

Focused Crawling using Asynchronous Cellular Learning Automata

S. Motiee

Computer Engineering and Information Technology
Department
Amirkabir University of Technology
Tehran Iran
motiee@aut.ac.ir

M. R. Meybodi

Computer Engineering and Information Technology
Department
Amirkabir University of Technology
Tehran Iran
mmeybodi@aut.ac.ir

Abstract: Web crawling is used to collect the web pages which will be indexed by a search engine. The search engine uses these crawled and indexed pages to answer users' queries. Since the volume of web pages is very high and it increases continuously, search engines can index a limited number of web pages. Therefore, in recent years, the focused crawler algorithms have been introduced which act selectively during crawling and collect the web pages related to a specific topic. In this paper, an asynchronous cellular learning automata based approach for focused crawling is proposed. The proposed approach is a combination of web structure and web usage mining techniques and is composed of two phases. In the first phase the relationship structure of pages is determined using asynchronous cellular learning automata, hyperlinks and users' behavior in visiting web pages, i.e. the related pages and their relevance degree are determined. In the second phase, the focused crawling is performed using the obtained relationship structure and the pages related to a specific topic are collected. Experimental results have shown the superiority of the proposed method (harvest rate and target recall) in comparison to Best First Crawler and its independency from initial set selection.

Keywords: Focused Crawling, Asynchronous Cellular Automata, Web Usage Data

[1]

[2] BestFirst

()

[3] SharkSearch

[4] PageRank

[5]

[6]

[7]

[8]

[9]

Spider

$$\begin{array}{llll} & & (\text{LA}) & A \\ & & Z^d & N=\{x_l, \dots, x_m\} \\ \beta & & \text{CLA} & F: \varphi^m \rightarrow \beta \end{array}$$

[13][12]

:(**ACLA**)

$$\begin{array}{ccccccc} \text{ACLA} & & & & & & \text{ACLA} \\ n & & d & & \text{ACLA} & & \\ & & & n & \rho & \text{CLA} & \text{CLA}=(Z^d,\varphi,A,N,F,\rho) \\ \rho_i & & LA_i & : & \text{ACLA} & & i & \text{LA} & \rho_i \\ & & \text{LA} & & \text{ACLA} & & & & \end{array}$$

[14].

[2-9]

$$\begin{array}{ccccccc} & &) & & (& &) \\ & & & & & & (\\ & & & & & & (\quad) \end{array}$$

$$\begin{array}{ccc} (\dots & &) \end{array}$$

(
Hub
Hub

Hub HITS HITS

$G=[0,w(n)-$ (ACLA)
 $n \qquad 2|\sqrt{n}| \quad w(n) \quad h(n) \qquad 1] \times [0,h(n)-1]$

ACLA

ACLA

()

()

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

()

$$f \qquad \beta$$

$$f(agent_i) = \max\{0, \frac{1}{9} \sum_{agent_j \in N(agent_i)} (1 - \frac{d(agent_i, agent_j)}{k})\} \quad ()$$

$$agent_j \quad agent_i \quad d(agent_i, agent_j) \quad k \quad agent_i \quad N(agent_i) \quad ()$$

$$d(agent_i, agent_j) = \sqrt{(s_{i,1} - s_{j,1})^2 + \dots + (s_{i,k} - s_{j,k})^2} \quad (1)$$

()

$$\left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \quad \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right)$$

$$\left(\begin{array}{c} \vdots \\ cell_i \\ \vdots \end{array} \right) \quad cell_i \quad \vdots \quad cell_i \quad \vdots \quad \left(\begin{array}{c} \vdots \\ cell_i \\ \vdots \end{array} \right).$$

$$\begin{array}{ccc} & & \vdots \\ & cell_i & cell_i & \bullet \\ (& cell_i & cell_i & \bullet \end{array}$$

$$L_{ReP} \quad . \quad () \quad () \quad agent_i \quad b \quad a \quad .$$
$$.(\quad) \quad agent_i \quad g \quad g \rightarrow agent_i$$
$$\quad path_i \quad Length(path_i)$$

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

$$b = \frac{\sum_{\forall \text{ cycle } i | (\text{agent}_i \text{ and } g \in N(\text{agent}_i)) \in \text{cycle}_i} \text{Length}(\text{cycle}_i)}{\sum_{\forall \text{ path } i | (\text{agent}_i \text{ and } g \in N(\text{agent}_i)) \in \text{path}_i} \text{Length}(\text{path}_i)} \quad ()$$

$$\begin{aligned} & \text{agent}_i \quad \text{agent}_i \quad () \\ & \quad \quad \quad \text{agent}_i \\ & \quad \quad \quad \text{agent}_i \\ & \quad \quad \quad c_2 \quad c_1 . \\ & \text{agent}_i \quad () \end{aligned}$$

$$\begin{aligned} & r(i, j) = f(\text{agent}_i) \frac{1}{d(\text{agent}_i, \text{agent}_j)} \quad () \\ & \quad \quad \quad j \quad i \quad \text{agent}_j \quad \text{agent}_i \quad r(i, j) \\ & \quad \quad \quad () \end{aligned}$$

//Relationship Structure Determination Algorithm

Define a 2D ACLA and initialize parameters

for each web page **do**

 assign an agent to a web page

 place agent randomly at cell

 equip agent with LA

end for

while (iterations exceeds a threshold)

 user_log: Array of [Number of Users][Users Path]

 /* user log, pages viewed by each user. Each row contains path of a user. */

for each cell **do** //traverse cells in row major method

 compute activation value for agent placed in current cell by equation (5)

if (agent's activation value > R) **then**

 select one action of agent's LA randomly

 move to an unoccupied neighbour cell based on selected action

 compute reward parameter by equation (8)

 compute penalty parameter by equation (9)

if (penalty parameter!= 0) **then**

 penalize(action) by equation (2)

else if (reward parameter!= 0) **then**

 reward(action) by equation (1)

end if

end if

end for

end while

Relationship Structure:= compute relevance degree of each agent and its neighbour by equation (10)

. ()

:

$$crawl_score(Page_j) = r(i, j) \times hub(i)$$

() $Page_j$ () $r(i, j)$ $hub(i)$ [15] HITS i Hub HITS

Lui [16]

[16]

$$m \quad n \quad () \quad n \quad m \quad s_{m,n} \quad ($$

()

[2] BestFirst

()

:

()

:()

	/
()	T_c /
	α_u
	λ /
$\Delta M'_t$	σ_m /
$\Delta M'_t$	μ_m /
	α_p
	σ_t /
	θ
	/

: ()

	c_l	,
	c_2	,
f	k	,
	B	,
	R	,

[17] BestFirst

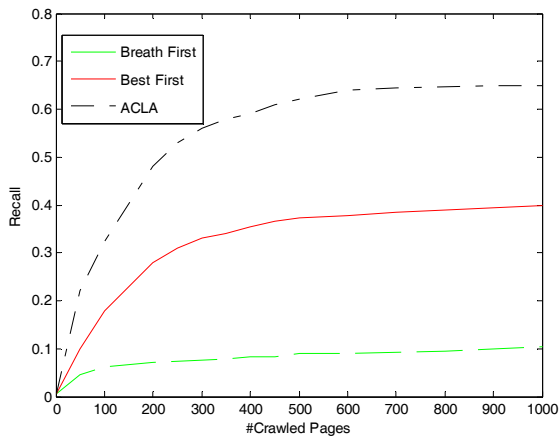
ACLA

()

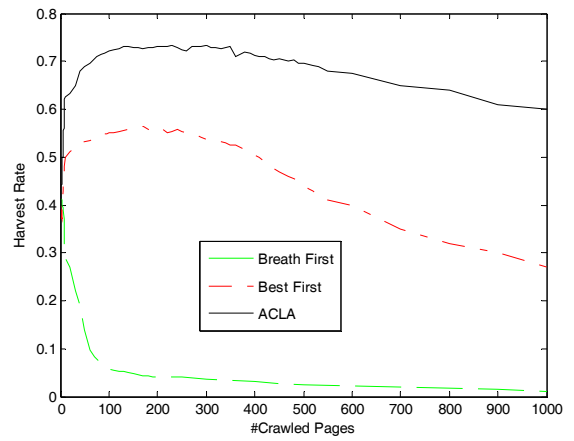
BestFirst BestFirst

A A A

ACLA ACLA BestFirst ()



. ()



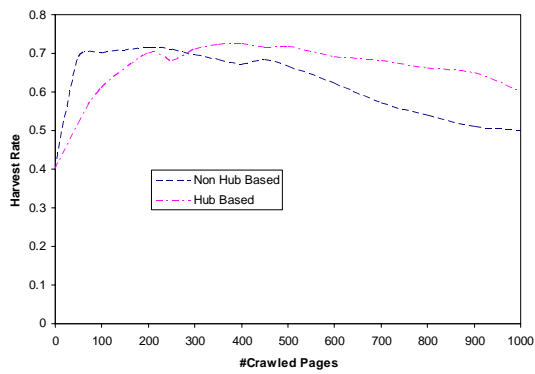
. ()

:

ACLA

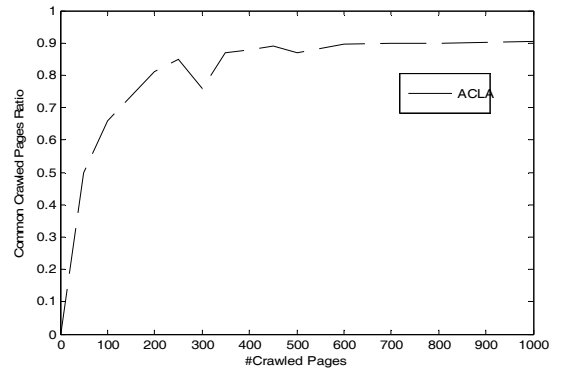
()

ACLA



Hub

. ()



. ()

:

Hub

Hub

Hub

Hub

()

ACLA

(

)

Hub

(Hub

)

()

$$\left(\begin{array}{c} \end{array} \right)$$

Hub

()

Best First

- [1] Gulli, A., Signorini, A., "The indexable web is more than 11.5 billion pages", *In Special interest Tracks and Posters of the 14th international Conference on World Wide Web*, Chiba, Japan, May 10-14, 2005.
- [2] Cho, J., Garacia-Molina, H., Page, L., "Efficient crawling through URL ordering ". *Comput Netw.* 30, pp. 161-172, 1998.
- [3] Hersovici, M., Javoci, M., Maarek, Y.S., Pelleg, D., Shtalham, M., "The shark-search algorithm—An application: Tailored Web site mapping", *Proceedings of the 7th International World-Wide Web Conference*, 1998.
- [4] Page, L., Brin, S., Motwani, R., Winograd, T., "The PageRank citation ranking: bringing order to the web ", *Stanford Publications*, 1998.
- [5] Chakrabarti, S., Van den Berg, M., Dom, B., "Focused Crawling: A New Approach to Topic-Specific Web Resource Discovery", *Proceedings of the 8th International WWW Conference*, Toronto, Canada, May 1999.
- [6] Diligenti, M., Coetzee, F., Lawrence, S., Giles, C., Gori, M., "Focused Crawling Using Context Graphs", *Proceedings of the 26th International Conference on Very Large Databases (VLDB 2000)*, Cairo, Egypt, September 2000.
- [7] Su, C., et al, "An Efficient Adaptive Focused Crawler Based on Ontology Learning ", *Proceedings of the Fifth International Conference on Hybrid Intelligent Systems (HIS'05)*, 2005.
- [8] Aggarwal, C., Al-Garawi, F., Yu, P., "Intelligent Crawling on the World Wide Web with Arbitrary Predicates", *Proceedings of the 10th International World Wide Web Conference*, Hong Kong, May 2001.
- [9] Pant, G., Srinivasan, P., "Learning to Crawl: Comparing Classification Schemes", *ACM Transactions on Information Systems*, Vol. 23, No. 4, pp. 430–462, October 2005.
- [10] Narendra, K. S. and Thathachar, M. A. L., *Learning Automata: An Introduction*, Prentice Hall, 1989.
- [11] 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.
- [12] Meybodi, M. R., Beigy, H., Taherkhani, M., "Cellular Learning Automata and Its Applications", *Journal of Science and Technology*, University of Sharif, No. 25, pp.54-77, Autumn/Winter 2003-2004.
- [13] Beigy, H, Meybodi, M. R., "Open Synchronous Cellular Learning Automata", *Advances in Complex Systems*, Vol. 10, No. 4, pp. 1-30, December 2007.
- [14] Beigy, H., Meybodi, M. R., "Asynchronous Cellular Learning Automata", *Automatica, Journal of International Federation of Automatic Control*, Vol. 44, No. 5, May 2008, to appear.
- [15] 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.
- [16] 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.
- [17] Menczer, F., Pant, G., Srinivasan, P., Ruiz, M., "Evaluating Topic-Driven Web Crawlers", *Proceedings of the 24th Annual International ACM/SIGIR Conference*, New Orleans, USA, 2001.

-
- ⁷ Sibling
 - ⁸ Bottomless Site
 - ⁹ Stationary
 - ¹⁰ Non-Stationary
 - ¹¹ Linear Reward-Penalty
 - ¹² Linear Reward epsilon Penalty
 - ¹³ Linear Reward Inaction
 - ¹⁴ Time Driven
 - ¹⁵ Step Driven
 - ¹⁶ Stagnation
 - ¹⁷ Harvest Rate
 - ¹⁸ Recall