



Mining users` navigation patterns for building web pages recommendation system

Farzad Soleymani Sabzchi^{1✉}, Shahram Jamali², Maryam Jafari³

(1) Department of Computer Engineering, Zanzan Branch, Islamic Azad University, Zanzan, Iran

(2) Computer Engineering Department, University of Mohaghegh Ardabili, Ardabil, Iran

(3) Department of Computer Engineering, Zanzan Branch, Islamic Azad University, Zanzan, Iran

f_soleymani63@yahoo.com; jamali@iust.ac.ir; m123_jafari@yahoo.com

Received: 2013/03/19; Accepted: 2013/05/02

Abstract

Due to the quick growth of the World Wide Web, retrieval of useful information from the Internet for a particular web user or a group of users becomes very difficult. Recommendation systems using web usage mining help providing an adaptive web environment for the web users. This paper presents a novel approach for page recommendation using fuzzy association rule mining algorithm. This method extracts previous users` access patterns and then employs them to recommend appropriate web pages for the active user. An illustrative example explains this method in details.

Keywords: Web recommender, Web usage mining, Fuzzy association rule mining, Pattern discovery

1. Introduction

Over the recent decades, fast growth of World Wide Web has led to huge amount of available information that is not simply accessible for the users. Recommender Systems (RS) help us to cope with confusion among the huge volume of information, and recommend pages that may fit our interests [1,2]. Web recommender systems are one of the applications obtained from Web Usage Mining (WUM) technique. WUM can elicit the web user's navigational behavior from secondary data, such as web server access logs, proxy server logs, browser logs, user profiles, user sessions, and user queries. Then it models the behavior as navigation patterns for designing the intelligent web sites [3,4,5]. Some examples of these web sites are recommending books, CDs and other products at Amazon.com [6], recommending movies by MovieLens [7], recommending books at LIBRA [8], and recommending electronic television program guides [9]. Several WUM techniques have been used to effectively develop recommendation systems. There have been efforts to use association rules [10], sequential patterns [11], and Markov models [12] in recommender systems.

Two effective factors in recommend pages to a web user are previous visited pages of the active user and the history of browsing the web pages of previous users. Hence, the main purpose of this paper is to recommend the appropriate pages based on discovered patterns from previous users` sessions (off-line) and active user`s navigational behavior (on-line). Patterns of previous users are extracted using Fuzzy Association Rule Mining

(FARM) technique in off-line part. Afterward, on-line part of proposed recommendersystem gets the active user's request, searches the antecedent of the rules for a match. In the case of finding a perfect match, recommends the pages in the consequence.

The rest of this paper is as follows. In the next section motivation of the paper is presented. In section 3 related works on web recommendation systems are reviewed. The proposed approach in this paper is described in section 4. An example to illustrate the proposed approach is given in section 5. Finally a conclusion of this work is presented in section 6.

2. Motivation of the research

Most of the web users complain about finding useful information on web sites. Web recommender systems predict the information needs of users and provide them with recommendations to facilitate their navigation. Recommender systems have been extensively explored in web mining. However, the quality of recommendations and the user satisfaction with such systems are still not optimal. All recommender systems based on web usage mining techniques have strengths and weaknesses. One of these weaknesses is lack of consideration of an appropriate measure to calculate the users' interest degree of pages. Applying fuzzy association rule mining to a recommender system leads to calculate users' interest of pages more accurately. Consequently generated recommendations will be more desirable.

3. Related Works

In this section some widely used techniques for implementation of recommender systems are described. The first technique clusters pages found from web server log files. A simple k-means algorithm clusters user sessions based on vector distances. This model is based on clusters of user sessions and uses the visiting time of pages and does not consider the visiting order of the pages [13]. In the click-stream tree model, the recommendations are generated using the recommender technique proposed in [14]. This technique uses the sequence of pages visited in a user session. In this method pairwise similarities between user sessions are calculated by a new similarity measure. Afterward the user sessions are clustered using a graph-based clustering algorithm. Nakagawa and Mobasher [15], proposed a hybrid recommender system that can intelligently switch among various recommendation models. Switching criterion is based on the degree of hyperlink connectivity and the neighborhood of a user's current location within the site to select the best recommendation model. This system uses three recommendation models are based on Association Rules, Sequential Patterns and Contiguous Sequential Pattern to generate navigation patterns of Web users.

In another work, the recommender system proposed by Mobasher et al., [16] is based on association rule mining method for web personalization. This approach present a data structure for storing the discovered frequent item sets from clickstream data. Presented recommendation algorithm applies this data structure to generate recommendations in real-time, without the need to extract all association rules from frequent item sets. This system produces recommendations based on matching the current user navigational behavior against patterns discovered through association rule mining.

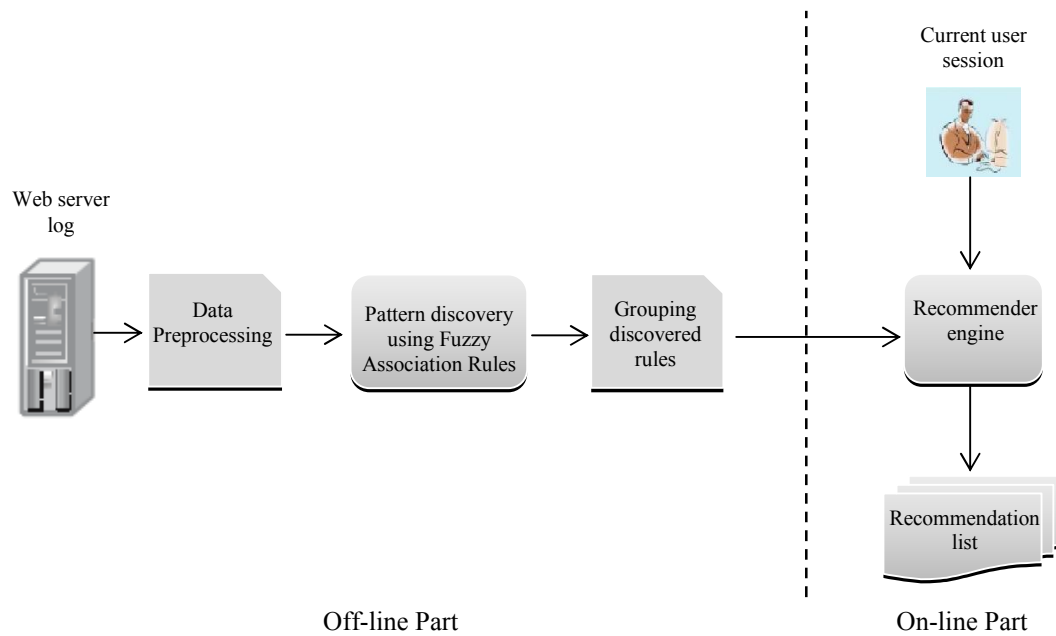


Figure 1. Recommender system architecture

4. Proposed Approach: Web Recommender System Based on FARM

The framework of usage based recommendation system shown in Fig. 1 which includes off-line part and on-line part. This system includes different phases such as preprocessing, pattern discovery, off-line part and on-line part, which are described in following subsections.

4.1 Preprocessing

Data preprocessing is the first phase of this proposed system. The results of data preprocessing directly impact the results of next steps. This phase consists of data cleaning, user and session identification and generating the database of log file. Usually, preprocessing is done by considering the visit time of a page.

A weight measure presented in [17] is used to approximate the interest degree of a web page for a user. Here, two concepts related to this measure are introduced: "*Frequency*" and "*Duration*". Frequency is the number of times that a page is accessed by different users. A parameter that must be considered to calculate the frequency of a page is the In-degree of that page (e.g. the number of incoming links to the page). The formula of "*Frequency*" is given in (1) which is normalized by the total number of visits of web pages in the session:

$$Frequency(P) = \frac{Numberofvisit(P)}{\sum_{Q \in T} Numberofvisit(Q)} \times \frac{1}{Indegree(P)} \quad (1)$$

"*Duration*" is defined as the time the user spent on a page. It is guessed that the longer time a user spends on a page means that the user is more interested in that page. "*Duration*" of a web page is shown in (2).

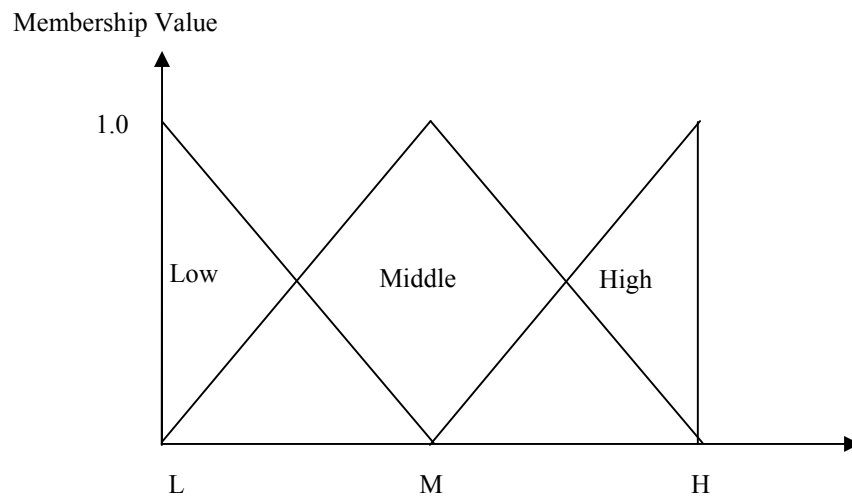


Figure 2. The triangle membership function used for partitioning interest degree of the page

$$Duration(P) = \frac{\frac{Total\ duration(P)}{Size(P)}}{\max_{Q \in T} \frac{Total\ Duration(P)}{Size(P)}} \quad (2)$$

"Frequency" and "Duration" are considered two strong indicators of users' interest. The harmonic mean of these two measures is used to represent the interest degree of a web page to a user in the session, shown as below:

$$Importance(P) = \frac{2 \times Frequency(P) \times Duration(P)}{Frequency(P) + Duration(P)} \quad (3)$$

Equation (3) guarantees that "Interest" of a page is high only when "Frequency" and "Duration" are both high.

4.2 Pattern Discovery

Pattern discovery is a phase which extracts the user behavioral patterns from the formatted data. In this phase, several data mining techniques are applied to obtain hidden patterns reflecting the typical behavior of users. Some important techniques for this phase are: path analysis, standard statistical analysis, clustering algorithms, association rules, classification algorithms, and sequential patterns. In this paper fuzzy partition method and fuzzy association rule mining technique are used to extract interested patterns.

4.2.1 Fuzzy partition method

The concepts of fuzzy sets and linguistic variables were proposed by Zadeh [18,19] and it is reasonable that different linguistic values can indicate diversity of degree of the attribute in many applications [20]. In the simple fuzzy partition method, K several linguistic values are defined in each quantitative attribute. K is also predefined before performing the method. Triangular and trapezoid membership functions are usually used for the linguistic values.

In this paper, the quantitative attribute "Interest degree" is partitioned by simple fuzzy partition method with triangle membership function. This attribute is divided into 3 linguistic values (*Low*, *Middle*, *High*) and its membership function is predefined in

Fig. 2. These membership functions can be represented by triangle fuzzy numbers. Limits of each linguistic variable can be displayed as below: $A_{Low}^{Interest} = (l, m, m)$ and $A_{Middle}^{Interest} = (l, m, h)$ and $A_{High}^{Interest} = (m, m, h)$.

Using formulas (4), (5) and (6), the attribute is partitioned into suitable linguistic value and obtained its membership in accordance with membership function, respectively.

$$\mu_{Low}^{Interest}(x) = \begin{cases} 1, & x < l \\ \frac{m-x}{m-l}, & l \leq x \leq m \\ 0, & x > m \end{cases} \quad (4)$$

$$\mu_{Middle}^{Interest}(x) = \begin{cases} \frac{x-l}{m-l}, & l < x \leq m \\ \frac{h-x}{h-m}, & m < x < h \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

$$\mu_{High}^{Interest}(x) = \begin{cases} 0, & x < m \\ \frac{x-m}{h-m}, & m \leq x < h \\ 1, & x \geq h \end{cases} \quad (6)$$

4.2.2 Fuzzy association rules

In general, the interesting or preference of each user is difficult to describe clearly. In order to generate recommendations, fuzzy association rules are used to provide behavioral patterns. Association rules indicate how much the degree of user interest in a page will effect on the next visits of pages in the web site. A rule consists of an antecedent and a consequent. The rule expresses that when the antecedent happens the consequent will also happen. The importance of each rule is measured using its support and confidence. The support of a rule is the number of instances that contain both the itemset of its antecedent and the itemset of its consequence. The conditional probability of the occurrence of the consequent given the antecedent is referred to as the confidence of the rule.

4.3 Off-line part of the Recommender System

The off-line part, models usage patterns from the web server access log data and builds a predictive model based on the extracted usage patterns. After generating frequent patterns by FP-Growth algorithm [21], the procedure of fuzzy association rules mining is executed. Steps that are performed in off-line part are as follows:

Off-line part algorithm:

Step 1. Preparation and preprocessing of web log data.

Step 2. Attributes will be partitioned into appropriate linguistic value using fuzzy partition method.

Step 3. Frequent patterns are extracted by using fuzzy FP-Growth algorithm

Step 4. All the fuzzy association rules are generated from available frequent patterns.

Step 5. The general fuzzy association rule is indicated by (7). The Fuzzy Support (FS) for all obtained fuzzy association rules are calculated using (8) as follows:

$$R : A_{s_1}^{x_1}, A_{s_2}^{x_2}, \dots, A_{s_k}^{x_k} \rightarrow A_{s_\alpha}^{x_\alpha} \quad (7)$$

$$FS(A_{s_1}^{x_1}, A_{s_2}^{x_2}, \dots, A_{s_k}^{x_k}) = \frac{1}{n} \sum_{j=1}^n [\mu_{s_1}^{x_1}(d_j^{x_1}) \times \mu_{s_2}^{x_2}(d_j^{x_2}) \times \dots \times \mu_{s_k}^{x_k}(d_j^{x_k})] \quad (8)$$

Where, x_i is the i th attribute, s_i is the term of degree in x_i , $A_{s_i}^{x_i}$ is the linguistic value of s_i in the i th attribute, $\mu_{s_i}^{x_i}$ is the membership function of s_i in the i th attribute, $d_j^{x_i}$ is the original value of the i th attribute in the j th tuple.

Step 6. If the FS value of fuzzy association rule is larger than or equal to the user-specified min FS, this rule will be accepted.

Step 7. Ineffective fuzzy association rules will be filtered out with fuzzy confidence. The Fuzzy Confidence (FC) of R is defined as follows:

$$FC(R) = FS(A_{s_1}^{x_1}, A_{s_2}^{x_2}, \dots, A_{s_k}^{x_k}, A_{s_\alpha}^{x_\alpha}) / FS(A_{s_1}^{x_1}, A_{s_2}^{x_2}, \dots, A_{s_k}^{x_k}) \quad (9)$$

When the FC(R) of rule R is less than the user-specified min FC, the rule R will be regarded as ineffective fuzzy association rule and filtered out.

Step 8. The useless fuzzy association rules are Pruned. If there exist two rules, denoted by R and S, possessing same consequence, and the antecedence of R is subset of S, then R is useless and can be removed.

Step 9. The interested discovered association rules are then grouped based on the number of items in their antecedent. The result is a group of rules with one item in the antecedent, a group with two items and so on.

4.4 On-line part of the Recommender System

The recommendation system takes the user's request (the active user session) and searches the items on the antecedent of each rule for finding matches. This process starts from the group of rules with one-item in the antecedent and then the next group with two items in the antecedent and so on. If matches are found, these rules are sorted according to their confidence values in descending order. Then the top N items of consequents from this list are recommended to the user.

When the user accesses two web pages, system searches the one-item group for the second web page and next, searches the two-items group for these two pages. When no match is found the system removes the user's browsing history and again the process starts from the beginning when a match is found. In the next section, an example to illustrate the method proposed in this paper is presented.

5. Illustrative Example

In this section, the function of the proposed recommendation system is illustrated with an example. This system uses a database of users' sessions. This database shown in Table 1 that consists of ten users' sessions and six web pages denoted P_1, P_2, P_3, P_4, P_5 and P_6 . Each item is represented by a tuple (page name: visit time). Features of these pages are shown in Table 2 that includes In-degree and the size of each page. The membership function used in this example is drawn in Fig. 2. In this membership function, $l = 0.15$, $m = 0.45$ and $h = 0.75$.

Table 1. The ten sessions used in the example

Session ID	Page Access Sequence
1	(P ₁ : 83), (P ₃ : 355), (P ₄ : 191), (P ₆ : 378)
2	(P ₂ : 207), (P ₃ : 267), (P ₄ : 455), (P ₅ : 407), (P ₆ : 411)
3	(P ₁ : 59), (P ₂ : 181), (P ₅ : 233), (P ₆ : 256)
4	(P ₁ : 67), (P ₂ : 391), (P ₃ : 271), (P ₄ : 241), (P ₅ : 291), (P ₆ : 501)
5	(P ₂ : 173), (P ₅ : 302), (P ₆ : 598)
6	(P ₁ : 120), (P ₃ : 107), (P ₄ : 219), (P ₅ : 391)
7	(P ₂ : 132), (P ₃ : 62), (P ₄ : 141), (P ₅ : 358), (P ₆ : 311)
8	(P ₁ : 95), (P ₂ : 283), (P ₃ : 211), (P ₆ : 279)
9	(P ₁ : 87), (P ₃ : 93), (P ₄ : 257), (P ₆ : 351)
10	(P ₁ : 139), (P ₂ : 98), (P ₃ : 148), (P ₅ : 183), (P ₆ : 393)

Table 2. Features of pages

Page	In-degree	Size
P ₁	1	293 KB
P ₂	2	382 KB
P ₃	2	170 KB
P ₄	1	131 KB
P ₅	3	111 KB
P ₆	3	209 KB

Table 3. The amount of interest to page in user's session

Session ID	Page Access Sequence
1	(P ₁ : 0.225), (P ₃ : 0.571), (P ₄ : 0.646), (P ₆ : 0.445)
2	(P ₂ : 0.207), (P ₃ : 0.414), (P ₄ : 0.734), (P ₅ : 0.734), (P ₆ : 0.385)
3	(P ₁ : 0.167), (P ₂ : 0.273), (P ₅ : 0.374), (P ₆ : 0.397)
4	(P ₁ : 0.156), (P ₂ : 0.368), (P ₃ : 0.482), (P ₄ : 0.646), (P ₅ : 0.52), (P ₆ : 0.452)
5	(P ₂ : 0.218), (P ₅ : 0.37), (P ₆ : 0.462)
6	(P ₁ : 0.199), (P ₃ : 0.247), (P ₄ : 0.53), (P ₅ : 0.374)
7	(P ₂ : 0.164), (P ₃ : 0.176), (P ₄ : 0.429), (P ₅ : 0.374), (P ₆ : 0.363)
8	(P ₁ : 0.361), (P ₂ : 0.429), (P ₃ : 0.559), (P ₆ : 0.462)
9	(P ₁ : 0.248), (P ₃ : 0.325), (P ₄ : 0.75), (P ₆ : 0.444)
10	(P ₁ : 0.371), (P ₂ : 0.186), (P ₃ : 0.429), (P ₅ : 0.364), (P ₆ : 0.462)

In this example, values are represented by three fuzzy regions: *Low*, *Middle* and *High*. Thus, three fuzzy membership values are produced for each page in a session according to the predefined membership functions (4), (5) and (6).

Here, assume that the predefined min FS and min FC are 0.20 and 0.60 respectively. The off-line part method can be described as follows:

Step 1,2: The results of these steps are presented in Table 3.

Step 3: using fuzzy FP-Growth algorithm, frequent patterns were extracted in Table 4.

Step 4,5: Fuzzy association rules from the results of step3 are discovered .

Step 6,7: Comparing FS and FC of the rules with predefined min FS and min FC, interested rules obtained and are shown in Table 5.

Step 8,9: By applying the results of Step 7, the pruned rules are grouped in Table 6.

Table 4. Conditional pattern base

Frequent patterns	Fuzzy values	Frequent Patterns	Fuzzy values
P ₄ .High	3.52	P ₂ .Low, P ₃ .Middle	2.98
P ₂ .Low	4.34	P ₂ .Low, P ₅ .Middle	4.27
P ₅ .Middle	4.51	P ₅ .Middle, P ₆ .Middle	3.76
P ₁ .Low	4.74	P ₃ .Middle, P ₅ .Middle	3.03
P ₃ .Middle	4.94	P ₁ .Low, P ₅ .Middle	2.98
P ₆ .Middle	8.15	P ₁ .Low, P ₆ .Middle	3.9
P ₄ .High, P ₆ .Middle	3.25	P ₁ .Low, P ₃ .Middle	3.8
P ₃ .Middle, P ₄ .High	3.52	P ₃ .Middle, P ₆ .Middle	4.62
P ₁ .Low, P ₄ .High	2.57	P ₆ .Middle, P ₃ .Middle, P ₁ .Low	2.96
P ₂ .Low, P ₆ .Middle	4.34	P ₆ .Middle, P ₃ .Middle, P ₄ .High	3.25
P ₆ .Middle, P ₁ .Low, P ₄ .High	2.3	P ₃ .Middle, P ₅ .Middle, P ₂ .Low	2.91
P ₃ .Middle, P ₁ .Low, P ₄ .High	2.57	P ₆ .Middle, P ₃ .Middle, P ₅ .Middle	2.28
P ₆ .Middle, P ₃ .Middle, P ₂ .Low	2.98	P ₃ .Middle, P ₁ .Low, P ₅ .Middle	2.23
P ₆ .Middle, P ₅ .Middle, P ₂ .Low	4.27	P ₆ .Middle, P ₃ .Middle, P ₁ .Low, P ₄ .High	2.3
		P ₆ .Middle, P ₃ .Middle, P ₅ .Middle, P ₂ .Low	2.91

Table 5. Fuzzy association rules

Association rules	FS	FC
P ₄ .High → P ₆ .Middle	0.3	0.86
P ₄ .High → P ₃ .Middle	0.25	0.7
P ₃ .Middle → P ₆ .Middle	0.43	0.87
P ₁ .Low → P ₆ .Middle	0.37	0.78
P ₅ .Middle → P ₆ .Middle	0.33	0.74
P ₂ .Low → P ₆ .Middle	0.37	0.85
P ₂ .Low → P ₅ .Middle	0.26	0.6
P ₂ .Low, P ₅ .Middle → P ₆ .Middle	0.22	0.87
P ₁ .Low, P ₃ .Middle → P ₆ .Middle	0.21	0.87
P ₃ .Middle, P ₄ .High → P ₆ .Middle	0.22	0.88

Table 6. Fuzzy association rules

Association rules	FS	FC	Group
P ₄ .High → P ₃ .Middle	0.25	0.7	group1
P ₂ .Low → P ₅ .Middle	0.26	0.6	group1
P ₂ .Low, P ₅ .Middle → P ₆ .Middle	0.22	0.87	group2
P ₁ .Low, P ₃ .Middle → P ₆ .Middle	0.21	0.87	group2
P ₃ .Middle, P ₄ .High → P ₆ .Middle	0.22	0.88	group2

The on-line part of the recommender system works as follow:

Assume three users visited these pages in this manner: user1 {P₂, P₅}, user2 {P₁}, user3 {P₂, P₃}. For user1, recommender system searches P₂ in the antecedents of rules in group1. Next it searches P₂, P₅ in the antecedents of rules in group2. Sorted rules according to FC are recommended to this user. These recommendations are P₆.Middle and P₅.Middle. For user2 who visited P₁, recommender system could find no appropriate page and waits until this user visit the second page without removing the user's browsing history. Recommender system could not find any matches in group2 for user3 with two visited pages P₂ and P₃. In this situation, the recommender system will drop the user's browsing history.

6. Conclusion

In this paper, a framework was proposed to generate recommendations for web users. For this purpose it employed web usage data. Our recommender system integrates the results from two distinct parts, off-line and on-line parts. It generates a single recommendation list for current user. Off-line part of the system discovers frequent patterns by identifying popular pages in user's navigation path based on fuzzy association rule mining. Next, on-line part uses these navigation patterns and current user's request to generate recommendations for the same current user. An illustrative example showed how the proposed method recommends pages based on the system history.

7. References

- [1] Campos L, Fernandez-Luna J, Huete J. A collaborative recommender system based on probabilistic inference from fuzzy observations. *Fuzzy Sets and Systems* 2008;159:1554–1576.
- [2] Goksedef M, Oguducu S. A Consensus Recommender for Web Users. Springer-Verlag Berlin Heidelberg 2007;287–299.
- [3] Buchner A, Mulvenna M. Discovering Internet marketing intelligence through online analytical web usage mining. *ACM SIGMOD Record* 1998;27:54–61.
- [4] Lin S, Chen M, Ho J, Huang Y. ACIRD: Intelligent Internet document organization and retrieval. *IEEE Transactions on Knowledge and Data Engineering* 2002;14:599–614.
- [5] Srivastava J, Cooley R, Deshpande, Tan P. Web usage mining: discovery and applications of usage patterns from web data. *ACM SIGKDD Explorations Newsletter* 2000;1:12–23.
- [6] Linden G, Smith B, York J. Amazon.com recommendations: item-to-item collaborative filtering. *IEEE Internet Computing* 2003;7:76–80.
- [7] Miller B, Albert I, Lam S, Konstan J, Riedl J. MovieLens unplugged: experiences with an occasionally connected recommender systems. *Intelligent User Interfaces* 2002;3:263–266.
- [8] Mooney R, Roy L. Content-based book recommending using learning for text categorization. *Digital Libraries* 2000;5:195–204.
- [9] O'Sullivan D, Smyth B, Wilson D, McDonald K, Smeaton A. Improving the quality of personalized electronic program guide, User modeling User-Adapted Interact 2004;14:5–35.
- [10] Nanopoulos A, Katsaros D, Manolopoulos Y. Effective prediction of web-user accesses: a data mining approach. Conference on Mining Log Data Across All Customer Touchpoints (WebKDD'01), San Francisco 2001.
- [11] Agrawal R, Srikant R. Mining sequential patterns. Eleventh International Conference on Data Engineering 1995;95:3-14.
- [12] Deshpande M, Karypis G. Selective Markov models for predicting web-page accesses. *ACM Transactions on Internet Technology (TOIT)* 2004;4:163-184.
- [13] Mobasher B, Dai H, Luo T, Nakagawa M. Discovery and evaluation of aggregate usage profiles for web personalization. *Data Mining and Knowledge Discovery* 2002;6:61–82.
- [14] Demir G, Uyar A, Gunduz S. Multiobjective evolutionary clustering of Web user sessions: a case study in Web page recommendation. *Soft Computing - A Fusion of Foundations, Methodologies and Applications* 2010;14:579-597.
- [15] Nakagawa M, Mobasher B. A hybrid web personalization model based on site connectivity. *WEBKDD Workshop (Washington, DC)* 2003;5:59-70.
- [16] Mobasher B, Dai H, Luo T, Nakagawa M. Effective personalization based on association rule discovery from web usage data. *Web information and data management* 2001;3: 9–15.
- [17] Liu H, Keselj V. Combined mining of Web server logs and web contents for classifying user navigation patterns and predicting users' future requests. *Data & Knowledge Engineering* 2007;61:304–330.
- [18] Zadeh L. The concept of a linguistic variable and its application to approximate reasoning. *Information Science* 1975;8:199-249.
- [19] Zadeh L. The concept of a linguistic variable and its application to approximate reasoning. *Information Science* 1975;8:301-357.

- [20] Hu Y, Hu J, Chen R, Tzeng G. Assessing weights of product attributes from fuzzy knowledge in a dynamic environment. *European Journal of Operational Research* 2004;154:125-143.
- [21] HAN J, PEI J, YIN Y, MAO R. Mining frequent patterns without candidate generation: A frequent-pattern tree approach. *Data Mining and Knowledge Discovery* 2004;8:53-87.