

Odkrivanje skupin vozlišč v velikih realnih omrežjih na osnovi izmenjave oznak

Lovro Šubelj

Doktorska disertacija

Mentor: izr. prof. dr. Marko Bajec
Fakulteta za računalništvo in informatiko
Univerza v Ljubljani

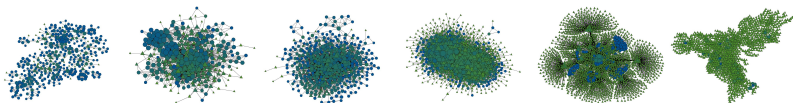
27. junij, 2013

- 1 Analiza omrežij
- 2 Odkrivanje skupin vozlišč
 - Osnovna izmenjava oznak
 - Uravnotežena izmenjava oznak
 - Napredna izmenjava oznak
 - Posplošena izmenjava oznak
- 3 Skupine v programskih omrežjih
- 4 Nadaljnje delo
- 5 Zaključek

Analiza omrežij

Različne vrste realnih omrežij (grafov).

socialna, informacijska, tehnološka, biološka itd.



Analiza omrežij: Newman (2008)

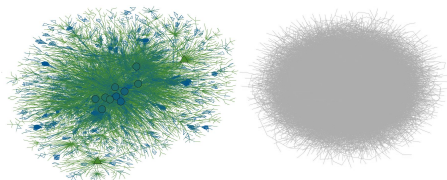
- preučevanje zgradbe omrežij
- razvoj potrebnih pristopov
- praktični primeri uporabe

Področje izjemno aktivno v številnih znanostih.

matematika, fizika, računalništvo, družboslovje, biologija itd.

Zgradba realnih omrežij

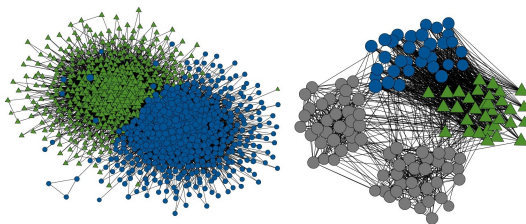
- globalne lastnosti omrežij (npr. mali svet)
- karakteristične skupine vozlišč (npr. skupnosti)
- pogosti vzorci vozlišč (npr. motivi, grafleti)
- lastnosti posameznih vozlišč (npr. zvezdišča)
- procesi nad omrežji (npr. širjenje)



Omejimo se na neusmerjena omrežja.

Karakteristične skupine vozlišč

- skupnosti (*community*) Girvan and Newman (2002)
(povezane) skupine tesno povezanih vozlišč
- moduli (*module*) Newman and Leicht (2007)
(nepovezane) skupine podobno povezanih vozlišč



Omejimo se na neprekrivajoče skupine vozlišč.

Vsebina

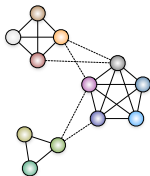
- 1 Analiza omrežij
- 2 Odkrivanje skupin vozlišč**
 - Osnovna izmenjava oznak
 - Uravnotežena izmenjava oznak
 - Napredna izmenjava oznak
 - Posplošena izmenjava oznak
- 3 Skupine v programskih omrežjih
- 4 Nadaljnje delo
- 5 Zaključek

Osnovna izmenjava oznak

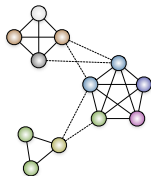
Osnovna izmenjava oznak (*label propagation*): Raghavan et al. (2007)

$$g_i = \operatorname{argmax}_g \sum_{v_j \in N_i} w_{ij} \delta(g_j, g)$$

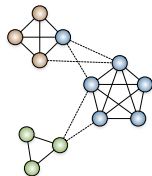
g_i posodablamo zaporedno v naključnem vrstnem redu.



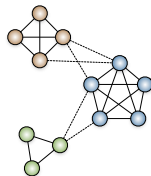
Začetno stanje



1. korak



2. korak



Končno stanje

g_i je oznaka skupine vozlišča v_i in w_{ij} utež na povezavi med v_i in v_j .

Analiza izmenjave oznak

Prednosti:

- brez predhodnega znanja (npr. število skupin)
- skoraj linearna časovna zahtevnost
- enostavna implementacija

Slabosti:

robustnost več razvrstitev že v manjših omrežjih Tibély and Kertész (2008)

natančnost slabša natančnost v omrežjih z nejasno zgradbo Leung et al. (2009)

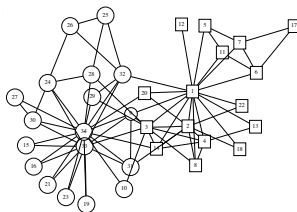
splošnost pristop omejen na odkrivanje skupnosti Šubelj and Bajec (2012c)

Vsebina

- 1 Analiza omrežij
- 2 **Odkrivanje skupin vozlišč**
 - Osnovna izmenjava oznak
 - **Uravnotežena izmenjava oznak**
 - Napredna izmenjava oznak
 - Posplošena izmenjava oznak
- 3 Skupine v programskih omrežjih
- 4 Nadaljnje delo
- 5 Zaključek

Robustnost izmenjave oznak

Osnovna izmenjava vrne > 500 razvrstitev v skupine. Tibély and Kertész (2008)

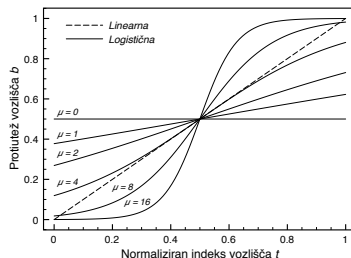


Vrstni red obravnave g_i se odraža kot preference vozlišč f_i . Šubelj and Bajec (2011a)

$$g_i = \operatorname{argmax}_g \sum_{v_j \in N_i} f_j \cdot w_{ij} \delta(g_j, g)$$

f_i je “moč” širjenja oznake vozlišča v_i . Leung et al. (2009)

Uravnotežena izmenjava oznak



Uravnotežena izmenjava (*balanced propagation*): Šubelj and Bajec (2011a)

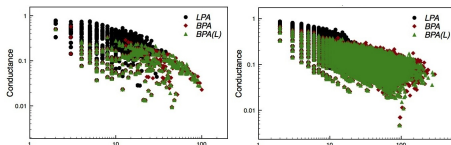
$$g_i = \operatorname{argmax}_g \sum_{v_j \in N_i} b_j \cdot w_{ij} \delta(g_j, g)$$

b_i je protiutež in $t_i \in (0, 1]$ normaliziran indeks vozlišča v_i .

Eksperimentalni rezultati

različnih razvrstitev v 1000 ponovitvah: Šubelj and Bajec (2011a)

	<i>karate</i>	<i>dolphins</i>	<i>books</i>	<i>football</i>	<i>jazz</i>	<i>elegans</i>
Osnovna izmenjava	184	525	269	414	63	707
Uravnotežena izmenjava	19	36	29	154	20	75



Robustnost se izboljša na račun časovne zahtevnosti. Šubelj and Bajec (2011c)

Znanstveni doprinos

Eur. Phys. J. B 40, 103–102 (2011)
DOI: 10.1140/epjb/e2011-10079-2

Regular Article

THE EUROPEAN
PHYSICAL JOURNAL B

Robust network community detection using balanced propagation

L. Šubelj^a and M. Bajec

University of Ljubljana, Faculty of Computer and Information Sciences, Ljubljana, Slovenia

Received 21 December 2010 / Received in final form 20 February 2011
Published online 4 May 2011 – © EDP Sciences, Società Italiana di Fisica, Springer-Verlag 2011

Abstract. Label propagation has proven to be an extremely fast method for deriving communities in large complex networks. Furthermore, due to its simplicity, it is also currently one of the most commonly adopted algorithms in the literature. Despite various subsequent advances, an important issue of the algorithm has not yet been properly addressed. Random (local) update rules within the algorithm strongly hamper its robustness, and consequently also the stability of the identified community structures. We note that an update rule can become an interesting propagation mechanism from various angles, and propose a balanced propagation that counteracts for the introduced weaknesses by utilizing such balance. We have evaluated the proposed approach on synthetic networks with identical partitions, and on several real-world networks with community structures. The results confirm that balanced propagation is significantly more robust than label propagation, when the performance of community detection is more important. Thus, balanced propagation retains high stability and algorithmic simplicity of label propagation, but improves on its stability and performance.

1 Introduction

Complex real-world networks can organize local structural modules (i.e., communities [1]) that are groups of nodes densely connected within and only loosely connected with the rest of the network. Communities may play important roles in different real-world systems – they can be related to functional modules in functional networks [2], or individuals with common interests in social networks [3]. Moreover, community structure also has a strong impact on dynamic processes taking place on such networks [4] and can thus provide an important insight into not only structural organization but also functional behavior of various real-world systems.

As a consequence, analysis of network community structure has been the focus of recent endeavor in different fields of science. There has also been a substantial number of community detection algorithms proposed in the literature over the last years [2, 4–12] (for a comprehensive survey see [13]). Nevertheless, due to its simplicity, only a small minority of these algorithms has been applied to large real-world networks with several millions of nodes, edges respectively.

A notable study towards this end was made by Hagmann et al. [5] who employed a simple label propagation to several significant communities in large real-world networks. Communities are identified by propagating (community) labels among nodes, thus, each node is assigned the label shared by most of its neighbors. Due to

very fast structural behavior of label propagation, densely connected sets of nodes form a community in some particular label after only a few iterations [7, 12]. The algorithm thus exhibits some linear complexity, which makes it applicable on networks with millions of nodes in a matter of minutes [1, 5]. The basic algorithm was further analyzed and refined by various authors [23, 15–20], when, due to its simplicity, label propagation is also currently one of the most commonly adopted algorithms in the literature.

Despite the above effects, an important issue of label propagation has not yet been properly addressed. The continuous convergence problem in some types of networks, Hagmann et al. [5] have proposed propagating labels among nodes (i.e., updating nodes' labels) in a random order. Although this updating strategy slows down convergence, the stability of the identified communities severely hampers the robustness of the algorithm, and consequently also the stability of the identified community structure. It has been noted that the algorithm reveals a large number of distinct community structures even in similar networks [7, 13, 16, 18], when these structures are also relatively different among themselves [14, 10]. Still, the stability of the algorithm can also be related to the significance of community structure in a network [10].

We argue that updating the nodes in some particular order can be more or less a better propagation mechanism [14] to the nodes that are updated at the beginning, and hence propagation preference to the nodes that are updated towards the end (and updating the nodes in a random order). The order of such updates thus governs the dynamics of the algorithm in a similar manner as

Generalized network community detection

Lovro Šubelj and Marko Bajec

University of Ljubljana, Faculty of Computer and Information Sciences,
EUROPEAN UNIVERSITY, Slovenia,
(lovro.subelj@fri.uni-lj.si)

Abstract. Community structure is largely regarded as an intrinsic property of complex real-world networks. However, several studies reveal that networks comprise even more sophisticated modules than classical ones, i.e., communities. These prototypical real-world networks can also be naturally partitioned according to common patterns of connections between the nodes. Recently, a propagation-based algorithm has been proposed for the detection of arbitrary network modules. We have advanced the label propagation algorithm by introducing a balanced update rule within the algorithm. The resulting algorithm is evaluated on various synthetic benchmark networks and real-world graphs. It is shown to be comparable to current state-of-the-art algorithms. However, in contrast to other approaches, it does not require some prior knowledge of the true community structure. To demonstrate its generality, we further employ the proposed algorithm for community detection in different synthetic and separate real-world networks, for generalized community detection and also predictive data clustering.

Keywords: link-driven community, link-pattern community, propagation, community detection, data clustering

1 Introduction

Over a decade of research in network analysis has revealed a number of common properties of complex real-world networks [23, 12]. Community structure [24] – the occurrence of coherent modules of nodes – is in particular interest as it provides an insight into not only structural organization but also functional behavior of various real-world systems [24, 5]. The analysis of communities has been the focus of many recent endeavors [23, 26], while community structure analysis is also considered as one of the most prominent areas of network science [2, 26].

However, most of the past work was constrained by community detection, and by higher density of links – real-world communities [24] (Fig. 1(a)). In contrast to the latter, recent studies reveal that networks comprise even more sophisticated modules than classical ones – communities [24] (Fig. 1(b)). In particular, real-world networks can also be naturally partitioned according to common patterns of connections among nodes – link-pattern communities [28, 14] (Fig. 1(c)). Link-pattern communities can be best related to observed functional modules in complex systems [27, 25], whereas they also capture a

^a e-mail: lovro.subelj@fri.uni-lj.si

Eur. Phys. J. B, 2011 (12 cit.)

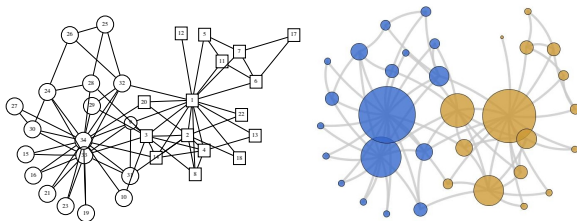
NEMO (ECML PKDD '11)

Vsebina

- 1 Analiza omrežij
- 2 **Odkrivanje skupin vozlišč**
 - Osnovna izmenjava oznak
 - Uravnotežena izmenjava oznak
 - **Napredna izmenjava oznak**
 - Posplošena izmenjava oznak
- 3 Skupine v programskih omrežjih
- 4 Nadaljnje delo
- 5 Zaključek

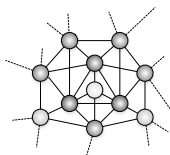
Natančnost izmenjave oznak

Natančnost izboljšamo z uporabo preferenc vozlišč f_i . Leung et al. (2009)

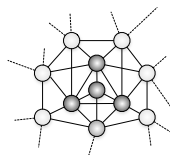


Osnovne lastnosti vozlišč niso primerne za f_i . Šubelj and Bajec (2011b)

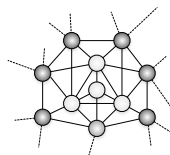
Zadržana in napadalna izmenjava oznak



Stopnje vozlišč



Zadržana



Napadalna

Zadržana izmenjava (*defensive propagation*): Šubelj and Bajec (2011b)

$$g_i = \operatorname{argmax}_g \sum_{v_j \in N_i} p_j \cdot w_{ij} \delta(g_j, g)$$

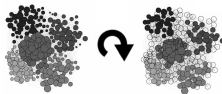
Napadalna izmenjava (*offensive propagation*):

$$g_i = \operatorname{argmax}_g \sum_{v_j \in N_i} (1 - p_j) \cdot w_{ij} \delta(g_j, g)$$

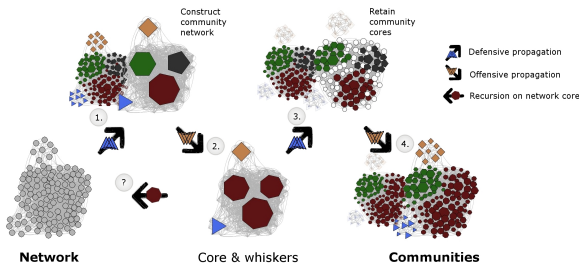
p_i je verjetnost pri naključnem sprehodu znotraj skupine g_i .

Napredna izmenjava oznak

Zadržana (napadalna) izmenjava doseže visok priklic (natančnost).



Napredna izmenjava (*diffusion propagation*): Šubelj and Bajec (2011b)



Eksperimentalni rezultati

Natančnost je primerljiva z najboljšimi pristopi. Šubelj and Bajec (2010)

Network	Description	Nodes	Edges	GMO	LPA	LPAD	LPAQ	LPAM	BDPA	DPA	No. CE ^c	T ^c
<i>karate</i>	Zachary's karate club [33]	34	78	0.381	0.416	0.402	0.399	0.420	0.419	0.420	0.02	
<i>dolphins</i>	Lusseau's bottlenose dolphins [38]	62	159		0.529	0.526	0.516	0.529	0.528	0.529	0.59	
<i>books</i>	Co-purchased political books [39]	105	441		0.526	0.519	0.522	0.527	0.527	0.527	0.46	
<i>football</i>	American football league [6]	115	616	0.556	0.606	0.606	0.604	0.605	0.606	0.606	0.37	
<i>elegans</i>	Metabolic network <i>C. elegans</i> [37]	453	2025	0.412	0.421	0.413	0.409	0.452	0.424	0.427^b	0.17	
<i>jazz</i>	Jazz musicians [40]	198	2742	0.439	0.443	0.443	0.445	0.445	0.444	0.444	0.00	
<i>netsci</i>	Network scientists [9]	1589	2742		0.902	0.947		0.907	0.960	1.00		
<i>yeast</i>	Yeast protein interactions [41]	2114	4480		0.694	0.799			0.725	0.824	1.04	
<i>emails</i>	Emails within a university [42]	1133	5451	0.503	0.557	0.560	0.537	0.582	0.555	0.562	0.01	
<i>power</i>	Western US power grid [32]	4941	6594		0.612	0.804			0.668	0.908	1.14	
<i>blogs</i>	Weblogs on politics [43]	1490	16718		0.426	0.426			0.426	0.426	1.00	
<i>pgp</i>	PGP web of trust [44]	10680	24340	0.849	0.754	0.844	0.726	0.884		0.869	1.08	
<i>asi</i>	Autonomous syst. of Internet [25]	22963	48436		0.511	0.591			0.528	0.600^b	1.02	0 s
<i>codmat³</i>	Cond. Matt. archive 2003 ^a [45]	27519	116181	0.661	0.616	0.683	0.582	0.755	0.634	0.735	1.00	1.5 s
<i>codmat⁵</i>	Cond. Matt. archive 2005 ^a [45]	36458	171736		0.586	0.643			0.608	0.683	1.00	
<i>kdd³</i>	KDD-Cup 2003 dataset [46]	27770	352285		0.624	0.630			0.619	0.617	1.00	3 s
<i>nec</i>	nec web overlay map [47]	75885	357317		0.693	0.738			0.703	0.767	1.03	
<i>epinions</i>	Epinions web of trust [48]	75879	508837		0.382	0.362			0.399	0.402	1.00	4.5 s
<i>amazon³</i>	Amazon co-purchasing 2003 [49]	262111	1.2M		0.682	0.749			0.701	0.857	1.01	20 s
<i>ndedu</i>	Webpages in nd.edu domain [50]	325729	1.5M		0.840	0.890			0.863	0.903	1.14	
<i>google</i>	Web graph of Google [3]	875713	4.3M		0.805	0.923			0.822	0.968	1.01	2.5 m
<i>nber</i>	NBER patents citations [51]	3.8M	16.5M		0.573	0.624			0.583	0.759	1.20	
<i>live</i>	Live Journal friendships [3]	4.8M	69.0M		0.538	0.539			0.557	0.693	1.00	44 m

Časovna zahtevnost blizu linearne $\mathcal{O}(m^{1,19})$. Šubelj and Bajec (2011b)

Znanstveni doprinos

Prispevek prejel fakultetno nagrado za raziskovalno delo.



Phys. Rev. E, 2011 (10 cit.)

ACNE (ECML PKDD '10)

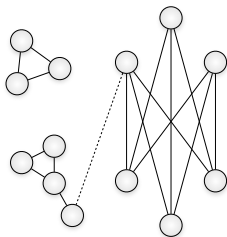
Vsebina

- 1 Analiza omrežij
- 2 Odkrivanje skupin vozlišč**
 - Osnovna izmenjava oznak
 - Uravnotežena izmenjava oznak
 - Napredna izmenjava oznak
 - Posplošena izmenjava oznak
- 3 Skupine v programskih omrežjih
- 4 Nadaljnje delo
- 5 Zaključek

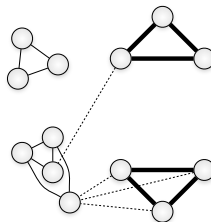
Splošnost izmenjave oznak

Pristop omejen na tesno povezana vozlišča (tj. skupnosti). Šubelj and Bajec (2012c)

Analogija med skupinami:



2 skupnosti, 2 modula



4 skupnosti

Oznake si izmenjujejo vozlišča na razdalji (največ) dva.

Posplošena izmenjava oznak

Posplošena izmenjava (*general propagation*): Šubelj and Bajec (2012c)

$$g_i = \operatorname{argmax}_g \left(\nu_g \sum_{v_j \in N_i} w_{ij} \delta(g_j, g) + (1 - \nu_g) \sum_{v_j \in N_i} w_{ij} / w_j \sum_{v_k \in N_j} w_{jk} \delta(g_k, g) \right)$$

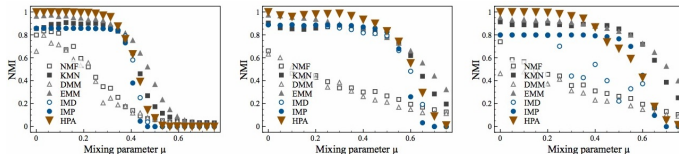
ν_g je blizu ena (nič) za skupnosti (module):

- lastnosti skupin Šubelj and Bajec (2012c)
- nakopičenost vozlišč Šubelj and Bajec (2011c)
- popravljena nakopičenost Šubelj and Bajec (2012a)

$\nu_g \in [0, 1]$ je parameter skupine vozlišč g .

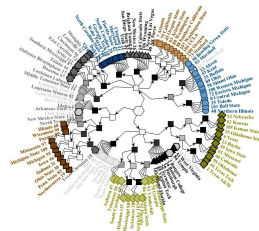
Eksperimentalni rezultati

Natančnost (vsaj) primerljiva z najboljšimi pristopi. Šubelj and Bajec (2012a)



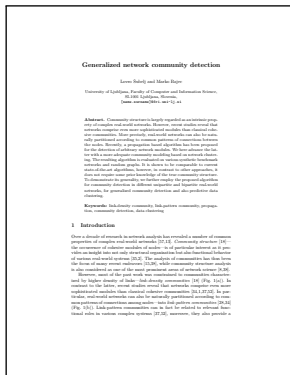
AUC pri napovedovanju povezav:

	$ N_i \cdot N_j $	$ N_i \cap N_j $	Infomap	Izmenjava
<i>football</i>	0.222	0.817	0.804	0.799
<i>politics</i>	0.646	0.862	0.763	0.762
<i>software</i>	0.779	0.826	0.724	0.766
<i>elegans</i>	0.812	0.920	0.631	0.641
<i>women</i>	0.564	0.290	0.578	0.699
<i>corporate</i>	0.456	0.481	0.649	0.748



Znanstveni doprinos

Prispevek izpostavljen s strani uredništva ter na straneh APS, ACM, IEEE ipd.



Eur. Phys. J. B, 2012 (3 cit.)

NEMO (ECML PKDD '11)

arXiv.org (2x)

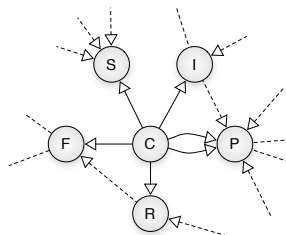
Vsebina

- 1 Analiza omrežij
- 2 Odkrivanje skupin vozlišč
 - Osnovna izmenjava oznak
 - Uravnotežena izmenjava oznak
 - Napredna izmenjava oznak
 - Posplošena izmenjava oznak
- 3 Skupine v programskih omrežjih**
- 4 Nadaljnje delo
- 5 Zaključek

Programska omrežja

Omrežja odvisnosti med razredi (*class dependency*): Šubelj and Bajec (2011d)

```
class C extends S implements I {  
    F field;  
    public C(P parameter) {  
        ...  
    }  
    public R function(P parameter) {  
        ...  
        return R;  
    }  
}
```

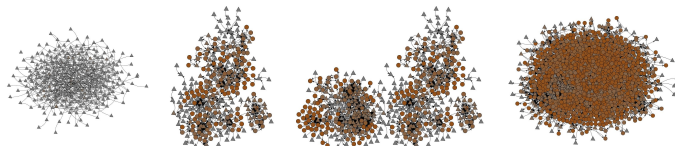


Podobne lastnosti kot druga realna omrežja. Valverde et al. (2002)

Skupine vozlišč še niso bile podrobno raziskane.

Analiza skupnosti vozlišč

Programska omrežja vsebujejo jasne skupnosti. Šubelj and Bajec (2011d)

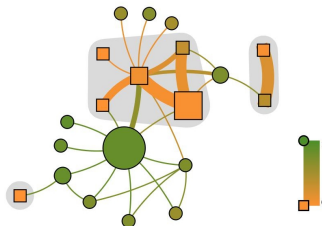


Skupnosti le delno sovpadajo s programskim paketi.

	Deljenje	Optimizacija	Izmenjava	Paketi*
<i>flamingo</i>	0.6466	0.6823	0.6485	0.2511
<i>colt</i>	-	0.6025	0.5599	-0.0332
<i>jung</i>	0.7210	0.7324	0.6874	0.3212
<i>org</i>	-	0.5599	0.5254	0.1830
<i>javax</i>	-	0.7667	0.7422	0.2907
<i>java</i>	-	0.4664	0.4132	0.2206

Analiza skupin vozlišč

Splošne skupine sovpadajo s programskim paketi. Šubelj and Bajec (2012c)



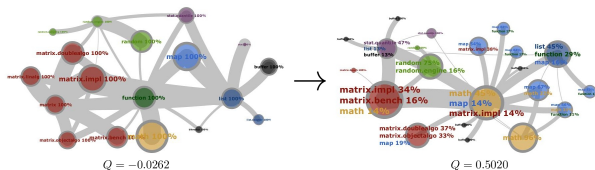
levo `jung.visualization.control.*Plugin` (100%)

sredina `jung.algorithms.layout.*Layout*` (54%), `...layout3d.*Layout` (23%);
`jung.graph.*(Graph|Multigraph|Tree)` (86%);
`jung.algorithms.filters.*Filter` (100%) itd.

desno `jung.io.graphml.parser.*Parser` (77%) in
`jung.io.graphml.*Metadata` (62%)

Uporaba skupin vozlišč

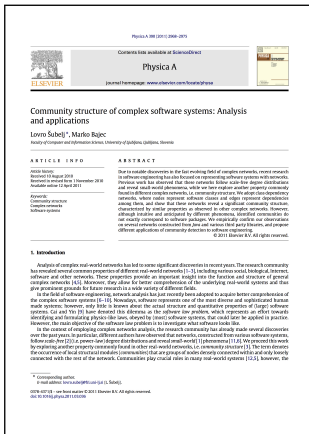
Reorganizacija paketov knjižnice (npr. modularno): Šubelj and Bajec (2011d)



Točnost napovedovanja paketov razredov: Šubelj and Bajec (2012b)

	I	I_{∞}	P	P_4	P_3	P_2	P_1
<i>flamingo</i>	2.65	4	0.566	←	0.572	0.793	1.000
<i>colt</i>	3.35	4	0.654	←	0.756	0.942	1.000
<i>jung</i>	2.97	4	0.617	←	0.663	0.857	1.000
<i>org</i>	3.50	7	0.616	0.616	0.714	0.989	1.000
<i>weka</i>	3.02	6	0.684	0.692	0.736	0.871	1.000
<i>javax</i>	3.11	5	0.626	0.631	0.816	0.982	1.000

Znanstveni doprinosi



Lovro Šubelj (FRI)



Odkrivanje skupin vozlišč v omrežjih



V pripravi

27. junij, 2013

23 / 27

Vsebina

- 1 Analiza omrežij
- 2 Odkrivanje skupin vozlišč
 - Osnovna izmenjava oznak
 - Uravnotežena izmenjava oznak
 - Napredna izmenjava oznak
 - Posplošena izmenjava oznak
- 3 Skupine v programskih omrežjih
- 4 Nadaljnje delo
- 5 Zaključek

Pregled izmenjave oznak

Contents	
1 Introduction	3
2 Label propagation	4
2.1 Synchronous propagation	4
2.2 Asynchronous propagation	4
2.3 Semi-supervised propagation	5
3 Label identification	6
3.1 General strategies	6
3.1.1 Standard propagation	6
3.1.2 Weighted propagation	6
3.1.3 Profound propagation	7
3.2 Performance-based strategies	8
3.2.1 Strength propagation	8
3.2.2 Degree propagation	8
3.2.3 Defusive propagation	9
3.2.4 Offensive propagation	9
3.2.5 Modularity propagation	10
3.2.6 Potts model propagation	10
3.3 Stability-based strategies	11
3.3.1 Momentum propagation	11
3.3.2 Controlled propagation	11
3.3.3 Attenuated propagation	12
3.3.4 Balanced propagation	12
3.4 Complexity-based strategies	14
3.4.1 Selective propagation	14
3.4.2 Pseudo propagation	14
4 Label ties	15
4.1 Random label	15
4.2 Label rotation	15
4.3 Label inclusion	15
4.4 Label priority	15
5 Propagation criteria	16
5.1 Label equilibrium	16
5.2 Label convergence	16
5.3 Label semi-convergence	16
5.4 Threshold convergence	16
6 Advanced propagation	17
6.1 Hierarchical propagation	17
6.2 Refining propagation	17
6.3 Hybrid propagation	17
6.4 Decisive propagation	17
6.5 Chain propagation	17
6.6 Parallel propagation	17
7 Other networks	18
8 Other groups	18
9 Applications	18

2

V pripravi



prof. dr. Xin Liu



dr. Steve Gregory

Analiza skupin vozlišč

Group extraction for real-world networks

Lovro Šubelj¹, Neil Blagus & Marko Bajec

University of Ljubljana, Faculty of Computer and Information Science, Slovenia

Background

Complex real-world networks contain characteristic groups of nodes with common linking pattern like densely linked communities [1]. These were the focus of most recent work and have diverse applications. However, many real-world networks also contain other groups of nodes that can be overlapping and other, whereas some parts of the networks reveal no significant groups.

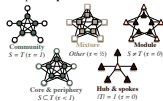
Group formalism

Let S be a group of nodes, T the linking pattern and τ the group parameter.

$$\rho(S, T) = \frac{|S \cap T|}{|S \cup T|}$$



Group examples



Group criterion

Let W be the group criterion, L the number of links and μ the (harmonic) mean size.

$$W(S, T) = \rho(S, T) (1 - \rho(S, T)) \left(\frac{L(S, T)}{|S|} - \frac{L(S, T^c)}{|S|} \right)$$

W is a local asymmetric criterion that favors the links between S and T , and penalizes for the links between S and T^c . (Note, however, that W disregards the links with both endpoints in S^c .) For $S = T$, W is consistent with a wide class of other models (e.g., stochastic blockmodel) [2].

Group extraction

A sequential extraction [2] of groups that can be overlapping, nested etc.

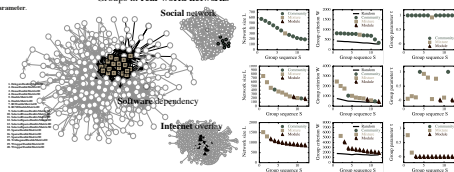
- (1) Find S and T that optimize criterion W (e.g., hub search).
- (2) Extract only the explained links between S and T (and isolated nodes).
- (3) Repeat until W is larger than expected in a random graph (by simulation).

Contributions

1. A simple formalism and criterion for general groups of nodes.
2. An adequate extraction procedure for statistically significant groups.
3. Characterization of the group structure of different real-world networks.

What are characteristic groups of nodes in real-world networks? Network (type) dependent.
What portion of network links is explained by the group structure? Between 60% and 90%.
What portion of network nodes is included in the group structure? More than 50%.

Groups in real-world networks



Network	Nodes	Links	Group	Community	Core	Mixture	Module	Background
			#	% Links	% Nodes	% Links	% Nodes	% Links
Author collaborator [3]	1580	2742	100	0.94	71% (47%)	0%	(0%)	15% (18%)
American football [1]	115	415	15	8.8	0.88 (83%)	9% (11%)	3% (7%)	25% (24%)
Leave-one search engine	1657	6988	123	12.1	0.55 (25%)	1% (2%)	20% (24%)	38% (34%)
Cdb computing [4]	227	963	15	10.3	0.41 (11%)	5% (6%)	69% (49%)	4% (6%)
Word adjacency [5]	112	425	4	11.2	0.28 (0%)	0%	34% (33%)	25% (10%)
Internet overlay [6]	707	1657	23	10.4	0.86 (1%)	12% (4%)	15% (7%)	34% (38%)
Southern action [6]	32	89	2	4.3	0.00 (0%)	0%	0%	80% (43%)

All extracted groups are statistically significant at 1% level.

References

- [1] Girvan, M. & Newman, M.E.J.: Community structure in social and biological networks. *P. Natl. Acad. Sci. USA* 99(12), 7821–7826 (2002).
- [2] Blaszko, T., Levine, R., & Shalizi, C.: Community structure in networks using the expectation-maximization. *Phys. Rev. E* 76(1), 016104 (2007).
- [3] Newman, M. E. J.: Finding community structure in networks using the expectation-maximization. *Phys. Rev. E* 76(1), 016104 (2007).
- [4] Eubank, J., & Ruge, M.: Community structure of complex software systems: Analysis and application. *Physica A* 369(1), 288–297 (2006).
- [5] Ledesma, J., Klemm, A., & Schuster, C.: Single-view time-dependent network analysis: detecting changes and predicting the future. *Proceedings of the ACM SIGMOD International Conference on Database Systems and Data Mining* (Chicago, IL, USA, 2005), pp. 177–187.
- [6] Davis, S., Gordon, B.R., & Graham, M.R.: Deep South (Chicago University Press, Chicago, 2003).



*Corresponding author: lovro.subelj@fri.uni-lj.si

NetSci '13

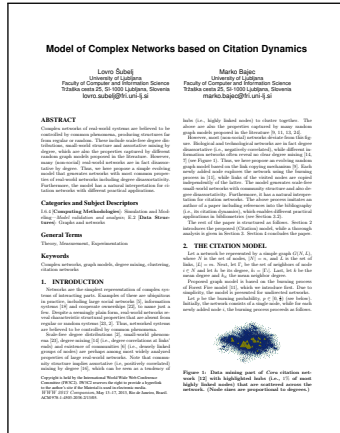
Analiza različnih omrežij

Zgradba programskih omrežij:

Modeli omrežij citiranj:



V pripravi



LSNA (WWW '13)

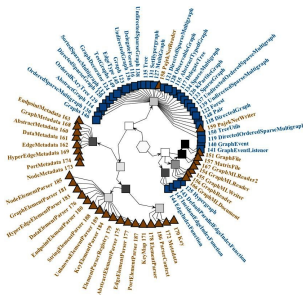
Vsebina

- 1 Analiza omrežij
- 2 Odkrivanje skupin vozlišč
 - Osnovna izmenjava oznak
 - Uravnotežena izmenjava oznak
 - Napredna izmenjava oznak
 - Posplošena izmenjava oznak
- 3 Skupine v programskih omrežjih
- 4 Nadaljnje delo
- 5 Zaključek**

Zaključek

Znanstveni doprinos:

- pristopi za odkrivanje skupin vozlišč v omrežjih
zahtevnost, splošnost, natančnost, robustnost, enostavnost, brez parametrov
- analiza in uporaba skupin v programskih omrežjih



Hvala za pozornost

Po zagovoru vabljeni v Laboratorij za podatkovne tehnologije.

- M. Girvan and M. E. J. Newman. Community structure in social and biological networks. *Proceedings of the National Academy of Sciences of United States of America*, 99(12): 7821–7826, 2002.
- I. X. Y. Leung, P. Hui, P. Liò, and J. Crowcroft. Towards real-time community detection in large networks. *Physical Review E*, 79(6):066107, 2009.
- M. E. J. Newman. The physics of networks. *Physics Today*, 61(11):33–38, 2008.
- M. E. J. Newman and E. A. Leicht. Mixture models and exploratory analysis in networks. *Proceedings of the National Academy of Sciences of United States of America*, 104(23): 9564, 2007.
- U. N. Raghavan, R. Albert, and S. Kumara. Near linear time algorithm to detect community structures in large-scale networks. *Physical Review E*, 76(3):036106, 2007.
- L. Šubelj and M. Bajec. Unfolding network communities by combining defensive and offensive label propagation. In *Proceedings of the ECML PKDD Workshop on the Analysis of Complex Networks*, pages 87–104, Barcelona, Spain, 2010.
- L. Šubelj and M. Bajec. Robust network community detection using balanced propagation. *European Physical Journal B*, 81(3):353–362, 2011a. doi: 10.1140/epjb/e2011-10979-2.
- L. Šubelj and M. Bajec. Unfolding communities in large complex networks: Combining defensive and offensive label propagation for core extraction. *Physical Review E*, 83(3):036103, 2011b. doi: 10.1103/PhysRevE.83.036103.
- L. Šubelj and M. Bajec. Generalized network community detection. In *Proceedings of the ECML PKDD Workshop on Finding Patterns of Human Behaviors in Network and Mobility Data*, pages 66–84, Athens, Greece, 2011c.
- L. Šubelj and M. Bajec. Community structure of complex software systems: Analysis and applications. *Physica A: Statistical Mechanics and its Applications*, 390(16):2968–2975, 2011d. doi: 10.1016/j.physa.2011.03.036.

- L. Šubelj and M. Bajec. Clustering assortativity, communities and functional modules in real-world networks. *e-print arXiv:12082518v1*, pages 1–21, 2012a.
- L. Šubelj and M. Bajec. Software systems through complex networks science: Review, analysis and applications. In *Proceedings of the KDD Workshop on Software Mining*, pages 9–16, Beijing, China, 2012b. doi: 10.1145/2384416.2384418.
- L. Šubelj and M. Bajec. Ubiquitousness of link-density and link-pattern communities in real-world networks. *European Physical Journal B*, 85(1):32, 2012c. doi: 10.1140/epjb/e2011-20448-7.
- G. Tibély and J. Kertész. On the equivalence of the label propagation method of community detection and a potts model approach. *Physica A: Statistical Mechanics and its Applications*, 387(19-20):4982–4984, 2008.
- S. Valverde, R. F. Cancho, and R. V. Solé. Scale-free networks from optimal design. *Europhysics Letters*, 60(4):512, 2002.