

A Trend-based Prediction System for Web User Behavior

^{a,c}Nien-Yi Jan, ^bShun-Chieh Lin, ^cNancy P. Lin, and ^cChung-I Chang

^aBusiness & Marketing Strategy Research Department
Telcommunication Lab., Chunghwa Telecom Co., Ltd

^bDepartment of Computer Science
National Chiao Tung University

^cDepartment of Computer Science and Information Engineering
Tamkang University

151 Ying-chuan Road, Tamsui
Taiwan 251, R.O.C.

yijan@cht.com.tw, <http://mail.tku.edu.tw/125502/csie/pclin.htm>

Abstract: - Since web applications make great progress, the latency of Internet owing to the network bandwidth becomes an urge problem in the cyber world. It is very important to deliberate on how to construct a prediction model to predict web users traveling path for adapting the website structure and improving the website performance. A trend based prediction model without extra information is proposed in this paper to generate prediction models with a sequence of pages for a proxy server prefetching the suitable pages. The trend similarity is the core of our proposed model which considers not only the page similarity but also position similarity. Two measures include page correctness rate and order correctness rate are proposed to evaluate accuracy of our prediction system.

Key-Words: - Trend similarity, Prediction system, Web mining, User behavior, Sequence mining

1 Introduction

As the World Wide Web is growing at a rapid rate, the web-based services and applications become more popular. Users are interested in seeking useful information on Internet and are highly requesting for the correctness in quickly response. How to adapt the structure of website to meet the needs of users becomes an important task for a website manager.

To build suitable models by analyzing user browsing logs to facilitate the perfecting in proxy servers by predicting the future traveling path. Proxy servers play a key role between users and websites, which could reduce the response time of user requests and save network bandwidth. [10][19][23] have proposed various effective caching and prefetching algorithms in proxy servers to improve the website performance through analyzing the user behaviors to cache. Web mining technology [3][9] is widely used on constructing prediction models to dig deeply into the interests of web users by analyzing their navigation behaviors. Therefore, the web administrator can adapt the website structure to provide the suitable information for users according to their characteristics. In this paper, we will focus on analyzing the user browsing behaviors to predict the probably request pages for users in the near future based upon the prediction models. Many distinctive methods are proposed to construct a prediction model such as association rules [23],

Markov chain [24], and longest repeating subsequences [15], most of them however need extra information, such as website topology and the similarity between pages, to improve the performance. Some researchers cluster similar browsing behaviors using clustering technology, most of them however focus on clustering similar browsing behaviors or similar pages. These prediction systems only predict which pages will be accessed but ignores the navigation sequences. In this paper, we will propose a trend based prediction system including two phases (prediction model constructing phase and predicting phase) to analyze user behaviors and predict the future traveling path based upon the trend similarity. Behavior prediction models can be constructed in the first phase using user browsing behavior sessions, where each model consists of a prefix pattern and a postfix pattern. The prefix pattern is used to help the system determine all similar behaviors and the postfix pattern is used to select the page candidates. In the second phase, the trend similarity of a new user browsing behavior can be compared with all prefix patterns. Only the postfix patterns whose the similarity is larger than a similarity threshold will be picked out to select suitable prediction page candidates. Besides selecting the page candidates, the ordering of these pages can be easily decided by an ordering grade considering the position of page candidates and the

trend similarity. Moreover, new prediction modes will be incrementally constructed for improving the accuracy and efficiency of prediction. Without predefining the pages similarity, we design trend similarity to consider the sequence of the predicting pages. The latest page browsing could be more useful for predicting the further browsing pages. Two measures including page correctness rate and order correctness rate are designed to evaluate the accuracy of our prediction system. Our proposed system can be applied on the e-learning system to prefetch the predicted materials on the device of student according to their environment or study performance.

2 Related Works

2.1 Web Mining

Due to the dramatic growth of the Internet, the number of transactions becomes larger and larger. It results in the difficulty of analyzing these complicated data using traditional analysis methods. Many researchers tend to discover the potential behaviors using data mining approaches to predict the behavior trend [18], to improve the performance of network usage, to support decision in EC, and to provide a guideline in designing websites, etc.

In [11], web mining can be divided into three categories according to different data types: (1) web content mining, (2) web usage mining, and (3) web structure mining. In 2000, Cooley considered user profile in web mining to improve the accuracy of analyzing the user behavior with similar profile [5].

Web content mining [4] is proposed to discover the meaningful information in the content of web pages. Unlike web content mining, web structure mining is proposed to discover the relationship or link between web pages according to the characteristic of web content. It can be used to adjust the structure of website after discovering the structure of web content. Also, understanding the structure between hyperlinks can help the web designer easily arrange the unique information flow to facilitate the query of meaningful information.

Since the browsing behaviors will be logged, these traveling paths of the user could be important for analyzing the interest of users. Based upon the analyzed results, several personalization of recommendation services can be provided to facilitate the management of a website. Web usage mining [1][9] is proposed to monitor the interaction between users and website for discovering useful information through analyzing the collected user navigation patterns [16]. M. Spiliopoulou [17] concluded three

goals in web usage mining according to different domains. (1) To predict the user behavior of a website can be applied to reduce the latency of network access. (2) To compare the realistic performance of the website with the ideal performance by adjusting the website to improve the efficiency of a website according to the log analysis results. (3) To adjust the website structure according to the user interest to improve the performance of website after understanding the personality of each user.

2.2 Web prediction

Many researchers proposed Markov Chain [24], self-similarity [6], clustering [20], and web mining [2][7][12] to predict the browsing behaviors of users. In [24], the Markov chain is proposed to predict the maximal probability of next browsing page by calculating each probability of each link according to the web structure and the historical browsing behaviors. The self-similarity methods are proposed to predict the user behaviors using statistics methodology since the browsing behaviors are heavy-tailed distribution [6]. [10][15][16][19] proposed web mining to analyze the user behaviors for predicting the future browsing behaviors. Nanopoulos focused on increasing the performance of prediction [13]. Yang used association rule mining to discover the co-occurrence web pages embedded into the request pages in advance for reducing the latency of network access [23]. Xie and Phoha proposed a belief function to cluster similar user behaviors for predicting reference [21].

Three purposes for predicting network user behaviors are focused on (1) making efficient recommendation and personalization for users, e.g., dynamic adapt the web site structure [14], (2) reducing the latency of browsing web sites using prefetching [10][19][23], and (3) detecting the anomaly behaviors and intrusions [8]. Since the limitation of space in storing the perfecting pages, it is important to predict which pages could be re-used when perfecting.

3 Prediction System Architecture

Due to the variousness and uniqueness of user behaviors, it is inappropriate to predict the browsing behaviors of current user according to the similarity comparison with single browsing cluster of older users. Therefore, a trend-based prediction model is proposed to predict the further traveling path by generating ordering browsing sequence.

3.1 Framework of Prediction

In order to reduce the latency of browsing websites, a trend based prediction method is proposed in our prediction system to download some relevant web pages in advance on local user's computer according to the previous browsing sequence. The prediction system is divided into two components which shown in Fig. 1. One is prediction model constructing phase, and the other is predicting phase.

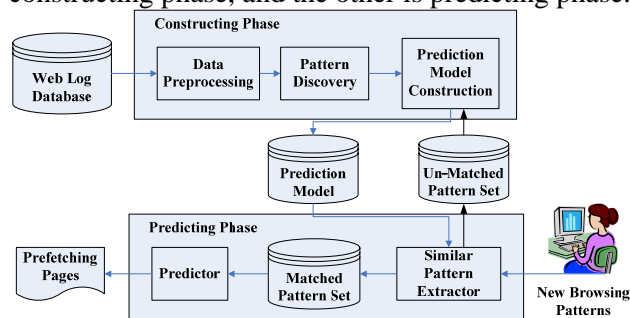


Fig. 1: Our prediction system architecture.

First, the prediction model of the website could be established in the constructing phase using web data mining approach after collecting enough web browsing log. In the predicting phase, the browsing behaviors of new users can be obtained to compare the similarity with the prediction model. Hence, the candidate pages could be prefetching to improve the browsing performance by applying on the replacement algorithm of proxy servers.

3.2 Prediction Model Constructing Phase

As mentioned above, the user behaviors are various, the prediction model constructing phase including data preprocessing, pattern discovery and prediction model construction steps is proposed to help experts discover useful common browsing patterns and then used to predict the further browsing sequences.

3.2.1 Prediction Model Definition

The browsing behavior can be treated as a sequence of browsing pages. The prediction model can be divided into prefix pattern and postfix pattern. Hence, the user browsing behavior can be predicted based upon the correlations between these patterns. The prediction model example can be shown as Fig. 2. In this example, the prefix pattern is {A0, B0, C3, D6}, which is length 4, and the postfix pattern is {D7, D8, D9, D16, D3, D9}, which is length 6. The model means that if a user browsed A0, B0, C3, D6 sequentially, the model will predict the D7, D8, D9,

D16, D3, D9 pages might be browsed after several minutes.

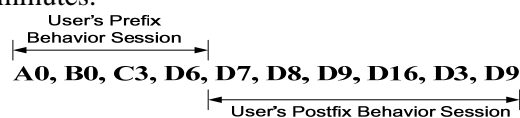


Fig. 2: User's behavior session example.

3.2.2 Data preprocessing

As we know, since the original web log records all user request pages, which often result in several log entries, e.g., graphics and scripts that are downloaded automatically. It is necessary to eliminate these irrelevant pages in order to obtain the actual user browsing behavior on the web server. Besides filtering the automatically generated pages, it is also important to distinct different behavior session for discovering different browsing behavior in a same user.

In order to obtain the actual browsing behavior, the embedded non-useful records, include multimedia files and the requests of aid agent provided by the web server, which are generated automatically in log database should be removed. Hence, retaining the useful records such as .asp or .html could improve the accuracy and practicability of prediction model to reduce the complexity of post processing.

The behavior session is defined as one browsing behavior of a user. For example, a user connects to one server from 10 AM to 12 AM on March 20 and then connects to the same server from 2 PM to 5 PM at the same day. Both browsing behavior should be treated as different behavior sessions due to the long period between these two browsing behaviors. Different sessions could be generated by different times from same users or different users. Therefore, based upon user login and logout period to distinct different session of a user would be more suitable for session identification. However, we cannot enforce a user to log into a popular server before browsing. Hence, Cooley's method [5] is used to distinguish user navigation patterns, and give a name Behavior Session according to IP and agent log during a period of time, e.g., 30 minutes.

After data preprocessing phase, the user behavior could be represented as a sequential of browsing pages. Each browsing sequence can be separate into two parts: prefix pattern and postfix pattern. The cut point can be chosen by the application.

3.2.3 Pattern Discovery

Although the user behaviors are various and distinctive, the browsing patterns of a structure website could be similar. These behavior patterns could be used to predict the further browsing behaviors of new users. Hence, it is important to discover the significant user patterns within web log database.

In order to obtain the significant patterns for predicting the user browsing patterns, a maximal closed itemsets algorithm [22] is chosen in this paper to discover the frequent prefix patterns using the part of prefix patterns dataset, i.e., satisfy the Support threshold.

3.3 Prediction Model Constructing

As mentioned above, the browsing behaviors are various, it means that users many browsing different postfix behavior patterns when they have the same prefix behavior patterns. On the other hand, they have the same postfix behavior patterns when they have similar prefix behavior patterns. Both mean that these users have similar browsing behaviors. Hence, it is reasonable to predict the postfix patterns through collecting all browsing behavior patterns of similar users. Two steps (including prediction page set determination and top-k prediction pages sorter) are designed in prediction model constructing phase.

3.3.1 Prediction page set determination Step

All pages in user's postfix behavior sessions whose occurrences reach Confidence threshold will be discovered as prediction set. To build the representative prediction model, all the page candidates could be discovered after setting a confidence value based upon each same prefix pattern sessions. For example, let the confidence value is set to 30%. There are 11 users with the prefix patterns {A0, B0, C10, D6} in Table 1. Therefore, we can obtain six pages {D2, D10, D29, D34, D41, D43}, which satisfies the confidence.

Since each prefix behavior session can lead to various number of postfix behavior sessions, the number of picked up prediction pages would be different. Therefore, the minimal length of prediction pages should be defined to avoid the useless prediction models. For example, if the minimal prediction length is set to 6, all prefix behavior sessions should be ignored where the number of page candidates is less than 6. Besides predicting the next pages, the ordering of page sequence are also important when building prediction model. Hence,

we proposed top-k prediction pages sorter to determine the ordering of prediction sequence. The length of Postfix behavior Patterns is equal to the number of page prediction set.

Table 1: The partial processed web log.

User ^o	Prefix behavior session ^o				Postfix behavior session ^o					
	1 ^o	2 ^o	3 ^o	4 ^o	1 ^o	2 ^o	3 ^o	4 ^o	5 ^o	6 ^o
U1 ^o	A0 ^o	B0 ^o	C10 ^o	D6 ^o	D10 ^o	D44 ^o	D3 ^o	D0 ^o	D29 ^o	D28 ^o
U2 ^o	A0 ^o	B0 ^o	C10 ^o	D6 ^o	D10 ^o	D34 ^o	D43 ^o	D41 ^o	D30 ^o	D19 ^o
U3 ^o	A0 ^o	B0 ^o	C10 ^o	D6 ^o	D23 ^o	D6 ^o	D34 ^o	D43 ^o	D41 ^o	D30 ^o
U4 ^o	A0 ^o	B0 ^o	C10 ^o	D6 ^o	D2 ^o	D9 ^o	D42 ^o	D28 ^o	C9 ^o	D29 ^o
U5 ^o	A0 ^o	B0 ^o	C10 ^o	D6 ^o	D34 ^o	D37 ^o	D40 ^o	D32 ^o	D0 ^o	D29 ^o
U6 ^o	A0 ^o	B0 ^o	C10 ^o	D6 ^o	D2 ^o	D9 ^o	D42 ^o	D19 ^o	D39 ^o	D15 ^o
U7 ^o	A0 ^o	B0 ^o	C10 ^o	D6 ^o	D34 ^o	D43 ^o	D41 ^o	D10 ^o	D0 ^o	D37 ^o
U8 ^o	A0 ^o	B0 ^o	C10 ^o	D6 ^o	D10 ^o	D4 ^o	D16 ^o	D32 ^o	D8 ^o	D12 ^o
U9 ^o	A0 ^o	B0 ^o	C10 ^o	D6 ^o	D2 ^o	D23 ^o	D36 ^o	D44 ^o	D3 ^o	D1 ^o
U10 ^o	A0 ^o	B0 ^o	C10 ^o	D6 ^o	D10 ^o	D34 ^o	D37 ^o	D38 ^o	D16 ^o	D29 ^o
U11 ^o	A0 ^o	B0 ^o	C10 ^o	D6 ^o	D2 ^o	C14 ^o	D18 ^o	D13 ^o	D43 ^o	D41 ^o
U12 ^o	A1 ^o	B0 ^o	C10 ^o	D6 ^o	D2 ^o	D9 ^o	D2 ^o	C14 ^o	D4 ^o	D16 ^o
U13 ^o	A1 ^o	B0 ^o	C10 ^o	D6 ^o	D23 ^o	D11 ^o	D31 ^o	D2 ^o	D41 ^o	D30 ^o
U14 ^o	A1 ^o	B0 ^o	C10 ^o	D6 ^o	D34 ^o	D6 ^o	D23 ^o	D36 ^o	D44 ^o	D27 ^o
U15 ^o	A1 ^o	B0 ^o	C10 ^o	D6 ^o	D23 ^o	D36 ^o	D39 ^o	D3 ^o	D0 ^o	D29 ^o
U16 ^o	A1 ^o	B0 ^o	C10 ^o	D6 ^o	D23 ^o	D11 ^o	D31 ^o	D2 ^o	D38 ^o	D4 ^o

3.3.2 Top-k Prediction Pages Sorter Step

Since the number of pages in prediction set might be too large and un-ordering, a top-k prediction pages sorter are designed to extract the representative and sorted prediction pages. In order to determine the ordering of page set, a weight vector is used to indicate that which page occurs in which position in each page candidate. The weight vector is set to increasing order. Equation (1) is used to calculate ordering grade to determine the ordering of page candidates.

$$O_{v_i}^j = (\sum S_{v_i}^j) / |PS_j| \quad (1)$$

where PS_j is the postfix behavior sessions with the same prefix behavior pattern j , v_i is the i^{th} page in a

page candidate set, $O_{v_i}^j$ is the ordering grade of v_i in PS_j , and $|PS_j|$ is the number of postfix page in PS_j .

In Equation (2), $S_{v_i}^j$ is the function of determining position weight before confirming the ordering of page candidates. The PenaltyWeight = $|PS_j| + 1$.

$$S_{v_i}^j = \begin{cases} PositionWeight & v_i \in PS_j \\ PenaltyWeight & v_i \notin PS_j \end{cases} \quad (2)$$

Table 2: The computing ordering score example.

Candidate	Order grade Computing	Sorting
D10	(1*4 + 4*1+7*6)/11	4.5
D34	(1*2+2*2+3*1+7*6)/11	4.6
D2	(1*4+7*7)/11	4.8
D29	(5*1+6*3+7*7)/11	6.5
D41	(3*1+4*1+5*1+6*1+7*7)/11	6.1
D43	(2*1+3*1+4*1+5*1+7*7)/11	5.7

Higher weight means higher probability to be clicked or downloaded by users according to the prefix pattern. For example, as mentioned above, {D2, D10, D29, D34, D41, D43} are the prediction page candidate shown in Table 1, the weight vector is set as {1, 2, 3, 4, 5, 6, 7} if the maximal length of postfix is 6. We can compute the ordering grade shown in Table 2 based upon Equations (1) and (2).

As we know, D10 occurs four times at the first position, once at the fourth position, and sixth do not occur at all prediction patterns. As shown in Table 2, we can obtain a predicting sequence D10, D34, D2, D43, D41 if we choose only top-5 postfix pages. The prediction model will be set to {A0, B0, C10, D6} → {D10, D34, D2, D43, D41}, where the first part is match pattern and the second part is probable pattern in the predicting phase.

3.4 Predicting Phase

The browsing behaviors of new visitors in the website can be predicted according to the prediction model. Since the browsing behavior might be similar with many prefix behavior sessions or not, the predictable patterns should be merged or the new prediction pattern should be constructed, respectively. The browsing behavior of a new visitor can be extracted to compare with all matching patterns in prediction model for predicting the future traveling path. The unmatched browsing patterns would be accumulated in unmatched pattern set to incremental constructing new predicting patterns to feedback into the prediction model.

3.4.1 Trend similarity

As we know, the user's navigation path is the time series. Later browsing behaviors would be more meaningful than the previous ones. Hence, a half decreasingly grading heuristic is proposed to implement the trend similarity to select all proper prediction patterns as Reference Patterns. The initial similarity is set to 0. To predict the browsing behaviors of new visitors with the generated browsing patterns, a pre-defined similarity threshold would be picked to filter some suitable prefix patterns to generate the probably patterns. Extracting the prediction patterns by trend similarity can get higher predicting accuracy.

$$Sim_{trend}(M_n^i, U_n^j) = \begin{cases} Sim_{trend}(M_{n-1}^i, U_{n-1}^j) + n, & \text{if } m_n = u_n \\ Sim_{trend}(M_{n-1}^i, U_{n-1}^j) / 2, & \text{if } m_n \neq u_n \end{cases} \quad (3)$$

where n is the length of the i^{th} comparing pattern. However, the similarity should be normalized by dividing the sum of position weight. For example, let

$M1 = \{ABCD\}$ be the Match Pattern and $U1 = \{ABCI\}$, $U2 = \{IBCD\}$ are user's navigation patterns. By equation (3), the similarity between M1 and U1 is 0.3, and the similarity between M1 and U2 is 0.9. Therefore, the prediction pattern that contains M1 is proper for U2 than U1. Hence, when a new user entered, the similarity between the prefix pattern of new pattern and prediction model will be calculated. Let's assume the similarity threshold is set to 0.5, the reference pattern set can be obtained when the similarity is large than the threshold. Therefore, we can obtain the reference pattern set is P1, P6, and P10. The partial prediction pattern models are shown in Table 3.

Table 3: The partial prediction pattern models

ρ	Match Pattern ρ				Probable Pattern ρ					Similarity ρ	
P1 ρ	A0 ρ	B0 ρ	C10 ρ	D6 ρ	D10 ρ	D17 ρ	D2 ρ	D29 ρ	D41 ρ	D43 ρ	0.61 ρ
P2 ρ	A1 ρ	B2 ρ	C3 ρ	D5 ρ	D22 ρ	D13 ρ	D9 ρ	D28 ρ	D13 ρ	D15 ρ	0 ρ
P3 ρ	A1 ρ	B4 ρ	C2 ρ	D6 ρ	D9 ρ	D5 ρ	D8 ρ	D17 ρ	D19 ρ	D11 ρ	0 ρ
P4 ρ	A0 ρ	B0 ρ	C4 ρ	D6 ρ	D34 ρ	D17 ρ	D18 ρ	D41 ρ	D38 ρ	D37 ρ	0.32 ρ
P5 ρ	A3 ρ	B5 ρ	C7 ρ	D9 ρ	D18 ρ	D36 ρ	D14 ρ	D35 ρ	D23 ρ	D40 ρ	0 ρ
P6 ρ	A4 ρ	B0 ρ	C10 ρ	D6 ρ	D22 ρ	D17 ρ	D9 ρ	D2 ρ	D43 ρ	D13 ρ	0.72 ρ
P7 ρ	A5 ρ	B5 ρ	C9 ρ	D24 ρ	D17 ρ	D20 ρ	D16 ρ	D39 ρ	D19 ρ	D21 ρ	0 ρ
P8 ρ	A0 ρ	B2 ρ	C10 ρ	D6 ρ	D10 ρ	D13 ρ	D3 ρ	D9 ρ	D29 ρ	D1 ρ	0.36 ρ
P9 ρ	A3 ρ	B9 ρ	C12 ρ	D6 ρ	D17 ρ	D9 ρ	D10 ρ	D15 ρ	D26 ρ	D29 ρ	0.22 ρ
P10 ρ	A0 ρ	B3 ρ	C10 ρ	D5 ρ	D10 ρ	D9 ρ	D17 ρ	D34 ρ	D41 ρ	D5 ρ	0.68 ρ

3.4.2 Sequence Prediction

At the previous step, all the postfix page candidates are chosen to decide the sequence of prediction pages. In Predictor, a prediction principal will be proposed to sort the page sequence and to decide the web user prospective traveling path. Similar to construct the prediction model, the ordering between page candidates should be ranked. Since the similarity between reference patterns and match pattern are usually different, the similarity should be considered as a merged weight when sorting these candidate pages. The grading function to sort the candidate pages is shown as function (4).

$$G_{v_i} = \sum (Sim_{Trend}(M^k) \times S_{v_i}^k) / \sum Sim_{Trend}(M^k) \quad (4)$$

where G_{v_i} is the ordering grade of predicting page v_i , $Sim_{Trend}(M^k)$ is the trend similarity between the match pattern M^k and the browsing pattern, $S_{v_i}^k$ is the position weight which is mentioned above. In this paper, we assume that the weight vector is {1, 2, 3, 4, 5, 6, 7} while the length of postfix pattern is 6.

For example, P1, P6, and P10 are the reference pattern set which are similar with the new user behavior. Therefore, the predicting page candidate set will be {D2, D5, D9, D10, D13, D17, D22, D29, D34, D41, D43}. Hence, the ordering grade of page D2 can be calculated using function (4).

$$G_{D2} = \frac{3 \times 0.61 + 4 \times 0.72 + 7 \times 0.68}{0.61 + 0.72 + 0.68} = 4.710$$

All the ordering grades of page candidates are shown as Table 4. Depending on the ordering grade result, the prediction sequence will be {D17, D10, D9, D2, D22, D41, D34, D43, D29, D13, D5}. As we know, the objective of constructing prediction model is to prefetch the suitable pages for a proxy server. Therefore, the necessary page according to the ordering grade can be prefetched.

Table 4: The ordering grades of all candidates.

Pages ^o	D2 ^o	D5 ^o	D9 ^o	D10 ^o	D13 ^o	D17 ^o	D22 ^o	D29 ^o	D34 ^o	D41 ^o	D43 ^o
Grade ^o	4.71 ^o	6.662 ^o	3.876 ^o	3.149 ^o	6.642 ^o	2.338 ^o	4.851 ^o	6.09 ^o	5.945 ^o	5.716 ^o	5.98 ^o

To construct the flexible prediction system, some space should be reserved to dynamically add new prediction patterns by integrating the prediction models. When the amount of new prediction patterns exceeds the reserved space, the prediction system will be offline for reconstruction to remove the out of date prediction patterns.

4 Experimental Result

To verify the accuracy of the prediction model, we established a small testing website with 80 web pages. Two measurements are proposed to evaluate the accuracy of the prediction models. The 10-fold evaluation is used to make a reasonable result.

4.1 Experiment Configuration

We simulated 20000 various user behavior sessions based upon website topology. The length of each session is set to 10 after data preprocessing, where the patterns will be ignored if the length is less than 10. In our experimental environment, 18000 user patterns are used to train the prediction model and the remaining 2000 patterns are used to test the accuracy of prediction models. 10-fold cross verification is used to enhance the confidence of evaluating our prediction model.

Predict: A4 B9 C18 D8 D42 D35 D14 D25 D5 D32
 Predict: A2 B6 C10 D23 D22 D3 D29 D43 D11 D19
 Predict: A2 B4 C10 D2 D40 D15 D36 D12 D4 D31
 Predict: A3 B10 C17 D8 D42 D32 D37 D25 D35 D5
 Predict: A3 B10 C13 D6 D9 D38 D16 D20 D18 D28
 Predict: A1 B1 C0 D3 D29 D22 D39 D43 D19 D23
 Predict: A2 B6 C10 D19 D22 D3 D29 D43 D23 D11
 Predict: A1 B1 C0 D0 D40 D4 D36 D31 D15 D12
 Predict: A3 B7 C15 D37 D5 D14 D25 D42 D8 D35

Fig. 3: Partial prediction patterns.

Lets assume the support is set to 0.002 to filter the maximal frequent prefix patterns and the confidence

is set to 50% to choose the candidate pages of postfix for each discovered prefix to construct prediction model. Figure 3 shows the partial prediction patterns in our experiment. In our experiment, each prediction model is divided into two parts. First is match pattern whose length is 4. Second is probable pattern whose length is 6.

4.2 Accuracy Measurement

In the prediction phase, each testing pattern is divided into two parts. First part is the browsed pages of this new user and the second part is the further browsing pages after verification. The similarity between matching pattern and the browsed pattern will be first calculated. The reference patterns will then be selected to generate the prediction result. To evaluate the accuracy of prediction result in our prediction system, two measuring criteria includes Page Correctness Rate (α) and Order Correctness Rate (β) are proposed to examine the predicting sequence.

The page correctness rate shown as function (5) is used to calculate the coverage of predicting sequence in the realistic further browsing patterns, where countA is the number of the same pages between testing pattern and predicting results, and l is the length of the testing pattern.

$$\alpha = \text{countA} / l \quad (5)$$

Since the ordering of a sequence can be treated as the set of multiple ordering pairs, the ratio of the same pairs of postfix patterns between predicting patterns and verification patterns can be calculated the correctness of the sequence ordering. The order

correctness rate is shown in function (6), where $\sum_{x=1}^{n-1} x$ is the total number of order-pairs, n is the number of page candidates, and countB is the same order pairs between testing pattern and predicting results.

$$\beta = \text{countB} / \sum_{x=1}^{n-1} x \quad (6)$$

For example, $T = \{A, B, C, D\}$ is one user actual browsing pattern, and $N = \{B, D, C, I\}$ is the predicting result by our prediction system. By Equation (5), the Page Correctness Rate (α) is $3/4 = 0.75$. Order-pairs included in T are $\{(A, B), (A, C), (A, D), (B, C), (B, D), (C, D)\}$. The order-pair (A, B) means that A comes before B and so on. In this example, the prediction system has predicted B, C, D pages. The actual order-pairs are (B, C), (B, D), (C, D). Only (B, D), (B, C) has been predicted and this results in that the Order Correctness Rate (β) is $2/3$.

4.3 Experimental Results

Table 5 shows the prediction correctness of the prediction system according to various similarity thresholds θ . The length of prefix pattern and the length of postfix pattern are 4 and 6, respectively.

Table 5: The prediction correctness with different similarity setting.

θ	0.0125	0.05	0.15	0.40
α	35%	46%	62%	68%
β	26%	34%	42%	46%

If the similarity thresholds are chosen 0.01, 0.05, 0.15, and 0.40, this means that there is at least one page are matched in first, second, third, and fourth, respectively. 2.1 pages can be predicted correctly in 6 predicted page set and the ordering correctness is 26% when similarity threshold is set to 0.0125. As shown in Table 5, the prediction correctness becomes high when the similarity threshold is setting high. Moreover, the latest browsing behavior is more important that the previous ones for predicting further browsing pages.

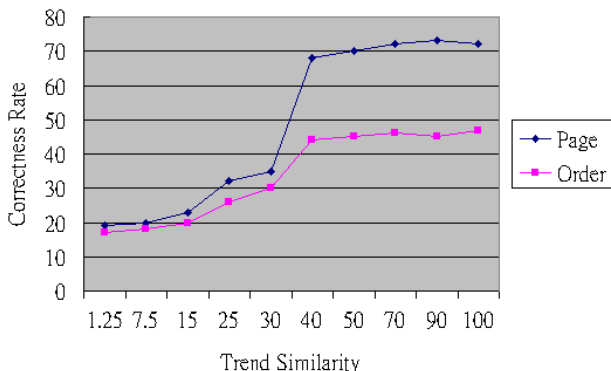


Fig. 4: The prediction accuracy between different trend similarities.

Fig. 4 shows the relationship between prediction correctness and the trend similarity. The result shows that the latest browsing behaviors can lead to more predictable pages for browsing. In Fig. 4, we observed that the rate of page correctness is less than 35% and the ordering correctness is less than 30% if the last page is unmatched. On the other hand, the page correctness increases to 73% and the order correctness increases to 46% when the last page is matched. This means the browsing trend are suitable for predicting further behaviors.

The last experiment is compared between our prediction system with the crisp prediction system according to the accuracy measurements where the similarity function is defined as trend similarity or

only page similarity. Our prediction method which integrates all models (higher than pre-defined similarity) and the crisp one chooses only one pattern with the highest similarity. Let's set the similarity threshold to 0.4 based upon the previous result.

Table 6: The accuracy comparison of our and crisp prediction models.

Model	Our Model		Crisp Model	
	Similarity	Trend	Page	Page
α	68%	52%	55%	46%
β	46%	37%	38%	32%

As shown in Table 6, the accuracy of the trend based similarity function is higher than the page similarity function in both models. Both page correct and page ordering accuracy of our model are better than the accuracy of crisp model. Therefore, our model should be useful in prefetching web pages for improving the performance of a proxy server. Choosing only one pattern to predict will lose the changeable behaviors of users.

5 Conclusion

In this paper, a trend based prediction system was proposed to predict the various web user behaviors. We design a trend similarity to select the proper prediction patterns for predicting a new user browsing behavior. Besides the previously constructing prediction models, our prediction patterns can be collected to learn the new prediction models to achieve the flexibility and adaptation of our prediction system. The experimental results represented that our model is better than page similarity based prediction system according to the page and order correctness rates. Our proposed model can be easily applied on the e-learning application to predict the user learning behaviors according to their learning environment.

Reference:

[1] S. Araya, M. Silva, and R. Weber, "A methodology for web usage mining and its application to target group identification," *Fuzzy Sets and Systems*, Vol. 148, No. 1, Nov., 2004, pp. 139-152.

[2] M. J. Asbagh and H. Abolhassani, "Web service usage mining: mining for executable sequences," in *Proc. of the 7th WSEAS Int'l Conf. on Applied Computer Science*, Vol. 7, Venice, Italy, 2007, pp. 266-271.

- [3] G. Castellano, A. M. Fanelli, M. A. Torsello, "LODAP: A LOG DATA Preprocessor for mining Web browsing patterns," in *Proc. of 6th WSEAS Int'l Conf. on Artificial Intelligence, Knowledge Engineering and Data Bases*, Corfu Island, Greece, 2007.
- [4] L. H. Chen, W. L. Chue, "Using web structure and summarisation techniques for web content mining," *Information Processing and Management*, Vol. 41, No. 5, Sep. 2005, pp. 1225-1242.
- [5] R. W. Cooley, "Web usage mining: discovery and application of interesting patterns from web data," Ph. D. dissertation, Dept. of Computer Science, University of Minnesota, May 2000.
- [6] M. Crovella and A. Bestavros, "Self-similarity in World Wide Web traffic: evidence and possible causes," *IEEE/ACM Trans. on Networking*, Vol. 5, Dec. 1997, pp. 835-846.
- [7] R. Ivancsy and I. Vajk, "Efficient sequential pattern mining algorithms," *WSEAS Trans. on Computers*, Vol. 4, 2005, pp. 96-101.
- [8] S. C. Lin, S. S. Tseng, and Y. T. Lin, "A new mechanism of mining network behavior," in *Proc. of the 6th Pacific-Asia Conf. on Advances in Knowledge Discovery and Data Mining*, 2002, pp. 218-223.
- [9] H. B. Liu and V. Kešelj, "Combined mining of web server logs and web contents for classifying user navigation patterns and predicting users' future requests," *Data & Knowledge Engineering*, Vol. 61, No. 2, May 2007, pp. 304-330.
- [10] W. W. Lou and H. J. Lu, "Efficient prediction of web accesses on a proxy server," in *Proc. of the 11th Int'l Conf. on Information and Knowledge Management*, 2002, pp. 169-176.
- [11] S. Madria, S. S. Bhowmick, W. K Ng, and E. P. Lim, "Research issues in Web data mining," in *Proc. of the 1st Int'l Conf. on Data Warehousing and Knowledge Discovery*, 1999, pp. 303-312.
- [12] J. Mamcenko, R. Kulvietiene, "Data mining technique for collaborative server log file analysis," *WSEAS Trans. on Information Science and Applications*, Vol. 2, Iss. 8, 2005, pp.1111-1115.
- [13] A. Nanopoulos, D. Katsaros, and Y. Manolopoulos, "Effective prediction of web-user accesses: A data mining approach," in *Proc. of the Workshop WEBKDD*, 2001.
- [14] M. Perkowitz and O. Etzioni, "Adaptive web sites," *Communications of the ACM*, Vol. 43, No. 8, August 2000, pp. 152-158.
- [15] J. Pitkow and P. Pirolli, "Mining longest repeating subsequences to predict world wide web surfing," in *Proc. of the USENIX Technical Conf.*, October 1999. pp. 139-150.
- [16] M. L. Shyu, S. C. Chen, and C. Haruechaiyasak. "Mining user access behavior on the WWW," in *Proc. of IEEE Int'l Conf. on Systems, Man, and Cybernetics*, Vol. 3, 2001, pp. 1717-1722.
- [17] M. Spiliopoulou, "Data mining for the web," in *Proc. of the 3rd European Conf. on Principles of Data Mining and Knowledge Discovery*, 1999, pp. 588-589.
- [18] V. S. Tseng and K. W. Lin, "Efficient mining and prediction of user behavior patterns in mobile web systems," *Information and Software Technology*, Vol. 48, No. 6, June 2006, pp. 357-369.
- [19] Y. H. Wu and A. L. P. Chen, "Prediction of web page accesses by proxy server log," *World Wide Web*, Vol. 5, No. 1, 2002, pp. 67-88.
- [20] J. T. Xiao and Y. C. Zhang, "Clustering of web users using session-based similarity measures," in *Proc. of the 2001 Int'l Conf. on Computer Networks and Mobile Computing*, 2001, pp. 223-228.
- [21] Y. J. Xie and V. V. Phoha, "Web user clustering from access log using belief function," in *Proc. of The 1st Int'l Conf. on Knowledge Capture*, 2001, pp. 202-208.
- [22] X. Yan, J. Han, and R. Afshar, "CloSpan: Mining closed sequential patterns in large database," in *Proc. of SIAM Int'l Conf. on Data Mining*, 2003, pp. 166-177.
- [23] Q. Yang, H. H. N. Zhang, and T. Y. Li, "Mining web logs for prediction models in WWW caching and prefetching," in *Proc. of the 7th ACM SIGKDD Int'l Conf. on Knowledge Discovery and Data Mining*, 2001, pp. 473-478.
- [24] J. H. Zhu, J. Hong, and J. G. Hughes, "Using Markov models for web site link prediction," in *Proc. of the 13th ACM Conf. on Hypertext and Hypermedia*, 2002, pp. 169-170.