



**HAL**  
open science

# HDM: A Client/Server/Engine Architecture for Real Time Web Usage Mining

Florent Massegia, Maguelonne Teisseire, Pascal Poncelet

► **To cite this version:**

Florent Massegia, Maguelonne Teisseire, Pascal Poncelet. HDM: A Client/Server/Engine Architecture for Real Time Web Usage Mining. Knowledge and Information Systems (KAIS), 2003, 5 (4), pp.439-465. 10.1007/s10115-003-0097-6 . lirmm-00191952

**HAL Id: lirmm-00191952**

**<https://hal-lirmm.ccsd.cnrs.fr/lirmm-00191952v1>**

Submitted on 26 Nov 2007

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# HDM : a Client/Server/Engine Architecture for Real Time Web Usage Mining

F. Masegla<sup>(1,2)</sup>

M. Teisseire<sup>(1)</sup>

P. Poncelet<sup>(3)</sup>

<sup>(1)</sup>LIRMM UMR CNRS 5506  
161, Rue Ada  
34392 Montpellier Cedex 5, France

<sup>(2)</sup>Laboratoire PRISM, Universit de Versailles  
45 Avenue des Etats-Unis  
78035 Versailles Cedex, France

<sup>(3)</sup>LGI2P - Ecole des Mines d'Ales - Site EERIE  
Parc Scientifique Georges Besse. 30035 Nimes Cedex 1, France

E-mail: {masegla, teisseire, poncelet}@lirimm.fr

## Abstract

The behaviour of a Web site's users may change so quickly that attempting to make predictions, according to the frequent patterns coming from the analysis of an access log file, becomes challenging. In order for the obsolescence of the behavioural patterns to become as null as possible, the ideal method would provide frequent patterns in real time, allowing the result to be available immediately. We propose, in this paper, a method allowing to find frequent behavioural patterns in real time, whatever the number of connected users is. Considering how fast the frequent behaviour patterns can change since the last analysis of the access log file, this result thus provides completely adapted navigation schemas for users behaviour predictions. Based on a distributed heuristic, our method also answers several tackled problems within the data mining framework: Discovering "interesting zones" (a great number of frequent patterns concentrated over a period of time, or the discovering of "super-frequent" patterns), discovering very long sequential patterns and interactive data mining ("on the fly" modification of the minimum support).

**Keywords :** real time, interactive data mining, zone mining, heuristic, distributed, long sequential patterns.

## 1 Introduction

With the growing popularity of the World Wide Web (Web), large volumes of data such as addresses of users or URLs which have been requested are automatically gathered by Web servers and collected in access log files. Analyzing a server access data can provide significant and useful information for performance enhancement, restructuring a Web site for increased effectiveness, and customer targeting in electronic commerce.

Discovering for relationships and global patterns that exist in large access log files is usually called *Web Usage Mining*. In this domain, the involved techniques mainly focus on the users behavioural patterns discovery from a Web server log file in order to extract relationships within collected data [14, 19, 5, 16, 6, 12]. These methods are designed to provide *association rules* or *sequential patterns*. Although association rules methods make it possible to discover very useful correlation for marketing needs, they are not optimal for analyzing a Web site users behaviours. Indeed, the absence of order between the frequent items avoids any attempt to know how the frequent requests are organized. To give a time order within the extracted patterns, the sequential patterns notion ([2, 17]) proposes

to refine the association rules problem. In this context, each fact is provided with a time stamp. Extracting sequential patterns from an access log file can provide the following kind of behaviours:

*60 % of users who visited /jdk1.1.6/docs/api/Package-java.io.html and /jdk1.1.6/docs/api/java.io.BufferedWriter.html, also visited, within the 2 following days, the URL /jdk1.1.6/docs/relnotes/deprecatedlist.html*

Once discovered, these patterns make up very useful knowledge for dynamically organizing a Web server hypertext structure. In that context, the active user access pattern could be matched with one or more sequential pattern discovered and navigational hints may be added to the pages accessed.

Indeed, the manager of a site for which the previous pattern applies can decide that if a client asked for the page `/jdk1.1.6/docs/api/Package-java.io.html` and then asks for the page `/jdk1.1.6/docs/api/java.io.BufferedWriter.html` it would thus be smart to use the knowledge: “this user might ask for the URL `/jdk1.1.6/docs/relnotes/deprecatedlist.html` with a confidence rate of about 60%.”

However, considering the increasing number of Web sites’ visitors, and the speed at which users’ behaviour may change from one day to another, the result obtained by analyzing an access log file may lose its relevance, if it is exploited too late. In order to avoid the obsolescence of the result, we propose a Heuristic based Distributed Miner (HDM) designed to get frequent behaviour patterns, answering the Web Usage Mining problem in real time.

In our opinion, extracting frequent behavioural patterns for the connected users, in real time and thanks to our architecture, has the following advantages over the classic approaches:

- As a matter of fact we can observe that, at this time, classic Web usage mining methods propose frequent behaviour patterns taking place over several days, weeks or even months (depending on the access log time range). The consequence is that if a frequent behaviour pattern is hidden among the large amount of data stored by the access log file, and if this pattern is frequent only when considering a specific time range and not the whole file, then this frequent pattern won’t be exhibited. Let us consider the pattern  $F1$  occurring in ten clients navigation sequences between June 12 and June 14, and nowhere else, in an access log file *Log* which takes place from May 01 to July 30. Let us consider that from June 12 to June 14, *Log* contains 15 clients, whereas it contains 100 clients, when considered entirely. Finally let us consider that the user’s minimum support is 60%. With such a minimum support, when considering the entire *Log* file,  $F1$  won’t be found as a frequent pattern because it occurs in ten clients sequences, thus having a threshold of 10%. Nevertheless, the owner of this Web site would surely be interested in knowing that from June 12 to June 14, 60% of the clients had a behaviour corresponding to  $F1$ . A classic access log analysis would not allow to find this pattern. This notion is presented and details are given in section 7.
- Once the result of an access log file analysis is obtained, the frequent behaviour patterns exhibited are taken into account for behaviour predictions attempts. When several different rules correspond to a user’s behaviour, they are both considered at the same level. Let us consider that one of these rules has been found thanks to the navigations of users close to the connected user, then this rule should supplant the other one, when trying to predict the connected user’s behaviour. The remaining question is to determine which users are close to another connected user. In our opinion, one of the best answers is: the other connected users. Since we can extract their frequent behaviour patterns in real time (thanks to our proposal), we are able to apply that point of view.
- As for most web usage mining processes based on an access log file analysis, the smaller the

number of clients is recorded in this file, the faster will be the data mining process. HDM is designed to provide the result in real time, whatever the number of clients is.

- A large number of methods are designed to tackle a major problem in frequent itemsets (or sequential patterns) mining: the size of the frequent patterns is so large, that most algorithms will enumerate a too large number of solutions before saturating the computer memory. HDM is designed to face that problem too, providing a satisfying result, whatever the largest frequent behaviour pattern length is.
- Finally, our proposal relies on the fact that, when a user is connected to a Web site, the computing power available on the machine this user navigates with, remains unused, until one exploits the architecture proposed in this paper. This computing power is insignificant when considered on its own, but the same computing power multiplied per thousands (or even hundreds of thousands), becomes strong enough in the condition of knowing how to use it.

Since our proposal is a heuristic based miner (summarized in [8]), our goal is to provide a result having the following characteristics:

Let us consider  $U$  the set of connected users at time  $t$ . When considering  $N_U$ , the navigation sequences of the users belonging to  $U$ , we can observe that these users share some behavioural patterns. Let us call *realResult*, the set of frequent behavioural patterns embedded in the navigation sequences of the users belonging to  $U$ , at time  $t$ . This “snapshot” (e.g. the association of  $U$ ,  $t$ ,  $N_U$ , and *realResult*) of a particular category, will be the center of our study. *realResult* is the result to obtain, e.g. the result that would be exhibited by an algorithm sequential pattern mining algorithm which would explore the whole set of solutions by working on  $N_U$ . Let us now consider *HDMresult* the result obtained at time  $t$ , by running the method presented in this paper. We want to minimize  $\sum_{i=0}^{size(HDMresult)} S_i / S_i \notin realResult$  (with  $S_i$  standing for a frequent sequence in *HDMresult*), as well as maximize  $\sum_{i=0}^{size(realResult)} R_i / R_i \in HDMresult$  (with  $R_i$  standing for a frequent sequence in *realResult*). In other words, we want to find most of the sequences occurring in *realResult* while avoiding the proposed result to be larger than it should (otherwise the set of all clients navigations would be considered as a good solution, which is obviously wrong). Finally, if that goal is reached, HDM will be able to find the frequent sequences on the users of  $U$  at time  $t$ , as well as a classic log mining algorithm would do on the log file limited to the users in  $U$  and to time  $t$ .

The rest of the paper is organized as follows. In section 2 we propose a presentation of the sequential pattern mining problem in large databases. Section 3 summarizes their usage by Web Usage Mining methods. Our contribution is presented in section 4. Section 5 gives performance evaluation and a discussion about our method is presented in section 6. We conclude and give future work opportunities in section 7.

## 2 Definitions

In this section we define the sequential pattern mining problem in large databases and give an illustration. The sequential pattern mining definitions are those given by [1] and [2].

In [1], the association rules mining problem is defined as follows:

**Definition 1** Let  $I = \{i_1, i_2, \dots, i_m\}$ , be a set of  $m$  literals (*items*). Let  $D = \{t_1, t_2, \dots, t_n\}$ , be a set of  $n$  transactions ; Associated with each transaction is a unique identifier called its *TID* and an *itemset*  $I$ .  $I$  is a  $k$ -*itemset* where  $k$  is the number of items in  $I$ . We say that a transaction  $T$  *contains*  $X$ , a set of some items in  $I$ , if  $X \subseteq T$ . The *support* of an itemset  $I$  is the fraction of transactions in  $D$  containing  $I$ :  $supp(I) = \|\{t \in D \mid I \subseteq t\}\| / \|\{t \in D\}\|$ . An *association rule* is an implication of the form  $I_1 \Rightarrow I_2$ ,

where  $I_1, I_2 \subset I$  and  $I_1 \cap I_2 = \emptyset$ . The rule  $I_1 \Rightarrow I_2$  holds in the transaction set  $D$  with confidence  $c$  if  $c\%$  of transactions in  $D$  that contain  $I_1$  also contain  $I_2$ . The rule  $r : I_1 \Rightarrow I_2$  has *support*  $s$  in the transaction set  $D$  if  $s\%$  of transactions in  $D$  contain  $I_1 \cup I_2$  (i.e.  $\text{supp}(r) = \text{supp}(I_1 \cup I_2)$ ).

Given two parameters specified by the user, *minsupp* and *minconfidence*, the problem of association rules mining in a database  $D$  aims at providing the set of frequent itemsets in  $D$ , i.e. all the itemsets having support greater or equal to *minsupp*. Association rules having confidence greater than *minconfidence* are then to be generated.

As this definition doesn't take time into consideration, the sequential patterns are defined in [2]:

**Definition 2** A *sequence* is an ordered list of itemsets denoted by  $\langle s_1 s_2 \dots s_n \rangle$  where  $s_j$  is an itemset. The *data-sequence* of a customer  $c$  is the sequence in  $D$  corresponding to customer  $c$ . A sequence  $\langle a_1 a_2 \dots a_n \rangle$  is a *subsequence* of another sequence  $\langle b_1 b_2 \dots b_m \rangle$  if there exist integers  $i_1 < i_2 < \dots < i_n$  such that  $a_1 \subseteq b_{i_1}, a_2 \subseteq b_{i_2}, \dots, a_n \subseteq b_{i_n}$ .

**Example 1** Let  $C$  be a client and  $S = \langle (3) (4\ 5) (8) \rangle$ , be the purchases of that client.  $S$  means that "C bought the item 3, then he bought 4 and 5 at the same moment (i.e. in the same transaction) and finally bought item 8".

**Definition 3** The *support* for a sequence  $s$ , also called  $\text{supp}(s)$ , is defined as the fraction of total data-sequences that contain  $s$ . If  $\text{supp}(s) \geq \text{minsupp}$ , with a minimum support value *minsupp* given by the user,  $s$  is considered as a *frequent sequential pattern*.

### 3 Access Log Files Analysis

In [12], we propose a Web Usage Mining method based on an Intrasite architecture. This section aims at explaining that work, in order to present in a synthetic way, the techniques involved when finding frequent behaviours for a Web site users.

#### 3.1 From Sequential Pattern Mining to Web Usage Mining

The general idea is similar to the principle proposed in [9]. It relies on three main steps. First of all, starting from a rough data file, a pre-processing is necessary to clean "useless" information. The second step starts from transformed data and applies data mining algorithms to find frequent itemsets or frequent sequential patterns. Finally, the third step aims at helping the user to analyze the results by providing a visualization and request tool.

Raw data are collected in access log files by Web servers. Each input in the log file illustrates a request from a client machine to the server (*http daemon*). Access log files format can differ, depending on the system hosting the Web site. For the rest of this presentation we will focus on three fields : client address, the URL asked for by the user and the time and date for that request. We illustrate these concepts with the access log file format given by the CERN and the NCSA [4], where a log input contains records made of 7 fields, separated by spaces [15]:

```
host user authuser [date:time] 'request' status bytes
```

Figure 1 is a sample, coming from one of the Lirmm<sup>1</sup> (our institution's) access log files:

Such an access log file will then be processed in two steps. First of all, the access log file is sorted by address and by transaction. Then each "noninteresting" data is pruned out from the file. During the sorting process, and in order to allow the knowledge discovery process to be more efficient, URLs and clients are mapped into integers. Each time and date is also translated into relative time, compared to the earliest time in the log file.

<sup>1</sup>Lirmm : Laboratoire d'Informatique, de Robotique et de Micro-electronique de Montpellier

```

193.49.105.224 - - [29/Nov/2000:18:02:26 +0200] "GET /index.html HTTP/1.0" 200 1159
193.49.105.224 - - [29/Nov/2000:18:02:27 +0200] "GET /PtitLirmm.gif HTTP/1.0" 200 1137
193.49.105.224 - - [29/Nov/2000:18:02:28 +0200] "GET /acceuil_fr.html HTTP/1.0" 200 1150
193.49.105.224 - - [29/Nov/2000:18:02:30 +0200] "GET /venir/venir.html HTTP/1.0" 200 1141
193.49.105.225 - - [29/Oct/2000:18:03:07 +0200] "GET /index.html HTTP/1.0" 200 1051
193.49.105.225 - - [16/Oct/2000:20:34:32 +0100] "GET /form/formation.html HTTP/1.0" 200 14617
193.49.105.225 - - [31/Oct/2000:01:17:40 +0200] "GET /form/formation.html#doct HTTP/1.0" 304 -
193.49.105.225 - - [31/Oct/2000:01:17:42 +0200] "GET /form/theses2000.html HTTP/1.0" 304 -
193.49.105.56 - - [22/Nov/2000:11:06:11 +0200] "GET /lirmm/bili/ HTTP/1.0" 200 4280
193.49.105.56 - - [22/Nov/2000:11:06:12 +0200] "GET /lirmm/bili/rev_fr.html HTTP/1.0" 200 2002
193.49.105.56 - - [07/Dec/2000:11:44:15 +0200] "GET /ress/ressources.html HTTP/1.0" 200 5003

```

Figure 1: Access log file example

**Definition 4** Let  $Log$  be a set of server access log entries.

An entry  $g$ ,  $g \in Log$ , is a tuple  $g = \langle ip_g, ([l_1^g.URL, l_1^g.time] \dots [l_m^g.URL, l_m^g.time]) \rangle$  such that for  $1 \leq k \leq m$ ,  $l_k^g.URL$  is the item asked for by the user  $g$  at time  $l_k^g.time$  and for all  $1 \leq j < k$ ,  $l_k^g.time > l_j^g.time$ .

The structure of a log file, as describe in definition 4, is close to the “Client-Time-Item” structure used by sequential pattern algorithms. In order to perform a frequent behaviour extraction on a log file, for each  $g$  in the log file, we first have to transform  $ip_g$  into a client number and for each record  $k$  in  $g$ ,  $l_k^g.time$  is transformed into a time number and  $l_k^g.URL$  is transformed into an item number. Figure 2 gives a file example obtained after that pre-processing. To each client corresponds a series of times and the URL asked for by the client at each time. For instance, the client 2 requested the URL “60” at time  $d4$ .

| Client | d1 | d2 | d3 | d4 | d5 |
|--------|----|----|----|----|----|
| 1      | 10 | 30 | 40 | 20 | 30 |
| 2      | 10 | 30 | 20 | 60 | 30 |
| 3      | 10 | 70 | 30 | 20 | 30 |

Figure 2: File obtained after a pre-processing

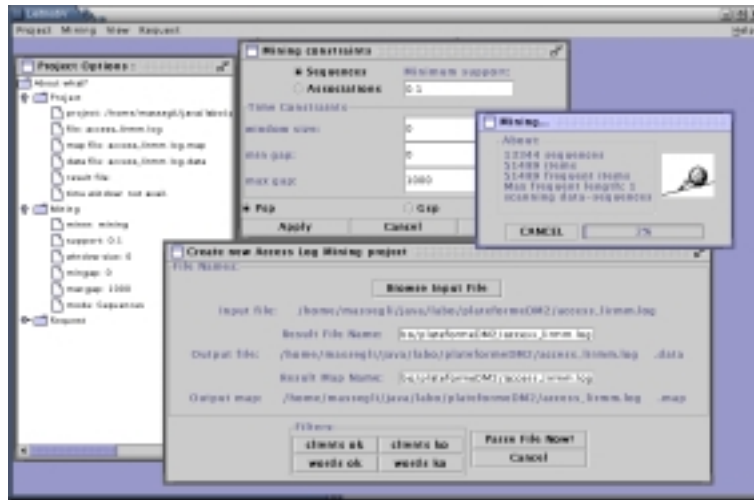


Figure 3: Leitmotiv while extracting sequential patterns

The goal is thus to find, according to definition 3 and thanks to a data mining step, the sequential patterns in that file that can be considered as frequent. The result can be, for instance  $\langle (10) (30) (20) (30) \rangle$  (with the file illustrated in figure 2 and a minimum support given by the user: 100%). Such a result, once mapped back into URLs, strengthens the discovery of a frequent behaviour, common to  $n$  users (with  $n$  the threshold given for the data mining process) and also gives the sequence of events composing that behaviour. Every step (data preprocessing, data mining, result presentation and manipulation...) is provided by our *LeitmotiV* soft, illustrated by figure 3.

### 3.2 Dynamic hypertext links modification

Jointly to the *LeitmotiV* system, we have developed a dynamic links generator using the rules coming from association rule or sequential pattern mining processes [13]. The generator is intended for recognizing a visitor according to his navigation through the pages of a server. When the navigation matches a rule given by *LeitmotiV*, the hypertext organization is dynamically modified. The functional architecture for the technique presented is depicted in figure 4.

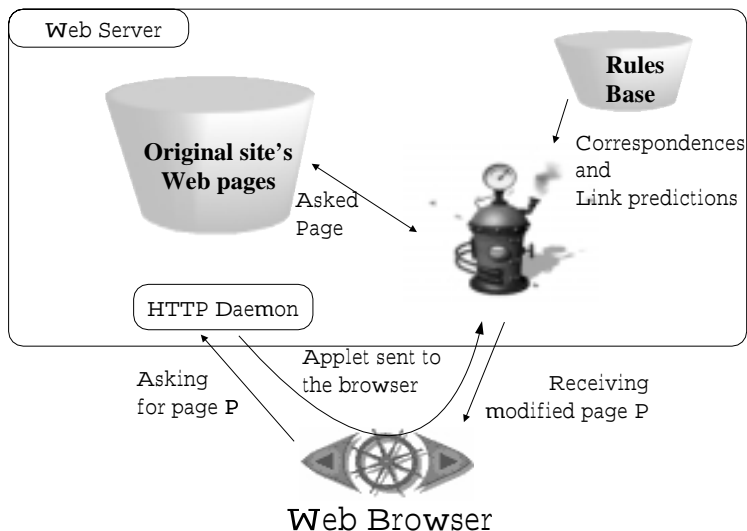


Figure 4: A dynamic links modification architecture, involving a web server and a web browser

Basically, the user's current behaviour can be compared to one or more sequential patterns and thus navigational hints can be added to the pages proved to be relevant for this category of users.

The Web server (*http daemon*) reacts to a customer request by returning an applet responsible for the connection to the *visitor manager module* in order to transmit visitor IP address and the required URL. At this time the client manager module is a java application running on the Web server site and using a client/server mechanism. When receiving an IP address and a required URL, the visitor manager examines the customer behaviour by using the *rule base* through the *correspondence module*. The latter checks if the customer behaviour, i.e. the client navigation, satisfies a rule previously extracted by the data mining process. When an input satisfies a rule in the correspondence module, the required page is modified by the *page manager* which dynamically adds links towards the consequent of the recognized rule. The applet then recovers the URL and displays the page in the navigator.

**Example 2** Among the several rules obtained when analyzing a Lirmm access log file, we noticed that 65% (*confidence*) of visitors asked for the page "accueil" and then "presentation" before they came back to "accueil". Then, such users asked for the page "formation" followed by "Ecole doctorale" and then "DEA informatique". That behaviour is summarized in figure 5.

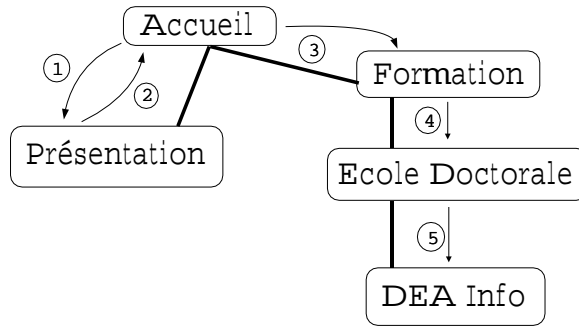


Figure 5: A navigation sequence example

Let us now consider a client, whose navigation sequence is  $\langle (accueil) (presentation) (accueil) (formation) \rangle$ . The beginning of this sequence is recognized in the rule base and matches the previous rule having a confidence of 65%. If that confidence is greater than the minimum confidence specified by the user, a link to the page "DEA Info" is added to the URL "formation", beside the one linking to "Ecole doctorale".

## 4 A better kind of knowledge

This section is devoted to the illustration of the kind of knowledge that can be found thanks to a real time web usage mining method. Let us consider the access log file illustrated at figure 6. In that file are recorded the transactions of five clients over two months. The sessions<sup>2</sup> are represented by the grey boxes. Let us now consider that we want to know the frequent behaviour, at time  $t$ , of the users that were connected at that same time (e.g.  $t$ ). We thus want to perform a data mining process over the users having a session  $s$  such that  $t \in [begin(s)..end(s)]$ . For instance, the user  $c5$  was in his session number 3 at time  $t$ , and the user  $c3$  was in his session number 4 at time  $t$ . The clients  $c1$ ,  $c3$  and  $c5$  have this characteristic.

If we consider the notion of snapshot given in the introduction, we have:

- $U = \{c1, c3, c5\}$
- $N_U =$ the rectangles in hash lines
- $t = t$ .

We thus want to know the frequent behaviours for the users  $c1$ ,  $c3$  and  $c5$  at time  $t$  (e.g. considering their transactions only before  $t$  (rectangles in hash lines)). The main interest of the result obtained with such a set of sequences will be to provide a knowledge dedicated to the connected users, as illustrated in the following exemple.

**Example 3** *Let us consider the access log file given in figure 7. This file gives the URL's requested by the users  $c1$ ,  $c2$ ,  $c3$ , and  $c4$  from time 1 to 7. Let us consider a first data mining process, with a support of 100%, on that log file. The frequent behaviour found will thus be:  $\langle (a) (g) (r) \rangle$ .*

*Let us now consider that we are at time 7, and we want to know the frequent behaviour of the users connected at time 7. As we can observe, the user  $c3$  has no request in the log file at time 7,  $c3$  is thus not considered as connected at that time 7. The connected users are thus  $c1$ ,  $c2$ , and  $c4$ . When*

<sup>2</sup>A session begins when the user logs on to the web site and it ends when the user leaves the site. Several tools are able to find the sessions in a log file, but in this article we use that concept for this illustration only (we will see later that our method does not need that concept since the connected users will be accessed in real time and not thanks to a session concept).



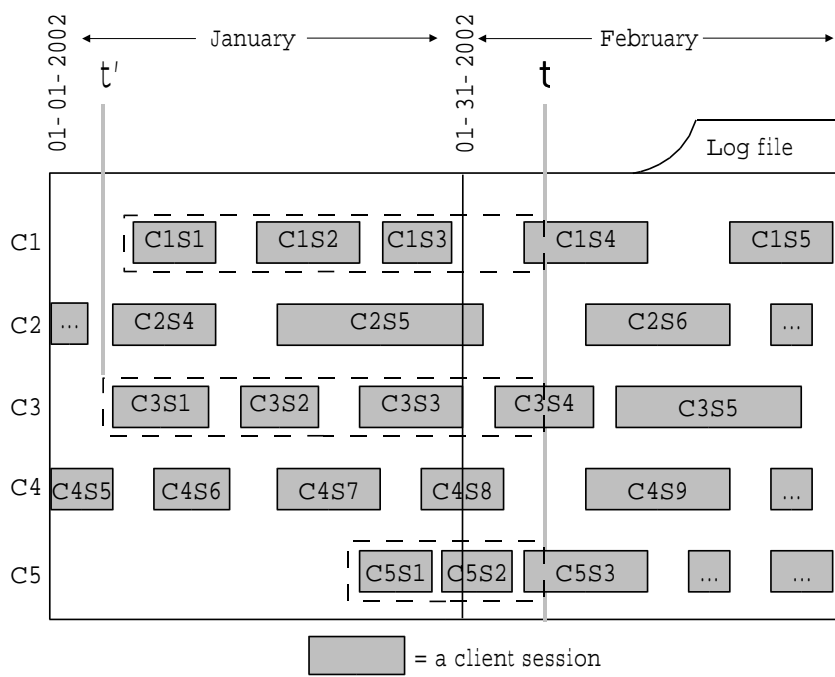


Figure 6: A log file with the sessions of 5 clients over two months

|    | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|----|---|---|---|---|---|---|---|
| c1 |   | a | v | g | p | r | w |
| c2 | a | b | g | n | r | o |   |
| c3 | a |   | v | g | r | t | w |
| c4 | a | v | x | g | r | u | w |

Figure 7: A log file, considered at time 7

considering only these users, we can observe that the frequent behaviour (with the same support of 100%) will be:  $\langle (a) (v) (g) (r) (w) \rangle$ . This frequent, obtained by mining only the connected users, is thus more interesting because with the same support it is longer and more significant.

### Discovering interesting periods

Let us now consider  $t'$  the minimum time of transactions occurring in the sequences of the users in  $U$  ( $t'$  is illustrated in figure 6). We can consider that the mining process on  $N_U$  will provide a result for a period taking place from  $t'$  to  $t$ . Depending on the threshold of that result, on the number of frequent sequences discovered, or any kind of characteristic for that result, we can decide whether that result is interesting or not. Nevertheless, since this result may not be discovered if the final user wants to analyze his log file from Jan 01 to Jan 31, this result might keep undiscovered. It thus can be interesting to know that from  $t'$  to  $t$  there is some particular result to look at.

The problem is thus to find a way to know whether a period is interesting or not. A classic mining process on an access log file  $Log$  taking place from  $d1$  to  $d2$ , would have to process that way:  
 foreach  $t \in [d1..d2]$  build a temporary log called  $LogTmp$  and use a data mining process on  $LogTmp$ .  
 If  $Log$  takes place over a period of time of one month, and if the time unit is the second, then there will be 2678400 possible choices for  $t$ . Enumerating all the possibilities in order to find an interesting  $t$  will thus be a bad solution.

On the other hand, if we are able to maintain the result in real time, for each  $t$ , then the question of mining interesting periods would be solved.

## 5 Metaheuristics and Data Mining

This section is not intended to be a presentation for meta heuristic methods. We mostly investigate here the way such methods can be used for Data Mining issues, and try to identify their faults.

Sequential pattern mining can rely on two different kinds of methods:

1. Algorithms, which explore the whole set of solutions, in order to provide a complete answer to the problem.
2. Heuristics, which guarantee a reliable result, approaching (and in a large number of cases, equal to) the global optimum with execution times largely reduced compared to algorithms described in the first point.

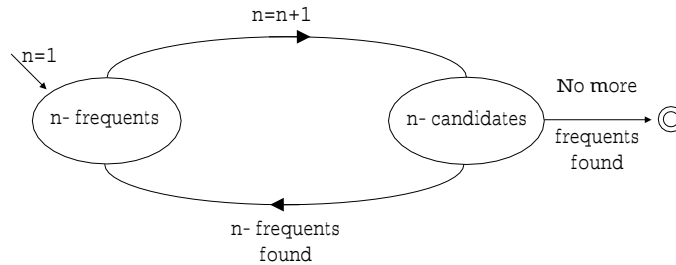


Figure 8: Generating-Pruning method

Algorithms categorized in the first point are numerous and based on methods like "Generating-Pruning" (cf figure 8), or on memory based methods which need to store all or part of the database in memory [11]. Unfortunately, when the problem becomes hard, such solutions might be unable to run in reasonable time (exploring the whole set of solutions, even when avoiding analysis of some subsets, is too CPU intensive).

Some methods then propose to provide a solution very similar (or even equal in numerous cases) to the global optimum [3, 18, 10]. Such methods correspond to the ones described in the second point and are intended to answer the two major difficulties one may have to face:

1. The result size avoids any enumeration of subsets of the result (what is done by a Generating-Pruning method).
2. The database size to analyze. When the database doesn't fit into the memory a memory based method will surely not be adapted.

There is, however, a third case, which also avoids using methods providing the optimal solution, and this case remains strongly related to the Web Usage Mining problems: fast growing data.

As a matter of fact, the behaviour of a Web site’s clients may change quickly, while its access log file is still being analyzed. Since an incremental algorithm cannot be considered as a real time method (even if executed as often as its execution time allows) we propose a stochastic algorithm for combinatorial optimisation problems, applied to the discovery of frequent sequential patterns in the Web Usage Mining framework. By “realtime”, we mean that the result can be updated every  $t$  seconds, with  $t$  the sum of :

- the time needed to contact a connected browser,
- the time needed to receive an answer of a connected browser,
- the time needed to compare  $x$  sequences to another sequence on the slowest distant machine.

We will give further details and a minimization of  $t$ , in a discussion about our architecture.

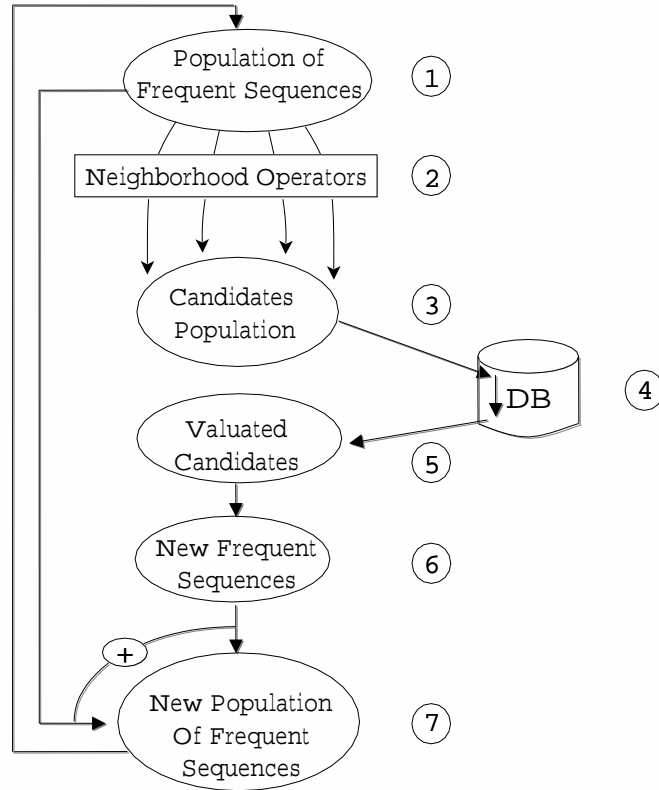


Figure 9: A data mining process based on a GA-like method

Figure 9 illustrates a possible method for solving the sequential pattern mining problem in databases using a meta heuristic.

This method is widely inspired from genetic algorithms (GA). Let us consider DB (step 4) as our access log file, using such a method for discovering frequent behavioural patterns from a Web site’s users collected navigations, is indeed possible. The starting point would be a population made of frequent (most asked) items (URLs) (step 1). Then, neighborhood operators propose a solution made of candidate sequences (steps 2 and 3) and those having the required threshold after a database scan (step 4) might take part in a finer solution (steps 5 and 6). Applying the neighborhood operators once again to the new solution will allow to refine the results and this process is repeated until no change enhances the solution (obviously, for Web Usage Mining, as data grows each time, each

change might enhance the solution and the solution will thus follow the clients behaviour on and on...).

That solution however presents an important issue: the number of scans to perform on the database in order to evaluate each proposed solution is too large. A Generating-Pruning based method should need  $k$  database scans (with  $k$  being the size of the largest frequent). A meta heuristic needs tens (even hundreds) of validations upon the database. This involves the two criteria that any heuristic has to face: exploration of the search space and exploitation of the proposed solutions.

For such a solution to be possible, the available computing power (needed to compare each candidate sequence with each sequence in the database) should increase at the same speed as the database size. That's the point of the solution proposed in this paper, and actually implemented at <http://www.kdd-tools.com>.

## 5.1 Computing Power

In order to avoid the problem related to the computing power, our proposal aims at evaluating each candidates set (population obtained after processing neighborhood operators) by each database entry: the connected machines. Each user connected to the Web site at any moment is using a browser and thus provides a local computing power (the one allowing him to navigate on the site). This computing power is available if one offers an appropriate architecture, designed to make the most of it. The database illustrated in figure 9, at step 4, can then be replaced by the set of all connected machines making that database up.

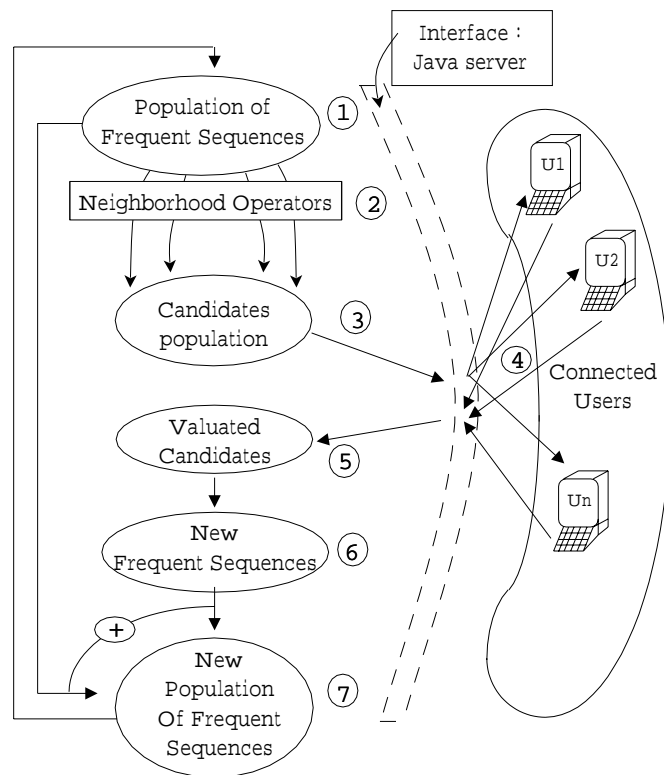


Figure 10: Our proposal, a client/server architecture where the client calculate at the server demand

This replacement is done by providing the heuristic with an interface, simulating a database scan by contacting each connected machine (fig 10). Each machine contacted this way will then compare its

own navigation sequence with each candidate that belongs to the candidates set just received. This comes down to comparing each data-sequence from the database (which groups all of the connected machines) with each candidate sequence, but allowing a satisfactory cost, considering the time required by a meta heuristic criterion. Let us consider the 4th step, in figure 9 on the previous page, the time complexity for this step was (likely)  $O(n * m * k * l)$ , with:

- $n$  the number of sequences in the access log file (database size)
- $m$  the number of candidate sequences (population size)
- $k$  the average size of sequences in the database
- $l$  the average size of candidate sequences to test.

We use an algorithm providing the longest sequence common to two sequences, with a time complexity:  $O(\text{firstSequenceLength} * \text{SecondSequenceLength})$  [7]. Our test for valuating candidates is mostly a rewriting of this algorithm.

On the other hand, considering figure 10, even if the time complexity remains the same, the calculations are distributed. Each calculation has a complexity, expressed as the following:  $O(m * k * l)$  with  $k$  the size of the navigation sequence on the contacted machine and  $l$  the average candidates size. The total number of calls to this calculation remains  $n$  but whatever the number of connected machines is, the delay<sup>3</sup> for a population valuation is  $\max(T1 + T2 + T3)$ , with :

- $T1$  : the time spent to send candidates to a machine  $n$
- $T2$  : the time spent to compare the sequence on  $n$  with the candidates
- $T3$  : the time spent to send results back from  $n$  to the server.

One of the main differences with usual data mining algorithms comes from the fact that we only work with the connected machines. All the machines corresponding to an entry in the access log file are not necessarily connected. This is one of the strengths of our method, which provides a result really adapted to the connected users because the results are validated on their category: other connected users. As we know, an access log file may take place over several days, weeks or even months. The frequent behaviour patterns coming from their analyze are long term behaviours, but such an analysis cannot provide short term behaviours corresponding to the most recent population (or even actual population).

## 5.2 Architecture

The method presented in the previous section, needs to be implemented as an architecture like the one described in figure 11. This architecture is implemented on the site '<http://www.kdd-tools.com>' and allowed us to test our proposal in a world wide experiment. When a machine accesses a page on the HTTP server, the server sends back the requested page with an embedded Java applet. This applet will then be responsible for contacting the Java server, in order to receive the calculations to run. At the beginning of the site's life cycle, the first connected client is the calculating engine, on which resides the stochastic algorithm. The engine will then regularly send to the server, the sequences to be distributed to the client machines.

The first problem to solve aims at maintaining pages thresholds, depending on user's requests, in real time. Example 4 gives an illustration of the method use for that purpose.

---

<sup>3</sup>A timeup has been implemented in our system, in case a machine cannot answer or the network has failures. This timeup can be set, and is generally less than 5 seconds.

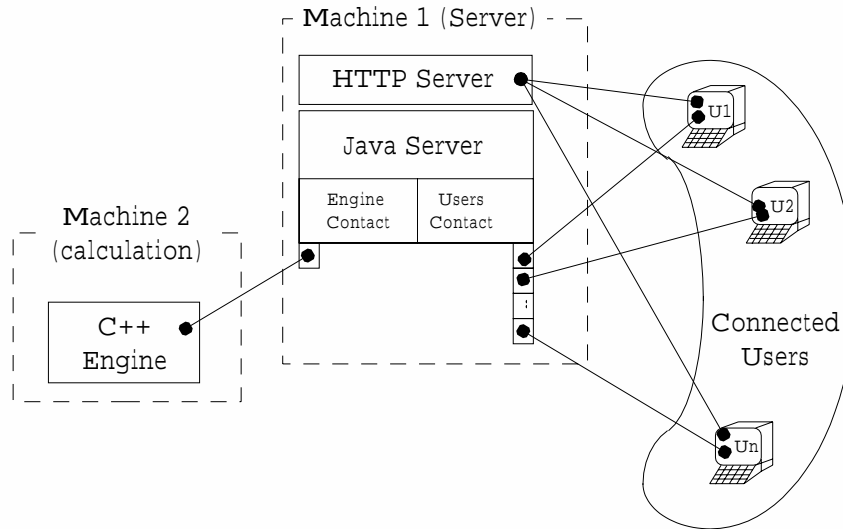


Figure 11: Connection between the candidates generator and the clients via the Web server

**Example 4** Let us consider that 2 users ( $U_1$  and  $U_2$ ) are already connected to the site. Their navigation sequences are  $\langle (1) (2) (1) (3) (5) \rangle$  for  $U_1$  and  $\langle (1) (2) (3) (4) (5) \rangle$  for  $U_2$ . The pages supports for this very moment are reported in the following table:

| Item    | 1    | 2    | 3    | 4   | 5    |
|---------|------|------|------|-----|------|
| Support | 100% | 100% | 100% | 50% | 100% |

This table is maintained by the Java server, and given regularly to the C++ engine which will use the frequent pages in order to provide candidates. Let us now consider a third user ( $U_3$ ) connecting to the Web site and asking for pages 1, 3 and then 2. Each requested page (if it does not appear in the user's sequence yet) will be reported to the Java server, in order to maintain the support of each page, depending on the number of connected users. Let us now consider that  $U_3$  asks for the page 3 again, this page is already in the navigation sequence of  $U_3$  so this action won't be reported to the Java server. Finally,  $U_3$  asks for the page 7, this is a new page for  $U_3$  so the Java server is contacted. The sequence of  $U_3$  is now  $\langle (1) (3) (2) (3) (7) \rangle$  and the new pages supports are:

| Item    | 1    | 2    | 3    | 4   | 5   | 7   |
|---------|------|------|------|-----|-----|-----|
| Support | 100% | 100% | 100% | 33% | 66% | 33% |

With a minimum support of 80%, the frequent behaviour pattern to find is  $\langle (1) (2) (3) \rangle$ . We don't consider 3 as occurring two times for  $U_3$  but only once, because of the definition of a frequent item. In fact, considering several occurrences of an item within a user's navigation sequence would give a biased answer since a client having asked  $\text{minSupport}$  times the item  $n$  while no other client asked for this item, would let the item  $n$  be frequent whereas it occurs in only one user's navigation sequence.

The architecture given in figure 12 does not only describe the way frequent pages are discovered in real time, but also gives an illustration of the operations performed during a user's navigation on the site:

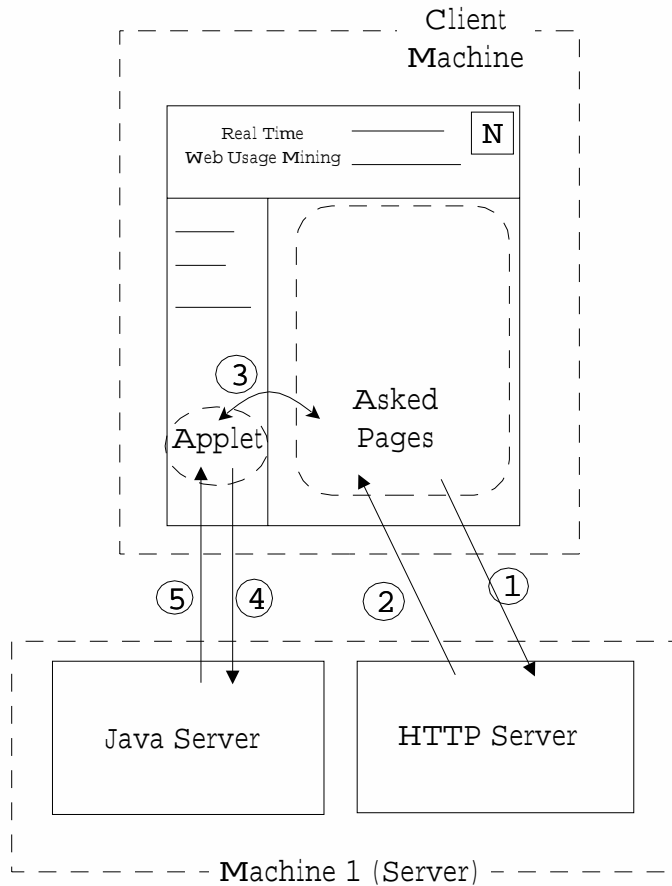


Figure 12: Connection between the candidates generator and the clients via the Web server

1. The browser asks for a page (the one needed by the user) by contacting the HTTP server.
2. The server sends the requested page, which, once arrived in the browser, will contact the applet in order to be added to the navigation sequence for this user (via a cookie).
3. if this page does not occur in the navigation sequence of the client, the applet, once contacted, will tell the Java server about this new page occurrence (in order to know its threshold).
4. The applet also sends back the results for the candidates comparison asked for by the engine.
5. The server regularly contact the applet to send the calculations asked for by the engine.

Let us give some details about this last point. As the engine knows the pages thresholds, it will use the frequent pages in order to propose candidates. These candidates, after being valuated by the connected users' machines, will become frequent or not, and frequent patterns found will be used with frequent pages in order to build new candidates and so on. The way such new candidates are proposed is developed in section 5.3. The C++ engine is thus based on algorithm 1. The criteria is based on the candidate ( $c$ ) size, the support of the candidate and an evaluation of the possibilities of accepting  $c$  even if it is a bad sequence. This feature allows HDM to follow the users behaviour patterns evolution since candidates are not valuated by a boolean function (e.g. "included or not") but by a function measuring the size of the candidate's subset included, compared to the total size of

the candidate, as described in section 5.3.

**function** *Hdm\_Engine*

Input: *JS* the Java server needed by the architecture

Output: *SP* the sequential patterns corresponding to the frequent behaviour patterns.

**while** *Web site is on-line* **do**

*pagesSupports*=getPagesSupports(*JS*); // Maintaining the pages supports

*candidates*=getValuation(*JS*); // Getting distributed calculations result

**for**  $c \in$  *candidates* **do**

**if** (*support*(*c*) > *minSupport* OR *criteria*) **then**

            insert(*c*, *SP*);

**endif**

**endfor**

*candidates*=neighborhoodOperators(*candidates*, *pagesSupport*);

    askForDistribution(*JS*, *candidates*);

**endwhile**

**end function** *Hdm\_Engine*

Algorithm 1: HDM engine

From a network point of view, the communications required by this method remain weak since the sequences are sent through a broadcast technique and an "on-the-fly" compression. This compression has a ratio of 20%, meaning that a 300 sequences population would cost approximately 8Kb, but takes only 1.5Kb on the bandwidth after being compressed.

### 5.3 Neighborhood operators and population valuation

The main idea of the HDM algorithm, is to propose candidates population thanks to previous frequent patterns and neighborhood operators, and then to send these candidates to the connected machines in order to know their threshold (or at least their distance from a frequent sequence). These two phases (neighborhood operators and candidates valuation) are explained in this section.

#### 5.3.1 Neighborhood operators

The neighborhood operators we used were validated thanks to an experiment performed on local databases (see section 6). We chose "Genetic-like" operators as well as operators based on sequential patterns properties. We present here some of the most efficient operators for the problem presented in this paper. When we talk about random sequence, we use a biased random such that sequences having a high threshold may be chosen before sequences having a low threshold.

##### *New frequent items*

When a new frequent item is occurring (after being requested by one or more users) it is used to generate all possible 2-candidate sequences with other frequent items. The candidates set generated is thus added to the global candidates set. Due to the number of candidate sequences to test, this operator only has a 15% ratio of accepted (e.g. frequent) sequences. This operator however remains essential since the frequent 2-sequences obtained are a base for other operators.

##### *Adding items*



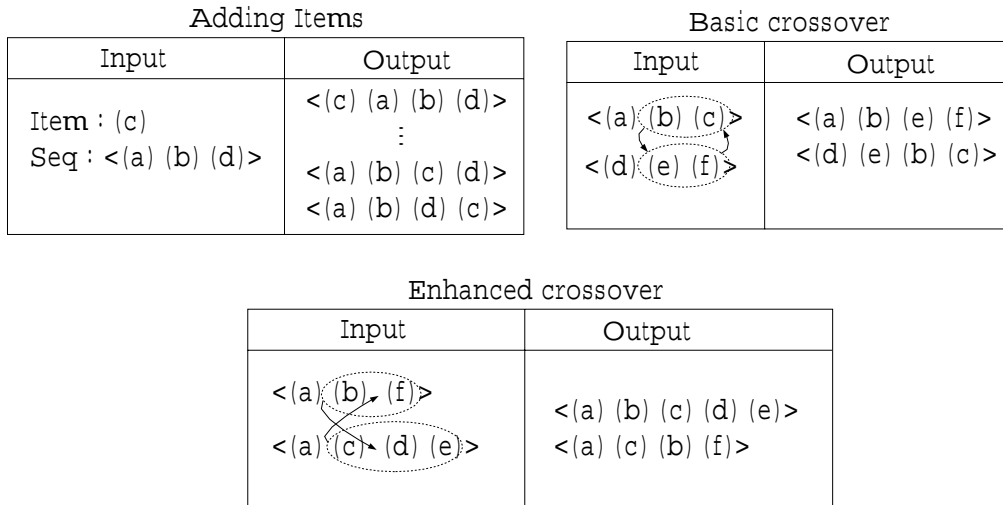


Figure 13: Some operators for finding frequent navigation sequences.

This operator aims at choosing a random item among frequent items and adding this item to a random sequence  $s$ , after each item in  $s$ . This operator generates  $length(s) + 1$  candidate sequences. For instance, with the sequence  $\langle (a) (b) (d) \rangle$  and the frequent item  $c$ , we will generate the candidate sequences  $\langle (c) (a) (b) (d) \rangle$ ,  $\langle (a) (c) (b) (d) \rangle$ ,  $\langle (a) (b) (c) (d) \rangle$  and finally  $\langle (a) (b) (d) (c) \rangle$ . This operator has a 20% ratio of accepted sequences, but the sequences found are necessary for the following operators.

#### *Basic crossover*

This operator (widely inspired from GA operators) uses two different random sequences and proposes two new candidates coming from their amalgamation. For instance, with the sequences  $\langle (a) (b) (c) \rangle$  and  $\langle (d) (e) (f) \rangle$ , we propose the candidates  $\langle (a) (b) (e) (f) \rangle$  and  $\langle (d) (e) (b) (c) \rangle$ . This operator has a good ratio (50%) thanks to frequent sequences embedded in the candidates generated by previous operators.

#### *Enhanced crossover*

Encouraged by the result obtained when running the previous operator, we developed a new operator, designed to be an enhancement of the basic crossover, and based on the frequent sequences properties. This operator aims at choosing two random sequences, and the crossover is not performed in the middle of each sequence, but at the end of the longest prefix common to the considered sequences. Let us consider two sequences  $\langle (a) (b) (f) \rangle$  and  $\langle (a) (c) (d) (e) \rangle$ . The longest common prefix of these two sequences is  $\langle (a) \rangle$ . The crossover will then start after the item following  $a$ , for each sequence. In our example, the two resulting candidate sequences are,  $\langle (a) (b) (c) (d) (e) \rangle$  and  $\langle (a) (c) (b) (f) \rangle$ . This operator has a success ratio of 35%.

#### *Final crossover*

An ultimate crossover operator was designed in order to improve the previous ones. This operator is based on the same principle as the enhanced crossover operator, but the second sequence is not randomly chosen. Indeed, the second sequence is chosen as being the one having the longest common

prefix with the first one. This operator has a ratio of 30%.

### *Sequences extension*

This operator is based on the following idea: frequent sequences are extended with new pages requested. The basic idea aims at adding new frequent items at the end of several random frequent sequences. This operator has a success ratio of 60%.

Figure 13 gives an illustration of some operators described in this section.

#### **5.3.2 Calculation performed by the client machine**

**function** *Client*

Input:  $CS$  the candidate sequences to evaluate and  $N$  the navigation sequence of the client.

Output:  $P$  the set of percentages assigned to each sequence.

**for**  $S \in CS$  **do**

**if** ( $S \subseteq N$ ) **then**

$P[S]=100+size(S)$ ; //  $S$  is included in  $N$ ,  $S$  is rewarded.

**endif**

**if** ( $size(S) \leq 2$ ) **then**

$P[S]=0$ ; //  $S$  has length 2 and is not included, it is not worth keeping it.

**endif**

  //Otherwise, give  $S$  a mark, and give short sequences a better mark

$P[S]=\frac{size(LCS(S,N))*100}{size(S)} - size(S)$ ;

**endfor**

**end function** *Client*

#### Algorithm 2: Client algorithm

When the candidates generation is complete, the C++ engine sends the candidates set to the Java server and this latter, through a broadcast technique, send it to the connected browsers. The calculation asked to the connected browser, is performed in the Java applet, loaded when the user logged on to the site, and compares each candidate to the navigation sequence of this user. The comparison aims at returning a percentage, representing the distance between the candidate and the navigation sequence. If the candidate is included in the sequence, the percentage should be 100% and this percentage will decrease when the amount of interferences (differences between the candidate and the navigation sequence) will increase. Furthermore, in order to obtain frequent sequences as long as possible, we use an algorithm that rewards long sequences if they are included in the navigation sequence. On the other hand, in order for the clients not to give a good mark to any long sequence, the algorithm has to avoid long, not included, sequences. To take all these parameters into account, the calculation performed by the client machine is described in algorithm 2.

## **6 Experiments**

We show here our experiments, designed to evaluate the result quality of the presented heuristic and also its ability to provide those results in real time. Before presenting our result, let us notice that our architecture has been validated by a test on a real web site. This web site (<http://www.kdd-tools.com>) has been built by our team, in order to be sure that the Java components and the communications will be working as expected. We thus asked to the research

|   |   |
|---|---|
| D | Number of customers (size of Database)                |
| C | Average number of transactions per Customer           |
| S | Average length of maximal potentially large Sequences |
| N | Number of items                                       |
| L | Size of the log file                                  |

Table 1: Parameters

| Dataset | D    | C  | S  | N   | L      | Server |
|---------|------|----|----|-----|--------|--------|
| F1      | 80K  | 18 | 53 | 2K  | 1,5 Gb | IIS    |
| F2      | 88K  | 32 | 58 | 2K  | 1,4 Gb | IIS    |
| L1      | 110K | 14 | 41 | 30K | 1 Gb   | Apache |
| L2      | 120K | 15 | 43 | 30K | 1,2 Gb | Apache |

Table 2: Synthetic datasets

members of the Lirimm to navigate through the web site for 15 minutes. This test has been repeated 3 times. The number of users connected during each test was about a hundred and each test was a success in validating the stability of the architecture and the quality of the result. Indeed, the result found by HDM at the end of each test was exactly the same as the result found on the log file with a classic log analysis. As the users were asked to keep connected for 15 minutes (e.g. a session of 15 minutes) at the same time, that was the expected result. If the users had not connected to the site simultaneously, the result found by HDM would have been expected different (e.g. corresponding to the last connected population, so we would not have been able to check the quality of the proposed frequent sequences in comparison to the whole population in the log file). Finally the experiment has been conducted in a world wide way since numerous access have been recorded in the log file for machines located in different countries (our server keeps connected at any time).

On the other hand, in order to validate HDM behaviour and result, we had to test it on local databases, where each data sequence was considered as a connected user. Our experiments were performed on files having characteristics described in tables 1 and 2.

The files F1 and F2 have been given by the web site <http://www.first-invest.com>, whereas the files L1 and L2 come from the Lirimm web site (<http://www.lirimm.fr>). FirstInvest is a financial web site, generating log files having size 1,5 Gb (approximately) every month. The Lirimm web site is less accessed, so we had to take two log files, covering a period of 5 months each, in order to get files having size 1 Gb. The number of items for the Lirimm access files is large because numerous pages use php and the parameters are written in the log files. FirstInvest uses a database to generate some pages, so its number of pages is large too. In both cases, the number of frequent items generally ranges from 50 to 70). The interesting point, when working on these two kind of files is to observe the behaviour of HDM when analyzing populations having different density and different navigation goals.

The experiments were performed on a 1,4 Ghz PC running Linux, with a local 60 Gb SCSI drive. The meta-heuristic and its simulated distribution have been written in C++.

For all those experiments we used an architecture simulator. This simulator had to be in contact with the engine (the meta-heuristic) and to work with a local database as the server would do with the connected users. To this end, the simulator considered each sequence in the database as a user and performed a scan over the database, in order to send to the engine the results of the comparison between the candidate set and each sequence in the database.

## 6.1 Quality measurement

In order to value the results provided by the HDM algorithm, we measured, for each proposed solution, the longest common sequence (LCS) between the sequences of the proposed solution and the real result to get (obtained by calling on the same data, the PSP algorithm, developed at the Lirmm by our team). Algorithm 3 is applied to each proposed population, in order to keep a trace of HDM behaviour. The goal of this algorithm is to provide the average quality of the sequences in the proposed solution. Basically, if all the sequences in the proposed solution are in the real result and if each sequence in the real result is in the proposed solution, then the proposed solution is considered as a result having a quality of 100% (e.g.  $proposedSolution \equiv realResult$ ). If a sequence in the proposed solution is only included (or not included at all) in a sequence of the real result, then the average quality will decrease. The last instructions aim at decreasing the global quality if a sequence in the real result is not in the proposed solution.

**function** *qualityMeasurement*

Input: *proposedSolution*, a solution to value. *PSP\_realResults* the real results to obtain (for comparison).

Output: *quality*, the quality percentage for proposed results compared to real results.

*sum*=0;

**for**  $s1 \in proposedSolution$  **do**

*localQuality*=0;

**for**  $s2 \in PSP\_realResults$  **do**

**if** ( $s1==s2$ ) **then** *localQuality*=100;

**else** *localQuality*=max(*localQuality*, (LCS( $s1,s2$ )/size( $s2$ ))\*100);

**endfor**

*sum* = *sum* + *localQuality*;

**endfor**

*quality*=*sum*/size(*proposedSolution*);

**for**  $s \in PSP\_realResults$  **do**

**if** ( $s \notin proposedSolution$ ) **then** *quality*=max(0,*quality*-1);

**endfor**

**return**(*quality*);

**end function** *qualityMeasurement*

Algorithm 3: Quality measurement algorithm

## 6.2 Validating the result quality

*Experiments on a stable behaviour:*

The goal of this first experiment is to validate the HDM algorithm, depending on the quality of the result provided. For that purpose we worked with the architecture simulator described above on the access log files F1, F2, L1 and L2. Basically, this experiment aims at making the snapshot last as long as necessary to validate the result quality. It comes down to test the architecture with a frozen population of connected users. We could thus know how much candidate sequences the engine has to test (starting from scratch) before providing a result having quality 100%.

Let us consider the first experiment, illustrated in Figure 14 (Log: F1). Each of the three graphs corresponds to a different size of the candidate set. Let us consider the graph corresponding to a candidate set having length 200. The first proposed solution has a quality of 76%. This is a good

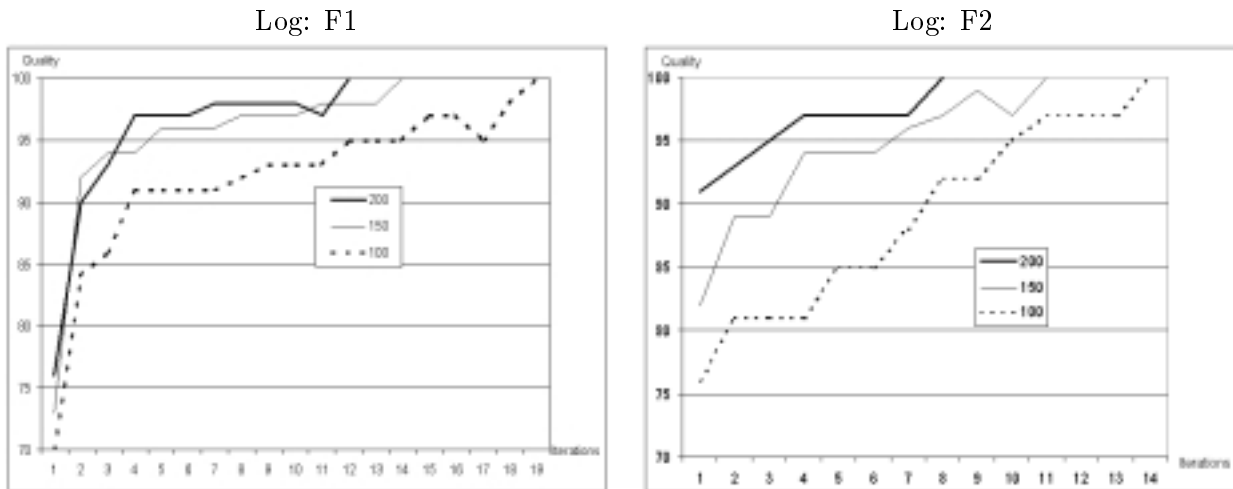


Figure 14: Results quality for each proposed solution

quality, because the engine already knows which are the frequent items, and the frequent sequences having length 2. This information is thus used to propose a first solution (based on a candidate set of 200 sequences). The engine then needs to propose 12 solutions before providing a result quality of 100%.

The other graphs stand for a candidate set having length 150 and 100. The main reason for being interested in varying the size of candidate set is that the bandwidth will be the limit for this size. As we can observe, for the logs F1 and F2, this size is rather important, since the difference of quality is 10% to 20% (as illustrated by the graphs). We thus provide a study about this difference, illustrated by Figure 15. This graph is the result of a second experiment conducted on F1. For each size of the candidate set, we report the number of iterations needed by the engine to provide a result having quality 100%. As we can observe in figure 14 and 15 for the log file F1, with a candidate set having size 200, the number of iterations is 12, with a candidate set having size 150, the number of iterations is 14 and with a candidate set having size 100, the number of iterations is 18. This graph clearly shows that the speed of the engine in terms of quality is depending on the number of sequences in the candidate population.

Let us now consider the third experiment, illustrated by Figure 16. The goal of this experiment is the same as the first, but the log files are L1 and L2, provided by the Lirmm web site. We can observe that the number of iterations before obtaining a 100% quality is larger (in comparison with F1 and F2). This is mainly due to a difference between the behaviours of the users for these two sites. Indeed, the access log file of the Lirmm does not hide a lot of frequent navigations, so we had to use a relatively weak support, and then the number of iterations before providing a result is thus larger. The size of the candidate set is not really important, since there are not much frequent behaviours in these files (meaning that 100 or 200 candidates are both enough to have a superset of the frequent sequences). We can observe, in figure 16, that HDM allows the results to degrade in order not to be trapped in a local optimum (with the file L1 and 100 candidates, from generation 15 to generation 18, the result quality decreases from 95% to 87%).

### 6.3 Validating the real time aspect of the approach

*Experiments on a variable behaviour:*

The goal of the first set of experiments was to work on a frozen population of users. In order to validate



Figure 15: Results quality for proposed populations, with frequent patterns having length 12

the real time aspect of our approach, we will now work on a changing population. This test aims at validating HDM ability at adapting to the behaviours, depending on the quality of a result following another one. To this end, we used two databases (DB1 and DB2) concealing distinct behaviours and we used HDM to find the frequent sequences. At first, only DB1 is considered to calculate the results, but with each generation proposed by HDM,  $x\%$  sequences from DB1 are replaced by  $x\%$  sequences from DB2. This process is repeated until DB1 is entirely replaced by DB2. The test then consists in estimating the quality of the results compared to those obtained by PSP on DB2, as soon as the replacement is achieved. The results of this test are reported in figure 17.

In this Figure, for the first graph, DB1 and DB2 are F1 and F2 (meaning that F1 has been replaced by F2). Let us consider  $x$ , the number of sequences from DB1 replaced by sequences from DB2 at each step ( $x$  thus stands for the speed at which the population will change). As we could expect, the faster DB1 is replaced by DB2, the worse the results get (for instance when F1 is replaced by F2 at the rhythm 1% per generation, results quality at the end of the replacement process is 100%, but if the rhythm is 10%, the quality decreases to 53%). This comes from the fact that HDM is not allowed to test a sufficient number of populations on DB2 (a rhythm of 10% means that HDM is allowed to test only 10 populations on the database), and furthermore DB2 is not revealed until the very end of the process. On the other hand, we can observe that with a progressive replacement (1% is indeed more realistic), the quality of the result ranges between 97% and 100%.

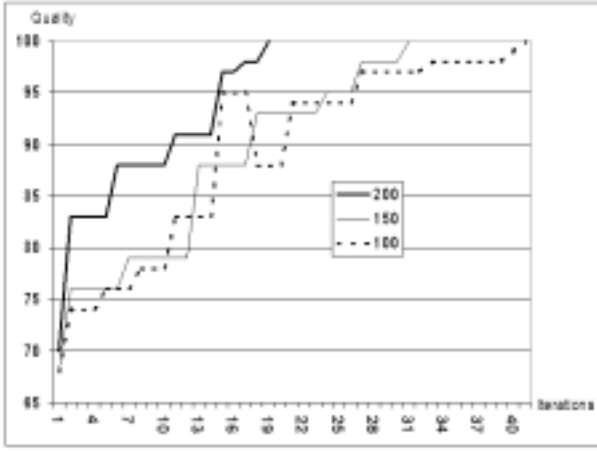
## 7 Discussion

The method we proposed in this paper, provides results having the following advantages:

*An immediate availability:*

As they are calculated in real time, the results available at time  $t$ , are, by nature, a representation of the behaviour corresponding to people connected at that same time  $t$ . The obtained results can then be used immediately, in two different ways:

Log: L1



Log: L2

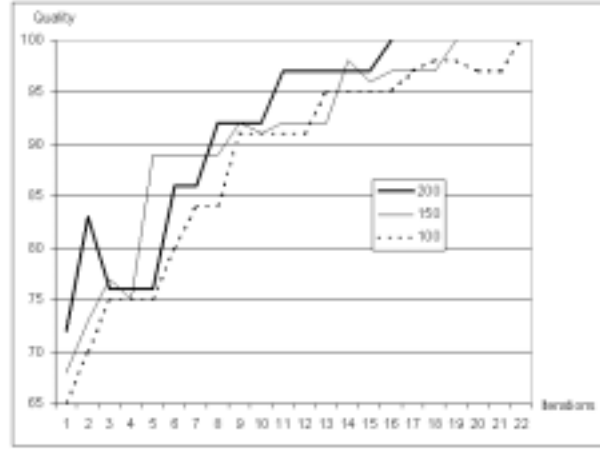


Figure 16: Results quality for each proposed solution

- As a data mining process result, considered as reliable and in which one can trust to take marketing or hypertext link modification decisions.
- As a pre-processing result (which doesn't cost a lot) designed to prepare the ground for a complete data mining process that can get help from these results in order to generate its candidates.

*A new kind of knowledge:*

Thanks to the concept of “snapshot”, given in the introduction, we can obtain a better kind of knowledge that can be discovered in a log file. Section 4 gave an example about the hidden knowledge that can be discovered thanks to our architecture. The question usually coming from the end user is of the kind: “*What is the frequent behaviour of my clients between June and August?*”. We are then within the framework of a data mining problem. However, at this time, and to the best of our knowledge, no method proposes to answer the question: “*Is there a period for which my database conceals more frequent behaviours? If yes, tell me which period it is*”. The given answer may take place over a different period than the one expected by the user, with, for instance, the following answer: “*The period from June 12 to Sept 2, conceals a frequent behaviour having a threshold of 72%*”. This behaviour may not have been discovered if the data mining process kept the period specified by the user. These zones can be discovered in real time, while HDM discovers behaviours, by recording frequent patterns for population having particular characteristics, such as:

- At least one frequent pattern having a very high threshold.
- Average threshold of the actual frequent patterns population is higher than the previous population threshold and the next population threshold (a top average threshold).
- The number of frequent patterns in the actual population is very high (heterogeneous frequent behaviour patterns) or very low.

Let us consider  $U$ , the set of connected users at time  $t$ , as given in the introduction. The result obtained by HDM on  $N_U$  is the same than the result obtained by a classic sequences extraction method applied to  $N_U$ . The problem of the classic extraction method is that it would have to try all the possible ways to obtain  $N_U$  (e.g. generating as many log files as there are time units in the log time range). For instance, with a log file starting the first of January and ending the 31st, the user would apply the following method:

$\forall i \in (1..2678400) \{ \log Tmp \leftarrow \text{rewriteLog}(\text{time} = i) ; \text{applyExtractingAlgorithm}(\log Tmp) \}$ .

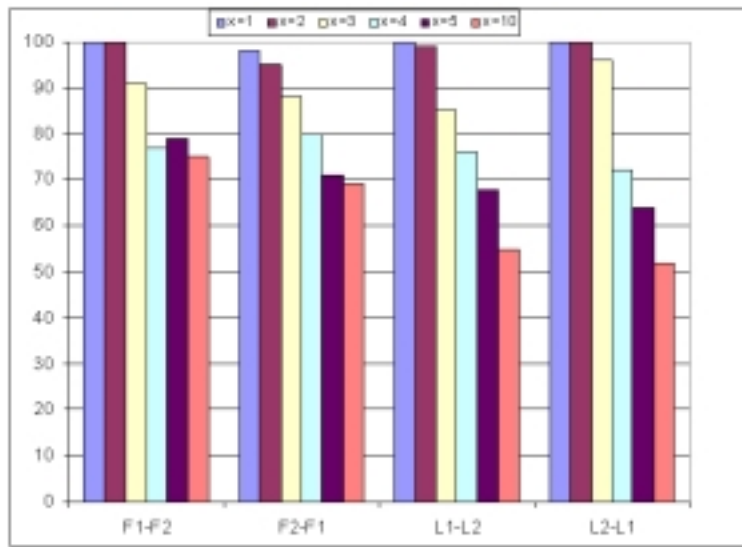


Figure 17: Results quality after a complete replacement of DB1 by DB2

Where  $rewriteLog(time\ t)$  is a function searching for each user which has been connected to the site at time  $t$ , and copy his navigation sequence (until  $t$ ) to a temporary log. This temporary log would then have to be processed by a data mining algorithm in order to decide whether the period given by  $t$  is interesting or not. The time needed to obtain all the possible  $N_U$  and all the corresponding frequent sequences is obviously too long, whereas by using our method we can obtain and update that “better knowledge” in real time.

*The results are dedicated to the connected users:*

Since the results are obtained in real time, the end user (i.e. the owner of the studied site) can use the results in an immediate way. That allows to target and to consider behaviours for users connected at this time, thanks to a global behaviour study over the users logged on.

*An inexhaustible computing power:*

Since HDM uses the computing power available on the connected machines, the problem provides a sufficient computing power to solve it, and the more complicated the problem gets, the more that computing power grows...

*Interactive data mining:*

At any moment, the Java server can accept a new user’s support, and send it to HDM immediately. HDM will then take this modification into account when merging the results of clients calculations and decide if a sequence is frequent or not. The next candidates population will then be proposed according to this new support. Since HDM runs continuously, as long as the site is on line, the support can be modified at any time.

*Mining long frequent sequential patterns:*

Since HDM does not enumerate an exhaustive list of the subsets of the frequent patterns to find, it does not depend on the size of the frequent sequences to find. We thus argue that our method can



find most (or even all) frequent patterns, whatever their size is. The experiments performed on local databases, allowed us to validate this feature while using HDM on an access log file concealing frequent patterns having size 40.

## 8 Conclusion and future work

In this paper, we proposed a real time method for Web Usage Mining designed to perform most calculations on the machines navigating on the Web site. Since the execution times of our method do not have to be evaluated, we centered our performance study on the quality of the result provided. Our experiments showed that the approach allows to find all the frequent sequences in most cases. With an aim of showing the feasibility of our approach, we presented the architecture of a Web site, designed to extract frequent behaviour patterns from users navigations. The main idea, based on the available computing power provided by connected browsers, has been evaluated during several test sessions, performed at the Lirmm. These tests, of about 15 minutes each, grouping a hundred users, allowed to find all the frequent sequential patterns exhibited after an access log analysis.

This method gives some future work opportunities, focusing on results quality.

We are working on *adapting the candidates population*, depending on the connected users. In a connected users categorization objective, it seems possible to maintain several populations at the same time (possibly managed by several engines, running simultaneously) and then to chose which population to send to which user, depending on different criteria:

- The time slots: depending on the slot, connected users may come from different countries. So, maintaining several candidates populations corresponding to users coming from the US, France or Japan, may allow to keep in memory the best behaviours for those categories of users and thus avoid rebuilding the population since we can get them evolving from their previous state.
- In order to extend that point of view, we can say that different categories of users will be exhibited from the extracted behaviours. Each behaviour can thus be refined if we send candidates trying to extend it, to only clients who correspond to this behaviour. That second point of view can be considered as an automation of the first one, since we don't chose the criteria that distinguish the users and we would let the data mining process make the distinction by itself.

Finally we are working on *studying new neighborhood operators*, in order to improve the heuristic side of the proposed approach.

## References

- [1] R. Agrawal, T. Imielinski, and A. Swami. Mining Association Rules between Sets of Items in Large Databases. In *Proceedings of the 1993 ACM SIGMOD Conference*, pages 207–216, Washington DC, USA, May 1993.
- [2] R. Agrawal and R. Srikant. Mining Sequential Patterns. In *Proceedings of the 11th International Conference on Data Engineering (ICDE'95)*, Tapei, Taiwan, March 1995.
- [3] R.J. Bayardo. Efficiently Mining Long Patterns from Databases. In *Proceedings of the 1998 ACM SIGMOD Conference*, pages 85–93, Almaden, USA, 1998.
- [4] World Wide Web Consortium. httpd-log files. In <http://lists.w3.org/Archives>, 1998.
- [5] R. Cooley, B. Mobasher, and J. Srivastava. Web Mining: Information and Pattern Discovery on the World Wide Web. In *Proceedings of the 9th IEEE International Conference on Tools with Artificial Intelligence (ICTAI'97)*, November 1997.

- [6] R.W. Cooley. Web Usage Mining: Discovery and Application of Interesting Patterns from Web Data. Technical report, University of Minnesota, 2000.
- [7] T. Cormen, C. Leiserson, and R. Rivest. *Introduction a l'Algorithmique*. ed. Dunod.
- [8] P. Poncelet F. Masegla, M. Teisseire. Real Time Web Usage Mining: a Heuristic based Distributed Miner. In *Proceedings of the 2nd International Conference on Web Information Systems Engineering (WISE'2001)*, Kyoto, Japan, December 2001.
- [9] U.M. Fayad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy, editors. *Advances in Knowledge Discovery and Data Mining*. AAAI Press, Menlo Park, CA, 1996.
- [10] D. Gunopulos, H. Manila, and S. Saluja. Discovering all Most Specific Sentences by Randomized Algorithms Extended Abstract. Technical report, Academy of Finland, 1997.
- [11] J. Han, J. Pei, and Y. Yin. Mining Frequent Patterns without Candidate Generation. In *Proceedings of the 2000 ACM SIGMOD Conference*, 2000.
- [12] F. Masegla, P. Poncelet, and R. Cicchetti. An Efficient Algorithm for Web Usage Mining. *Networking and Information Systems Journal*, 2(5-6), December 1999.
- [13] F. Masegla, P. Poncelet, and M. Teisseire. Usinf Data Mining Techniques on Web Access Logs to Dynamically Improve Hypertext Structure. *ACM SigWeb Letters*, 8:13–19, October 1999.
- [14] B. Mobasher, N. Jain, E. Han, and J. Srivastava. Web Mining: Pattern Discovery from World Wide Web Transactions. Technical Report TR-96-050, Department of Computer Science, University of Minnesota, 1996.
- [15] C. Neuss and J. Vromas. *Applications CGI en Perl pour les Webmasters*. Thomson Publishing, 1996.
- [16] M. Spiliopoulou and L.C. Faulstich. WUM: A Tool for Web Utilization Analysis. In *Proceedings of the EDBT Workshop WebDB'98*, Valencia, Spain, March 1998.
- [17] R. Srikant and R. Agrawal. Mining Sequential Patterns: Generalizations and Performance Improvements. In *Proceedings of the 5th International Conference on Extending Database Technology (EDBT'96)*, pages 3–17, Avignon, France, September 1996.
- [18] H. Toivonen. Sampling Large Databases for Association Rules. In *Proceedings of the 22nd International Conference on Very Large Databases (VLDB'96)*, September 1996.
- [19] O. Zaïane, M. Xin, and J. Han. Discovering Web Access Patterns and Trends by Applying OLAP and Data Mining Technology on Web Logs. In *Proceedings on Advances in Digital Libraries Conference (ADL'98)*, Santa Barbara, CA, April 1998.