学術俯瞰システムの使い方 How to Use Academic Landscape System



2015-01-21

Sustainability 2007 and 2013



Kajikawa et al., 2014. Sustainability science: the changing landscape of sustainability research.

<u>Agenda</u>

学術俯瞰システムの概要

Overview of Academic Landscape System
論文のデータの入手

How to Get your Dataset of Papers

学術俯瞰システムの操作方法

How to Operate the System

「学術俯瞰システム」とは

学術俯瞰システムは、大規模な書誌情報をテキストマイニングとネットワーク分析により 自動分析し、学術分野の俯瞰を可能にするシステムです。

The Academic Landscape System enables us to **overview science and technology** with text-mining and network analysis of huge bibliographic information.



<u>Background</u>:知識の爆発の現状(光触媒の例)

80年代は日本が3-4割。現在は日本が1割、中国が4割。

専門家も全てを把握しきれないほど、 情報が増えている。



学術俯瞰システムのしくみ The Concept of Citation Network

Paper A PHYSICAL REVIEW E 68, 066133 (2004)

Fast algorithm for detecting community structure in networks

M. E. J. Navman test of Physics and Center for the Study of Complex Journal, University of Michigan, Ann Arbor, Michigan 48109-1120, USA (Received 22 September 2003, revised mananeity in revised 22 March 2004, published 18 June 2004)

Many networks display community structure—groups of vertices within which connections are dense but between which they are sparser—and sensitive computer algorithms have in recent years been developed for detecting this structure. These algorithms, however, are computationally demanding, which limits their appli-

cation to small networks. Here we describe an algorithm which gives excellent omputer-generated and real-world networks and is much faster, typically th

previous algorithms. We give several example applications, including one to a than 50 000 physicists.

DOI: 10.1103/PhysRevE.69.066133

L INTRODUCTION

There has in recent years been a surge of interest within the physics community in the properties of networks of many kinds, including the Internet, the world wide web, citation for example, while the world wide web numbers in the billions [14

hous [17]. In this paper, therefore, we propose another algorithm for detecting community structure. The algorithm operates on different principles from that of Girvan and Newman (GN), but, as we will show, gives qualitatively similar results. The worst-case running time of the algorithm is $O((m^+n)n)$, or $O(n^2)$ on a sparse graph. In practice, it runs to completion on current computers in reasonable times for networks of up to a million or so vertices, bringing within reach the study of

FAST ALGORITHM FOR DETECTING COMMUNITY.

mitias in m

networks, transportation networks, software call graphs, email networks, food webs, and social and biochemical networks [1-4]. One property that has attracted particular attention is that of "community structure": the vertices in net-works are often found to cluster into tightly knit groups with a high density of within-group edges and a lower density of between-group edges. Girvan and Newman [5,6] proposed a computer algorithm based on the iterative removal of edges with high "betweenness" scores that appears to identify such structure with some sensitivity, and this algorithm has been employed by a number of authors in the study of such diemployed by a minose of sumori in the study of such di-twens systems as networks of enail messages, social net-works of animals, collaborations of jazz musicians, meta-bolic networks, and guess networks [5–11]. As pointed out by Newman and Gravan [6], the principal disability surges of the algorithm is the high computational demands it makes. In its simplest and fastest form, it runs in worst-case time $O(m^2n)$ on a network with m edges and n vertices, or $O(n^3)$ on a sparse network (one for which m scales with n in the limit of large n, which covers essentially all networks of current sciantific interest, with the possible exception of food webs). With typical computer resources available at the time of writing, this limits the algorithm's use to networks of a few thousand vertices at most, and substantially less than this for interactive applications. Increasingly, however, there is inter-est in the study of much larger networks; citation and col-laboration networks can contain millions of vertices [12,13].

Let a he the fra

69.066133-1

If a particular divi edges than would

Paper B

Paper C

1530.3755/3004/60/6/066133/5/822 50

FIG. 5. Cumulative distribution function of the sizes of commu

nities found in one of the subnetworks of the physics collaboration

graph, as described in the text. The dotted line represents the slope

the plot would have if the distribution followed a power law with

-1.6, although this conclusion should be treated with caution

as there is significant deviation from a perfect power law

[20]. Narrowing our focus still further to the particular one of mathematical and the state of the state o

these communities that contains the present author, we find

the structure shown in the right panel of Fig. 4. Feeding this one last time through the algorithm, it breaks agart into com-munities that correspond closely to individual institutional

research groups, the anthor's group appearing in the comer of the figure, highlighted by the dashed box. One could pur-

PACS number(a): 89.75.H been considered int

Our algorithm any network, the G produces some divis gardless of whether sion. To test whet define a quality for Let e_{ij} be one-hal that connect vertices the total fraction of : will be the diago fraction of edges the half). Then $\Sigma_{c,c,i}$ is groups. All other e value of this sum uities is go order 1. On its own if we not all vertic trivial and not par

A more useful a culate the sum $\Sigma_i e$ would take if edge gives a score of a single community. ings

to vertices in group by noting that a = together at rando connect vertices 1 larity to be

nent -1.6

R. Albert and A.-L. Barabiai, Rev. Mod. Phys. 74, 47 (2002). S. N. Dorogovtsev and J. F. F. Mendes, Evolution of Networks: From Biological Nets to the Internet and WWW (Oxford University Press, Oxford, 2003).
 M. E. J. Newman, SIAM Rev. 45, 167 (2003).

- [5] M. Girvan and M. E. J. Newman, Proc. Natl. Acad. Sci. U.S.A. 99, 7821 (2002)
- [6] M. E. J. Newman and M. Girvan, Phys. Rev. E 69, 026113 (2004) 171 D. Wilkinson and B. A. Huberman, Proc. Natl. Acad. Sci.
- U.S.A. 101, 5241 (2004). [8] P. Holme, M. Huss, and H. Jeong, Bioinformatics 19, 532
- (2003)[9] R. Guimerk, L. Danon, A. Diaz-Guilera, F. Giralt, and A. Arenas, Phys. Rev. E 68, 065103 (2003).
- [10] J. R. Tyler, D. M. Wilkinson, and B. A. Huberman, in Proceedings of the First International Conference on Con and Technologies, edited by M. Huysman, F. Wenger, and V.
- Wulf (Kluwer, Dordrecht, 2003). [11] P. Gleiser and L. Danon, Adv. Complex Syst. 6, 565 (2003).
- S. Redner, Eur. Phys. J. B 4, 131 (1998).
 M. E. J. Newman, Proc. Natl. Acad. Sci. U.S.A. 98, 404 (2001)

PHYSICAL REVIEW E 69, 066133 (2004

sue this line of analysis further, identifying individual groups, iteratively breaking them down, and looking, for example, at the patterns of collaboration between them, but we leave this for later studies

IV. CONCLUSIONS

In this paper we have described an algorithm for extracting community structure from networks, which has a consid-erable speed advantage over previous algorithms, running to completion in a time that scales as the square of the network size. This allows us to study much larger systems than has previously been possible. Among other examples, we have applied the algorithm to a network of collaborations between more than 50 000 physicists, and found that the resulting community structure corresponds closely to the traditional divisions between specialties and research groups in the

We believe that our method will not only allow for the extension of community structure analysis to some of the very large networks that are now being studied for the first time, but will also provide a useful tool for visualizing and understanding the structure of these networks, whose dounting size has hitherto made many of their structural properties obscime

ACKNOWLEDGMENTS

The author thanks Leon Danon, Pablo Gleiser, David Lusseau, and Douglas White for providing network data used in the examples. This work was supported in part by the Na-tional Science Foundation under Grant No. DMS-0234188.

S. H. Strogatz, Nature (London) 410, 268 (2001). [14] J. Kleinberg and S. Lawrence, Science 294, 1849 (2001).

- [15] B. Evenitt, Cluster Analysis (John Wiley, New York, 1974).
 - [16] J. Scott, Social Network Anahois: A Handbook, 2nd ed. (Sage London, 2000).
 - [17] W. W. Zachary, J. Anthropol. Res. 33, 452 (1977) [18] D. Lusseau, Proc. R. Soc. London, Ser. B 270, S186 (2003).
 - [19] The criterion for deciding correct classification is as follows We find the largest set of vertices that are grouped together by the algorithm in each of the four known communities. If the ithm puts two or more of these sets in the same group then all vertices in those sets are considered incorrectly classified. Otherwise, they are considered correctly classified. All other vertices not in the largest sets are considered incorrectly classified. This criterion is quite harsh-there are cases in which one might consider some of the vertices to have been

identified correctly, where this method would not. Even with this harsh definition, however, our algorithm performs well, and a laser definition would only make its performance more [20] This power law is different from the one observed in an email

network by Guimerà et al. [9]. They studied the histogram of community sizes over all levels of the dendrogram; we are looking only at the single level corresponding to the maximum value of O.

Namely, relations between papers are shown as a networks.



\downarrow More papers...





学術俯瞰システムのしくみ

The Flow of Analysis



Type of Citations





- direct-citation
- co-citation
- bibliographic coupling

Direct citation

- Edge between two documents when *paper C* cites *paper A* directly
- An earlier paper is cited by a new paper

Co-citation

- Edge between two documents cited by the same paper(s)
- Considering the number of citations

Bibliographic coupling

- Edge between two documents citing the same paper(s)
- Considering the number of citations

	Direct vs Co-citation vs Bibliographic
Nodes	Direct > Bibliographic > Co-citation
Edges	Bibliographic > Co-citation > Direct
Year	Direct = Bibliographic > Co-citation

<u>Agenda</u>

学術俯瞰システムの概要
 Overview of Academic Landscape System

 論文のデータの入手
 How to Get your Dataset of Papers

 学術俯瞰システムの操作方法
 How to Operate the System

最新データのアップロード

論文の場合(学内LAN)

Web of Scienceから最新データを入手する

- 東京大学 GACoS, or
- ・ "Web of Science"で検索
- 特許の場合
 - Thomson Innovationのデータに対応しています
 - 入手困難な場合は、イノベーション政策研究センター にご相談ください



(1) Search "Web of Science" within UT LAN, and



最新データのアップロード

(2) Search on Web of Science

Web of Science TM InCites TM Journal Citation Reports	Essential Science Indicators SM	EndNote ®	Sign In 🔻 Hel	lp English 🔫
WEB OF SCIENCE™		Select "Web of	Science Core	e Collection"
Search Web of Science [™] Core Colle	tion 🔽	My Tools	s 🔻 Search History	Marked List
)	Welcome to the new	Web of Science! View a	ı brief tutorial.
Basic Search 🔽 ‴stem cell*″ + Add Anoth	Topic Pr Field Reset Form Your	Search favorite conditio	Click her improve	re for tips to e your search.
All years From 1900 v to 2015 v				11

最新データのアップロード

(2) Search on Web of Science

Web of Science ™	InCites TM	Journal Citation Reports®	Essential Science Indicators SM	EndNote®		Sign In 🔻 Help	English 🔻
WEB C	<u>Q&A</u> Q: 名言 A: 「ce にしま	詞の複数形やヌ Ⅱ*」「identif*」の ∵す。	舌用した動詞も様 りように、複数の	資素したい パターン t	。 バある文字ヵ	いら、*(ア)	スタリスク)
Basic Searcl	Q:「er A:「" e	nergy policy」の energy policy "_	ように、2語の連 」とダブルクォーラ	続状態を テーション・	保って検索し で囲みます。	たい。	
TIME SPAN All years	Q:検 A:Wo (a) (b)	素結果が10万(Sでは10万件ま 10万件に納ま 期間で分割し 俯瞰システム	牛を超える。 でしかアクセスの るように、クエリー て、各々10万件。 では、約20万件の	の許可を与 一の選定を までに結果 のデータセ	きえていませ とやり直す。 是を減らして ?ットを分析し	ん。 ダウンロー た実績が	-ドする。 あります。
From 190	00 🗸 to	2015 🗸					12



(3) Select "Save to Other File Formats"

Web of Science ™ InCites [®] Journal Cit	ation Reports® Essential Science Indicators SM EndNote®	Sign In 🔻 Help 🛛 English 🔫			
WEB OF SCIENCE [™]					
Back to Search		My Tools 👻 Search History Marked List			
Results: 52,997 (from Web of Science Core Collection)	Sort by: Publication Date newest to oldest 🗸				
You searched for: TOPIC: ("stem cell*")More	Select Page	Add to Marked List Analyze Results Citation Report feature not available. [?]			
Refine Results	1. Systematic analysis by tangential flow f Save to Other File Formats By: Maurer, Elizabeth I Save to RefWorks NANOTOXICOLOGY Volume: 8 Issue: 7 Pages: 718-7 NOV 2014	dissolution Times Cited: 0 ations (from Web of Science J.; et al. Core Collection) 727 Published:			
Search within results for	Full Text View Abstract Image: Straight of the	ealthy Times Cited: 0 (from Web of Science			

最新データのアップロード

Update and Upload the Latest Data



« collider_accelerator > collider_accelerator_2015					
ファイルをすべて展開					
名前	種類	サイズ	圧縮率		
 2013_000500.tsv	TSV ファイル	4,760 KB	66%		
2013_000851.tsv	TSV ファイル	3,710 KB	65%		
🗎 2014_000500.tsv	ブラマ ファイル	5,919 KB	64%		
2014_00083	ファイル	2,470 KB	65%		
2015_00	PTIL	86 KB	65%		
		、/ ┌ \	~		

最新データのアップロード

\Rightarrow (5) Compress into a zip file, and upload it.



俯瞰システムでの分析に失敗するとき (データセットに起因するFAQ)



- A."Web of Science Core Collection"を選 んでいますか?
 - 「横断検索」(All Databases)は、引用
 関係データが含まれません

Q.「データセットに記法がおかしいファイ ルが含まれています」とエラーになる



A.壊れていないファイルからコピー&
 ペーストして修復してください
 – WoS内部で既に壊れていますので、
 再ダウンロードに効果はありません¹⁶

Web of Scienceの論文数 (2015-01-02)

2,500,000	
2,000,000	
1,500,000	
1,000,000	
500,000	
0	
	1900 1905 1915 1915 1915 1925 1925 1925 1925 192

注意事項

Precautions

- アクセス権限 Access Authority
 - 論文本文は権限のない人に渡さないようにしてください。

Don't give papers to unauthorized people.

- アクセス権限のないデータを利用する場合 if you are not authorized:
 - アクセス権限のある方と共同研究を行ってください。

your research must be a collaborative research with people who is authorized.

	坂田·森研所属者 IPRC所属者	東京大学所属者	その他
データベースから抽出	0	×	×
Web of Scienceから取得してアップロード	0	0	×
Thomson Innovation (特許データベース)	0	×	×

 * 論文等で公表される場合、または、企業で使用される場合は、 権限者(坂田先生・森先生)に共同研究を相談してください。
 If you publish your paper using the system, your paper must be a collaborative research. 授業や研究の文献調査など、権限内の非公開での利用は問合せ不要です。
 You have already been permitted in private use.

アカウントの作成

How to make your account

- アカウントを作成します。Staffs make your account.
 - 必要な情報 Please e-mail your data to us.
 - アカウント名 account
 - 半角アルファベット、半角数字、@/./+/-/_ で30文字以下。
 - You can use 30 or less English characters, which are alphabets, number, @/./+/-/_.
 - アカウント名に関しては、他の利用者から見える場合があります。
 Other users can see your account. Don't use your secret data.
 - ・ パスワード Password
 - 連絡先メールアドレス E-mail address
- URL of the System
 - <u>http://academic-landscape.com/</u>