第19回年次大会予稿

引用論文の分散値を重み付けとして考慮したページランクアルゴリ ズムによる主要論文の抽出

Extraction of key literature based on PageRank algorithm considering variance values of cited literatures as weighting

大槻明1*, 川上あゆみ1, 林剛2, 川村雅義2

Akira OTSUKI^{1*}, Ayumi KAWAKAMI¹, Takeshi HAYASHI², Masayoshi KAWAMURA²

1 お茶の水女子大学

Ochanomizu University

〒112-8610 東京都文京区大塚2丁目1番1号

E-mail: otsuki.akira@ocha.ac.jp

2 東京大学

The University of Tokyo

〒113-8656 東京都文京区弥生2-11-16東京大学工学部9号館320号室

*連絡先著者 Corresponding Author

学術俯瞰の分野における最近の研究動向は、参考文献の引用分析により実現するサイテーションマップが主流であり、ネットワーク構築やクラスタ化までの自動化はなされているが、各クラスタがどのような集団であるかの意味付けまでの自動化はなされておらず、専門家が手動で分析している現状である。ゆえに、各クラスタの自動解釈を最終的な目的として、本発表では各クラスタの主要論文の自動抽出を目指す。具体的には、論文をノード、引用をエッジとする有向グラフと考え、各ノードに発表年数を持たせたうえで、あるノードに入るエッジの元ノードの発表年数の分散を調べることでそれぞれの重要度の計算を試みる。そして、それらの重要度を基に、時間軸を持つ可視化グラフの構築を目指す。

Even though Citation Map has been in the mainstream of recent study trend in a field of academic landscape achieving the stage of automated clustering, the reality is that experts manually analyze semantic attachment about what kind of group each cluster is.

Therefore, we try to achieve an automated extraction of key literature of each cluster in the report, by setting automated interpretation of each cluster as the final purpose of the study.

キーワード: 学術俯瞰, 引用分析, データベース, ネットワーク分析

Keyword1, Science highangle, Citation analysis, Database, Network analysis

1 はじめに

学術俯瞰の分野における最近の研究動向は、引用ネットワーク分析が主流であり、自動クラスタ化まで実現されている.しかし、クラスタリングにより同定された各領域の特定や主要論文の自動抽出までは実現されてはいない.ゆえに、本研究ではクラスタリングにより同定された各領域の主要論文を自動で特定する手法について研究する.具体的には、引用される側の論文の発表年数の分散を調べ、その分散値をページランクアルゴリズムに適応することにより各論文の重要度を算出する.また、この重要度を基に、時間軸を持つ可視化グラフの構築を目指す.

2 先行研究

学術俯瞰の分野において、Small[1]は、被引用数は上位 1%の論文からなる共引用ネットワークを分析し、科学分野で成長している領域を追跡する方法を提案した。また、松尾[2]は、図1のとおり、引用ネットワークの構築、最大連結成分の取得、クラスタリング、可視化を行うことで学術論文引用ネットワークを分析した。しかし、クラスタリングにより同定された各領域の特定や主要論文、主要研究者の抽出について自動化はなされておらず、この部分は専門家が手動で分析しているのが現状である。

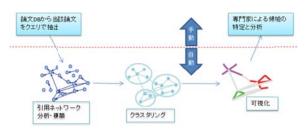


図1 ネットワーク分析を応用した学術俯瞰の手順

3 提案手法

前節の課題を解決するために, 本研究では 各領域の主要論文の自動抽出について考え る. 引用件数が同じでも, 「一時期に大量に 引用された」場合や、「長期間少しずつ引用 されている」場合などが考えられるため、従 来の引用分析だけでは、それぞれの重要度を 計算する事が難しい. ゆえに, 本研究では, 上述のそれぞれの「場合」に対し、論文をノ ード、引用をエッジとする有向グラフと考え, 各ノードに発表年数を持たせたうえで,ある ノードに入るエッジの元ノードの発表年数 の分散を調べることでそれぞれの重要度の 計算を試みる. そして、それらの重要度を基 に、時間軸を持つ可視化グラフの構築を目指 す. 以下に全体の流れを記し, 次節からその 詳細について述べる.

- 1. 論文 DB からキーワード (クエリ) 検索により論文数を絞る
- 2. 引用論文を中心にリスト化する.
- 3. 上記 2.のリストから論文発表年数の 分散分析を行うことによって、各引用 論文に重み付けを行う
- 4. 上記 3. の重み付けページランクアルゴリズムに適応して各引用論文の重要度を算出する.
- 重要度(ノード・エッジ)を基に可視化

3.2 論文DBからキーワード(クエリ)検索による論文数の絞り込み

本研究では、論文DBとしてSCOPUSを採用した. 「clustering」というクエリを用いて論文数を絞った結果、87、399件の論文数に絞り込まれた.

3.3 引用論文を中心にリスト化

前節のリストは、各論文がどの論文を引用 しているかという並び順になっているが、そ れを各引用論文がいつ、どのような論文に引 用されているのかといった並び順に再リス ト化する.

3.4 論文発表年数の分散分析を行うことによる各引用論文の重み付け

下記1)~3) により各引用論文の重み付けを行う.

1) ヒストグラムの最大値を抽出 最も引用数が多い年度を次の関数で抽 出し、MaxYearに格納する.

 $MaxYear=max\{y(x)/y(x):=y年に参照された$ 回数} (1)

2) 引用期間の特定

年度の古い年度から1年度毎に調べ、最初に見つかったMaxYearの10%以上の引用数の年度を引用がされ始めた開始年度としStartYearに格納する.そして10%以下の引用数の年度になった時点で、その年度をlast yearに格納する.そして論文が引用され始めてから引用されなくなった年度までの期間を次の式で求める.

Period := (LastYear + 1) - StartYear (2)

また、図2のように、ヒストグラムの山が複数存在する場合は、この作業を繰り返しそれぞれPeriodO、1…nに格納する.

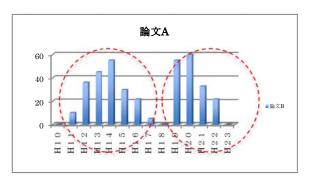


図2 ヒストグラムの形が正規分布から外れているケース

3) ヒストグラムの分散 (標準偏差) の算出 当該論文を引用する論文の発表年数の分 散 (標準偏差) を調べることで当該論文がど のくらいの期間にわたって引用されている かについて調べる. なお, 標準偏差の一般的 な求め方は次のように表され, 求められた標 準偏差値はVarianceに格納する.

$$Variance = \frac{\sum (x-x^{-})^{2}}{(n-1)}$$
 (3)

なお、図2のようにヒストグラムの形が正規分布から明らかに外れているようなケース(山がいくつもある様な場合)は、Period0、1…nの分散(標準偏差)を算出し、それらの平均値をVarianceに格納する。そして、Varianceを引用論文の重み付けの値として利用する。

3.5 各引用論文の重要度の算出

PageRankアルゴリズム[3]は、ハイパーリンク構造のような相互参照関係があるときに、どのページがもっとも「重要」であるかを定量的に算出する手法である.本研究では、このアルゴリズムを利用し、各引用論文の重要度の算出する.なお、この重要度の算出は次のように表される.

- (1)各論文は、固有の得点を持っている。 各引用もまた、固有の得点を持っている。
- (2) ある論文X に対して,
 - X の得点を P とする。
 - Xが他論文から引用されている得点を それぞれ Variance, …, Variance, と する。
 - Xが他論文を引用している得点をそれぞれの, …, 0 とする。

このとき、次が成り立つものとする。 $Variance_1+\cdots+Variance_n=P$

$$O_1 = \cdots = O_m = \frac{p}{m} \left(= \frac{\sum_{i=1}^n \text{Variance}_i}{m} \right)$$

すなわち、各論文に「流れ出す」引用の 得点の総和と、各論文から「流れ込む」引 用の得点の総和が等しくなるようにして、 その総和をその論文の得点と考え,この得 点が高いほど、その論文は重要であると考 える.そして,各論文から「流れ込む」引 用の得点計算にVarianceの値を適応することにより,各領域における主要論文の特定 を目指す.従来のアルゴリズムでは,「流れ 込む」引用が複数個あった場合,得点は均 等に割り振られていたが、本研究では Varianceの値が高いものにより多く「流れ 込む」と考え計算することで,被引用年次 を反映した重要度を計算する.

3.6 重要度を基に可視化

前節で導出した重要度を元に引用ネットワークとして可視化したものが図3である. 各ノードには論文名を表示しており、前節の重要度が高いほど、より大きなノードとして表現される. なお、本ツールのことをHiAc(Highangle of Academic) と呼ぶ.

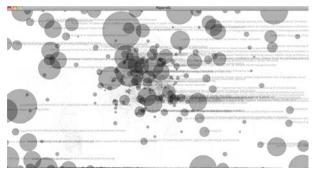


図3 重要度に基づき可視化した例

4 評価実験

4.1 専門家が手動で抽出した主要論文 との比較検証

立堀[4]らが2004年に発表した研究動向調査報告は、IBMの論文データベースを用いて、1999年以降のソフトウエア・アーキテクチャ研究分野の動向を調査し、主要な51論文を手動で抽出したものである。本評価実験では、この専門家が手動で抽出した主要論文をどこまで自動で抽出できるかについて検証する.

立堀らにおける51論文の抽出方法は,GoogleScholarを用いて年あたり引用数を求め,その上位40論文を抽出している.この40論文に加えて,立堀らが特に重要と考えた国際会議に絞ったうえで,ソフトウエア・アーキテクチャに関する11の論文を抽出している.また,立堀らは,定量的な研究動向評価のために,独自の分類方式を採用している,具体的には,ソフトウエア開発プロセスにおいて,ソフトウエア・アーキテクチャの果たす役割を,次の5つの役割のどれに着目しているかによって51論文を分類しており,それを図にしたものが図4である.

・[R]アーキテクチャへの要件(Requirement) ⇒様々な利害関係者のシステムへの要件 をアーキテクチャに反映する.

- - ⇒アーキテクチャの設計は、メタモデルに 基づいてアーキテクトが行う.
- ・[C] アーキテクチャの用いた利害関係者と のコミュニケーション(Communication) ⇒全ての要件を満たすことができない場 合,要件を調整するために,アーキテク チャを用いて利害関係者と交渉する.
- ・[S] アーキテクチャとシステム間の同期 (Synchronization)
 - ⇒抽象的なアーキテクチャは、実際に動く システムの実装に落とさなければなら ない. 逆に、システムに変更があった場 合、それはアーキテクチャにも反映され るべきである.

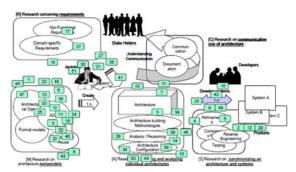


図4 51論文をソフトウエア・アーキテクチャの5つの役割に分類したイメージ

立堀らの報告に対して、HiAcで同様に主要論文を抽出したものが図5である。HiAcにおける論文抽出について述べる。まず、論文DBとして SCOPUSを用いた。そして、クエリとして「Software Architecture」を、また、年代として「1999年~2004年」をそれぞれ設定して論文を抽出した。なお、表1の「SCOPUS」欄のとおり、51論文のうち、〇のついた23論文以外は、そもそもSCOPUSには論文が存在しなかったため、対象外としている。また、表1の論文番号を図4及び図5上

でも表記している. さらに、図4における5つの役割は、定量的な研究動向評価のために立堀らが独自に設定した分類方式であるため、HiAcの場合(図5)では、手動でクラスタの微調整を行った.

図5から、引用論文の分散値を重み付けとして考慮したページランクアルゴリズム(以下「本アルゴリズム」という)によって可視化した結果、23論文はどれも主要な論文として抽出されていた.特に、論文番号41-51は、立堀らが主要な論文を手動で抽出したものであり、その中でSCOPUSに存在したものは、論文番号44-48であるが、それら全てを主要論文として抽出できていた.

このように、専門家が抽出した論文を本アルゴリズムによって分析することにより、主要な論文として自動で抽出することが可能となる。ただ、[R]のクラスタにおける緑のノードなど、立堀らが抽出していない論文もHiAcでは主要な論文として抽出していたが、これらの論文がどのような意味を持つものかについては、今後の課題として検討していきたい。

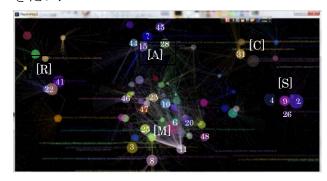


図5 51論文をHiAcで抽出したイメージ

表1 立堀らが抽出した51論文とSCOPUSで 抽出した19論文との対比表

N	C	Y	Author	Title	S
o	О	e			o
	n	a			P
	f	r			U
					S
1	W	2	João Pedro Sousa,	Aura: An Architectural	×
	I	0			^\

	С	0	D :10 1	E 1 C II	
	S A	2	David Garlan	Framework for User Mobility in Ubiquitous Computing Environments	
2	I	2	Aldrich, J.,	ArchJava: Connecting	0
	C	0	Chambers, C.,	software architecture to	
	S E	2	Notkin, D.	implementation	
3	I	2	Mehta, Nikunj R.,	Towards a taxonomy of	0
	C S	0	Medvidovic,	software connectors	0
	E	0	Nenad, Phadke,		
			Sandeep		
4	I	2	Batory, D.,	Scaling step-wise	0
	C S	0	Sarvela, J.N.,	refinement	
	E	3	Rauschmayer, A.		
5	I	1	Nenad	A Language and	×
	C S	9	Medvidovic,	Environment for	
	Е	9	David S.	Architecture-Based	
			Rosenblum,	Software	
			Richard N. Taylor		
6	I C	2	Dashofy, E.M.,	An infrastructure for the	0
	S	0	Van Der Hoek,	rapid development of	
	Е	2	A., Taylor, R.N.	XML-based architecture	
				description languages	
7	I C	1 9	Kazman, Rick,	Experience with	0
	S	9	Barbacci, Mario,	performing architecture	
	Е	9	Klein, Mark,	tradeoff analysis	
			Carriere,		
			S.Jeromy, Woods, Steven G.		
8	I	1	Bowman, Ivan T.,	Linux as a case study: Its	_
	C	9	Holt, Richard C.,	extracted software	0
	S E	9	Brewster, Neil V.	architecture	
9	I	1	Keller, Rudolf K.,	Pattern-based	_
	C	9	Schauer,	reverse-engineering of	0
	S E	9	Reinhard,	design components	
	Li.		Robitaille,		
1	U	2	Aler, R., Borrajo,	Reconciling the needs of	0
0	M L	0	D., Camacho, D.,	architectural description	0
	L	0	Sierra-Alonso, A.	with object-modeling	
			· ·	notations	
1	W	1	Kruchten, P.,	Describing software	0
1	I C	9	Selic, B.,	architecture with UML	
	S	9	Kozaczynski, W.		
1	A E	2	Jonathan Aldrich,	Architectural Reasoning in	
2	C	0	Craig Chambers	ArchJava	×
	0	0 2	Craig Chambers	Alciijava	
	O P	2			
1	W	2	Eric M.	A Highly-Extensible	×
3	I C	0	Dashofy,André	XML-Based Architecture	
	S	1	van der		
	Α		Hoek,Richard N.		
			Taylor		
1 4	I C	1	Jean-Marc	A Systematic Approach to	×
4	s	9	DeBaud,Klaus	Derive the Scope	
Ļ.	Е	9	Schmid	D 1 - 1' - 1'	
1 5	I C	1	Bosch, Jan	Product-line architectures	0
,	S	9		in industry: A case study	
1	E	9	A: 4 C	D-1	
6	C	0	Anita Sarma,	Palantír: Raising	×
	S	0	Zahra Noroozi,	Awareness among	
	Е	3	and André van der Hoek	Configuration Management Workspaces	
1	W	2	Shang-Wen	Using Architectural Style	
7	I	0	Cheng, David	as a Basis for System	×
	C S	0 2	Garlan, Bradley	Self-repair	
	A	2	Schmerl,	5011-10pan	
1	U	1	Fiadeiro, J.L.,	Interconnecting objects via	×
8	M	9	Andrade, L.F.	contracts	^
	L	9			
1	P	2	Jan Bosch, Gert	Variability Issues in	×
9	F E	0	Florijn, Danny	Software Product Line	
		1	Greefhorst, Juha		
	_	_			

2	T	1	Kuusela,	11. CC 4 1 1C	-
0	I C	1 9	Dashofy, Eric M., Medvidovic,	Using off-the-shelf	
	S	9	Nenad, Taylor,	middleware to implement connectors in distributed	
	Е	9	Richard N.	software architectures	
2	S	2	Martin L. Griss	Implementing Product-Line	
1	P L	0	Transmi E. Grigo	Features By Composing	
	C	0		Component Aspects	
2	F	2	Uchitel, S.,	Detecting implied scenarios	
2	S E	0	Kramer, J.,	in message sequence chart	
		1	Magee, J.	specifications	
2 3	I C	2	Fielding, Roy T.,	Principled design of the	
	S	0	Taylor, Richard	modern web architecture	
2	E W	0	N. Nenad	Aggaging the Cuitability of	H
4	I	9	Medvidovic and	Assessing the Suitability of a Standard Design Method	
	C S	9	David S.	for Modeling Software	
	A		Rosenblum	Architectures	
2	I	1	Di Nitto,	Exploiting ADLs to specify	
5	C S	9	Elisabetta,	architectural styles induced	
	E	9	Rosenblum, David	by middleware	
				infrastructures	L
6	0	2	Riehle, D.,	The architecture of a UML	
-	P	0	Fraleigh, S.,	virtual machine	
	S L	1	Bucka-Lassen, D., Omorogbe, N.		
	Α		_		L
2 7	W I	1 9	Mark H. Klein,	Attribute-Based	
′	C	9	Rick Kazman,	Architecture Styles	ĺ
	S	9	Len Bass, Jeromy Carriere, Mario		1
	Λ		Barbacci		
2	W	1	Jeff Magee,Jeff	Analyzing the behaviour of	
8	I	9	Kramer, Dimitra	distributed software	
	C S	9	Giannakopoulou	architectures: A case study	
_	A E	2		-	L
2	C	2	Aldrich, J.,	Language Support for Connector Abstractions	
	О	0	Sazawal, V., Chambers, C.,	Connector Abstractions	
	O P	3	Notkin, D.		
3	G	2	Sandeep	Generators for Synthesis of	T
0	P C	0	Neema,Ted	QoS Adaptation in	
	E	2	Bapty,Jeff	Distributed Real-Time	
			Gray, Aniruddha	Embedded Systems	
2	т	2	S. Gokhale	0 (6) 4	L
1	I C	2	Kazman, R.,	Quantifying the costs and benefits of architectural	
	S	0	Asundi, J., Klein, M.	decisions	
3	E W	2	Bridget	A Compositional Approach	
2	I	0	Spitznagel, and	for Constructing	
	C S	0	David Garlan	Connectors	
	Α				
3	E C	2	Marcus Fontoura,	UML-F: A Modeling	
-	О	0	Wolfgang Pree,	Language for	
	O P	0	Bernhard Rumpe	Object-Oriented Frameworks	ĺ
3	F	1	Pascal Fradet,	Consistency Checking for	H
4	S	9	Daniel Le	Multiple View Software	
	Е	9	Métayer and	Architectures	
		L	Michaël Périn		
3 5	R	2	Jon G. Hall	Relating Software	
5	Е	0	Michael Jackson	Requirements and	
		2	Robin C. Laney	Architectures using	
2		_	Bashar Nuseibeh	Problem Frames	L
6	S P	2	Jan Bosch	Maturity and Evolution in	
_	L	0		Software Product Lines:	
	С	2		Approaches, Artefacts and Organization	1
3	R	1	John Grundy	Aspect-oriented	H
7	E	9	Joini Grandy	Requirements Engineering	
		9		for Component-based	
				Software	1

				Cyatama	
3	F	2	Nimo Vo1-	Systems Deadlesk detection in	
8	S	0	Nima Kaveh,	Deadlock detection in	×
"	Ē	0	Wolfgang	distribution object systems	
	-	1	Emmerich		
3	F S	1 9	Michel	Algebraic software	×
	E	9	Wermelinger, José	architecture reconfiguration	
		9	Luiz Fiadeiro		
4 0	I C	2	Spitznagel, B.,	A compositional	0
U	S	0	Garlan, D.	formalization of connector	
	Е	3		wrappers	
4	F S	2	Uchitel, S.,	System architecture: The	×
1	E	0	Chatley, R.,	context for scenario-based	
	_	4	Kramer, J.,	model synthesis	
			Magee, J.		
4	F	2	Zhang, X., Young,	Refining code-design	×
2	S E	0	M., Lasseter,	mapping with flow analysis	
	£	4	J.H.E.F.		
4	I	2	Hasselbring, W.,	The Dublo architecture	×
3	C S	0	Reussner, R.,	pattern for smooth	
	E	4	Jaekel, H	migration of business	
			ĺ	information systems: An	
				experience report	
4	I	2	Matinlassi, M.	Comparison of software	0
4	C S	0		product line architecture	
	E	0		design methods: COPA,	
	L	7		FAST, FORM, KobrA and	
				OADA	
4	I	2	Caporuscio, M.,	Compositional verification	0
5	C	0	Inverardi, P.,	of middleware-based	0
	S E	0	Pelliccione, P.	software architecture	
	Е	*	i cinecione, i .	descriptions	
4	I	2	Grechanik, M.,	Design of large-scale	0
6	C	0	Batory, D., Perry,	polylingual systems	O
	S E	0 4	D.E.	poryiniguai systems	
4	I	2	François, A.R.J.	A hybrid architectural style	
7	C	0	1 1011y015, A.IV.J.	for distributed parallel	0
	S E	0		processing of generic data	
	E	4		streams	
4	I	2	Khare, R., Taylor,	Extending the	_
8	C	0	R.N	REpresentational State	0
	S	0	13.13	Transfer (REST)	
	Е	4			
				architectural style for	
4	I	2	Hong Von Dori'd	decentralized systems	.
9	C	0	Hong Yan, David	DiscoTect: A System for	×
	S	0	Garlan, Bradley	Discovering Architectures	
5	E	2	Schmerl,	from Running Systems	
0	C	0	Bas van der	Polyphony in Architecture	×
	S	0	Raadt, Jasper		
	Е	4	Soetendal, Michiel		
			Perdeck,Hans van		
Ļ		_	Vliet		
5	I C	2	Ian Gorton,	Architecting in the Face of	×
1	C S	0	Jereme Haack	Uncertainty: An Experience	
Ш	Е	4		Report	

5 むすび

本論文では、学術論文引用ネットワーク分析において、クラスタリングにより同定された各領域における主要論文の自動抽出について試みた、具体的には、引用論文の発表年

数の分散について分析し、その結果をページ ランクアルゴリズムに応用することにより 各論文の重要度を算出した. そして, 立堀ら が2004年に発表したソフトウエア・アーキテ クチャの研究分野の動向調査報告と比較検 証することにより,専門家が手動で抽出した 主要論文をどこまで自動で抽出できるかに ついて検証した. その結果, 対象論文はすべ て本アルゴリズムにおいても主要な論文と して抽出できていた. つまり, 専門家が抽出 した論文を本アルゴリズムによって分析す ることにより, 主要な論文として自動で抽出 することができた. 今後は、本アルゴリズム で自動抽出された論文のさらなる分析を行 い,専門家が抽出した論文との比較検証をし ていきたいと考える.

参考文献

- [1] Small, H. (2006). Tracking and predicting growth areas in science. Scientometrics, 68, 595-610
- [2] 松尾豊. (2008), 学術俯瞰とウェブからの情報抽出, 「イノベーション政策及び政策分析手法に関する国際共同研究」成果報告書No. 4, pp43-59
- [3] Lawrence Page, Sergey Brin, Rajeev Motwani, Terry Winograd, 'The PageRank Citation Ranking: Bringing Order to the Web', 1998
- [4] 立堀道昭,丸山宏,小林真, Daniel Yellin,吉田尚志,川井奈央,:ソフトウエア・アーキテクチャ研究動向の調査報告概要,情報処理学会研究報告,2005