P</i> , VQR weighted kappa	0.54 (17.29)**

²⁹ Landis and Koch (1977) characterize the range of values 0–0.20 as “slight agreement”, 0.21–0.40 as “fair agreement”, 0.41–0.60 as “moderate agreement”, 0.61–0.80 as “substantial agreement”, and 0.81–1 as “almost perfect agreement”. These guidelines are somewhat arbitrary and by no means universally accepted. Fleiss (1981) for instance characterizes kappas over 0.75 as “excellent”, 0.40 to 0.75 as “fair to good”, and below 0.40 as “poor”. Kappa has also been shown to increase with the number of classes (only 4 in our case). Since the most common scales to subjectively assess the value of kappa mention “adequate” and “fair to good”, these are the terms that we use in the paper to convey the meaning of the statistic when commenting the estimated kappas.

<i>K</i> values	Description
Landis and Koch (1977)	
<0.00	Poor
0.00–0.20	Slight
0.21–0.40	Fair
0.41–0.60	Moderate
0.61–0.80	Substantial
0.81–1.00	Almost perfect
Altman (1991)	
<0.20	Poor
0.21–0.40	Fair
0.41–0.60	Moderate
0.61–0.80	Good
0.81–1.00	Very good
Fleiss et al. (2003)	
<0.40	Poor
0.40–0.75	Fair to good
>0.75	Excellent
George and Mallery (2003)	
<0.51	Unacceptable
0.51–0.60	Poor
0.61–0.70	Questionable
0.71–0.80	Acceptable
0.81–0.90	Good
0.91–1.00	Excellent
Stemler and Tsai (2008)	
<0.50	Unacceptable
>0.50	Acceptable
	unacceptable

Table 13Kappa statistic for the amount of agreement between F and P scores.

	Total sample (1)	Economics (2)	History (3)	Management (4)	Statistics (5)
F and P , linear weight kappa	0.54 (18.11)**	0.56 (11.94)**	0.32 (2.95)**	0.49 (8.91)**	0.55 (9.41)**
F and P , VQR weighted kappa	0.54 (17.29)**	0.56 (11.53)**	0.29 (2.56)**	0.50 (8.37)**	0.55 (9.18)**

Note: The table reports the kappa statistic and the associated z-value in parenthesis for the total sample and by research sub-area.

* Indicates significance at the 5% level.

** Indicates significance at the 1% level.

«The second row in Table 13 reports the “VQR weighted” kappa. The resulting statistic is quite similar to the linearly weighted kappa, indicating **fair to good agreement** for the total sample (**0.54**) and for Economics, Management and Statistics, and **poor agreement for History (0.29)**.»

Therefore:

“the agencies that run these evaluations could feel confident about using bibliometric evaluations and interpret the results as highly correlated with what they would obtain if they performed informed peer review” (Bertocchi et al. 2015)

Is this true?

Do they agree? Bibliometric evaluation versus informed peer review in the Italian research assessment exercise

Alberto Baccini¹  · Giuseppe De Nicolao²

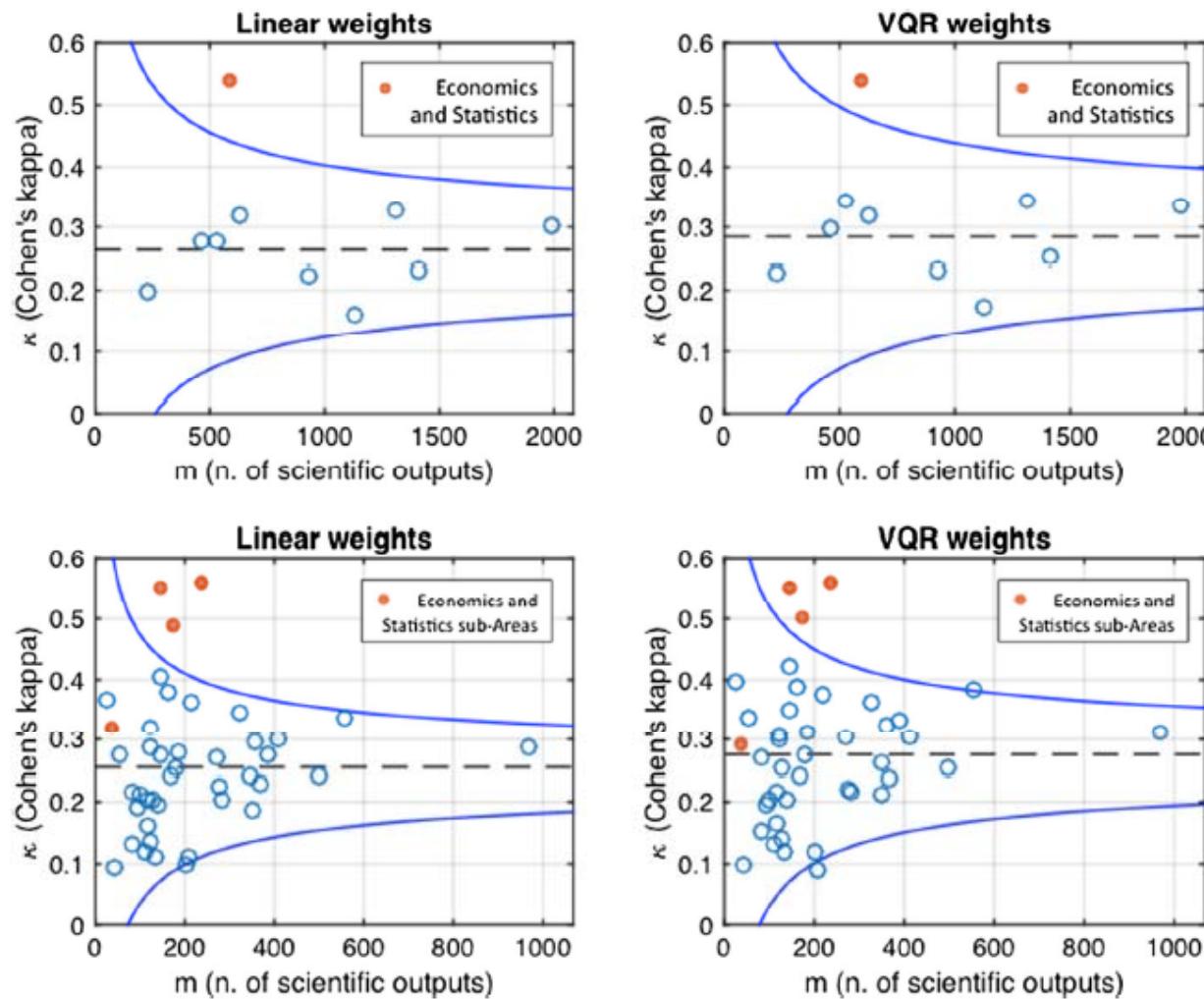
Abstract During the Italian research assessment exercise, the national agency ANVUR performed an experiment to assess agreement between grades attributed to journal articles by informed peer review (IR) and by bibliometrics. A sample of articles was evaluated by using both methods and agreement was analyzed by weighted Cohen's kappas. ANVUR presented results as indicating an overall “good” or “more than adequate” agreement. This paper re-examines the experiment results according to the available statistical guidelines for interpreting kappa values, by showing that the degree of agreement (always in the range 0.09–0.42) has to be interpreted, for all research fields, as unacceptable, poor or, in a few cases, as, at most, fair. The only notable exception, confirmed also by a statistical meta-analysis, was a moderate agreement for economics and statistics (Area 13) and its sub-fields. We show that the experiment protocol adopted in Area 13 was substantially modified with respect to all the other research fields, to the point that results for economics and statistics have to be considered as fatally flawed. The evidence of a poor agreement supports the conclusion that IR and bibliometrics do not produce similar results, and that the adoption of both methods in the Italian research assessment possibly introduced systematic and unknown biases in its final results. The conclusion reached by ANVUR must be reversed: the available evidence does not justify at all the joint use of IR and bibliometrics within the same research assessment exercise.

E negli altri GEV come va?

Table 2 Weighted kappas values for Areas and sub-areas

	Sample	Linear weighted kappas	VQR weighted kappas			
Area 1: Mathematics and informatics	631	0.3176	0.3173	Area 6: Medicine	1984	0.303
Informatics	164	0.3794	0.3896	Experimental medicine	347	0.2407
Mathematics	121	0.3218	0.3102	Clinical medicine	968	0.2883
Analysis and probability	179	0.2551	0.2755	Surgical sciences	554	0.3368
Applied mathematics	167	0.2426	0.2403	Public health	115	0.2023
Area 2: Physics	1412	0.2302	0.2515	Area 7: Agricultural and veterinary sciences	532	0.2776
Experimental physics	139	0.1957	0.2049	Agricultural sciences	387	0.2741
Theoretical physics	499	0.2428	0.2559	Veterinary	145	0.2747
Physics of matter	349	0.1862	0.2099	Area 8: Civil engineering and architecture	225	0.1994
Nuclear and sub-nuclear physics	45	0.0951	0.1001	Structural engineering	99	0.2106
Astronomy and astrophysics	270	0.2708	0.3048	Infrastructure engineering	126	0.2037
Geophysics	28	0.3671	0.3975	Area 9: Industrial and information engineering	1130	0.1615
Applied physics, teaching and history	82	0.2153	0.2715	Mechanical engineering	125	0.1355
Area 3: Chemistry	927	0.2246	0.2296	Industrial engineering	81	0.1325
Analitical chemistry	276	0.2261	0.2192	Nuclear engineering	117	0.1606
Inorganic and industrial chemistry	283	0.2024	0.2158	Chemical engineering	201	0.0996
Organic and pharmaceutical chemistry	368	0.2304	0.2368	Electronic engineering	210	0.1105
Area 4: Earth sciences	458	0.2776	0.2985	Telecommunication engineering	135	0.1117
Geochemistry etc.	123	0.287	0.2996	Bio-engineering	110	0.1214
Structural geology	96	0.1891	0.1932	Informatics	145	0.4052
Applied geology	56	0.2736	0.3375	Infrastructure engineering	6	na
Geophysics	183	0.277	0.3125	Area 13: Economics and statistics	590	0.54
Area 5: Biology	1310	0.3287	0.3453	Economics	235	0.56
Integrated biology	325	0.3451	0.3648	Economic history	37	0.32
Morfo-functional sciences	216	0.3629	0.3775	Management	175	0.49
Biochemistry and molecular biology	410	0.2998	0.304	Statistics	143	0.55
Genetics and pharmacology	359	0.296	0.3248	All areas	9199	0.32
						0.38

Source: (ANVUR 2013). Final Report; Appendix B; Appendix A of each Area Report. All data



**Cohen's
kappa for
Economy and
Statistics:
a statistical
anomaly?**

Fig. 2 Funnel plots: a point with coordinates (m , κ) represents a (sub)-area having m evaluated products and whose Cohen's kappa is κ . Cohen's kappas for Area 13 (full circles) are compared to the mean kappa (dashed) and 95 % prediction limits (continuous), based on kappas collected in the other nine areas (open circles). *Top* The kappas refer to the 10 areas. *Bottom* The kappas refer to the sub-areas. *Left* Linearly-weighted kappas are considered. *Right* VQR-weighted kappas are considered

Baccini e De Nicolao: Area 13, “a fatally flawed experiment”

- random sampling took into account authors' requests to be evaluated by peer review;
- the referees might have known that they were part of the experiment;
- the referees might have known the precise merit class in which each article was classified by using bibliometrics;
- the synthesis of the two referee's judgments was defined by a Consensus Group composed by (at least) two panel members;
- the panel members forming the Consensus Groups knew that their final judgment would be used for the experiment;
- at least 53 % of the IR evaluations was not expressed by referees, but directly by the Area 13 panelists.

For these reasons, results reached for Area 13 have to be considered as fatally flawed by virtue of the protocol modifications introduced by the area panel

Baccini e De Nicolao: Area 13, “a fatally flawed experiment”

- Nel campionamento casuale sono entrati prodotti (quanti?) per i quali gli autori hanno richiesto di essere valutati con IR (quindi non è casuale);
- I referee avrebbero potuto sapere che facevano parte dell'esperimento;
- I referee avrebbero potuto conoscere facilmente l'esatta classe di merito in cui ciascun articolo è stato classificato utilizzando la bibliometria;
- la sintesi dei due giudizi dei referee è stata definita da un Gruppo di Consenso composto da (almeno) due membri del panel;
- i membri dei gruppi di consenso sapevano che il loro giudizio finale sarebbe stato utilizzato per l'esperimento;
- almeno il 53 % delle valutazioni IR non è stato espresso dai referee, ma direttamente dai membri del panel.

Per questi motivi, i risultati raggiunti per l'area 13 devono essere considerati fatalmente errati in virtù delle modifiche del protocollo introdotte dal panel di Area 13

Comment to: Do they agree? Bibliometric evaluation versus informed peer review in the Italian research assessment exercise

Graziella Bertocchi¹ · Alfonso Gambardella² ·
Tullio Jappelli³ · Carmela Anna Nappi⁴ · Franco Peracchi⁵

Many of the points raised by Baccini and De Nicolao (henceforth BD) were already addressed in the RP paper. Other points are either incorrect or not supported by evidence.

Reply to the comment of Bertocchi et al.

Alberto Baccini¹  · Giuseppe De Nicolao²

Bertocchi et al.’s comment dismiss our explanation and suggest that the difference was due to “differences in the evaluation processes between Area 13 and other areas”. In addition, they state that all our five claims about Area 13 experiment protocol “are either incorrect or not based on any evidence”. Based on textual evidence drawn from ANVUR official reports, we show that: (1) none of the four differences listed by Bertocchi et al. is peculiar of Area 13; (2) their five arguments contesting our claims about the experiment protocol are all contradicted by official records of the experiment itself.

5. Concordanza o fallacia statistica?

Evaluating scientific research in Italy: The 2004–10 research evaluation exercise

Alessio Ancaiani¹, Alberto F. Anfossi^{1,2}, Anna Barbara^{1,3},
Sergio Benedetto¹, Brigida Blasi¹, Valentina Carletti¹, Tindaro Cicero¹,
Alberto Ciolfi¹, Filippo Costa^{1,4}, Giovanna Colizza¹,
Marco Costantini^{1,3}, Fabio di Cristina¹, Antonio Ferrara¹,
Rosa M. Lacatena¹, Marco Malgarini^{1,*}, Irene Mazzotta¹,
Carmela A. Nappi¹, Sandra Romagnosi¹ and Serena Sileoni¹

¹*Agenzia Nazionale di Valutazione del Sistema Universitario e della Ricerca (ANVUR), Via Ippolito Nievo 35 - 00153 Rome, Italy*, ²*Compagnia di San Paolo Sistema Torino, Piazza Bernini 5, IT-10138 Turin, Italy*, ³*Gabriele D'Annunzio Chieti-Pescara University Via dei Vestini, 31 - 66013 Chieti Scalo, Italy* and ⁴*Department of Information Engineering, Pisa University, Via Caruso 16 - 56122 Pisa, Italy*

*Corresponding author. Email: marco.malgarini@anvur.it

Table 2. K-Cohen statistic

Area	F e P, linear weights	F e P, VQR weights
Mathematics and Computer Sciences	0.3176 (10.25)*	0.3173 (0.74)*
Physics	0.2302 (14.26)*	0.2515 (15.10)*
Chemistry	0.2246 (10.67)*	0.2296 (10.42)*
Earth Sciences	0.2776 (8.72)*	0.2985 (8.50)*
Biology	0.3287 (16.38)*	0.3453 (15.67)*
Medicine	0.3024 (19.18)*	0.3351 (19.04)*
Agricultural and Veterinary Sciences	0.2776 (10.83)*	0.3437 (11.57)*
Civil engineering and Architecture	0.1994 (5.03)*	0.2261 (5.10)*
Industrial and Information Engineering	0.1615 (10.56)*	0.1710 (10.91)*
Economic and Statistics	0.54 (18.11)*	0.6104 (17.27)*
Total	0.3152 (44.48)*	0.3441 (44.55)*

* indicates significance at 1% level.

«K is always **statistically different from zero**, showing that there is a **fundamental agreement** among the two distributions which **may not be attributed to mere chance**, regardless of the weight used to calculate the differences among the two distributions. The value of K ranges from 0.16 to 0.61 depending on the area and weights, being on average equal to 0.32, a value that is usually considered as '**poor to fair**' in the literature (Landis and Koch 1977).»

Therefore:

“results of the analysis relative to the degree of concordance and systematic difference may be considered to validate the general approach of combining peer review and bibliometric methods” (Ancaiani et al. 2015)

Is this true?

The significance fallacy

Kühberger et al. BMC Research Notes (2015) 8:84
DOI 10.1186/s13104-015-1020-4

BMC Research Notes

RESEARCH ARTICLE Open Access

The significance fallacy in inferential statistics

Anton Kühberger^{1*}, Astrid Fritz², Eva Lermer³ and Thomas Scherndl¹

Abstract
Background: Statistical significance is an important concept in empirical science. However the meaning of the term varies widely. We investigate into the intuitive understanding of the notion of significance.
Methods: We described the results of two different experiments published in a major psychological journal to a sample of students of psychology, labeling the findings as 'significant' versus 'non-significant.' Participants were asked to estimate the effect sizes and sample sizes of the original studies.
Results: Labeling the results of a study as significant was associated with estimations of a big effect, but was largely unrelated to sample size. Similarly, non-significant results were estimated as near zero in effect size.
Conclusions: After considerable training in statistics, students largely equate statistical significance with medium to large effect sizes, rather than with large sample sizes. The data show that students assume that statistical significance is due to real effects, rather than to 'statistical tricks' (e.g., increasing sample size).
Keywords: Statistical significance, Practical significance, Effect size, NHST, Sample size

Background
There is continuing debate on the usefulness and validity of the method of Null Hypothesis Significance Testing (NHST, e.g., [1-3]). Several journals edited special issues on this topic (e.g., *Journal of Experimental Education* in 1993; *Psychological Science* in 1997; *Research in the Schools* in 1998) that culminated in the question: What is beyond the significance test ritual (*Journal of Psychology* in 2009)?
The debate has led to an increased awareness of the problems associated with NHST, and these problems are linked to what has been referred to as a 'crisis of confidence' [4]. Among the dominant recommendations for NHST is reporting of effect size as a supplement to the *p* value [5]. Accordingly, not only the statistical significance of a result should be valued but also the effect size of the study (e.g., [1,6-12]). This should prevent readers from holding the false belief that significant results are automatically big and important, or otherwise, that not significant means 'no effect at all'. Although these misinterpretations are often avoided by reporting effect size, significance means no effect, are often referred to in the literature (e.g., [3,13-17]) their empirical basis is weak.

* Correspondence: Anton.kuehberger@sbg.ac.at
¹Department of Psychology and Centre of Cognitive Neuroscience, University of Salzburg, Hellbrunnerstr. 34, 5020 Salzburg, Austria
Full list of author information is available at the end of the article

© 2015 Kühberger et al. licensee BioMed Central. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly credited. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated.

BioMed Central

the false belief
that [statistically]
significant results
are automatically
big and important

Una nozione insegnata in tutti i corsi di statistica di base: la differenza tra *statistical* e *practical* significance

Researchers want to test a new medication that claims to raise IQs to genius levels (175+). In the population, the average IQ is 100. A sample of 40 individuals has a mean IQ of 104 with a standard deviation of 15.

$H_0: \mu = 100$

$H_a: \mu \neq 100$

A t-test is performed and the null hypothesis is rejected. It is concluded that the medication raises IQ.

But the medication claimed to raise IQs to genius (175+) levels. Even though we found statistical significance, the medication does not meet the practical value it claimed to. It lacks practical significance.

Statistical Versus Practical Significance

Practical Significance

- Practical (or clinical) significance asks the larger question about differences
 - Are the differences between samples big enough to have real meaning?
- Although men and women undoubtedly have different IQs, is that difference large enough to have some practical implications

For a reasonably large sample size n , this μ would lead to an x value near 101, so we would not want this sample evidence to argue strongly for rejection of H_0 when $\mu = 101$ is observed.

For various sample sizes, Table 8.1 records both the P -value when $x = 101$ and also the probability of not rejecting H_0 at level .01 when $\mu = 101$.

An illustration of the Effect of Sample Size on Products and Defects

Products produced by the machine (following the repair)

Defective

$H_0: p = .20$

$H_a: p < .20$

64 defective

Conditions are met

$\hat{p} = 64/400 = .16$

$z = -2$

$p\text{-value} = .023$

Conclusion: H_0 can be rejected.

Statistical vs. Practical Significance

Statistical significance vs. Practical Significance

Statistical significance can be argued through the interpretation of the P -value.

A statistically significant result has a P -value of less than .05 (see previous slide).

Practical significance can be argued in relation to the effect size. It depends on the study's context and norms.

An example where practical significance is of greater importance could be in medication, where 2mg could make a huge difference in effects on a patient, but the P -value might suggest otherwise. In a different context, 1mg of sugar per kilogram may not be of practical significance.

Further examples concerning when practical significance is or is not important can be found in the Causation, Chapter 8, page 12.

Statistical vs. Clinical or Practical Significance

Wit G Hopkins
Auckland University of Technology
Auckland, NZ

- Statistical significance
 - P values and null hypotheses
- Confidence limits
 - Precision of estimation
- Clinical or practical significance
 - Probabilities of benefit and harm
 - Examples

Statistical vs. Practical Significance

Experimental significance

Statistical significance

Practical significance

"Are the effects observed due to the changing levels of the factors real or due to random chance?"

"If a factor does have a statistically significant effect, then is the degree of effect practically useful given what we are trying to accomplish?"

Statistical vs. Practical Significance

Example: Suppose a weight loss program recruits 10,000 people for a randomized experiment.

- A difference in average weight loss of only 0.5 lbs could be found to be statistically significant
- Suppose the experiment lasted for a year. Is a loss of $\frac{1}{2}$ a pound practically significant?

Not everything can be quantified!

The one I PEA to go to the gym

Now I have to eat a lot more

Statistical vs. Practical Significance

Example: Suppose a weight loss program recruits 10,000 people for a randomized experiment.

- A difference in average weight loss of only 0.5 lbs could be found to be statistically significant
- Suppose the experiment lasted for a year. Is a loss of $\frac{1}{2}$ a pound practically significant?

Statistical vs. Practical Significance

- Statistical significance (e.g., $p < 0.05$) does not imply practical relevance
- Results should be both: (1) statistically and (2) practically significant in order to influence policy
- Example: A drug may induce a statistically significant reduction in blood pressure. However, if this reduction is 1 mmHg in your systolic BP, then it is not a useful (practical and clinically relevant) drug.

Una citazione riferita proprio alla kappa di Cohen

Statistical significance “is generally of little practical value, since a relatively low value of kappa can yield a significant result. In other words, a value such as $k = 0.41$ (in spite of the fact that is statistically significant) may be deemed by a researcher to be too low a level of reliability (i.e. degree of agreement) to be utilized within a practical context” (Sheskin 2003).



“the results reported by Ancaiani et al. **do not support a good concordance between peer review and bibliometrics.** [...] On the basis of these data, the conclusion that it is possible to use both technique as interchangeable in a research assessment exercise appears to be **unsound.**” (Baccini and De Nicolao 2017)

6. Dati chiusi, concordanza non replicabile

Dal 2014 abbiamo tentato di replicare l'esperimento

- ANVUR non fornisce i dati necessari
(mail 10/2/2014 a Presidente Fantoni)

 lunedì 10/02/2014 11:21
Alberto Baccini <alberto.baccini@unisi.it>
Richiesta dati VQR
A 'Presidenza@anvur.org'

Gentile presidente,
sto tentando di riprodurre i risultati ANVUR relativi alla concordanza tra risultati bibliometrici e IR (Appendice B del rapporto finale e appendici A dei rapporti di Area).
Le informazioni disponibili pubblicamente non permettono di raggiungere tale fine e neanche di ricalcolare gli indici di concordanza.
Sono pertanto a chiedere di avere accesso alle informazioni elencate in calce a questa mail, che al momento sono utilizzate da membri GEV e collaboratori ANVUR in pubblicazioni scientifiche.
Chiederei inoltre di conoscere in dettaglio gli algoritmi di sintesi utilizzati dai GEV 1-9 per la sintesi dei punteggi dei revisori cui si fa riferimento nei rapporti di area, ma che non sono pubblicati in quanto tali.
Sono a disposizione per ogni ulteriore chiarimento in merito alla mia richiesta,
Cordiali saluti,

Alberto Baccini

Descrizione dei dati

Per ciascun articolo che è stato utilizzato nella sintesi di concordanza.

Identificativo dell'articolo
Area
SSD
Valutazione bibliometrica dell'articolo
Identificativo del revisore P1 (basta un codice univoco del revisore, salvaguardando l'anonymato)
Se il revisore P1 è membro del GEV
Punteggio attribuito da P1 a criterio rilevanza
Punteggio attribuito da P1 a criterio originalità/innovazione
Punteggio attribuito da P1 a criterio internazionalizzazione
Valutazione di sintesi del revisore P1
Identificativo del revisore P2 (basta un codice univoco del revisore, salvaguardando l'anonymato)
Se il revisore P2 è membro del GEV
Punteggio attribuito da P2 a criterio rilevanza
Punteggio attribuito da P2 a criterio originalità/innovazione
Punteggio attribuito da P2 a criterio internazionalizzazione
Valutazione di sintesi del revisore P2
Identificativo del revisore P3 (basta un codice univoco del revisore, salvaguardando l'anonymato)
Se il revisore P3 è membro del GEV
Punteggio attribuito da P3 a criterio rilevanza
Punteggio attribuito da P3 a criterio originalità/innovazione
Punteggio attribuito da P3 a criterio internazionalizzazione
Valutazione di sintesi del revisore P3
Valutazione di sintesi dei giudizi dei revisori

prof. alberto baccini
dipartimento di economia politica e statistica
via p.a. mattioli 10
53100 siena
tel. +39 0577 235233
fax +39 0577 235235
<http://www.econ-pol.unisi.it/baccini>

Evaluating scientific research in Italy: The 2004–10 research evaluation exercise

Alessio Ancaiani¹, Alberto F. Anfossi^{1,2}, Anna Barbara^{1,3},
Sergio Benedetto¹, Brigida Blasi¹, Valentina Carletti¹, Tindaro Cicero¹,
Alberto Ciolfi¹, Filippo Costa^{1,4}, Giovanna Colizza¹,
Marco Costantini^{1,3}, Fabio di Cristina¹, Antonio Ferrara¹,
Rosa M. Lacatena¹, Marco Malgarini^{1,*}, Irene Mazzotta¹,
Carmela A. Nappi¹, Sandra Romagnosi¹ and Serena Sileoni¹

¹*Agenzia Nazionale di Valutazione del Sistema Universitario e della Ricerca (ANVUR), Via Ippolito Nievo 35 - 00153 Rome, Italy*, ²*Compagnia di San Paolo Sistema Torino, Piazza Bernini 5, IT-10138 Turin, Italy*, ³*Gabriele D'Annunzio Chieti-Pescara University Via dei Vestini, 31 - 66013 Chieti Scalo, Italy* and ⁴*Department of Information Engineering, Pisa University, Via Caruso 16 - 56122 Pisa, Italy*

*Corresponding author. Email: marco.malgarini@anvur.it

Research Evaluation, 26(4), 2017, 353–357

doi: 10.1093/reseval/rvx013

Advance Access Publication Date: 27 April 2017

Article

OXFORD

A letter on Ancaiani et al. 'Evaluating scientific research in Italy: the 2004-10 research evaluation exercise'

Alberto Baccini¹ and Giuseppe De Nicolao²

¹Department of Economics and Statistics, University of Siena, Piazza San Francesco 7, Siena, 53100, Italy, and

²Department of Electrical, Computer and Biomedical Engineering, University of Pavia, Pavia, Italy

*Corresponding author. Email: alberto.baccini@unisi.it

A letter on Ancaiani et al. 'Evaluating scientific research in Italy: the 2004–10 research evaluation exercise'

Alberto Baccini¹ and Giuseppe De Nicolao²

This letter documents some problems in Ancaiani et al. (2015). Namely the evaluation of concordance, based on Cohen's kappa, reported by Ancaiani et al. was not computed on the whole random sample of 9,199 articles, but on a subset of 7,597 articles. The kappas relative to the whole random sample were in the range 0.07–0.15, indicating an unacceptable agreement between peer review and bibliometrics. The subset was obtained by non-random exclusion of all articles for which bibliometrics produced an uncertain classification; these raw data were not disclosed, so that concordance analysis is not reproducible. The VQR-weighted kappa for Area 13 reported by Ancaiani et al. is higher than that reported by Area 13 panel and confirmed by Bertocchi et al. (2015), a difference explained by the use, under the same name, of two different set of weights. Two values of kappa reported by Ancaiani et al. differ from the corresponding ones published in the official report. Results reported by Ancaiani et al. do not support a good concordance between peer review and bibliometrics. As a consequence, the use of both techniques introduced systematic distortions in the final results of the Italian research assessment exercise. The conclusion that it is possible to use both technique as interchangeable in a research assessment exercise appears to be unsound, by being based on a misinterpretation of the statistical significance of kappa values.

Protocolli dell'esperimento

Protocollo 5X5

Protocollo 4X4

		Bibliometric Indicator			
		1	2	3	4
n. of citations	1	A	IR	IR	IR
	2	A	B	C	D
	3	A	B	C	D
	4	IR	IR	IR	D

		Bibliometric Indicator			
		1	2	3	4
n. of citations	1	A	IR	IR	IR
	2	A	B	C	D
	3	A	B	C	D
	4	IR	IR	IR	D

Protocollo 5X5 vs. protocollo 4X4

Table 1. Agreement between informed peer review and bibliometrics

Areas	Whole sample 5×5 protocol		Reduced sample 4×4 protocol ^a		
	N	Unweighted kappa	N	Linear-weighted kappa	VQR-weighted kappa
Area 1 Mathematics and Informatics	631	0.13	438	0.32	0.32
Area 2 Physics	1,412	0.12	1,212	0.23	0.25
Area 3 Chemistry	927	0.14	778	0.22	0.23
Area 4 Earth Sciences	458	0.12	377	0.28	0.3
Area 5 Biology	1,310	0.15	1,058	0.33	0.35
Area 6 Medicine	1,984	0.14	1,602	0.30	0.34
Area 7 Agricultural and Veterinary Sciences	532	0.12	425	0.28	0.34
Area 8a Civil Engineering	225	0.07	198	0.20	0.23
Area 9 Industrial and Information Engineering	1,130	0.10	919	0.16	0.17
Area 13 Economics and Statistics	590	0.37	590	0.54	0.61
All areas	9,199	0.16	7,597	0.32	0.38

^aData drawn from ANVUR report. Appendix B. Not reproducible.

All other data, our elaboration from ANVUR publicly available raw data. Appendix B of ANVUR report.

R, psych package ver. 1.6.6 <https://cran.r-project.org/web/packages/psych/psych.pdf>.

Protocollo 5X5 vs. protocollo 4X4

Table 1. Agreement between informed peer review and bibliometrics

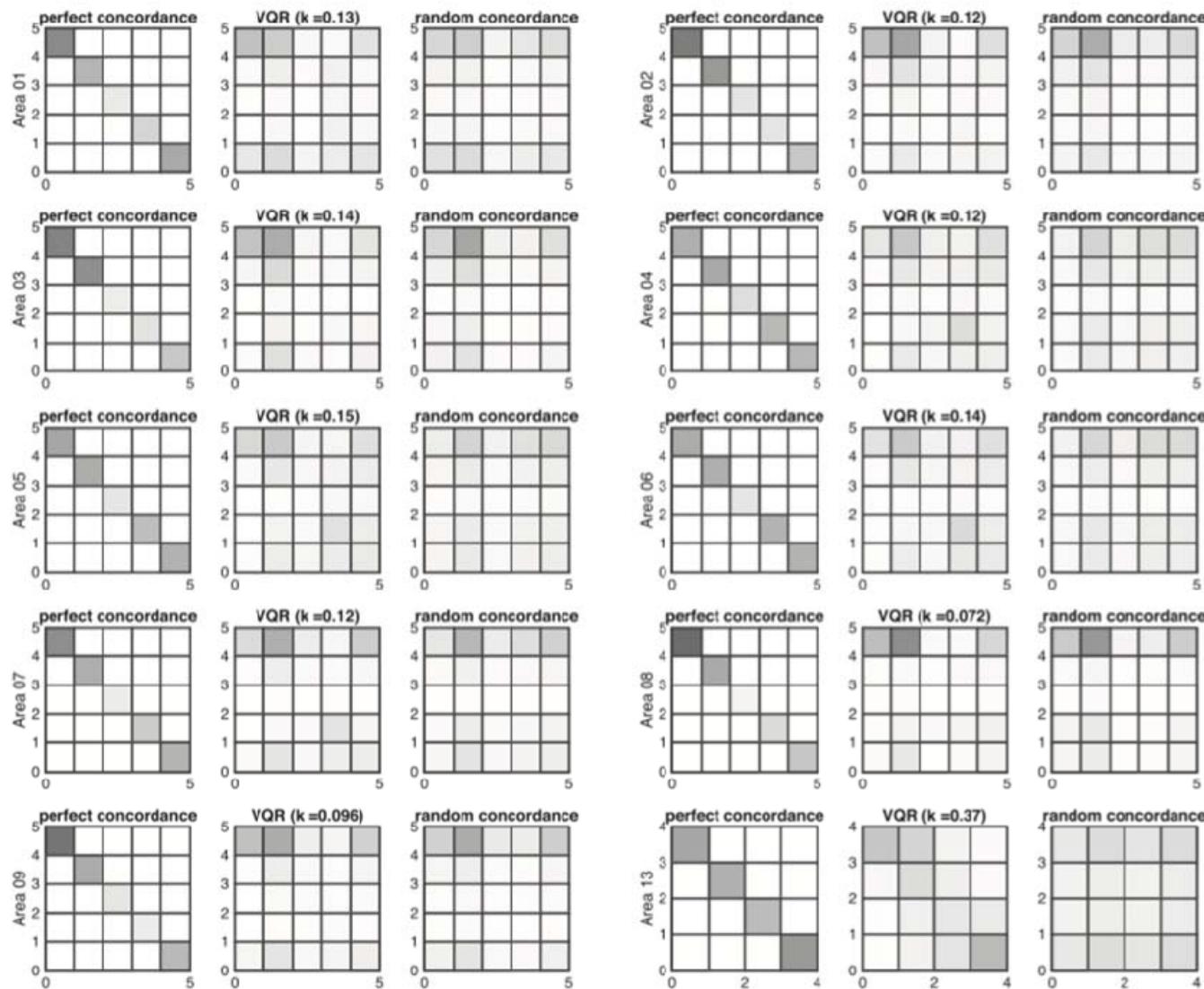
Areas	Whole sample 5 × 5 protocol		Reduced sample 4 × 4 protocol ^a		
	N	Unweighted kappa	N	Linear-weighted kappa	VQR-weighted kappa
Area 1 Mathematics and Informatics	631	0.13	438	0.32	0.32
Area 2 Physics	1,412	0.12	1,212	0.23	0.25
Area 3 Chemistry	927	0.14	778	0.22	0.23
Area 4 Earth Sciences	458	0.12	377	0.28	0.3
Area 5 Biology	1,310	0.15	1,058	0.33	0.35
Area 6 Medicine	1,984	0.14	1,602	0.30	0.34
Area 7 Agricultural and Veterinary Sciences	532	0.12	425	0.28	0.34
Area 8a Civil Engineering	225	0.07	198	0.20	0.23
Area 9 Industrial and Information Engineering	1,130	0.10	919	0.16	0.17
Area 13 Economics and Statistics	590	0.37	590	0.54	0.61
All areas	9,199	0.16	7,597	0.32	0.38

^aData drawn from ANVUR report. Appendix B. Not reproducible.

All other data, our elaboration from ANVUR publicly available raw data. Appendix B of ANVUR report.
R, psyc package ver. 1.6.6 <https://cran.r-project.org/web/packages/psych/psych.pdf>.

**valori bassi
di kappa non
pubblicati da
ANVUR**

**BIBLIOMETRIC EVALUATION
(Excellent, Good, Acceptable, Limited, IR)**



**PEER REVIEW EVALUATION
(Excellent, Good, Acceptable, Limited, IP)**

Ancaiani et al. 2015

Table 2. K-Cohen statistic

Area	F e P, linear weights	F e P, VQR weights	P1 e P2, linear weights	P1 e P2, VQR weights
Mathematics and Computer Sciences	0.3176 (10.25)*	0.3173 (0.74)*	0.3595 (10.22)*	0.3516 (9.82)*
Physics	0.2302 (14.26)*	0.2515 (15.10)*	0.23317 (11.65)*	0.2271 (11.33)*
Chemistry	0.2246 (10.67)*	0.2296 (10.42)*	0.2501 (10.02)*	0.2381 (9.60)*
Earth Sciences	0.2776 (8.72)*	0.2985 (8.50)*	0.2500 (6.72)*	0.2548 (6.48)*
Biology	0.3287 (16.38)*	0.3453 (15.67)*	0.2750 (12.13)*	0.2717 (11.39)*
Medicine	0.3024 (19.18)*	0.3351 (19.04)*	0.2460 (13.48)*	0.2356 (12.22)*
Agricultural and Veterinary Sciences	0.2776 (10.83)*	0.3437 (11.57)*	0.1570 (4.60)*	0.2656 (12.22)*
Civil engineering and Architecture	0.1994 (5.03)*	0.2261 (5.10)*	0.2029 (4.07)*	0.1943 (3.85)*
Industrial and Information Engineering	0.1615 (10.56)*	0.1710 (10.91)*	0.1935 (8.30)*	0.1818 (7.77)*
Economic and Statistics	0.54 (18.11)*	0.6104 (17.27)*	0.40 (12.93)*	0.4599 (12.94)*
Total	0.3152 (44.48)*	0.3441 (44.55)*	0.2853 (34.63)*	0.2816 (32.86)*

460 G. Bertocchi et al. / Research Policy 44 (2015) ·

Table 13
Kappa statistic for the amount of agreement between F and P scores.

APPA TAB.6

	Total sample		Economics
	(1)	(2)	
F and P, linear weight kappa	0.54 (18.11)	0.56 (11.94)	
F and P, VQR weighted kappa	0.54 (17.29)	0.56 (11.53)	
P1 and P2, equal weights	0.40 (12.93)	0.44 (9.06)	
P1 and P2, VQR weights	0.39 (12.06)	0.42 (8.28)	

Note: The table reports the kappa statistic and the associated z-value in parenthesis for the total sample

* Indicates significance at the 5% level.

** Indicates significance at the 1% level.

Errore nei dati o altro?

Altro: ci sono due sistemi di pesi chiamati nello stesso modo

Table 3. VQR weights. Matrix used by ANVUR and Ancaiani et al

		Informed peer review			
		A	B	C	D
Bibliometrics	A	1	0.8	0.5	0
	B	0.8	1	0.8	0.5
	C	0.5	0.8	1	0.8
	D	0	0.5	0.8	1

Note: This matrix attributed to agreement, one-class, two-class, and three-class disagreement weights modeled on the basis of the score (1, 0.8, 0.5, and 0) associated to the four categories in which papers are classified (A, B, C, and D). For example, consider two papers: a paper classified as A by bibliometrics and classified as B by peer review; and a second paper classified B by bibliometrics and C by peer review. Both have a one-class disagreement and a weight of 0.8, which appears arbitrary. In fact, in the former case, the score error is $1.0 - 0.8 = 0.2$, while in the latter one, it is $0.8 - 0.5 = 0.3$.

Table 4. VQR weights. Matrix used by Area 13 panel

		Informed peer review			
		A	B	C	D
Bibliometrics	A	1	0.8	0.5	0
	B	0.8	1	0.7	0.2
	C	0.5	0.7	1	0.5
	D	0	0.2	0.5	1

Note: This matrix attributed to agreement, one-class, two-class, and three-class disagreement weights modeled on the basis of the difference between the scores associated to the four categories in which papers are classified (A, B, C, and D). For example, consider two papers: a paper classified as A (Score 1) by bibliometrics and classified as B (Score 0.8) by peer review; and a second paper classified B (Score 0.8) by bibliometrics and C (Score 0.5) by peer review. Both have a one-class disagreement; the difference between the two scores for the first paper is 0.2, and the weight is $1 - 0.2 = 0.8$; for the second paper, the difference between the two scores is 0.3, and the weight is $1 - 0.3 = 0.7$.

Altri dati che non quadrano. Perché?

Furthermore two values reported in Table 2 of Ancaiani et al. differ from the corresponding ones published in the ANVUR report ([ANVUR 2013](#): Appendix B, p. 22). Namely, the value $k = 0.3441$ for the agreement between peer review and bibliometrics for all areas reported by Ancaiani et al. differs from $k = 0.38$ published in the ANVUR report (Table 1), and the value $k = 0.2816$ for the agreement between two reviewers for all areas differs from $k = 0.33$ published in the ANVUR report (Table 2). We were not able to explain these discrepancies, given that the result cannot be replicated due to the aforementioned unavailability of raw data for the 4×4 protocol.

Reply to the letter on Ancaiani et al. 'Evaluating Scientific research in Italy: The 2004–10 research evaluation exercise'

Sergio Benedetto^{1,*}, Tindaro Cicero², Marco Malgarini² and Carmen Nappi²

¹Politecnico di Torino, Dipartimento di Elettronica e telecomunicazioni, Corso Castelfidardo 39, Turin, Italy and

²ANVUR, Via Lovolito Nievo 35, 00153, Rome

Abstract

Baccini and De Nicolao (2017) provide some criticism on the results showed in Ancaiani et al (2015) concerning the Italian Evaluation exercise (VQR in the Italian acronym). In this reply we provide ample evidence that the issues raised do not weaken the main results previously presented in any substantial way.

RT. A Journal on Research Policy & Evaluation 1 (2017)
Submitted on 10 July 2017, published on 22 July 2017, open for comments

Doi: 10.13130/2282-5398/8872



Errors and secret data in the Italian research assessment exercise. A comment to a reply

Alberto Baccini*, Giuseppe De Nicolao**

Errori inspiegabili nella replica

Table 1. Sampling distribution

Area	Number of bibliometric articles (population of reference)	Number of articles in the full sample	Number of articles in the subsample
1	6,758	631	438
2	15,029	1,412	1,212
3	10,127	927	778
4	5,083	458	377
5	14,043	1,310	1,058
6	21,191	1,984	1,603
7	6,284	532	425
8	2,460	225	198
9	12,349	1,130	919
13	5,681	590	590
Total	99,005	9,199	7,598

Table 2. Bibliometric distribution in the sample and in the whole population

Evaluation class	Population	%	Sample	%
A	4,7583	48.1	4,419	48.0
B	15,739	15.9	1,457	15.8
C	5,180	5.2	479	5.2
D	1,486	13.6	1,242	13.5
IR	17,010	17.2	1,602	17.4

Population: 86.998

ERROR:
47.583?

ERROR
7,597

ERROR

7. Conclusioni



Brenno, il capo dei Galli che avevano occupato Roma, aveva acconsentito ad andarsene in cambio di un tributo in oro. I romani, pesandolo, gli avevano fatto notare che i pesi della bilancia erano truccati; Brenno aveva replicato gettando sul piatto anche la spada e dicendo, appunto, "Guai ai vinti!", per comunicare che l'unica legge che riconosceva era quella del più forte.

http://btfp.sp.unipi.it/dida/kant_7/ar01s10.xhtml#ftn.idp2831664

Messa sul piatto la spada, ANVUR usa la bilancia

- Perché questo straordinario sforzo di disseminazione è stato prodotto da studiosi che lavorano per l'ANVUR?
- Probabilmente perché la pubblicazione su riviste scientifiche rappresenta una giustificazione ex-post del doppio sistema di valutazione, senza precedenti, sviluppato e applicato dall'ANVUR.
- La metodologia e i risultati della valutazione della ricerca sono giustificati ex post da documenti scritti dagli studiosi che hanno sviluppato e applicato la metodologia adottata dal governo italiano.
- I risultati di questi lavori non possono essere replicati perché i dati non sono stati messi a disposizione di studiosi diversi da quelli che lavorano per l'ANVUR.

Politica vaccinale

- Governo prescrive un nuovo vaccino obbligatorio in conformità con la raccomandazione di un rapporto emesso da un'agenzia come la Food and Drug Administration.
- Un paio d'anni dopo l'adozione obbligatoria, le riviste accademiche pubblicano articoli, scritti da membri della commissione della FDA che ha pubblicato il rapporto.
- Questi articoli riproducono, senza dichiararlo, i contenuti e le conclusioni del rapporto della FDA, fornendo così una giustificazione scientifica de facto - anche se ex post - del rapporto stesso.
- Quando ricercatori indipendenti richiedono dati per replicare i risultati, l'agenzia non risponde o, in alternativa, nega i dati che si asseriscono riservati.

Inquinamento della letteratura e integrità delle riviste

Scientometrics
DOI 10.1007/s11192-017-2384-0



Do social sciences and humanities behave like life and hard sciences?

Andrea Bonacorsi^{1,2} · Cinzia Daraio³ · Stefano Fantoni⁴ ·
Viola Folli⁵ · Marco Leonetti^{5,7} · Giancarlo Ruocco^{5,6}

Research Policy 46 (2017) 911–924

Contents lists available at ScienceDirect

Research Policy

journal homepage: www.elsevier.com/locate/respol

Gender effects in research evaluation

Tullio Jappelli^{a,*}, Carmela Anna Nappi^b, Roberto Torrini^c

^a University of Naples Federico II, Italy
^b Anvur, Italy
^c Bank of Italy, Italy

Journal of Informetrics 10 (2016) 224–237

Contents lists available at ScienceDirect

Journal of Informetrics

journal homepage: www.elsevier.com/locate/joi

Nondeterministic ranking of university departments[☆]

Andrea Bonacorsi^a, Tindaro Cicero^{b,*}

^a DESTEC, School of Engineering University of Pisa Largo, Lucio Lazzarino 2, 56125 Pisa, Italy
^b ANVUR Italian Agency for the Evaluation of Universities and Research Institutes, Via Ippolito Nievo 35, 00153 Rome, Italy

**Distributed or Concentrated Research Excellence?
Evidence From a Large-Scale Research
Assessment Exercise**

Andrea Bonacorsi
DESTEC Department, School of Engineering, University of Pisa, Largo Lucio Lazzarino 2, Pisa 56125, Italy;
Italian Agency for the Evaluation of Universities and Research Institutes (ANVUR), Via Ippolito Nievo 35,
Rome 00153, Italy. E-mail: a.bonacorsi@gmail.com

Tindaro Cicero
Italian Agency for the Evaluation of Universities and Research Institutes (ANVUR), Via Ippolito Nievo 35,
Rome 00153, Italy. E-mail: tindaro.cicero@anvur.it

F1000Research 2015, 4:196 Last updated: 09 SEP 2015



RESEARCH ARTICLE
**Journal ratings as predictors of articles quality in Arts,
Humanities and Social Sciences: an analysis based on the
Italian Research Evaluation Exercise [version 1; referees: 3
approved]**

Andrea Bonacorsi, Tindaro Cicero, Antonio Ferrara, Marco Malgarini
ANVUR, Via Ippolito Nievo 35, Rome, 00153, Italy

About

Research

Commentary

Education

Events

Donate

BLOG • VIDEOS • COLLECTIONS

How Pseudoscientific Rankings Are Distorting Research

By **Alberto Baccini and Giuseppe De Nicolao**

JAN 18, 2018 | PUBLIC & PRIVATE INSTITUTIONS | INSTITUTIONS, POLICY & POLITICS

