mmeybodi@aut.ac.ir                                    forsati@mrl.ir

:

.

.

"              "  "                          "

.              .

.

*HITS*                      .

.              .

.

HITS                                                          :

# An Efficient Algorithm based on Web Usage Data and Structure of the Site for Web Page Recommendation

**R. Forsati**

Electrical and Computer Engineering and Information
Technology Department
Islamic Azad University
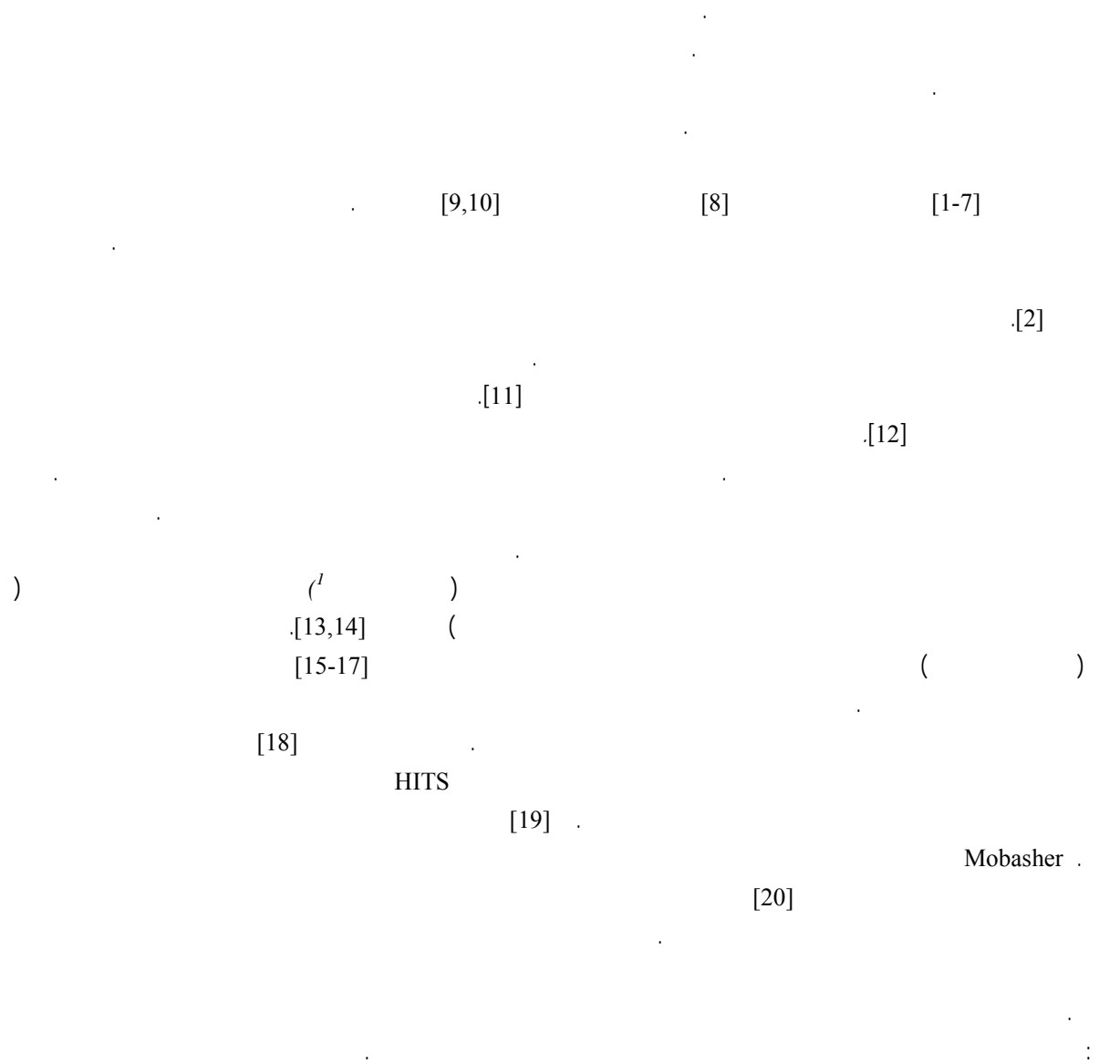Qazvin, Iran
forsati@mrl.ir

**M. R. Meybodi**

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

**Abstract:** One of the challenging tasks in improving web personalization algorithms is the simultaneous use of user's log data and underlying site's link information. In this paper after introducing a novel method for weighting web pages; an algorithm that simultaneously uses both web usage logs and underlying site's link information is proposed. The duration time of visited page and frequency of visiting is employed to measure the weight of each page which correctly determines the user's interest and importance of each page. Also the in-degree of each page is used to benefit the pages with low in-degree with high frequency and duration time. The proposed algorithm solves two main problems in web personalization. The first problem is the recommendation of recently generated pages which are not visited yet and the second problem is that the quality of system significantly decreases with increasing the number of recommended pages. In the proposed algorithms, the first recommending page is generated by using a novel weighted association rule mining algorithm. Then this page is expanded using HITS algorithm with pages which

belong to the same cluster as the recommended pages belong. For clustering the web pages, an algorithm based on Learning Automata and graph partitioning is presented. This method gives chance to pages which are not visited in any session in the log. Simulation results on real data set reveals that the proposed algorithm improves the quality of recommended pages significantly and solves the above mentioned problems.

**Keywords:** Web mining, Weighted association rule miming, Learning automata, HITS algorithm
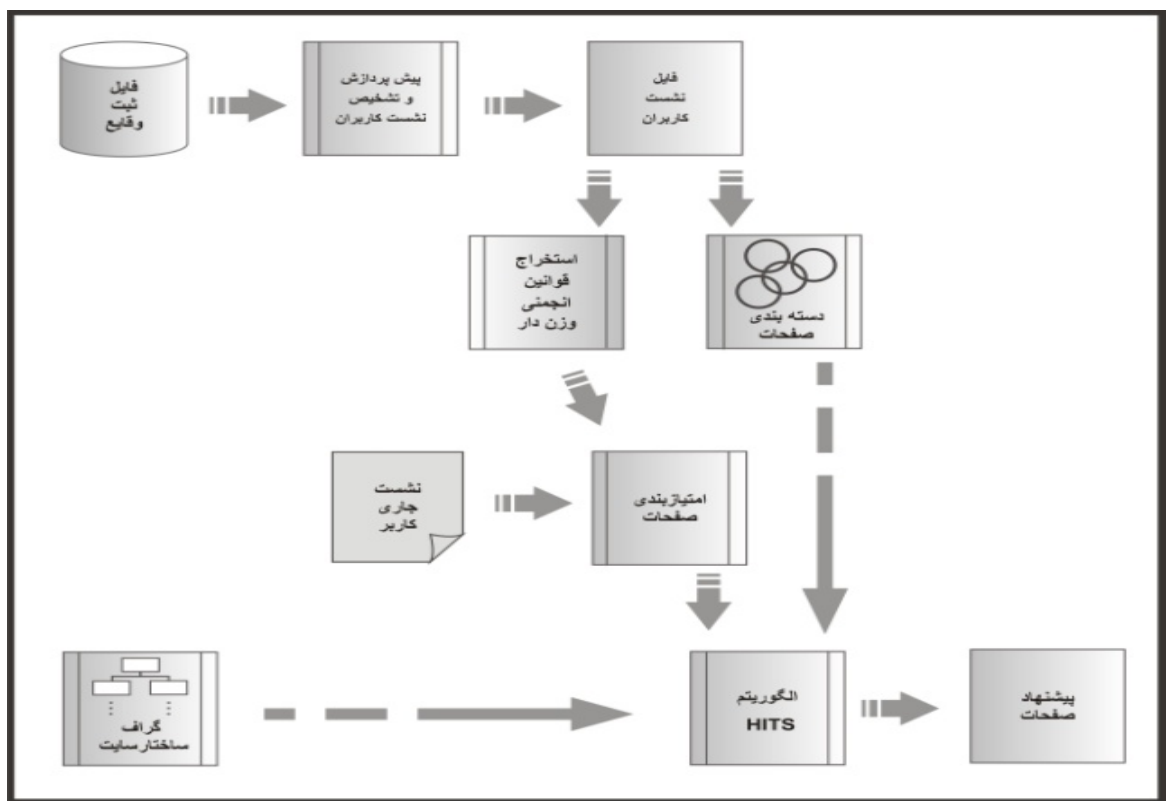
.

.

.

.

[1-7]                     [8]                     [9,10]           .

.

.[2]

.[11]

.[12]

.                               .

.

)                               (         )

.[13,14]         (

[15-17]                                                         (           )

.

[18]                           .

HITS

[19]   .

Mobasher  .

[20]

.

.

:

---

:

:

•
•
•

[2,3,21]

1

:

.                                            $P$

$P = \{p_1, p_2, ..., p_m\}$                          URL          .          $T$

$T = \{t_1, t_2, ..., t_n\}$                    $t_i \in T$

$t_i = \{(p_1, w_1), (p_2, w_2), ..., (p_m, w_m)\}$          $P$          .          $t_i$          $m$

$w_i$          $p_i$          $t_i$          .

"          "          "          "

:

.

[22]

.                                            .

.

.                "          "          "          "

.

$d_p(P)$                          $f_p(P)$          .

٤

$$f_p(P) = \frac{Visit(P)}{\sum_{Q \in T} Visit(Q)} * \frac{1}{In\deg ree(p)} \qquad (\ )$$

$$d_p(P) = \frac{\dfrac{Duration(P)}{Size(P)}}{\max_{Q \in T}(\dfrac{Duration(P)}{Size(P)}} \qquad (\ )$$

$\alpha$      .

"      " "      "      .

.      $\alpha$

$$W(p) = \frac{\alpha * f_p(P) * d_p(P)}{f_p(P) + d_p(P)} \qquad (\ )$$

.

( )      . [23]

.      $T = \{t_1, t_2, ..., t_{n}\}$      .

.      $I = \{i_1, i_2, ..., i_n\}$

.      $X, Y$      $X \Rightarrow Y, where X \subset I, Y \subset I, X \cap Y = \phi$

.[24]      $Y$ ,      $X$

.

$$r = \{(p_i, w_i), (p_j, w_j), ...(p_K, w_k) \Rightarrow (p_m, w_m), ...(p_n, w_n) \, و \, S \, و \, C\} \in R$$

.      $w_i$      $C$      $S$

.      ( )

.

.

.

---

[5] Support
[6] Confidence

$$Confidence = \frac{\sup port(X \cup Y)}{\sup port(X)}$$

( )

*rs*

.                                                                                        .

.                                  $p_1 \rightarrow p_2 \rightarrow p_3 \rightarrow ... \rightarrow p_k$

*rw*

.

.                                                                                        .

[25-27]

.                                                                  .

.

.                                                                            .

)                                              (     )                                    (          )

.                                                                                              (

$S =< w_1^{\ s}, w^s_{\ 2},..., w_n^{\ s} >$                    (*n*                    ) *S*

.                                                      $R =< w^r_{\ 1}, w^r_{\ 2},..., w_m^{\ r} >$

$$similarity(S,R) = \frac{\sum_k w_k^{\ r}.w_k^{\ s}}{\sqrt{\sum_k (w^r_{\ k})^2 \times \sum_k (w^s_{\ k})^2}}$$

( )

.                                        (                              *URL*) *u*

$rank(S,u) = Confidence(R) * similarity(S,R)$                                      ( )

*2*                                        .

.                                        (                    )      *N*

.                      *HITS*                                              $N+1$

٦

*HITS*

.

.

[28]

)　　　　　　(　　　　　　　　　　　)　　　　　　　　　　　　　　　.

.　　　　　(

.　　　　　　　　　　　　　[29]

.

.

.[30]

*n*

.　　　　　　　　　　　　*P*　　　　　　　　　　　　　　*n-1*　　　31

.

*j*　　*i*　　　　　　　　　　　　　　.　　　　　　　　　　　　　　　　　　*i*

*P*　　　　　　　　　　　　.　　　*i*　　　*j*

.　　　　[21]

.　　　　　*j*　*i*　　　　　　　*p_{ij}*　　*i*　　　*j*

.　　　$(p_{ij} \neq p_{ji})$　*P*

*s_{ij}*　.　　　*S*　　　　　　　*P^{T}*　　　　　*P*

.　　*j*　*i*

$$S = P \cdot P^{T}$$

( )

$$s_{ij} = \sum_{k} p_{ik} p_{kj}$$

.　　　*S*

*hMeTis*　　　　.

.　　　　　　　.

.

# *HITS*

---

[7] Synonym
[8] Homonym
[9] Minable
Transpose [10]

*HITS*

[28]

)　　　　　　(　　　　　　　　　　　)　　　　　　　　　　　　　　　.

(

[29]

.[30]

*n*

*P*　　　　　　　　　　　　　　*n-1*　　　31

*j*　　*i*　　　　　　　　　　　　　　　　　　　　　　　　　　　　*i*

*P*　　　　　　　　　　　　　　　*i*　　　*j*

[21]

*j*　*i*　　　　　　　*p_{ij}*　　*i*　　　*j*

$(p_{ij} \neq p_{ji})$　*P*

*s_{ij}*　　　　　*S*　　　　　　　*P^{T}*　　　　　*P*

*j*　*i*

$$S = P \cdot P^{T}$$

( )

$$s_{ij} = \sum_{k} p_{ik} p_{kj}$$

*S*

*hMeTis*

# *HITS*

---

[7] Synonym
[8] Homonym
[9] Minable
Transpose [10]

$N+1$ 	.

.	.	*HITS* [32]

*HITS*	.

.

.[33]

.	*rt*	.

*rt* .	*rs*	*rt*	.	*rt*

.

.

.

*Lui*	.	[34]	.

.	*CTI DePaul*	.
.[35]	*CTI DePaul*

.

.	*rt*	*rw* :

*rw*	.	.

.	*rw*	*rw* + *rt*

.	$rp = \{x_{rw+1}, x_{rw+2}, ..., x_{rw+|rs|}\}$	.	[36]

:

$$Pr\,ecision(rs, rp) = \frac{|rs \cap rp|}{|rs|}$$	( )

$$Coverage(rs, rp) = \frac{|rs \cap rp|}{|rp|}$$	( )

.	$Pr\,ecision$

.

$rw + 1$	.

.

$Coverage$	.

.

.

٨

.

[11] *AR* . AR

. K ( )

d I d+1 I

. . rw

rw+1

( )

.

AR .

rw . AR .

.

. AR



:

*HITS*

1. *J. Srivastava, R. Cooley, M. Deshpande, and P. N. Tan, Web Usage Mining: Discovery and Applications of Usage Patterns from Web Data ,SIGKDD Explorations, Vol. 1, No. 2, 2000, pp.12 23.*

2. *B. Mobasher, R. Cooley, and J. Srivastava, Automatic Personalization Based on Web Usage Mining, Communications of the ACM, vol. 43, no.8, 2000, pp. 142 151.*

3. *B. Mobasher, H. Dai, T. Luo, and M. Nakagawa, Discovery and Evaluation of Aggregate Usage Profiles for Web Personalization, Data Mining and Knowledge Discovery, vol. 6, no. 1, 2002, pp. 61-82.*

4. *G. Pierrakos, C. Paliouras, C. Papatheodorou and C. D. Spyropoulos, Web Usage Mining as a Tool for Personalization: A Survey, User Modeling and User-Adapted Interaction, vol. 13, no. 4, 2003, pp. 311-372.*

5. *Kazienko, P., Adamski, M.: AdROSA - Adaptive Personalization of Web Advertising. Information Sciences 177(11), 2269 2295, 2007.*

6. *P. Kazienko, M. Kiewra, Personalized Recommendation of Web Pages. In: Nguyen, T. (ed.) Intelligent Technologies for Inconsistent Knowledge Processing. Advanced Knowledge International, International, Adelaide, South Australia, ch. 10, 2004, pp. 163 183.*

7. *P. Kazienko, Filtering of Web Recommendation Lists Using Positive and Negative Usage Patterns , Springer-Verlag Berlin Heidelberg, 2007.*

8. *F. Heylighen and J. Bollen, Hebbian Algorithms for a Digital Library Recommendation System , Proceedings of the International Conference on Parallel Processing Workshops (ICPPW 02), 2002, pp. 439-446.*

9. *J. Zhu, J. Hong, and  J. G. Hughes,  Mining Conceptual Link Hierarchies from Web Log Files for Adaptive Web Site Navigation, ACM Transactions on internet Technology (TOIT), in press, 2003.*

10. *J. Zhu, J. Hong, and  J. G. Hughes,  Using Markov Chains for Link Prediction in Adaptive Web Sites, Proc. Of Soft-Ware: First International Conference on Computing in an Imperfect World, Lecture Notes in Computer Science, Springer, Belfast, April 2002, ,pp.60-78.*

11. *B. Mobasher, H. Dai, T. Luo, M. Nakagawa, Effective Personalization based on Association Rule Discovery from Web Usage Data ,Proceedings of the 3$^{rd}$ ACM Workshop on Web Information and Data Management, 2001.*

12. *H. Dai, B. Mobasher, Integrating Semantic Knowledge with Web Usage mining for Personalization, 2004.*

13. *B. Mobasher, H. Dai, Y. Sun, J. Zhu, Integrating Web Usage and Content Mining for More Effective Personalization", Proceeding of the EC-WEB Conference, 2003.*

14. *B. Mobasher, H. Dai, T. Luo, M. Nakagawa, "Using Sequential and Non-sequential Patterns for Predictive Web Usage Mining Tasks", Proceedings of the IEEE International Conference on Data Mining, Maebashi City, Japan, 2002.*

15. *M. Richardson, P. Domingos, The Intelligent Surfer: Probabilistic Combination of Link and Content Informationin PageRank , in Neural Information Processing System, 2002.*

16. *T. Haveliwala, Topic-Sensitive PageRank , Proceeding of WWW Conference, Hawaii, 2002.*

17. *M. S. Aktas, M. A. Nacar, F. Menczer, Personalizing PageRank Based on Domain Profiles , Proceeding of WEBKDD Workshop, Seattle, 2004.*

18. *J. Wang, Z. Chen, L. Tao, W. Ma, L. Wenyin, Ranking User s Relevance to a Topic through Link Analysis on Web Logs , Proceeding of the WIDM 02, 2002.*

19. *J. Borges, M. Levene, Data Mining of User Navigation Patterns , in Revised Papers from the International Workshop on Web Usage Analysis and User Profiling, 2000, pp. 92-111.*

20. *M. Nakagawa, B. Mobasher, A Hybrid Web Personalization Model Based on Site Connectivity , in Proceeding of the 5<sup>th</sup> WEBKDD Workshop, Washington DC, 2003.*

"                                                        "                                    .

.

22. *P.K. Chan, A Non-invasive Learning Approach to Building Web User Profiles, in: Workshop on Web usage Aanalysis and User Profiling, Fifth International Conference on Knowledge Discovery and Data Mining, San Diego, 1999.*

23. *Demiriz, A. 2002. Enhancing Product Recommender Systems on Sparse Binary Data, to be published in the Journal of Data Mining and Knowledge Discovery, 2003.*

24. *A. Agrawal and R. Srikant, "Fast Algorithms for Mining Association Rules in Large Databases", In Proceedings of the 20th International Conference on Very Large Data Bases (VLDB'94), Santiago, Chile, 1994.*

25. *B. Mobasher," Web Usage Mining and Personalization", In Practical Handbook of Internet Computing. Munindar P. Singh (ed.), CRC Press, 2005.*

26. *A. Demiriz, "Enhancing Product Recommender Systems on Sparse Binary Data", accepted to be published in the Journal of Data Mining and Knowledge Discovery, 2003.*

27. *B. Mobasher, H. Dai, T. Luo, and M. Nakagawa: Improving the Effectiveness of Collaborative Filtering on Anonymous Web Usage Data. In Proceedings of the IJCAI  Workshop on Intelligent Techniques for Web Personalization (ITWP01), Seattle, WA, 2001.*

28. *M. Junichiro, M. Yutaka, I. Mitsuru, F. Boi, "Keyword Extraction from the Web for FOAF Metadata," Proceeding of 1st International Workshop on Friend of a Friend, Social Networking and the Semantic Web, Galway, Ireland, 2004*

29. *J. Manuel Barrueco Cruz, T. Krichel, "Automated Extraction of Citation Data in a Distributed Digital Library," Proceedings of the 2nd International Workshop on New Developments in Digital Libraries, 2002.*

"                                                        "                                    .

.

31. *M .A. L. Thathachar, R. Harita Bhaskar, "Learning Automata with Changing Number of Actions," IEEE Transactions on Systems Man and Cybernetics, vol. 17, no. 6, Nov. 1987.*

32. *J. Kleinberg, "Authoritative sources in a hyper-linked environment". Proc. ACM-SIAM Symposium on Discrete Algorithms, 1998. Also appears as IBM Research Report RJ 10076(91892), May 1997.*

33. *O. R. Za¨ıane, J. Li, R. Hayward, "Mission-Based Navigational Behavior Modeling for Web Recommender Systems," Springer-Verlag Berlin Heidelberg, 2006.*

34. *J. Mori, Y. Matsuo, M. Ishizuka, B. Faltings, "Keyword Extraction from the Web for FOAF Metadata", Proceeding of 1st International Workshop on Friend of a Friend, Social Networking and the Semantic Web, Galway, Ireland, 2004.*

35. *J. Manuel Barrueco Cruz, T. Krichel, "Automated Extraction of Citation Data in a Distributed Digital Library," Proceedings of the 2nd International Workshop on New Developments in Digital Libraries, 2002.*

36. *T. Haveliwala, "Topic-Sensitive PageRank", in Proceedings of the 11thInternational Conference on World Wide Web, New York: ACM Press, 2002,, pp. 517–526.*