

袁培国

引文在评价研究工作上的作用和引文分析中应注意的问题

ABSTRACT Whether and how much a paper has been cited indicates its effects on later research work. The author thinks that there are some important factors in analysis as follows: sampling quantity, sampling period, article type, retrieval method, disciplinary difference, difference of attitudes, co-authors, authors with the same name, etc. 26 refs

KEY WORDS Citation Scientific research Evaluation

CLASS NUMBER G257.5

在科学研究特别是基础研究成果的评价方面,研究论文的出版数量曾起过并继续起着重大作用。但在质量评价方面,长期以来一直是靠专家的定性分析,既耗时又昂贵,并且一些主观因素的影响也难以排除。当前科技飞速发展,学科专业越分越细,同时跨学科、多学科研究方兴未艾,论文数量迅猛增加,论文数量翻番的周期越来越短。以论文数量和专家的定性分析来评价科学研究工作越来越难以满足需要。在1945~1988年世界被引用最多的400篇论文中,作者应邀在ISI的“引文经典”(citation classics)栏介绍有关情况的205篇论文里,在出版前审稿过程中碰到问题的就有22篇之多,占10.7%^[1]。正是在这样的情况下,美国科学情报研究所(ISI)公司的引文索引为学术研究的质量评价提供了全新的工具。以引文统计分析来评价科研工作的方法在国际上为越来越多的人所接受。大量研究表明:通过引文数据分析(引文数、引文影响因子、相对影响等)得出的结果与著名、杰出等定性评价和多项指标的定量、定性分析综合评价结果有着高度的相关性。

论文发表以后,是否被引用和引文量的多少,说明它们对后来研究工作的影响和作用。出版发行的大量论文中,被引用的只是其

中一部分,而长期大量被引用的更是很少一部分。1981~1990年ISI标引的全世界872万多篇论文中,在1981~1993年期间从未被引用过的就达48.8%;1981~1993年巴西被ISI标引的论文的作者共52808人,同期从未被引用过的作者达44.1%,而被引用过10次以上的作者仅占17.5%^[2]。1945~1988年被SCI收录的被引用文献近3273万件,件均被引5.3次,44年间仅被引用过1次的占56%左右,被引用100次以上的为45%,被引用1000次以上的为40%^[3]。

根据SCI的统计,1961年被引用过的25.6万名第一作者中,在1962和1963年获诺贝尔奖者的人均被引次数、被引论文数和人均每篇论文被引次数分别为所有被引第一作者人均的30倍、17倍和近2倍。对1~15年不同时期里被引用最多的50名到1000名科学家的统计分析发现:按被引次数从高到低排序,世界科学家的前1%甚至更少些的作者中,相当大一部分是已获得或以后可能获得诺贝尔奖的人;对他们中的获诺贝尔奖者和非获奖者的人均被引次数和论文篇均被引次数比较后发现:前者明显高于后者;在发表论文总数上,前者甚至少于后者^[4]。

1961~1975年被引用最多的249名第一

作者中,到 1977 年已有 42 人获诺贝尔奖(1975 年时为 38 人,到 1990 年增加到 51 人),任一或数个国家科学院院士者 151 人^[5]。1965~1978 年期间发表的论文同期被引用最多的 1000 名科学家(包括第一作者和非第一作者)中,获诺贝尔奖者 35 人(到 1990 年增加到 61 人),任一或多国科学院院士者 373 人。对诺贝尔奖获得者和院士们的被引次数分别统计表明:前者明显高于后者^[6]。研究还发现,美国科学家获国家科学基金会拨款与科学家自身 9 个变量的互关性上,论著被引次数是相关性最强的^[7]。分析表明:按作者的被引次数从高到低排序,取样数目越大,获诺贝尔奖和任院士者人数越多,但百分比却变小。也就是说,排名越前者,其成就和影响可能越大,获奖或担任荣誉职务的机率就越高。如果按学科专业排序,情况更加明显。

《美国新闻与世界报道》(US News & World Report)每年两次分别以多项定性和定量指标对美国高校和各学科研究生院进行综合评价排序。ISI 1994 年根据 1981~1993 年期间学科专业的引文相对影响(即某单位在某学科专业的论文篇均被引次数与全世界同一学科专业论文篇均被引次数的比值)排出 21 个学科专业美国的前十名研究大学^[8]。在引文相对影响排名中,同一院校有 6 个以上学科专业名列前 10 名的大学共 10 所,它们在前一种多项指标排名榜上 1989~1993 年几乎全部名列前 20 名;其中 9 个以上学科专业名列前 10 名的 5 所大学几乎均排在综合评价排名榜的前 6 名。引文相对影响的前 10 名与相关学科研究生院的多项指标综合评价排名的相关性更强:按引文相对影响排名前 10 名的经济学/商学和法学院校在 1990~1994 年的商学和法学研究生院排名榜上均列前 18 名;化学、生物/生化和计算机科学的前 10 名分别有 6 或 7 名列入 1993 年相关学科综合评价前 8 名研究生院中。

研究表明:通过引文数据的积累分析,可

以说明个人、刊物、学科专业、院校研究机构、国家或地区研究工作的相对影响。通过简单的量化客观运算可以产生出与繁杂的主观定性选择排序及多项定量和定性相结合的分析结论高度相关的结果。但引文数据分析也和其他量化指标一样,有其自身的局限性。因此,在具体应用时,应区别不同情况,分析对待。

首先,引文数据取样必须保证一定的量。一般说,取样群体越大,时间跨度越长,其可靠性越强。在个体级的特定作者或特定刊物上,其局限性往往尤为明显,一定要结合其他情况全面分析。例如:1950~1977 年获诺贝尔物理、化学和医学奖的 162 人,在 1961~1975 年间人均被引 2877 次,最高为 18888 次,最低仅 79 次。被引数最低的是 1963 年 3 位诺贝尔物理学奖得主之一、德国的 J·H·D·Jensen。他只发表过 14 篇论文,均用德文出版,且都在 1961 年 SCI 创刊以前发表的,此外,他 1949 年论述原子核结构的论文在发表后作为被引论文不久即被他人的论著所取代^[9]。这些原因造成他的论文很少被引用,但不能抹杀他开创性的研究成果。

其次,要注意不同学科专业的区别。由于不同学科专业的发展历史、从业人员的多少、参阅和引用文献的习惯和多少等因素的影响,引文数差别很大。在基础工程方面,一篇被引用 30~40 次的论文就可能是经典,而在生命科学方面,被引用 30~40 次的论文可能有数百万篇,只是极普通的论文^[10]。在 1981~1993 年期间,根据 SCI 的统计,分子生物学、免疫学方面的所有论文篇均被引 14~15 次,而数学、计算机科学和农业科学则只有 3~3.5 次,工程方面则不足 3 次^[11]。在不加学科专业区分的全面引文统计和引文影响(篇均被引次数)排序中,排名前者基本为生命科学、医学方面的科学家和论文所把持。例如:1981~1990 年期间发表的论文在同期被引用最多的 100 名科学家中,除 2 名物理学 1 名化学科学家外,几乎全部都是医学、生命科

学方面的科学家^[12]; 1945~ 1988 年期间全世界被引用最多的 400 篇论文中, 生命科学和医学方面的论文占近 4/5^[13]。

同一学科的不同专业方向的论文和作者的被引情况也有很大差别。如 1965~ 1978 年被引用最多的 1000 名科学家中, 理论化学的作者人均被引 5227 次, 分析化学的作者人均被引 2822 次, 物理化学、有机金属化学、无机化学和有机化学的作者人均被引次数在 3600~ 3800 次左右^[14]。与此相应, 一些小专业的杰出科学家甚至诺贝尔奖获得者也很难在不分学科专业的引文排名榜前数百名甚至上千名的名单中出现。如射电天文学家、诺贝尔奖获得者 A · A. Penzias 和 R · W. Wilson 在 1965~ 1978 年被引最多的 1000 名科学家中就未出现。但在射电天文学方面, 他绝对排在被引最多的前 5 名中^[15]。

第三, 取样的时期、论文种类的界定, 甚至检索手段和策略都会对引文统计带来重大影响。例如: 同样都从 SSCI 为统计源对社会学方面论文的被引情况进行的两次统计分析结果就大为不同。D · Hamilton 统计 1984 年发表的社会学论文到 1988 年底被引用过的仅为 22.6%; D · M. Boff 和 L · L. Hargens 对 1974 年发表的社会学论文进行统计, 到 1975 年底被引用过的就达 45.4%, 到 1980 年底 85.4% 的论文都被引用过^[16]。据分析, 产生差别的原因主要是前者利用计算机检索, 对源文献和引文间的注录差别无法进行纠正; 而后者在查检时, 经核对纠正了此类错误, 查出前者漏检的引文。其次, 前者用的源文献不仅包括后者用的“论文”和“研究快讯”, 而且还包括 SSCI 收录的“书评”、“信函”、“订正”等并非真正的“研究论文”, 从而扩大了源文献的数量。

根据 Gerfield 的估计, 被引用最多的 300 名作者, 在统计取样基础年份每变一年时, 人名变化约为 7.5%^[17]。

第四, 由于事物发展和人们认识过程不

同, 论文发表后的被引时机和时间长短亦不同, 科学上的发现被人们承认、重视、接受一般都有一个过程, 过程长短不同, 从发展看总体上这一过程越来越短。在科学史上, 有些发现要多少年甚至几十年后才得到承认的例子并不是很个别的, 遗传学家 Mendel 就是一例^[18]。也有一些发现一出现就立即引起广泛注意, 很快成为“热点”, 从而被大量引用。Garfield 曾把被大量引用的论文的被引用情况归纳为 5 类: 高速飞弹型、流星型、迟开之花、双峰型和持久型^[10]。

当然, 在科学史上也有一些划时代的发现在发表后不久便被融汇为普遍承认的常识, 本应很高的被引次数却没有持续下去。爱因斯坦的相对论, Watson 和 Crick 有关 DNA 双螺旋结构的论文就是典型。

第五, 多作者论文被引次数的计算问题。随着多学科、跨学科研究的迅速发展, 多作者论文大增, 有些论文作者达数百人。随之而来的是作者个人被引次数的计算问题。据统计, 以第一作者身份的被引次数和以各种(包括第一和非第一)作者身份的被引次数间差距很大。如 1961~ 1975 年期间, 诺贝尔奖获得者 Glaser, SCC Ting, A · R. Prokhorov 以第一作者身份被引次数分别为 101, 170 和 146; 而把非第一作者身份的被引数计算在内, 则分别为 343, 303 和 1031^[20]。1965~ 1978 年被引最多的 1000 名科学家, 按 38 个学科专业分类, 各类中的作者人均以第一作者身份被引数超过以非第一作者身份被引数的只有天体物理学和神经生物学两类, 特别在核医学方面, 人均以第一作者身份被引数(403)不足以非第一作者身份被引数(4051)的 1/10。这 1000 名科学家中, 包括斯坦福线性加速器中心的 23 位合作者, 他们的一篇合作论文被引 626 次, 如果没有这篇论文, 他们中的 10 人就不会出现在这 1000 名科学家中。

第六, 同姓名作者的区分问题。由于西文参考文献中的名字普遍采用首字母简写, 许

多刊物也不给出作者全名。因此,在引文索引中同一姓名往往代表着多个甚至十几、数十位不同的作者。在 1965~1978 年被引用最多的 1000 名科学家中经姓名甄别后方定下来的达 15%^[21]。在 1961~1975 年被引用最多的 250 名科学家中, K·Alder 是德国诺贝尔化学奖获得者,在核对他的被引论文时发现其中有 1958 年他过世后发表的,经查核,那是另一位被大量引用的瑞士物理学家 K·Alder 的。E·Fischers 仅在世界科学家名人录中就有 7 位,甚至有两位 Emil Fischers 分别为德国诺贝尔化学奖获得者和瑞士生理学家^[22]。因此,在进行个人被引次数统计时要注意鉴别。这种情况在中国人以西文发表文章时更为突出,在以西文发表文章的中国人名鉴别中还要注意的一点是复名有时取 1 个首字母简写,有时取 2 个字母简写,即复名中每个字的首字母。不考虑这一点往往会漏检被引用次数。

此外,在对被大量引用论文的看法上,有些被引科学家,其中包括相当一批诺贝尔奖获得者,就不认为被引用最多的论文是他最重要的。Heine Fraenkel-Conrat 及其合作者都认为每人至少可推出比被引最多的那篇更重要的另外 10 篇论文;哪篇论文被大量引用并不是出于它的精髓,而是出于它里面谈的方法^[23]。相当普遍的看法认为被大量引用的论文中,方法方面的论文占很大优势。到 1990 年被引用超过 20 万次的 Oliver Lowry 1951 年的论文就是例证。对论文的引用也有多种不同的情况,有全面引证,有仅引用其中的字词、数据。P·Pichappon 曾把被引和引用他人论文的论文间的关系分为四类^[24]。因此,在分析引文数据时,还应考虑引用强度,这就需要结合进行内容分析。

D·J·D·Price 认为引文多少与交流的有效性有很大关系。国际现在通用引文统计源 SCI 是以英文刊物为主的检索刊物,因此在引文统计上英语论文占绝对统治地位。此外,某

些刊物在世界上以多种语言出版发行,因此这些刊物上的论文引文就可能重复计算。

在引文统计中,自引和同一研究群体内的互引也是一个重要影响因素。据分析,国际上自引一般占被引数的 13% 左右^[25]。我国科学家被 1989~1993 年 SCI 收录的论文,1994 年被引数中自引占 17.5%,被 1990~1994 年 SCI 收录的论文,在 1995 年的被引数中自引占 19.2%;两者的被引论文篇数中,自引分别占 24.7% 和 27.0%^[26]。

尽管以引文评价科研工作还可提出那样那样的局限性,但它毕竟是我们可以得到的世界科学界对科学家及其论著看法的最精确指标^[27]。只要进行认真比较、斟酌,适当地结合专家的分析、解释和同行评审,并通过统计学的权重、控制和补偿等措施,它会在研究工作质量评审中发挥最佳的量化指标作用。

参考文献

- 1 Juan Miguel Campanario. Have Referees Rejected Some of the Most-Cited Articles of All Times Journal of the American Society for Information Science 1996, 47(4)
- 2 J. Leta and L. DeMeis A Profile of Science in Brazil Scientometrics 1996, 35(1): 33~44
- 3 E. Garfield The Most-Cited Papers of All Times, SCI 1945~1988 Part 1A. Current Contents 1990, (7)
- 4, 12 E. Garfield and A. Welljam's Dorof. Of Nobel Class: A Citation Perspective on High Impact Research Authors Current Contents 1992, (33), (35)
- 5, 9, 20, 22, E. Garfield The 250 Most-Cited Authors, 1961~1975 Part 1, 2 Current Contents 1977, (49), (50)
- 6, 14, 15, 21. E. Garfield, The 1000 Most-Cited Contemporary Scientists, 1965~1978 Part 2A~2D. Current Contents 1982, (9), (21), (22), (24)
- 7 S. Cole, et al Peer Review and the Support of Science Scientific American 1977, 237: 34~41

- 8 C. King, et al, America's Best Research University? Stanford Soars in Top Ten Tournament Science Watch. 1994, 5(9)
- C. King, et al Harvard Hot in Physical Sciences as Top Ten Tournament Terminates Science Watch. 1994, 5(10)
- 10, 23, 27 E. Garfield, Citation Classics-Four years of the Haman Side of Science Current Contents 1981, (22)
- 11 C. King, et al, Scotland Harrests Bumper Crop of Highly Cited Studies on Agriculture Science Watch. 1994, 5(4)
- S. M. Itton. South African Science A nem ic Serious Action Surely Advisable Science Watch. 1995, 6(3)
- 13, 19 E. Garfield, The Most-Cited Papers of All Times, SCI 1945~ 1988 Part1-4 Current Contents 1990, (7), (8), (26), (34); 1991, (21)
- 16 D. M. Boff and L. L. Hargens Are Sociologists' Publications Uncited? Citation Rates of Journal Articles, Chapters, and Books Current Contents 1993, (5)
- 17 E. Garfield, The 1000 Contemporary Scientists Most-Cited 1965~ 1978 Part1. Current Contents 1981, (41)
- 18 Don Schauder, Electronic Publishing of Professional Articles: Attitudes of Academics and Implications of the Scholarly Communication Industry. Journal of the American Society for Information Science 1994, 45(2)
- 24 P. Pichappon, Levels of Citation Relations Between Papers Journal of the American Society for Information Science 1996, 47(8)
- 25 E. Garfield, The Impact Factor. Current Contents 1994, (25)
- 26 中国科技信息研究所 中国科技论文统计与分析(年度研究报告). 1994, 1995
- 袁培国 通信地址: 南京市汉口路 22 号, 南京大学图书馆. 邮编: 210093.
- (来稿时间: 1997. 8. 21. 编发者: 李万健)

(上接第 52 页) 始, 凡订《广州日报》一年或一年以上的社会各界人士, 只需交纳 45 元的借书押金及 5 元的手续费, 均可办理该中心的借书证, 免费上机检索该报数据库, 上网浏览电子报刊和借阅图书报刊资料。迄今已有 200 多人办理了借书证, 其中以大专以上学历的人居多。可见, 不但新闻工作者需要地方报纸信息, 其他工作在各条战线的人们也有着同样的需求。尽管《广州日报》图书阅览中心的做法距离地方报纸信息数据库商品化的要求尚远, 但毕竟开了向社会提供电子版地方报纸信息服务的先河, 其经验十分值得推广和借鉴。

笔者设想, 各省市、地区地方报社联合建设本省市、本地区的报纸信息数据库, 并通过信息网络与本省市的其他专业报社, 其他各省市的报纸信息网络联通, 组成全国报纸信息网络, 进入全国的信息国道, 走向国际互联网络。在这样的信息网络的环境下, 地方报纸

信息数据库向社会开放, 提供有偿服务, 将会使地方报纸信息数据库兼得社会效益和经济效益, 走上良性循环的发展道路。

参考文献

- 1 宋明亮 我国报纸信息数据库开发的现状与对策 中国图书馆学报, 1995(1): 60~ 65
- 2 李献线 关于建设和发展我国新闻数据库的思考 中山大学硕士学位论文, 1996: 20
- 3 曹树金, 罗春荣 论信息网络时代图书馆的数据库建设和服务选择 高校文献信息学刊, 1994(3~ 4): 24~ 27
- 4 Long, A. Full-text newspaper retrieval is hard to manage: fact or fiction. Proceedings of the 9th National Online Meeting: Medford, 1988: 213~ 216

潘燕桃 中山大学信息管理系硕士毕业, 现为该系讲师。通讯地址: 广州市。邮编 510275。

(来稿时间: 1997. 8. 21. 编发者: 翟凤岐)