

# A Scalable Behavior Model for Temporal Prediction of Web User Access Sequences

Enrique Frías-Martínez and Vijay Karamcheti

Computer Science Department, Courant Institute of Mathematical Sciences  
New York University  
715 Broadway, New York, NY 10003 USA  
{frias, vijayk}@cs.nyu.edu

**Abstract.** One of the important Internet challenges in coming years will be the introduction of intelligent services and the creation of a more personalized environment for users. Analysis of Web server logs has been used in recent years to model the behavior of web users in order to provide intelligent services. In this paper we propose a model for predicting sequences of user accesses which is distinguished by two elements: scalability and sequentiality. Scalability, in this context, means that the proposed model can be adapted to the characteristics of the server to more efficiently capture its behavior. The concept of sequentiality in our model consists of three elements: (1) preservation of the sequence of the click stream in the antecedent, (2) preservation of the sequence of the click stream in the consequent and (3) a measure of the time gap between the antecedent and the consequent using the number of user clicks.

## 1 Introduction

Since Etzioni [12] first proposed the term “Web Mining” a lot of research has been done in this area. One topic that has received a lot of attention is modeling the behavior of web users, in other words, being able to predict the requests of users. Having a model of user behavior has made possible the implementation of a great variety of intelligent services. Examples of these services include: redesigning of web sites [24], personalization of e-commerce sites [27], recommendation of pages [18], construction of web pages in real-time [23], adaptation of web pages for wireless devices [4], improvement of web search engines, and prefetching [11][19][30][31].

The main techniques traditionally used for modeling user’s patterns are clustering and association rules. These two approaches produce systems which lack two important characteristics of web user access: sequentiality and temporality. In this context sequentiality implies reflecting the order of the requests of the user, and temporality refers to being able to capture when the predicted actions are actually going to happen.

In this paper we present a model that constructs sequential association rules to capture the sequentiality and temporality in which web pages are visited.

In order to do this, the rules constructed by the model preserve (1) the sequence of the click stream of the antecedent and (2) the sequence of the click stream of the con-

sequent. There is a third significant element in our model: the rules also reflect the distance between the antecedent and the consequent measured by the number of user clicks to go from one to another.

The concept of distance between the antecedent and the consequent is very important for the prediction system because it allows the rules to express not only what pages are going to be accessed but also precisely when they are going to be accessed. This is especially useful for prefetching applications or for recommendation systems. Additionally, the concept of distance can be used as a measure of the rules's quality. For example if we want to redesign web pages for wireless devices by finding shortcuts, the distance of a rule can be used to measure how useful that shortcut is.

Up to now, the algorithms developed for prediction have not taken into account the characteristics of the specific web server they are trying to model. Our model introduces the property of scalability. This means that the prediction system can be adapted, depending on the size of the server, in order to more efficiently capture the behavior of its users. The model offers a balance between the storage space needed and the efficiency of the prediction system. The balance will be mainly determined by the application that is going to be developed using the prediction system.

The proposed model can be used as a black box for any of the applications previously mentioned.

The organization of the paper is as follows. Section 2 summarizes the motivation and prior work done in this area. Section 3 presents the Scalable Sequential Behavior Model. Section 4 implements some examples of the Scalable Sequential Behavior Model and analyzes the results. Finally, in Section 5 we conclude the article with the conclusions and a discussion of future work.

## **2 Motivation and Related Work**

The main techniques used for pattern discovery are clustering and association rules [25].

Clustering, applied in the context of web mining, is a technique that makes it possible to group similar browsing patterns or to divide the web pages of a site into groups that are accessed together. This information can be used in the recommendation process of a page [18] or by search engines to display related pages along with their results. Also, in this context, clustering has a distinguishable characteristic: it is done with non-numerical data. This implies that, usually, the clustering techniques applied are relational. This means that we have numerical values representing the degrees to which two objects of the data set are related. Some examples of relational clustering applied to web mining are [18] and [14]. Several authors have also considered the inherent fuzziness of the data presented in the web mining problem and have developed relational fuzzy clustering algorithms [15].

Association rule discovery aims at discovering all frequent patterns among transactions. The problem was originally introduced by Agrawal et al. [1] and is based on detecting frequent itemsets in a market basket. In the context of web usage mining, association rules refer to sets of pages that are accessed together. Usually these rules

should have a minimum support and confidence to be valid. The Apriori algorithm [1] is widely accepted to solve this problem. Association rules can be used to re-structure a web site [24], to find shortcuts, an application especially useful for wireless devices [4], or to prefetch web pages to reduce the final latency [11].

The data used to obtain frequent patterns in a web mining problem has a very important characteristic: it is sequential. The user accesses a set of pages in a given order and it is very important to capture this order in the final model obtained. Unfortunately, the two previous methods lack any kind of representation of this order. Clustering identifies groups of pages that are accessed together without storing any information about the sequence. Association rules indicate groups that are presented together.

Some authors have already dealt with the problem of capturing sequentiality in association rules for web mining. The approach taken in [9] considers sequences of each session to produce rules. [21] presents the PPM algorithm, which also preserves the order of access and basically uses a Markov prediction model. The main limitation of most of these approaches is that those algorithms only detect patterns that correspond to consecutive sequences.

In this paper we present a model that is able to detect patterns produced by non consecutive sequences and to preserve the order in which those web pages are visited. The model expresses those patterns using rules.

This ability to detect patterns constructed with non consecutive sequences introduces the possibility of measuring the distance between the antecedent and the consequent of a rule. Some algorithms, like [19], are designed to detect non consecutive sequences, but there is no indication of the distance between them. In our model the distance between the antecedent and the consequent is measured in terms of the number of user clicks to go from one to the other. To date no model deals with the concept of distance between the antecedent and the consequent of a rule. This concept is very important for any application (such as a recommendation system) that attempts to infer characteristics of a web site because it provides information about when the pages are going to be visited. This concept is different from finding association rules that have explicit temporal information, as done in [3], or from looking for rules that give temporal relations between different sessions of the same user, as done in [17]. Our method gives a temporal relation within the same session between the antecedent and the consequent by measuring the distance between them.

The models and algorithms developed up to date for user access prediction ([2],[4],[26], etc.) apply the same approach to all servers regardless of their size. In our approach, we consider that the same model can not be applied to any kind of servers. In other words, that the characteristics of the server play an important role in efficiently capturing user behavior. In order to capture user's behavior, the model we propose is scalable. This means that it can be adapted to the inherent characteristics (number of web pages, number of users, architecture of the site, etc.) of each server. The scalability introduced in the model makes it possible to obtain a trade-off between the number of rules and the prediction accuracy of the prediction system. We also give indications of how to obtain this balance for different types of servers.

### 3 Scalable Sequential Behavior Model

In this section we begin by formalizing the preparation of the data. Then, we present the Clustering of Users, which is used to introduce scalability. Next, the concept of Sequential Association Rule is introduced and the definition of Scalable Sequential Behavior Model is given.

#### 3.1 Preparing the Data

The syntax of the log file that contains all requests that a site has processed is specified in the CERN Common Log Format [7]. Basically an entry consists of (1) the user's IP address, (2) the remote logname of the user, (3) the access date and time, (4) the request method (GET, POST ...), (5) the URL of the page accessed, (6) the protocol (HTTP 1.0, HTTP 1.1,...), (7) the return code and (8) the number of bytes transmitted. This set of logs contains enough information to reveal the set of sessions served by the site. The W3C Web Characterization Activity (WCA) [29] defines a user session as the click-stream of page views for a single user across the entire Web. In our context, the information we really want to obtain is a server session, defined as the click-stream in a user session for a particular web server. A click-stream is defined as a sequential series of page view requests of a user. Also, the W3C defines a user as a single individual that is accessing one or more servers from a browser. Our system defines a compiler that transforms a set of logs  $L$  expressed as,

$$L = \{L_1, \dots, L_{|L|}\} \quad (1)$$

$$L_i = (IP_i, LOGNAME_i, TIME_i, METHOD_i, URL_i, PROT_i, CODE_i, BYTES_i), \forall i / i = 1 \dots |L| .$$

Into a set of sessions  $S$ ,

$$S = \{S_1, \dots, S_{|S|}\}, \quad (2)$$

where  $|L|$  is the number of logs of  $L$  and  $|S|$  is the number of sessions of  $S$ . Each session is defined by a tuple  $((IP, LOGNAME), PAGES)$ , where  $(IP, LOGNAME)$  identifies the user of the session and  $PAGES$  the set of pages requested,

$$S_i = ((IP_i, LOGNAME_i), PAGES_i) \text{ with } PAGES_i = \{url_{i,1}, \dots, url_{i,p_i}\}, i = 1 \dots |S|, \quad (3)$$

where  $p_i$  is the number of pages requested by user  $(IP_i, LOGNAME_i)$  in session  $S_i$ . The set of URLs that form a session satisfy the requirement that the time elapsed between two consecutive requests is smaller than  $\Delta t$ . The value we have used is 30 minutes, based on the results of [6] and [25]. Figure 1 presents the algorithm of the compiler. The filters implemented by our algorithm delete all entry logs that do not refer to a URL or that indicate an error. Also, sessions of length one or sessions three times as long as the average length of the set of sessions  $S$  are erased. This is done to eliminate the noise that random accesses or search engines would introduce to the model.

```

Input :  $L, \Delta t, |L|$ 
Output :  $S, |S|$ 
function Compiler( $L, \Delta t, |L|$ )
  for each  $L_i$  of  $L$ 
    if  $METHOD_i$  is  $GET$  and  $URL_i$  is  $WEB\_PAGE$  then
      if  $\exists S_k \in OPEN\_SESSIONS$  with  $IP_k = IP_i$  And  $LOGNAME_k = LOGNAME_i$  then
        if  $(TIME_i - END\_TIME(S_k)) < \Delta t$  then
           $S_k = ((IP_k, LOGNAME_k), PAGES_k \cup URL_i)$ 
        else
          CLOSE SESSION ( $S_k$ )
          OPEN NEW SESSION( $(IP_i, LOGNAME_i), (URL_i)$ )
        end if
      else
        OPEN NEW SESSION( $(IP_i, LOGNAME_i), (URL_i)$ )
      end if
    end for
end for

```

**Fig. 1.** Algorithm of the compiler implemented.

Finally, the average length of the set of sessions  $S$  is defined as,

$$N = \frac{\sum_{i=1}^{|S|} p_i}{|S|}. \quad (4)$$

### 3.2 Clustering the Users

The set of different users of a system is expressed by,

$$USERS = \{(IP_1, LOGNAME_1), \dots, (IP_1, LOGNAME_{1,h(1)}), \dots, (IP_t, LOGNAME_{t,1}), \dots, (IP_t, LOGNAME_{t,h(t)})\} \quad (5)$$

Where  $t$  is the number of different  $IPs$  of  $S$ , and  $h(i)$   $i=1, \dots, t$  is the number of different users of each  $IP$ . The total number of users of the system can be expressed as:

$$\sum_{k=1}^t h(k) \quad (6)$$

The set of different users is going to be clustered according to a cluster policy  $P$  in order to create a scalable model.

The purpose of these clusters is to group the users that have the same behavior and to allow a trade-off between the size and the personalization capabilities provided by the model. Given  $p$  the number of clusters defined by  $P$ , the set of clusters of users  $CU$  can be expressed as:

$$\begin{aligned}
P &= \{P_1, \dots, P_p\} & (7) \\
CU &= \{CU_1, \dots, CU_p\} \\
CU_i &= \{(IP_s, LOGNAME_{s,d}), \dots, (IP_e, LOGNAME_{e,r})\}, i=1, \dots, p \text{ and } 1 \leq s \leq e \leq t \\
\text{where } (IP_l, LOGNAME_{l,m}) \in CU_i &\Leftrightarrow (IP_l, LOGNAME_{l,m}) \text{ is } P_i \\
CU_i \cap CU_k &= \emptyset, \forall i, k \text{ with } i \neq k \\
\bigcup_{i=1}^p CU_i &= USERS
\end{aligned}$$

For example, given the web server of a university department a possible clustering policy is  $P = \{Professor, Graduate\_Students, Students, Others\}$ .

The classification of users is going to be used to cluster the set of sessions  $S$ . The set of clustered sessions  $CS$ , can be expressed as:

$$\begin{aligned}
CS &= \{CS_1, \dots, CS_p\} & (8) \\
CS_i &= \bigcup_k (IP_k, LOGNAME_k) \in CU_i, \forall i=1, \dots, p, \forall k=1, \dots, |s| \\
\bigcup_{i=1}^p CS_i &= S \\
CS_i \cap CS_j &= \emptyset, \forall i, j \text{ with } i \neq j
\end{aligned}$$

Each  $CS_i$  groups the sessions of the users that are part of the same  $CU_i$  cluster.

### 3.3 Sequential Association Rules

The concept of Sequential Association Rule (SAR) is based on the N-Gram concept. In the context of web mining, an N-Gram of a session  $S_i$  is defined as any subset of N consecutive URLs of that session.

A Sequential Association Rule (SAR), given  $|A|$  the length of the sequence of URLs of the antecedent,  $|C|$  the length of the sequence of URLs of the consequent, and  $n$  the distance between the antecedent and the consequent, is defined as:

$$\begin{aligned}
A &\rightarrow C & (9) \\
A &:= url_j, \dots, url_{|A|+j} \\
C &:= url_l, \dots, url_{|C|+l} \\
l &= |A| + j + n + 1, \quad n \in \mathbb{N}^+
\end{aligned}$$

A SAR expresses the following relation: if the last click stream of length  $|A|$  of a user is  $A$  then in  $n$  clicks the set of URLs  $C$  will be requested.

Each SAR is constructed from an N-Gram obtained from a session. This means that for the SAR of the definition some session  $S_k$  of a given  $CS_i$  contains an N-Gram, with  $N = |A| + |C| + n$ , that satisfies,

$$S_k = (\dots, url_{k,j}, \dots, url_{k,j+|A|}, url_{k,j+|A|+1}, \dots, url_{k,j+|A|+n}, url_{k,l}, \dots, url_{k,l+|C|}, \dots), \quad (10)$$

with  $l = j + |A| + n + 1$ .

Each SAR has a degree of support and confidence associated with it. The support of a rule is defined as the fraction of strings in the set of sessions of  $CS_i$  where the rule successfully applies. The support of a rule is given by,

$$\theta(A?_1 \dots ?_n C) = \frac{|S_k \in CS_i / (A?_1 \dots ?_n C) \in S_k|}{|CS_i|}. \quad (11)$$

Where  $?_1 \dots ?_n$  represents the set of any  $n$  pages in the session, and  $A?_1 \dots ?_n C \in S_k$  is defined as an N-Gram of  $S_k$ . This is the way to model the distance between the antecedent and the consequent when obtaining the support. The confidence of a rule is defined as the fraction of times for which if the antecedent  $A$  is satisfied, the consequent  $C$  is also true in  $n$  clicks. The confidence of the rule is defined as,

$$\sigma(A?_1 \dots ?_n C) = \frac{\theta(A?_1 \dots ?_n C)}{\theta(A)}. \quad (12)$$

$SR(CS_i)_{|A|,n,|C|}$ , for the set of sessions grouped by  $CS_i$ , is defined as a set of tuples  $(SAR, Counter)$ , where each tuple has a SAR with an antecedent of length  $|A|$ , a consequent with length  $|C|$  and a distance  $n$  between the antecedent and the consequent. The *Counter* of each tuple indicates the number of times that the correspondent rule occurs in  $S$ . Figure 2 shows the algorithm used to obtain  $SR(CS_i)_{|A|,n,|C|}$ .

```

Input :  $|A|, n, |C|, CS_i$ 
Output :  $SR(CS_i)_{|A|,n,|C|}$ 
function Obtain_SR( $|A|, n, |C|, CS_i$ )
   $SR(CS_i)_{|A|,n,|C|} = \emptyset$ 
  for each  $S_k$  of  $CS_i$ 
    for  $j = 1$  to  $p_k$ 
      if  $j + |A| + n + |C| \leq p_k$ 
         $SAR = url_{k,j}, \dots, url_{k,j+|A|} \xrightarrow{n} url_{k,j+|A|+n}, \dots, url_{k,j+|A|+n+|C|}$ 
        if  $(SAR, ?) \in SR(CS_i)_{|A|,n,|C|}$  then
           $Counter((SAR, ?)) = Counter((SAR, ?)) + 1$ 
        else
           $SR(CS_i)_{|A|,n,|C|} = SR(CS_i)_{|A|,n,|C|} \cup (SAR, 1)$ 
        end if
      end if
    end for
  end for

```

**Fig. 2.** Algorithm to obtain  $SR(CS_i)_{|A|,n,|C|}$

In order to preserve the relevant information, only those SARs which have support and confidence bigger than a given threshold are considered. With  $\theta'$  the threshold of the support and  $\sigma'$  the threshold of the confidence, we will talk about the set of rules  $SR(CS_i)_{|A|,n,|C|, \theta', \sigma'}$  as the rules of  $SR(CS_i)_{|A|,n,|C|}$  with a support bigger than  $\theta'$  and a confidence bigger than  $\sigma'$ .

$SR(CS_i)_{|A|,n,|C|,\theta',\sigma'}$  is defined as a set of elements  $(SAR, \theta_{SAR}, \sigma_{SAR})$ , where  $SAR$  is a Sequential Association Rule, and  $\theta_{SAR}$  and  $\sigma_{SAR}$  the support and confidence associated with it, satisfying  $\theta_{SAR} > \theta'$  and  $\sigma_{SAR} > \sigma'$ .  $SR(CS_i)_{|A|,n,|C|}$  contains enough information to obtain  $\theta_{SAR}$  and  $\sigma_{SAR}$ . The set of  $SARs$  of  $SR(CS_i)_{|A|,n,|C|,\theta',\sigma'}$  is a subset of  $SR(CS_i)_{|A|,n,|C|}$ .

In order to optimize the storage and access to the rules that define  $SR(CS_i)_{|A|,n,|C|,\theta',\sigma'}$ , we group the rules with the same antecedent, storing for each set of rules the consequent and the degree of support and confidence associated with it. The structure of  $SR(CS_i)_{|A|,n,|C|,\theta',\sigma'}$  is:

$$SR(CS_i)_{|A|,n,|C|,\theta',\sigma'}(A) = \begin{cases} ((C_{1,1}, \theta_{1,1}, \sigma_{1,1}), \dots, (C_{1,l_1}, \theta_{1,l_1}, \sigma_{1,l_1})), & \text{if } A = A_1 \\ \dots \\ ((C_{k,1}, \theta_{k,1}, \sigma_{k,1}), \dots, (C_{k,l_k}, \theta_{k,l_k}, \sigma_{k,l_k})), & \text{if } A = A_k \\ \emptyset, & \text{otherwise} \end{cases} \quad (13)$$

Where  $A$  is a click stream of length  $|A|$ ,  $A_i, i=1\dots k$ , with  $|A_i|=|A|$ , is the set of antecedents of the rules of  $SR(CS_i)_{|A|,n,|C|,\theta',\sigma'}$ ,  $C_{i,j}, i=1\dots k, j=1\dots l_k$ , with  $|C_{i,j}|=|C|$ , the set of consequents for each antecedent  $A_i$ , and  $l_k$  the number of consequents of each antecedent.

### 3.3 Sequential Behavior Model

A Sequential Behavior Model is defined by a tuple  $\{RU, \Phi\}$  where  $RU$  is a set of rules and  $\Phi$  is the decision policy function.

$RU$  is defined as:

$$RU((IP, LOGNAME), A) = \begin{cases} SR(CS_1)_{|A|,n,|C|,\theta',\sigma'}(A), & \text{if } (IP, LOGNAME) \in CU_1 \\ \dots \\ SR(CS_p)_{|A|,n,|C|,\theta',\sigma'}(A), & \text{if } (IP, LOGNAME) \in CU_p \end{cases} \quad (14)$$

The function  $\Phi$  is defined as:

$$\Phi((C_{i,1}, \theta_{i,1}, \sigma_{i,1}), \dots, (C_{i,l_i}, \theta_{i,l_i}, \sigma_{i,l_i})) = \{C_{i,j}, \dots, C_{i,b}\} \\ \text{with } 1 \leq i \leq k, 1 \leq j \leq l_k, 1 \leq b \leq l_k \quad (15)$$

Where the independent variable of  $\Phi$  is the set of consequents obtained from  $RU$  for a given user and a given antecedent (click stream), and the dependent variable is the consequent or set of consequents predicted. Some examples of policy functions include,

$$\Phi((C_{i,1}, \theta_{i,1}, \sigma_{i,1}), \dots, (C_{i,l_i}, \theta_{i,l_i}, \sigma_{i,l_i})) = \{C_{i,j}\} \\ \text{with } \sigma_{i,j} \geq \sigma_{i,m} \forall m / m = 1\dots l_i, m \neq j \quad (16)$$

$$\Phi((C_{i,1}, \theta_{i,1}, \sigma_{i,1}), \dots, (C_{i,l_i}, \theta_{i,l_i}, \sigma_{i,l_i})) = \{C_{i,j}\} \quad (17)$$

with  $\theta_{i,j} \geq \theta_{i,m} \forall m / m = 1 \dots l_i, m \neq j$

$$\Phi((C_{i,1}, \theta_{i,1}, \sigma_{i,1}), \dots, (C_{i,l_i}, \theta_{i,l_i}, \sigma_{i,l_i})) = \{C_{i,j}, C_{i,b}\} \quad (18)$$

with  $\sigma_{i,j} \geq \sigma_{i,b} \geq \sigma_{i,m} \forall m / m = 1 \dots l_i, m \neq j, m \neq b$

The first example predicts the consequent with the biggest confidence, the second one predicts the one with the biggest support, and the third one picks the two consequents with the highest confidence. The policy function can be defined depending on the specific application and on the characteristics of the web log used.

The tuple  $(RU, \Phi)$  defines the on-line execution of the prediction system. Given  $A$  the click stream of the last  $|A|$  pages requested by the user  $(IP, LOGNAME)$ , the set of predicted pages will be given by,

$$\Phi(RU((IP, LOGNAME), A)) \quad (19)$$

### 3.4 Profiles of the Sequential Behavior Model

The web servers found in the Internet have very different sets of characteristics. These characteristics include: the number of users; the number of web pages that the server has; the architecture of the web server and the number of links per page.

The prediction systems developed to date ([2],[4],[11],[26], etc.) do not consider the different characteristics of the web servers in order to more efficiently model their behavior. Nevertheless, the Model proposed in this paper considers that the size and the characteristics of the server are essential to capturing the behavior of their users. In other words, the same model cannot be used to generate a prediction system for a university department server and to capture the behavior of an e-commerce site.

The proposed Sequential Behavior Model is scalable because it can be adapted to the characteristics of the server that is going to be modeled. The Model can be presented for different profiles: global, personalized and cluster.

#### 3.4.1 Global Sequential Behavior Model (GSBM)

The Global Sequential Behavior Model considers only one cluster of users. This means that the system will also consider just one cluster of the set of sessions. Formally,

$$\begin{aligned} p &= 1 & (20) \\ CU &= \{CU_1\} \\ P &= \{P_1\}; P_1 = \text{"Any user of the system"} \\ CS &= \{CS_1\} \end{aligned}$$

The Global Sequential Behavior Model will be defined by the tuple  $(RU, \Phi)$ , where  $RU$  is expressed as:

$$RU((IP, LOGNAME), A) = \{SR(CS_1)_{|A_i, n_i, C_i, \theta_i, \sigma_i}(A), \text{if } (IP, LOGNAME) \in CU_1\} \quad (21)$$

The Global Model is small because it has only one set of rules, but also has a limited personalization capability.

### 3.4.2 Clustered Sequential Behavior Model (CSBM)

The Clustered Model considers a set of  $p$  clusters. These clusters should be designed to group the set of users that have a common behavior. Each set of users will end up having a set of rules that describe their behavior in  $RU$ .

This profile produces bigger models than the Global approximation but also increases the personalization capabilities. The size and the prediction efficiency will depend on the number of clusters that the system has defined.

### 3.4.3 Personalized Sequential Behavior Model (PSBM)

The Personalized Model is a particularization of the Clustered Model. In this case each one of the users of a system will have its own cluster, which means that each user will have its own set of rules to describe his behavior. The Model produced by this profile is bigger than the other two but it also has the maximum possible personalization capability.

This profile is very interesting considering that an increasing number of sites identify its users (using password or cookies) to provide them with personalized services. Having a model that describes the behavior of each user will make it possible to provide a higher degree of adaptation and personalization.

The main criticism that arises from this idea is that too much space is needed to store each user's set of personalized sequential association rules. Nevertheless, e-sites already possess a lot of information about each user, ranging from name, address, or credit cards numbers to layouts and preferences, as can be seen in [5],[16] and [28]. The inclusion of a personal set of association rules, as will be shown, will not cause a large increase in the amount of needed storage.

It is not necessary for all the users of a server that implements a Personalized Sequential Behavior Model to have a personal set of rules. A more useful approach considers that, given the set of users of a server, a smaller subset is typically responsible for the largest part of the system's load. Each member of this set of frequent users will have a personal cluster. The non-frequent users, will be grouped in another cluster and will have their own set of rules. Formally, being  $g$  the number of frequent users of the system,

$$\begin{aligned} p &= g + 1 & (22) \\ CU &= \{CU_1, \dots, CU_g, CU_{g+1}\} \\ P &= \{P_1, \dots, P_g, P_{g+1}\} \\ CS &= \{CS_1, \dots, CS_g, CS_{g+1}\} \end{aligned}$$

The number of clusters is  $g+1$ , where the last cluster groups the non-frequent users of the system. The clustering policy is defined as  $P_1=$ Frequent User #1, ...,  $P_g=$ Frequent User #g,  $P_{g+1}=$ All non-frequent users.  $CS_j$  groups the set of sessions of the frequent user #1 and  $CS_{g+1}$  groups all the sessions of the non-frequent users.

The Personalized Sequential Behavior Model will be defined by the tuple  $(RU, \Phi)$ , where RU is expressed as:

$$RU((IP, LOGNAME), A) = \begin{cases} SR(CS_1)_{|A|,n|C|,\theta',\sigma'}(A), \text{ if } (IP, LOGNAME) \text{ is Frequent\_User\_}\#1 & (23) \\ \dots \\ SR(CS_g)_{|A|,n|C|,\theta',\sigma'}(A), \text{ if } (IP, LOGNAME) \text{ is Frequent\_User\_}\#g \\ SR(CS_{g+1})_{|A|,n|C|,\theta',\sigma'}(A), \text{ if } (IP, LOGNAME) \in CU_{g+1} \end{cases}$$

#### 3.4.4 General Rules for Application of the Profiles

The different profiles give the designer of the prediction system the capability of choosing a trade-off between the memory needed by the set of rules and the personalization capabilities provided. Each one of the profiles is useful in different contexts:

- The Global Sequential Behavior Model is able to efficiently capture the user behavior of a simple site. We consider a simple site as one that has a small number of users and/or a small number of web pages and/or a simple architecture (normally a tree architecture). A typical example of this kind of site is the web server of a university department.

- The Clustered Sequential Behavior Model is capable of modeling the behavior of more complex systems. A complex system is defined as a system with a high number of users and/or a high number of web pages and/or a highly interconnected structure. The clustering policy  $P$  will have to be defined according to the intelligent service that is going to be implemented with the Behavior Model.

- The Personalized Sequential Behavior Model is recommended for highly complex sites and/or for sites that need individual information for each user.

The next section presents the implementation of different Sequential Behavior Models for different web logs.

## 4. Examples of Sequential Behavior Models

### 4.1 Characteristics of the logs used

We have selected three sets of logs in order to cover different types of servers, ranging from small sites with few users to complex commercial sites with a large number of users.

The first log is from the Computer Science Department (CS) of the Polytechnic University of Madrid [10]. The training set is for September 2001. The test set has been defined using log data for 1<sup>st</sup> October 2001. This is an example of a small site, with a simple tree architecture and a small number of visitors.

The second log is from the NASA Kennedy Space Center server [20]. It contains 3,461,612 requests collected over the months of July and August 1995. This is an example of a medium site server, with a complex architecture and a medium number of users. The training set considered was July 1995, and the test set has been defined using log data for the 1<sup>st</sup> of August 1995.

The third site is ClarkNet [8], a commercial Internet site provider, which contains 3,328,587 requests over a period of two weeks. This is an example of a large commercial site, with a highly complex architecture and a high number of visitors. The training set has been defined from 28 August 1995 to 9 Sep 1995, and the test set is the log data for 10<sup>th</sup> Sep 1995.

These training sets are the inputs to the algorithms and filters presented in Section 3.1. Table 1 presents some characteristics of the training sets including the processing time of the algorithm presented in Table 1 for each log. Although the filtering process eliminates a lot of sessions, we keep the relevant part of the requests. The processing time is given for an implementation in LISP of the compiler, running with Linux, on a Pentium III 450MHz.

**Table 1.** Training Log Characteristics.

	CS	NASA	CLARKNET
Size	27M	160M	308.6M
Dates	Sep. 2001	Jul. 1995	28 Aug.-9 Sep. 1995
Processing Time	12 min	4 h 35 min	8 h 50 min
# of sessions before filtering	3,499	124,666	224,935
# of sessions after filtering	839	57,875	83,011
# of requests before filtering	6,244	352,844	511,536
# of requests after filtering	3,504	215,223	283,844

The characteristics of the test logs defined for each one of the training sets are given in Table 2. The processing time indicates the time needed to obtain the set of sessions from the logs using the algorithm presented in Figure 1. The number of sessions and number of pages requested shown in Table 2 are post-filtering.

**Table 2.** Test Log Characteristics.

	CS	NASA	CLARKNET
Size	170K	6.8M	19M
Date	Oct. 1 <sup>st</sup> 2001	Aug. 1 <sup>st</sup> 1995	Sep. 10 <sup>th</sup> 1995
# of sessions	53	2753	5725
# of Pages requested	138	9521	18233
Average length of Session	2.6	3.4	3.9
Processing Time	25 sec.	2 min 40 sec	6 min 10 sec

## 4.2 Implementation of Global Sequential Behavior Models

We have obtained the Global Sequential Behavior Model of each log defined by  $SR(CS_l)_{l,l,l}$ . After that, a threshold of 1% for the support and of 5% for the confidence has been applied, obtaining  $SR(CS_l)_{l,l,l,1,5}$ . Other values of support and confidence have been tested, but the best prediction rate is obtained with 1% and 5%. The characteristics of these  $SR$ s are presented in Table 3. The processing time indicates the time needed to process each one of the training sets using the algorithm presented in Figure 2.

**Table 3.** Characteristics of the SR obtained.

	CS	NASA	CLARKNET
Processing Time	50 sec	3 min 20 sec	7 min 40 sec
# of SAR of SR before applying thresholds	392	15,644	49,245
# of SAR of SR after applying thresholds	94	116	166

Figure 3 presents the accuracy of the prediction system for CS, NASA and Clark-Net logs, with  $RU$  defined by  $SR(CS_l)_{l,l,l,1,5}$  and different  $\Phi$  functions. The results are obtained comparing the actual next page of the session with the page or set of pages predicted by the system. The first set of columns represents the results for the function  $\Phi$  defined as:

$$\Phi_1((C_{i,1}, \theta_{i,1}, \sigma_{i,1}), \dots, (C_{i,l_i}, \theta_{i,l_i}, \sigma_{i,l_i})) = \{C_{i,j}\} \quad (24)$$

*with  $\sigma_{i,j} \geq \sigma_{i,m} \forall m / m = 1 \dots l_i, m \neq j$*

In other words, only the consequent with the highest confidence is given as the prediction. The rest of the functions  $\Phi$  defined for each set of columns are, in order:

$$\Phi_2((C_{i,1}, \theta_{i,1}, \sigma_{i,1}), \dots, (C_{i,l_i}, \theta_{i,l_i}, \sigma_{i,l_i})) = \{C_{i,j}, C_{i,b}\} \quad (25)$$

*with  $\sigma_{i,j} \geq \sigma_{i,b} \geq \sigma_{i,m} \forall m / m = 1 \dots l_i, m \neq j, m \neq b$*

$$\Phi_3((C_{i,1}, \theta_{i,1}, \sigma_{i,1}), \dots, (C_{i,l_i}, \theta_{i,l_i}, \sigma_{i,l_i})) = \{C_{i,j}, C_{i,b}, C_{i,h}\} \quad (26)$$

*with  $\sigma_{i,j} \geq \sigma_{i,b} \geq \sigma_{i,h} \geq \sigma_{i,m} \forall m / m = 1 \dots l_i, m \neq j, m \neq b, m \neq h$*

$$\Phi_4((C_{i,1}, \theta_{i,1}, \sigma_{i,1}), \dots, (C_{i,l_i}, \theta_{i,l_i}, \sigma_{i,l_i})) = \{C_{i,1}, \dots, C_{i,l_i}\} \quad (27)$$

Function  $\Phi_2$  gives as the set of predicted pages the two consequents with higher confidence,  $\Phi_3$  the set of three consequents with higher confidence, and  $\Phi_4$  all the predictions that have a support and a confidence bigger than the thresholds used.

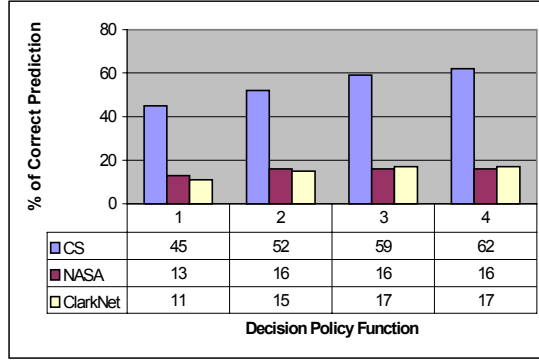


Fig. 3. Results of the model for each set of logs.

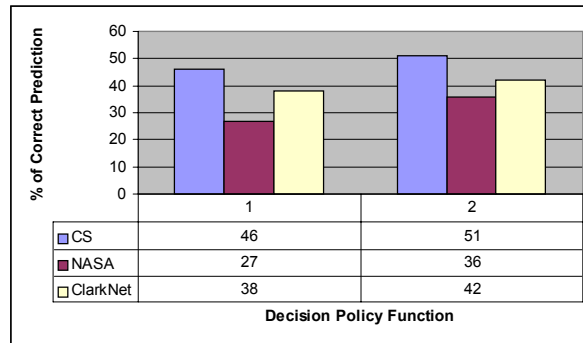
### 4.3 Implementation of Clustered Sequential Behavior Models based on IP

In this section we present a CSBM where the clustering policy used is defined using the IP address. In this case, the set of users that have the same 16 first bits of their IP equal are going to be part of the same cluster. The objective of this policy is to cluster the users that share a common behavior, and this is based on the fact that the users that come from the same subnet have a greater probability of having similar behavior than two users that have completely different IPs. Table 4 presents the characteristics of the clusters obtained for each set of logs.

Table 4. Characteristics of the clusters of each log.

	CS	NASA	CLARKNET
# of Clusters	236	6,667	8,510
# of different IPs	840	37,821	58,035
Processing Time to obtain the Clusters	20 sec.	2 min. 10 sec.	3 min 40 sec.
Average # of IP per cluster	3.5	5.6	6.8
Processing Time to obtain $SR(CS)_{1,1,1,1,5}, i=1, \dots, p$	15 sec.	1 min. 05 sec.	1 min. 30 sec.
Average # of SAR per cluster	4.6	6.3	7.2

For each cluster we have obtained  $SR(CS)_{1,1,1,1,5}, i=1, \dots, p$ , with  $p=236$  for the CS log,  $p=6,667$  for the NASA log, and  $p=8,510$  for the ClarkNet log. Figure 4 presents the accuracy of the prediction system for each log using the decision policy functions  $\Phi_1$  and  $\Phi_2$  previously defined.



**Fig. 4.** Results of the CSBM for each log.

In order to generate an efficient CSBM the training log should contain elements of all the possible clusters that the set of visitors of the system may form. This allows the recommendation system to give a recommendation for any possible visitor. If the training log lacks some of these users, the prediction system will not be able to give a recommendation for them (which will count as an incorrect prediction).

In the testing logs defined for CS, NASA and ClarkNet only 2%, 5% and 6% of the sessions correspond to a user that is not included in any cluster. For those users, the system's prediction is obviously wrong. When the percentage of sessions of users that are not part of any of the clusters generated with the training log is significant compared with the total amount of visits, a GSBM can be applied for those users.

The solution in this case is similar to the one presented in the PSBM, where a cluster groups all the possible visitors of the system that are not part of any of the other clusters.

#### 4.5 Implementation of Personalized Sequential Behavior Models

We are going to develop a PSBM for NASA and ClarkNet logs. As presented in section 3.4.3 PSBMs are especially suitable for systems where a core of the users is responsible for the best part of the load. In this case, we define a frequent user as the one that has more than two visits in the training log. Also, because the logs lack from any kind of information about each individual user, we will identify an IP as an user, bearing in mind that a single IP can be used by a group users.

Table 5 shows the characteristics of our training logs.

**Table 5.** Some characteristics of the Training Logs.

	NASA	CLARKNET
# of Sessions	124,666	224,935
# of Users	37,821	58,035
# of Users with just one visit	30,875	49,085

As can be seen in Table 5, in the case of the NASA log, 18% of the users are responsible of 75% of the visits. Similar behavior is also observed for the ClarkNet log, where 15% of the users are responsible for 78% of the visits. These sites possess the

perfect characteristics for the implementation of a PSBM. For each cluster we have obtained  $SR(CS_{i,1,1,1,1,5})$ ,  $i=1, \dots, g+1$  where  $g=6,946$  for the NASA log, and  $g=8,950$  for ClarkNet. Table 6 presents the characteristics of the  $RUs$  constructed.

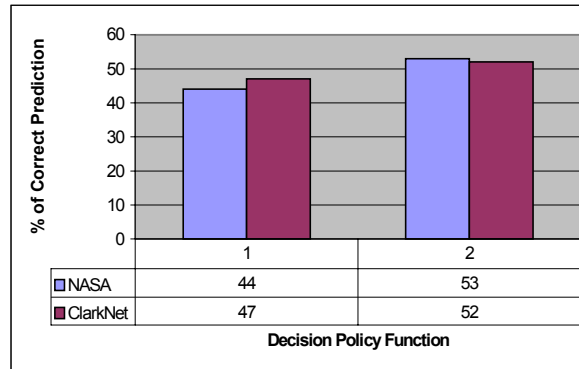
**Table 6.** Characteristics of the Personal SR constructed.

	NASA	CLARKNET
$g$	6946	8950
Average # of rules per $SR(CS_i)$ , $i=1, \dots, g$	10.3	9.7
Total Processing Time	1 min 4 sec	2 min 30 sec

Figure 5 presents the prediction accuracy of the set of PSBM constructed for the NASA log and the ClarkNet log. The test logs in this case have been modified to contain only the set of visits of the users that have a personal  $SR$ . This allows us to obtain the prediction accuracy of the set of personal  $SRs$ ,  $SR(CS_i)$ ,  $i=1, \dots, g$ .

Each set of columns presents the percentage of pages correctly predicted for each test log using two different policy functions.  $\Phi_1$  gives the prediction accuracy using the consequent with the highest confidence. In both cases, NASA and ClarkNet, the system achieves at least a 44% correct prediction.  $\Phi_2$  considers the two consequents with highest confidence, and in that case the correct prediction is over 50%.

The global prediction rate will be given by the prediction rate of  $SR(CS_{g+1})$ , that models the non-frequent users, and by the prediction rate of the set of personal  $SRs$ . In the case of the NASA logs, 75% of the visits are generated by users that have a personal  $SR$ . In the rest of the load, 25%, the prediction accuracy will be given by  $SR(CS_{g+1})$ . Using  $\Phi_1$  as the decision policy function gives a total prediction accuracy of 0.36 ( $0.75*0.44+0.25*0.13$ ). Using  $\Phi_2$  the prediction accuracy goes up to 0.43 ( $0.75*0.53+0.25*0.16$ ). In the case of the ClarkNet log, when using  $\Phi_1$  the final accuracy is 0.33 and when using  $\Phi_2$  it is 0.42.



**Fig. 5.** Prediction Accuracy for the NASA log.

## 4.6 Results Analysis

As can be seen in Figure 3, the percentage of correct prediction using the GSBM for each log is, in the worst case, 45% for CS, 13% for NASA and 15% for ClarkNet.

The result for CS is satisfactory in the sense that the prediction system obtained allows us to implement efficient intelligent services. Nevertheless, the results for NASA and ClarkNet are not satisfactory. The cause of this difference in the prediction system is that the sites are actually very different. In CS we have a site with a well-defined tree architecture, a low number of visitors and a low number of links per page. The other two sites have a large number of pages, a higher number of links per page and a higher number of visitors. More importantly, they have highly interconnected architectures. Based on these results we speculate that in those sites, user behaviors cannot be captured properly with a Global Sequential Behavior Model because access order is not a global property in complex sites. The inability of the GSBM to efficiently capture global behavior in sites with complex structures and high numbers of visitors, suggests that a more local approach should be taken.

The results of the CSBM based on IP (Figure 4) prove that a more local approach makes it possible to obtain a better prediction rate in complex sites. In this case, the NASA prediction rate, in the worst case, goes up to 27% (14% more than the GSBM) and ClarkNet up to 36% (20% more than the GSBM). Nevertheless a more local approach in the case of a simple site, like CS, does not mean necessarily a better prediction rate. As we can see, the prediction rate for CS stays the same. With these results, we speculate that a Global Sequential Behavior Model can efficiently capture the behavior of a site which has clear patterns, which occurs mainly in sites that have a tree architecture and a low number of visitors.

Although the CSBM increases the prediction rate of complex sites, this prediction rate is not necessarily enough for any intelligent application that is going to be implemented using the proposed model. For the cases where the prediction rate of CSBM is not enough, the PSBM will provide the highest prediction rate possible.

The results of the PSBM show the prediction rate for the lowest level approach. In this case, the NASA prediction rate goes up to 44% (29% more than the GSBM and 17% more than the CSBM based on IP), and the ClarkNet goes up to 47% (36% more than the GSBM and 9% more than the CSBM based on IP). And this is for the worst case, because when using more complex decision policy functions the successful prediction rate goes up to 53% for NASA and 52% for ClarkNet. We consider these results completely satisfactory, in the sense that they allow the implementation of any kind of intelligent service.

These results show that PSBM produces for highly complex sites the same successful prediction rates as the GSBM produces for simple sites. The PSBM is able to efficiently capture the behavior of complex highly interconnected sites with an acceptable increase in memory usage. This increase in memory is less significant when one takes into account that most sites already have a lot information about each user for other personalization purposes.

The results also prove that the proposed model can be adapted to the different characteristics of the site that is going to be modeled. In other words, the model is scalable. The results of this section have been used to give the rules of section 3.4.4.

## 5 Conclusions and Future Work

We have considered the problem of modeling web user behavior. For that purpose, a Scalable Sequential Behavior Model has been designed.

One of the distinguishable characteristics of our model is its scalability. The necessity of a scalable model comes from the fact that the set of Internet servers have very different characteristics. Our model is able to adapt itself to the characteristics of the server in order to obtain the maximum possible efficiency. Our Scalable Model allows a trade-off between the number of rules and the prediction accuracy of the Model obtained.

Also, the model is able to capture the inherent sequentiality of web visits. The model is constructed using a set of Sequential Association Rules which reflect the order of the set of URLs of the antecedent and the consequent, and also the distance between the antecedent and the consequent measured in the number  $n$  of clicks between the two click-streams. This distance makes it possible to design systems that not only predict which pages are going to be visited but also when they are going to be visited.

To the best of our knowledge, our model is the only one that is scalable and that incorporates a distance metric in its rules.

The Model has been designed as a black box to be used in any application that is based in a prediction system

A deeper study of the impact of the parameters of the model on the prediction rate is needed. We are especially interested in studying the results depending on the parameter  $n$ , the distance in clicks between the antecedent and the consequent. The possibility of knowing not only what pages are going to be requested but also when are going to be requested can improve the efficiency of a lot of services (prefetching, recommendation systems, etc.). It will also be interesting to study the parameter  $C$ , especially study the optimum value  $|C|$  for a given set of sessions. Having consequents with  $|C| > 1$  is useful for some applications such as the construction of web pages in real time or prefetching.

We plan to apply our prediction model to the NYU HOME page. NYU HOME is a site that allows its users, the NYU community, to personalize their channel contents, ranging from e-mail, to news or weather forecasts. The site is slower than sites serving static content, because the pages have to be generated in real-time. As can be seen, this page has the ideal characteristics for the implementation of a PSBM: it has a large set of frequent users and the system already has personal information about each user. We will decrease the latency time for each user by pregenerating the pages using a PSBM.

## Acknowledgements

This work has been partially supported by the Spanish Department of Research grant EX2001-46775835.

## References

1. Agrawal, R., Imielinski, T., Swami, A.: Mining Association Rules between Sets of Items in Large Databases. In: Proceedings of the 1993 ACM SIGMOD Conference, Washington DC, USA (1993)
2. Albrecht, D., Zukerman, I., Nicholson, A.: Pre-sending documents on the WWW: A comparative study. In: IJCAI99-Proceedings of the Sixteenth International Joint Conference on Artificial Intelligence (1999)
3. Ale, J.M., Rossi, G.H.: An Approach to Discovering Temporal Association Rules. In Proceedings of SAC'00, Marh 19-21, Como, Italy, (2000) 294-300
4. Anderson, C. R., Domingos, P., Weld, D.S.: Adaptive Web Navigation for Wireless Devices. In Proceedings of the Seventeenth International Joint Conference on Artificial Intelligence (IJCAI-01) (2001)
5. Belkin, N.: Helping People Find What They Don't Know. In Communications of the ACM, Vol. 43, No. 8, (2000) 58-61
6. Catledge, L., Pitkow, J.: Characterizing browsing behaviors on the world wide web. In Computer Networks and ISDN Systems, Vol. 27, No. 6 (1995)
7. CERN Common Log Format, <http://www.w3.org/Daemon/User/Config/Logging.html>
8. ClarkNet Internet Provider Log, <http://www.web-caching.com/traces-logs.html>
9. Chen, M.S., Park, J.S., Yu, P.S.: Efficient Data Mining for Path Traversal Patterns. In IEEE Transactions on Knowledge and Data Engineering, Vol. 10, No. 2 (1998) 209-221
10. Departamento de Tecnologia Fotonica Log (Univerisdad Politecnica de Madrid), <http://www.dtf.fi.upm.es>
11. Duchamp, D.: Prefetching Hyperlinks. In: Proc. Second USENIX Symp. on Internet Technologies and Systems, USENIX, Boulder, CO (1999) 127-138
12. Etzioni, O.: The World Wide Web: Quagmire or gold mine. In: Communications of the ACM, Volume 39, No. 11 (1996) 65-68
13. Freedman, D.: Markov Chains. In: Holden-Day Series in Probability and Statistics (1971)
14. Joshi, A., Joshi, K., Krishnapuram, R.: On mining Web Access Logs. In: Proceedings of the ACM-SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery, (2000) 63-69
15. Krishnapuram, R., Joshi, A., Nasraoui, O., Yi, L.: Low-Complexity Fuzzy Relational Clustering Algorithms for Web Mining. In IEEE Transactions on Fuzzy Systems, Vol. 9 (4), (2001) 595-608
16. Manber, U., Patel, A., Robinson, J.: Experience with Personalization on Yahoo!. In: Communications of the ACM, Vol. 43, No. 8 (2000) 35-39
17. Maseglier, F., Poncelet, P., Cicchetti, R.: An efficient algorithm for Web usage mining. In Networking and Information Systems Journal, Vol. X, n° X (2000) 1-X
18. Mobasher, B., Cooley, R.: Automatic Personalization Based on Web Usage Mining. In Communications of the ACM, Vol 43:8 (2000)142-151
19. Nanopoulos, A., Katsaros, D., Manolopoulos, Y.: Effective Prediction of Web-user Accesses: A Data Mining Approach. In: Proceeding of the WEBKDD 2001 Workshop, San Francisco, CA (2001)
20. NASA Kennedy Space Center Log, <http://www.web-caching.com/traces-logs.html>
21. Palpanas, T., Mendelzon, A.: Web Prefetching using Partial Match Prediction. In Proceedings of the 4<sup>th</sup> Web Caching Workshop (1999)
22. Shi, W., Wright, R., Collins, E., Karamcheti, V.: Workload Characterization of a Personalized Web Site and Its Implications for Dynamic Content Caching. In: Technical Report TR2002-826, Department of Computer Science, New York University (2002).
23. Spiliopoulou, M., Pohle, C., Faulstich, L.: Improving the Effectiveness of a Web Site with Web Usage Mining. In: Proceedings of WEBKDD99 (1999) 142-162

24. Srikant, R., Yang, Y.: Mining Web Logs to Improve Website Organization. In: Proceedings of WWW10, Hong Kong (2001) 430-437
25. Srivastava, J., Cooley, R., Deshpande, M., Tan, P.: Web Usage Mining: Discovery and Applications of Usage Patterns from Web Data. In SIGKDD Explorations, ACM SIGKDD (2000)
26. Su, Z., Yang, Q., Zhang, H.: A prediction system for multimedia pre-fetching on the Internet. In: Proceedings of the ACM Multimedia Conference 2000, ACM (2000)
27. VanderMeer, D., Dutta, K., Datta, A.: Enabling Scalable Online Personalization on the Web. In Proceedings of EC'00 (2000) 185-196
28. Wells, N., Wolfers, J.: Finance with a Personalized Touch. In: Communications of the ACM, Vol. 43:8 (2000) 31-34
29. World wide web committee web usage characterization activity. <http://www.w3.org/WCA>
30. Yang, Q., Tian-Yi, I., Zhang, H.: Mining High-Quality Cases for Hypertext Prediction and Prefetching. In D.W. Aha and I. Watson (Eds.): ICCBR 2002, LNAI 2080, Springer-Verlag Berlin Heidelberg (2001) 744-755
31. Yang, Q., Zhang, H., Li, I., Lu, Y.: Mining Web Logs to Improve Web Caching and Prefetching. In: N. Zhong et al (Eds.): WI 2001, LNAI 2198, Springer-Verlag Berlin Heidelberg (2001) 483-492