

mmeybodi@aut.ac.ir, motiee@ce.aut.ac.ir

HITS

HITS

HITS

Identification of Web Communities using Cellular Learning Automata

S. Motiee M. R. Meybodi

Computer Engineering and Information Technology Department
Amirkabir University of Technology
Tehran Iran

Abstract: A collection of web pages which are about a common topic and are created by individuals or any kind of associations that have a common interest on that specific topic is called a web community. Since at present, the size of the web is about 3 billion pages and it is still growing very fast, identification of web communities has become an increasingly hard task. In this paper, a cellular learning automata based approach for identification of web communities is proposed. The proposed approach is a combination of web structure, web usage and web content mining techniques. The proposed approach determines the related pages and their relevance degree using the cellular learning automata and the users' behavior in visiting the web pages. Then, a HITS based algorithm is applied on the structure obtained from the previous step to identify the web communities related to arbitrary subjects. The web community obtained using this method is not dependent on a specific web graph structure. To evaluate the proposed approach, it is implemented and the results are compared with the results of two existing methods: HITS and a complete bipartite graph based method. Experimental results have shown the superiority of the proposed method.

Keywords: Web Community, Cellular Learning Automata, HITS Algorithm, Web Usage Data

$K_{3,3}$

[2]

[8] [6]

[2]

[2]

[3]

(RPA)

RPA

[3]

Hub
Authority

HITS

Authority

Hub

[5]

HITS

Authority Hub [7] HITS
[8] [6] [5] [4]

HITS

[4]

Kumer

Trawling

[5]

$$\begin{array}{l} d \\ : \\ CLA=(Z^d,\varphi,A,N,F) \\ d \qquad \qquad Z^d \end{array}$$

$$\begin{array}{llll} & \varphi & E=\{\alpha,\beta,c\} & \\ (LA) & A & \beta=\{\beta_1,\beta_2,...,\beta_r\} & \alpha=\{\alpha_1,\alpha_2,...,\alpha_r\} \\ & & c_i & c=\{c_1,c_2,...,c_r\} \end{array}$$

$$Z^d \qquad \qquad N=\{x_l,...,x_m\}$$

$$\begin{array}{llll} & & c_i & \alpha_i \\ \beta & CLA & F\colon \varphi^m \rightarrow \beta & \end{array}$$

$$\begin{array}{l} \beta \\ (\end{array}$$

$$\begin{array}{l} \{\alpha,\beta,p,T\} \\ \alpha=\{\alpha_1,\alpha_2,...,\alpha_r\} \\ \beta=\{\beta_1,\beta_2,...,\beta_r\} \\ p=\{p_1,p_2,...,p_r\} \\ p(n+1)=T[\alpha(n),\beta(n),p(n)] \quad T \end{array}$$

$$\begin{array}{ll} \alpha_i & \mathbf{n} \\ \beta(n) & \end{array}$$

$$\begin{array}{l} p_i(n+1)=\ p_i(n)+a(1-\beta(n)). \\ (1-p_i(n))-b.\beta(n).p_i(n) \\ p_j(n+1)=\ p_j(n)+a(1-\beta(n)). \qquad \text{if } j \neq i \qquad (\end{array}$$

$$\mathbf{HITS}$$

$$p_j(n)+\frac{b.\beta(n)}{r-1}-b.\beta(n).p_j(n)$$

$$\mathbf{HITS}$$

$$\begin{array}{llll} b & a & b & a \\ a & b & L_{R-P} & \\ .[9] & L_{R-I} & b & L_{R\&P} \end{array}$$

$$\mathbf{HITS}$$

$$[1]$$

HITS

HITS

Hub Authority

Authority Authority

Authority Hub Hub

Hub

HITS

Hub

Authority

Hub

Authority

Hub

Authority

Authority Hub

N

$Aut[A]$ Authority

N A

$Hub[A]$ Hub

$Hub[A]$

Hub Aut

N A

$$Aut[A] = \sum_{(B,A) \in N} H[B] \quad ()$$

N A

$$Hub[A] = \sum_{(B,A) \in N} A[B] \quad ()$$

Aut Hub

Authority Hub

Hub

Authority

HITS

HITS

[10]

$$a = c_1 \frac{\sum_{\forall a \in N(agent_i) \text{ and } (a \rightarrow agent_i \text{ or } agent_i \rightarrow a)} 1}{\sum_{\forall a \in N(agent_i)} 1} + c_2 \sum_{\forall path_i | a \text{ and } N(agent_i) \in path_i} \frac{1}{Length(path_i)} \quad ()$$

c_1 و c_2 نیز ضرایبی هستند که میزان اهمیت هر یک از دو عامل موثر در محاسبه پاداش را تعیین می کنند. حاصل جمع این دو ضریب برابر با یک می باشد.

$N(agent)$

$$b = \frac{\sum_{\forall cycle_i | a \text{ and } N(agent_i) \in cycle_i} Length(cycle_i)}{\sum_{\forall path_i | a \text{ and } N(agent_i) \in path_i} Length(path_i)} \quad ()$$

$f(agent_i)$

$$p_a(agent_i)$$

R

p_a

() ()

$$L_{R-I}$$

$$d(agent_i, agent_j) =$$

$$\sqrt{(s_{i,1} - s_{j,1})^2 + \dots + (s_{i,k} - s_{j,k})^2}$$

k .

$$n$$

m

$$S_{m,n}$$

$$f(agent_i) = \max\{0, \frac{1}{(2s_x + 1) \times (2s_y + 1)}\}$$

()

$$r(i, j) = f(agent_i) \frac{1}{d(agent_i, agent_j)}$$

()

$$\sum_{agent_j \in N(agent_i)} (1 - \frac{d(agent_i, agent_j)}{k}) \quad (1)$$

$$f$$

$$S_y \quad S_x$$

$agent_i$

$N(agent_i)$ (

k

$$d(agent_i, agent_j)$$

$$p_a(agent_i) = \frac{\beta^2}{\beta^2 + f(agent_i)^2} \quad ()$$

β_{45}

	l
(\quad)	T_c
α_u	
λ	l
ΔM^v_t	σ_m
ΔM^v_t	μ_m
α_p	
σ_t	l
θ	
	l

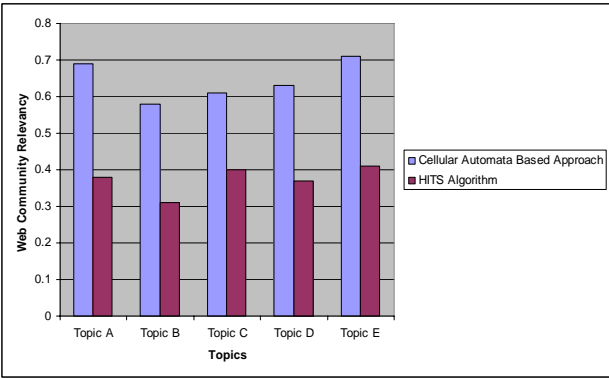
() ()
 - [5] HITS
 () ()
 ()

c_1	$,$
c_2	$,$
f	k
B	$,$

Hub : Authority Hub
 Authority
 :
 $Authority(i) = \sum_{\forall j \rightarrow i} r(j,i) \times Hub(j)$
 $Hub(i) = \sum_{\forall i \rightarrow j} r(i,j) \times Authority(j)$
 $j \quad i \quad r(i,j)$

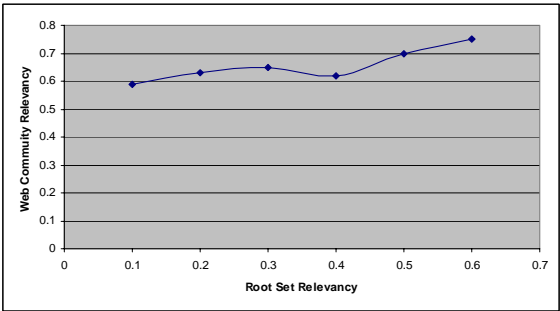
Authority Hub
 :
 Authority Hub
 Authority Hub
 :
 Hub
 Authority Hub
 Authority Hub
 Authority Hub
 Hub
 Hub
 Hub
 Lui [11]

: ()
 [11]
).
 (.



HITS

()



()

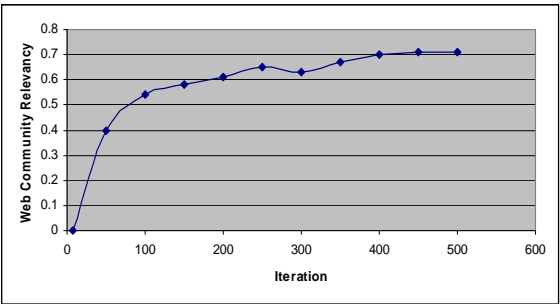
[5]

()

[5]

()

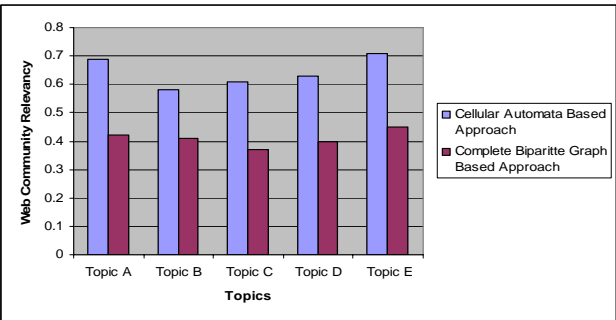
[5]



()

HITS

:



()

HITS

()

HITS

HITS

)

(

HITS

Authority Hub

- [12] Beigy, H., Meybodi, M. R., “*Open Synchronous Cellular Learning Automata*”, Proceedings of the 8th world Multi-conference on Systemics, Cybernetics and Informatics (SCI2004), pp. 9-15, Orlando, Florida, USA. July 18-21, 2004.

¹ Web Mining

² Web Structure Mining

³ Hyperlink

⁴ Hyperlink Analysis

⁵ Related Page Algorithm

⁶ Complete Bipartite Graph

⁷ Stationary

⁸ Non-Stationary

⁹ Linear Reward-Penalty

¹⁰ Linear Reward epsilon Penalty

¹¹ Linear Reward Inaction

¹² Root Set

¹³ Base Set

HITS

- [1] Beigy, H. and Meybodi, M. R., “*A Mathematical Framework for Cellular Learning Automata*”, Advances on Complex Systems, Vol. 7, Nos. 3-4, pp. 295-320, 2004.
- [2] Toyoda, M., Kitsuregawa, M., “*Creating a Web Community Chart for Navigating Related Communities*”, In Proc. Hypertext 2001, pp.103-112, 2001.
- [3] Gibson, D., Kleinberg, J. M., Raghavan, P., “*Inferring Web Communities from Link Topology*”, In Proc. of the 9th ACM Conference on Hypertext and Hypermedia. Pittsburgh, PA, pp. 225-234, 1998.
- [4] Kumar, R., Raghavan, P., Rajagopalan, S., Tomkins, A., “*Trawling the Web for Emerging Cyber-Communities*”, Proc. of the 8th WWW Conference, 1999.
- [5] Imafuji, N., Kitsuregawa, M., “*Effects of Maximum Flow Algorithm on Identifying Web Community*”, Proc. of the 4th international Workshop on Web information and Data Management (McLean, Virginia, USA, November 08 - 08, 2002). WIDM '02. ACM Press, New York, NY, pp. 43-48, 2002.
- [6] Flake, G., Lawrence, S., Giles, C.L., “*Efficient Identification of Web Communities*”, the 6th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. Boston, MA, pp. 150-160, 2000.
- [7] Kleinberg, J., “*Authoritative Sources in a Hyper-linked Environment*”, Proc. Of ACM-SIAM Symposium on Discrete Algorithms, 1998. Also appears as IBM Research Report RJ 10076(91892) May 1997.
- [8] Flake, G. W., Lawrence, S., Giles, C. L., Coetzee, F. M., “*Self-Organization and Identification of Web Communities*”, IEEE Computer, Vol. 35, No. 3, pp. 66-71, 2002.
- [9] Narendra, K. S. and Thathachar, M. A. L., *Learning Automata: An Introduction*, Prentice Hall, 1989.
- [10] Chen, X. Xu, and Chen, Y., “*A Novel Ant Clustering Algorithm Based on Cellular Automata*”, Proc. IEEE/WIC/ACM International Conference on Intelligent Agent Technology (IAT 04), 2004.
- [11] Liu, J., Zhang, S. and Yang, J., “*Characterizing Web Usage Regularities with Information Foraging Agents*,” IEEE Transactions on Knowledge and Data Engineering, Vol. 16, No. 4, pp. 566-584, 2004.