

关键词共引分析的科学计量方法研究

黄文彬 王冰璐 (北京大学信息管理系 北京 100871)

步 一 (印第安纳大学信息学、计算机与工程学院 印第安纳布鲁明顿 47408)

闵 超 (南京大学信息管理学院 江苏 210023)

摘 要 文献共引分析、作者共引分析和期刊共引分析等共引分析方法经常作为绘制知识图谱的重要手段。本研究以引文关键词为基础扩展传统共引分析方法,提出关键词共引分析(Keyword Co-citation Analysis,KCA)方法作为绘制知识图谱的新路线。该方法主要包括数据集遴选、原始关键词共引矩阵构建、相关矩阵转化和知识图谱绘制及其评估四个主要步骤。网络分析和MDS测度分析结果显示,KCA方法所绘制的知识图谱能够清楚描绘出领域的知识结构,且聚类性能较好。

关键词 引文关键词 共引分析 引文分析 文献计量学 科学计量学

A Study on Scientometrics of Co-citation Analysis of Keywords

Huang Wenbin Wang Binglu (Department of Information Management, Peking University, Beijing, 100871)

Bu Yi (School of Informatics, Computing, and Engineering, Indiana University, Bloomington, Indiana, 47408)

Min Chao (School of Information Management, Nanjing University, Jiangsu, 210023)

Abstract Co-citation analysis mainly includes document co-citation analysis, author co-citation analysis, and journal co-citation analysis. They are frequently-used methods in mapping knowledge domains. As an alternative exploration in mapping knowledge domains, this paper expands traditional co-citation methods into Keyword Co-citation Analysis (KCA). The basic steps of KCA include data collection, raw keyword co-citation matrix construction, correlation matrix transformation, and knowledge domain mapping and results interpretation. The results based on network analysis and MDS-measurement show that KCA has a good performance in mapping knowledge domains and depicting scientific intellectual structures; meanwhile, it has satisfying clustering results.

Keywords citation keyword, co-citation analysis, citation analysis, bibliometrics, scientometrics

1 引言

知识图谱主要用于展示科学知识的发展过程与可视化领域结构关系^[1],并已被许多不同领域研究人员应用于探索领域的学科发展状况和科学知识结构(scientific intellectual structures)^[2,3]。目前,知识图谱的研究视角包含以文献调研和综述分析为主的传统研究范式、以科学社会学为基础的研究范式、基于引文分析的文献计量范式和基于复杂网络理论的社会网络分析研

究范式^[4]。其中,在基于引文分析的文献计量的知识图谱绘制研究中,引用关系(citation)、共引关系(co-citation)、耦合关系(bibliographic coupling)、合著关系(co-authorship)和共词关系(co-word)是五种常被使用的学术网络关系^[4-7]。在上述几种文献计量关系当中,共引关系主要描绘了两个书目要素(如文献、期刊、作者等)共同被其他书目要素引用的关系。由于数据简单可获取,且具有一定的动态性和展望性,共引关系是知识图谱绘制中较常使用的一种学术关系^[4,8]。共引分析(co-

citation analysis)通过计算书目要素两两间的共引数量(频次)得到原始共引矩阵(raw co-citation matrix),并通过一系列转化、数据分析和可视化手段,绘制出某领域的知识图谱。常见的共引分析方法主要有文献共引分析(Document Co-citation Analysis, DCA, 1973年由Small H提出)^[9]、作者共引分析(Author Co-citation Analysis, ACA, 1981年由White H D和Griffith B C共同提出)^[10]和期刊共引分析(Journal Co-citation Analysis, JCA, 1991年由McCain K W提出)^[11]。通过共引分析可以发现某研究领域中处于研究前沿且备受关注、多次被引的书目要素,找寻作者的研究路径和研究偏好,进而促进学术合作和学术交流^[8,12]。特别是ACA,该方法已被广泛地应用在许多领域^[13-15],作为评估该领域发展现况和科学知识结构分析的参考^[12];有关该方法的改良也层出不穷^[16-20]。

情报学家Morris和Martens^[7]列举出文献计量学中包括论文、论文作者、发表论文期刊、论文索引词、参考文献、参考文献作者和发表参考文献期刊在内的七个基本书目要素,如图1所示。在共引分析方法中,对参考文献、参考文献作者和发表参考文献期刊的分析恰好对应了上段所述的DCA、ACA和JCA。然而,除了上述书目要素外,参考文献索引词(关键词)也可作为计量书目的重要指标,因为它可以鲜明直观地表述该参考文献论述的主题,是该参考文献内容与主题的最直观反映。分别包含有关键词A和B的两篇文献被同时引用(共引),在本文中称作关键词A和B被共引。

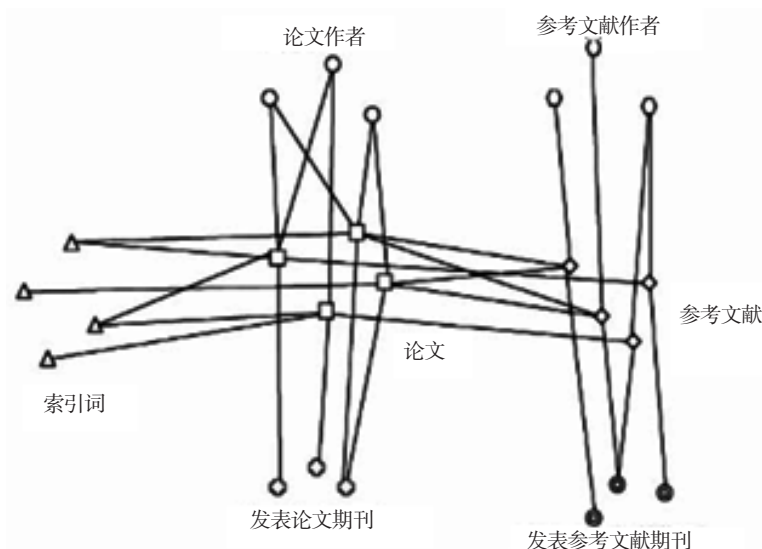


图1 文献计量学中七个基本书目要素^[4,7]

本文认为,两个关键词被共引,体现了这两个关键词在内容或主题上有一定的联系;两个关键词被共引的次数越多,它们之间的这种联系越强烈。这一基本假设与ACA的基本假设^[10,21]极为类似。基于此,本文尝试从参考文献关键词的角度扩展传统共引分析方法,提出关键词共引分析(Keyword Co-citation Analysis, KCA)方法,并进行了实证研究。实证研究显示,与其他传统共引分析方法相比,KCA绘制出的知识图谱更为直观,其聚类效果也较为令人满意。

2 关键词共引分析方法(Keyword Co-citation Analysis, KCA)

KCA方法依序进行数据集遴选、原始关键词共引矩阵构建、相关矩阵转化和知识图谱绘制及其评估四个主要基本步骤进行,如图2所示。

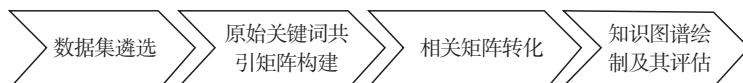


图2 关键词共引分析方法的步骤

2.1 数据集遴选

数据集遴选的方式与ACA相似,可选择较大的宏观学科,也可选择微观学科或者一个专题^[8,22]。数据库的选择上,可以选择Web of Science或Elsevier Scopus等综合型数据库,也可选择Digital Bibliography and Library Project (DBLP)等专指收录的数据库^[8,22]。接着对数据库以某领域学科索引词进行检索,并以人工采集或机器获取的方式抽取参考文献内的关键词,并利用自然语言处理技术对关键词进行相似归并(如同义词归并、单复数及其他同源或相似形式的归并)。

2.2 原始关键词共引矩阵构建

数据集中被归并处理后的两个关键词被其他文献同时引用1次,该关键词对的共引频次相应加1,依照这种方法可以计算出数据集中所有关键词之间的共引频次。然后,构建一个对称矩阵(原始关键词共引矩阵),其中矩阵的行列标识均为数据集中的关键词。关于该矩阵的对角线元素处理,可以采用ACA常用的五种处理方式^[8]。简单起见,可以全部置为空缺值或0。

设在步骤2.1中遴选出的数据集为 Ω ,其中包含了 λ 个经过归并处理的关键词,记

作 $k_1, k_2, \dots, k_\lambda$ 。根据传统共引分析的定义,假如关键词 k_i 和 k_q 所在的文献被其他文献共引了 $m_{i,q}$ ($1 \leq i, q \leq \lambda$)次。这样,数据集 Ω 所包含的全部 λ 个关键词两两之间的共引次数可构成一个原始关键词共引矩阵 M 。 M 是一个 $\lambda \times \lambda$ 的对称矩阵,其行(列)标识为关键词 $k_1, k_2, \dots, k_\lambda$,其元素由 $m_{i,q}$ 组成,代表了相应关键词对的共引次数。

2.3 相关矩阵转化

原始关键词共引矩阵一般较为稀疏,需要进行必要的缩减,即将非零元素数过少的行列进行删除处理^[8,21]。设经过缩减后的对称矩阵为 M' ,其行列数均为 n ,矩阵元素记为 m_{ij}' ($1 \leq i, j \leq n$)。此外,与ACA类似,经缩减后的共引矩阵,还需要进行标准化处理,以便于进一步核查分析对象间的相关程度,使KCA分析结果更易于解释和接受^[8,12]。Pearson相关系数法、余弦系数、Jaccard系数、欧氏距离和Chi-square计算等不同方法均可应用于KCA的相关矩阵转化中。关于这些相似测度的优缺点及应用范围,请参照参考文献[8]和[23]。本文使用的Pearson相关系数进行矩阵转化。设 M' 经转化后成为 $n \times n$ 的矩阵 M'' ,该矩阵的行列标识均和 M' 相同。其中, M'' 的元素 $m_{i,j}''$ 的计算方式如下:

$$m_{i,j}'' = \frac{n \sum_{p=1}^n m_{ip}' m_{jp}' - (\sum_{p=1}^n m_{ip}') (\sum_{p=1}^n m_{jp}')}{\sqrt{n \sum_{p=1}^n m_{ip}'^2 - (\sum_{p=1}^n m_{ip}')^2} \sqrt{n \sum_{p=1}^n m_{jp}'^2 - (\sum_{p=1}^n m_{jp}')^2}}$$

2.4 知识图谱绘制及其评估

知识图谱绘制可以采用SPSS多维尺度分析(Multi-Dimensional Scaling, MDS)或Gephi^[24]等一系列可视化工具。在KCA的可视化知识图谱中,节点表示经过处理后的关键词,边表示关键词之间经过相关转化后的共引频次,节点间距离表示关键词之间的紧密程度(共引频次越高,越紧密,距离越小)。知识图谱的聚类可以使用人工聚类或者使用布局算法(如Yifan Hu^[25]或ForceAtlas2^[26]算法等)进行聚类。此外,研究人员也可以应用聚类分析(clustering analysis)或因子分析(factor analysis)等技术来辅助对知识图谱的解读。

除了通过领域专家对绘制后的知识图谱进行定性的评估外,从定量层面对知识图谱进行评估也是必要的。MDS测度分析(MDS-measurement)是常用定量评估ACA的手段之一^[17],主要通过计算知识图谱中类内节点距离和与类间节点距离之和的商值 σ 来做算法的效能评估。 σ 值越小说明知识图谱的聚类效果越好,从而表明知识图谱在定量层面上有较好的绘制结果——同类节点较为聚拢,不同类节点较为分散。由于MDS测度分析对于知识图谱的绘制工具和计量单位具有较

好的普遍适用性,即使用不同的绘制工具和不同的距离测度对于知识图谱的MDS测度值 σ 没有显著影响,KCA也可以类似地利用MDS测度分析进行知识图谱的定量评估。

3 实证结果及其分析讨论

3.1 数据及其处理

Journal of the Association for Information Science and Technology(即《美国信息科学与技术学会期刊》,下文简称JASIST)于1950年创刊,其在2016年《期刊印证报告》(*Journal Citation Report*, JCR)^[27]公布的影响因子高达2.322,是图书情报学领域的国际顶级期刊。本文从美国科学情报所(Institute for Scientific Information, ISI)的Web of Science(WoS)数据库中选取2012~2015年4年间JASIST出版的866篇原文和与之对应的38910篇引文,包含每篇文章的题名、作者、发表时间、卷期号、引文第一作者和引文发文年份、引文DOI等相关信息。随后,本文在WoS数据库中利用引文的DOI信息进行引文关键词采集,每篇文章保留至多5个关键词,并对关键词进行编码。具体做法包括:(1)对大小写不同的关键词和特殊符号进行合并,例如“HINDEX”和“h-index”;(2)对单复数和动名词等同源关键词进行合并,例如“citation”和“citations”、“ranking”和“rank”;(3)词频统计,并选取其中排名最高的200个关键词。经过一系列人工消歧工作,我们最终在200个关键词中遴选了100个该领域关键词用以绘制关键词共引分析知识图谱。表1展示了词频超过50的关键词及其出现频次。

表1 KCA中的高频关键词及其频次

序号	关键词	词频	序号	关键词	词频
1	citation	265	7	search	63
2	science	176	8	impact	62
3	information	152	9	indicator	61
4	relevance	138	10	bibliometrics	60
5	communication	75	11	co-citation	55
6	web	65	12	journal	51

3.2 知识图谱分析

下页图3展示了利用本文提出的方法根据所遴选出的100个关键词所绘制的知识图谱。节点代表遴选出的高频关键词,节点的大小与该节点的加权重(即某个高频关键词与其他所有高频关键词的共引值之和)成正比,边代表由该边相连的两个节点(即两个关键词)存在共引关系,边的粗细与该边相连的两个节点所代表关键词的共引强度成正比。此外,节点间的位置越近,代表节点所指示的关键词之间在共引关系中的

KCA 具有较好的聚类效果。虽然本实证研究中的 σ 值不及一些经过改进的ACA方法 σ 值小,但相对这些方法中,KCA需要的输入信息仅有引文关键词一项,输入量小、运算快速。因此,从MDS测度分析的定量角度看,使用KCA方法绘制的知识图谱聚类效果较好,意即类内各点较聚拢、类间各点较离散。

表2 已有文献中的MDS测度分析结果值参考

文献	使用的共引分析方法	MDS测度值(σ)
本文	1 KCA方法	11.74
文献[16]	2.1 传统ACA方法	13.97
	2.2 结合DCA思想的ACA方法	10.32
	2.3 结合引文、原文时间信息的ACA方法	10.35
	2.4 结合DCA思想和引文、原文发表时间信息的ACA方法	7.22
文献[17]	3.1 传统ACA方法	12.18
	3.2 结合引文发表时间信息的ACA方法	9.96
	3.3 结合引文发表期刊信息的ACA方法	10.47
	3.4 结合引文关键词信息的ACA方法	9.91
	3.5 结合引文发表时间和引文发表期刊信息的ACA方法	9.25
	3.6 结合引文发表时间和引文关键词信息的ACA方法	8.83
	3.7 结合引文发表期刊和引文关键词信息的ACA方法	9.08
	3.8 结合引文发表时间、期刊和引文关键词信息的ACA方法	8.12

与传统共引分析相比,KCA在绘制知识图谱方面的优势有三:第一,图谱解读与理解方便直接。在使用ACA或者DCA进行知识图谱绘制过程中,绘制出的可视化图形中节点表示的是作者或者文献,不仅需要研究人员查阅大量文献进行图谱含义解读,还不利于读者直观理解图谱;但KCA生成的知识图谱,其节点表示的是参考文献关键词(索引词),节点的含义可被研究人员和读者直接理解,方便直接。第二,不需要进行作者姓名消歧。传统ACA为了保证分析的精确性,在原始共引矩阵生成之前,需要花费大量的精力进行作者姓名消歧,其高昂的计算代价使得很多研究人员在使用之时望而却步;而KCA只需要进行关键词的归并,这在目前的自然语言处理中属于较为基本的工作。第三,KCA的研究思想有助于推动领域本体研究的发展。高频关键词本身便代表了学科的研究热点和主要内涵,共引关系所构建的词间关系一定程度上代表了词与词之间的相似度,可以进一步研究其具体关系含义,即上下位类、同义和参见关系等,从而帮助学科更好地认识其内在的体系构建,做好本体的研究和进一步挖掘。

4 结语

本文以JASIST期刊2012~2015年的学术论文作为数据集,通过参考文献的共引关系计算关键词的使用频率与引用关系,提出了关键词共引分析方法

(KCA)绘制领域知识图谱,从不同视角挖掘关键词之间的潜在关系。实证结果显示,以关键词为研究主体的KCA研究能够一定程度上展现该领域的知识图谱。另一方面KCA所反映的词间相似度可以为后续领域本体的研究基础,将关键词共引分析中的“关键词”概念进行扩展,利用话题建模(topic modeling,例如LDA模型^[37])等方法代替关键词进行共引分析。例如,使用LDA模型的作者-会议-话题(Author-Conference-Topic,ACT)算法^[38]计算书目要素在一定数量话题下的概率分布,形成书目要素的话题分布向量;通过计算向量之间相似度得到书目要素之间的相似度。ACT算法可以用于文献、作者、期刊等各个书目要素,作为参考文献关键词的替代与补充。此外,自然语言处理领域中的主题抽取(topic extraction)技术^[39]也可应用到共引分析中进行更为精确的知识图谱绘制。

参考文献

[1] 刘则渊,陈悦,侯海燕. 科学知识图谱:方法与应用[M]. 北京:人民出版社,2008:3-5.

[2] 薛晓芳. 知识可视化理论、方法和工具及军事医学应用研究[D]. 北京:中国人民解放军军事医学科学院,2014.

[3] 辛伟,雷二庆,常晓等. 知识图谱在军事心理学研究中的应用:基于ISI Web of Science数据库的CiteSpace分析[J]. 心理科学进展,2014,22(2):334-347.

[4] 赵丹群. 试论科学知识图谱的文献计量学研究范式[J]. 图书情报工作,2012,56(6):107-110.

[5] Yan E, Ding Y. Scholarly network similarities: how bibliographic coupling networks, citation networks, co-citation networks, topical networks, coauthorship networks, and co-word networks related to each other[J]. Journal of the American Society for Information Science and Technology, 2012, 63(7):1313-1326.

[6] Boyack K W, Klavans R. Cocitation analysis, bibliographic coupling, and direct citation: which citation approach represents the research front most accurately? [J]. Journal of the American Society for Information Science and Technology, 2010, 61(12):2389-2404.

[7] Morris S A, Van der Veer Martens B. Mapping research specialties [J]. Annual Review of Information Science and Technology, 2008, 42(1):213-295.

[8] 步一,刘天祯,赵丹群等. 国外作者共引分析研究评述[J]. 情报杂志,2015,34(12):48-53.

[9] Small H. Co-citation in the scientific literature: a new measure of the relationship between two documents [J]. Journal of the American Society for Information Science, 1973, 24(4):265-269.

[10] White H D, Griffith B C. Author cocitation: a literature measure of intellectual structure [J]. Journal of the American Society for Information Science, 1981, 32(3):163-171.

- [11] McCain K W. Mapping economics through the journal literature: an experiment in journal cocitation analysis [J]. Journal of the American Society for Information Science, 1991, 42(4): 290-296.
- [12] 步一, 刘天祯, 黄文彬. 优化传统作者共引分析的研究初探——综合引文发表时间信息的作者共引分析方法 [J]. 图书情报知识, 2015, 32(6):89-97.
- [13] Ma R. Discovering and analyzing the intellectual structure and its evolution of LIS in China, 1998-2007 [J]. Scientometrics, 2012, 93(3):645-659.
- [14] Zhao R, Chen B. Applying author co-citation analysis to user interaction analysis: a case study on instant messaging groups [J]. Scientometrics, 2014, 101(2):985-997.
- [15] Chen C. Visualizing semantic spaces and author co-citation networks in digital libraries [J]. Information Processing and Management, 1999, 35(3):401-420.
- [16] 黄文彬, 步一, 王冰璐. 作者共引分析方法的扩展与效能改进研究 [J]. 图书情报知识, 2017, 36(2):75-82.
- [17] Bu Y, Liu T, Huang W-B. MACA: a modified author co-citation analysis method combined with general descriptive meta-data of citations [J]. Scientometrics, 2016, 108(1):143-166.
- [18] Persson O. All author citations versus first author citations [J]. Scientometrics, 2001, 50(2):339-344.
- [19] Zhao D, Strotmann A. Counting first, last, or all authors in citation analysis: a comprehensive comparison in the highly collaborative stem cell research field [J]. Journal of the American Society for Information Science and Technology, 2011, 62(4): 654-676.
- [20] Jeong Y-K, Song M, Ding Y. Content-based author cocitation analysis [J]. Journal of Informetrics, 2014, 8(8):197-211.
- [21] McCain K W. Mapping authors in intellectual space: a technical overview [J]. Journal of the American Society for Information Science, 1990, 41(6):433-443.
- [22] Eom S. Author Co-Citation Analysis: Quantitative Methods for Mapping the Intellectual Structure of An Academic Discipline [M]. Hershey, NY: Information Science Reference, 2008.
- [23] Mëgnigbëto E. Controversies arising from which similarity measures can be used in co-citation analysis [J]. Malaysian Journal of Library and Information Science, 2013, 18(2):25-31.
- [24] Bastian M, Heymann S, Jacomy M. Gephi: an open source software for exploring and manipulating networks [C]. Proceeding of the 3rd International Conference on Web and Social Media (ICWSM), 2009:361-362.
- [25] 胡一凡. 最优结构图:从复杂关系中发现规律 [M]. 上海: 上海科学技术出版社, 2007.
- [26] Jacomy M, Venturini T, Heymann S, et al. ForceAtlas2: a continuous graph layout algorithm for handy network visualization [J]. PLoS ONE, 2014, 9(6):e0098679.
- [27] <http://thomsonreuters.com/en/products-services/scholarly-scientific-research/research-management-and-evaluation/journal-citation-reports.html>.
- [28] Newman M E. Modularity and community structure in networks [J]. Proceedings of the National Academy of Sciences of the United States of America, 2006, 3(23):8577-8582.
- [29] White H D, McCain K W. Visualizing a discipline: an author cocitation analysis of information science 1972-1995 [J]. Journal of the American Society for Information Science, 1998, 49(4):327-335.
- [30] Yu L. Understanding information inequality: making sense of the literature of the information and digital divides [J]. Journal of Librarianship and Information Science, 2006, 38:229-252.
- [31] Hearst M A. TileBars: visualization of term distribution information in full text information access [C]. Proceedings of the SIGCHI Conference on human factors in computing systems ACM Press/Addison-Wesley Publishing Co., 1995:59-66.
- [32] Hirsch J E. An index to quantify an individual's scientific research output [J]. Proceeding of the National Academy of Sciences of the United States of America, 2005, 102:16569-16572.
- [33] Yan E, Ding Y. Discovering author impact: a PageRank perspective [J]. Information Processing and Management, 2011, 47(1):125-134.
- [34] Ding Y. Scientific collaboration and endorsement: network analysis of coauthorship and citation networks [J]. Journal of Informetrics, 2011, 5(1):187-203.
- [35] Ebadi A, Schiffauerova A. How to receive more funding for your research? Get connected to the right people! [J] PLoS ONE, 2015, 10(7):e0133061.
- [36] <http://www.scienceofteams.org/>.
- [37] Hastie T, Tibshirani R. Discriminant analysis by Gaussian Mixtures [J]. Journal of the Royal Statistical Society, 1996, 58(1): 155-176.
- [38] Tang J, Jin R, Zhang J. A topic modeling approach and its integration into the random walk framework for academic search [C]. Proceeding of the 2008 Eighth IEEE International Conference on Data Mining IEEE, 2008:1055-1060.
- [39] Chen K Y, Luesukprasert L, Chou S C T. Hot topic extraction based on timeline analysis and multidimensional sentence modeling [J]. IEEE Transactions on Knowledge and Data Engineering, 2007, 19(8):1016-1025.

[作者简介] 黄文彬, 男, 1977年生, 北京大学信息管理系副教授。

王冰璐, 女, 1995年生, 北京大学信息管理系本科生。

步一, 男, 1994年生, 印第安纳大学信息学、计算机与工程学院博士研究生(通讯作者)。

闵超, 男, 1990年生, 南京大学信息管理学院博士研究生。

收稿日期: 2017-06-23