

Web Navigation Prediction Using Multiple Evidence Combination and Domain Knowledge

Mamoun A. Awad and Latifur R. Khan

Abstract—Predicting users' future requests in the World Wide Web can be applied effectively in many important applications, such as web search, latency reduction, and personalization systems. Such application has traditional tradeoffs between modeling complexity and prediction accuracy. In this paper, we study several hybrid models that combine different classification techniques, namely, Markov models, artificial neural networks (ANNs), and the All- K th-Markov model, to resolve prediction using Dempster's rule. Such fusion overcomes the inability of the Markov model in predicting beyond the training data, as well as boosts the accuracy of ANN, particularly, when dealing with a large number of classes. We also employ a reduction technique, which uses domain knowledge, to reduce the number of classifiers to improve the predictive accuracy and the prediction time of ANNs. We demonstrate the effectiveness of our hybrid models by comparing our results with widely used techniques, namely, the Markov model, the All- K th-Markov model, and association rule mining, based on a benchmark data set.

Index Terms—Artificial neural networks (ANNs), association rule mining (ARM), Dempster's rule, Markov model, N-gram.

I. INTRODUCTION

WEB PREDICTION is the problem of predicting the next web page that a user might visit after surfing in a website. The importance of web prediction originates from the fact that various applications, such as latency reduction, web search, and recommendation systems, can be made more effective through the use and the improvement of web prediction.

One of the early applications of web prediction is the latency of viewing of web documents [6]. Traditional solutions are based on caching and prefetching [2], [3], [9]. Other advanced intelligent methods [10], [11] acquire knowledge from surfers' previous path history and utilize that in prediction. Pandey *et al.* [10] present an intelligent prefetching method based on a proxy server using association rule mining (ARM) to generate association rules that are later used to predict future requests.

World Wide Web (WWW) prediction can also improve search engines. The entire structure of the WWW can be pictured as a connected graph, where each node corresponds to a website, and surfers navigate from one node to another. The distribution of the visits over all WWW pages can be computed

and used in reweighting and reranking results. In such scenario, we consider the surfer path information to be more important than the keywords that were entered by the user [14].

Another application of web prediction is recommendation systems, in which we try to find the top k users having the same interests or tastes to a target user record. ARM is a well-known model that is used in recommendation systems. Mobasher *et al.* [7] propose the frequent item set graph to match an active user session with frequent item sets and predict the next page that the user is likely to visit. Prediction for the active session is based on the confidence of the corresponding association rule.

Other prediction models that are widely used in WWW predictions and its related applications include k nearest neighbors (NN), artificial neural network (ANN), fuzzy interference, and Markov model. Joachims *et al.* [18] propose the k NN-based recommender WebWatcher. The WebWatcher is a learning tour guide agent that extracts knowledge from user's previous clicks and from the hypertext structure. Nasraoui and Krishnapuram [19] propose a web recommendation system using fuzzy inference. Clustering is applied to group profiles using hierarchical unsupervised niche clustering. Context-sensitive uniform resource locator (URL) association is inferred using a fuzzy approximate-reasoning-based engine. Levene and Loizou [1] compute the information gain from the navigation trail to construct a Markov chain model to analyze user navigation pattern through the web. The main contribution of [1] is that they present a mechanism to estimate the navigation trail. Pitkow and Pirolli [5] explore pattern extraction and pattern matching based on the Markov model that predicts future surfing paths. Longest repeating subsequences (LRS) is proposed to reduce the model complexity (not predictive accuracy) by focusing on significant surfing patterns.

Our work is related to the path-based prediction model using the N-gram model [14] and the LRS model [5]. However, our approach differs from them in the following ways: First, only one path-based prediction technique is used when combining different N-gram models. Second, the main focus of LRS is to reduce the modeling complexity by reducing the data set. Third, all these models are probabilistic, i.e., it depends on the frequencies of patterns/occurrences in the training set. Therefore, our model can predict some values that Markov models cannot (i.e., our model can predict for some unobserved values). In the work of Nasraoui and Krishnapuram [19], the focus is to use a set of URL predictors by creating a neural network for each profile independently with a separate training set. Their goal is to overcome the high complexity of the architecture and training in case that one neural network is used. In our work, we not only use a set of predictors but also fuse them

Manuscript received October 17, 2007; revised February 21, 2007. This paper was recommended by Associate Editor M. Kamel.

M. A. Awad was with the University of Texas at Dallas, Richardson, TX 75083-0688 USA. He is now with the United Arab Emirates University, Al Ain, UAE (e-mail: mamoun.awad@uaeu.ac.ae).

L. R. Khan is with the Department of Computer Science, University of Texas at Dallas, Richardson, TX 75083-0688 USA (e-mail: lkhan@utdallas.edu).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TSMCA.2007.904781

92 in a hybrid model for prediction. Our goal is to improve the
93 accuracy using different prediction techniques, namely, ANN,
94 the Markov model, the All- K th model, and using different
95 N-gram models.

96 One important subtlety of web prediction is that web pre-
97 diction is a multiclass problem of a large number of classes
98 (11 700 classes in our experiments). Here, we define a class (or
99 a label) as a unique identifier that represents a web page in a
100 web site. Most multiclass techniques, such as one-versus-one
101 and one-versus-all, are based on generalizing binary classifiers,
102 and prediction is resolved by checking against all these binary
103 classifiers. As a result of that, prediction accuracy is very
104 low [4], because the prediction model has many conflicting
105 outcomes from the classifiers.

106 There are several problems with the current state-of-the-art
107 solutions. First, models such as Markov and ARM models are
108 unable to generalize beyond training data [5]. This is because
109 prediction using ARM and LRS pattern extraction is done based
110 on choosing the path of the highest probability in the training
111 set; hence, any new surfing path is misclassified, because it has
112 zero probability. Second, prediction using ARM suffers from
113 well-known limitations including scalability and efficiency [7],
114 [13]. Finally, many of the previous methods have ignored
115 domain knowledge as a means to improving prediction.

116 In this paper, we present a new approach to improving the ac-
117 curacy in web prediction. Our approach is based on generating
118 a hybrid prediction model by fusing two different classification
119 models. We use four classification models, namely: 1) ANNs;
120 2) ARM; 3) Markov model; and 4) *All- K th-model*. ARM and
121 Markov model are powerful techniques for predicting seen data,
122 i.e., already observed data; however, they cannot predict beyond
123 training data (see Section III-A). On the other hand, the All-
124 K th model and ANN are powerful techniques that can predict
125 beyond training data. In other words, the ANN and All- K th
126 models can predict some values that the Markov model and
127 ARM cannot. We combine the All- K th model with ANN by
128 fusing their outcomes using Dempster's rule.

129 Nonetheless, when dealing with a large number of classes or
130 when there is a possibility that one instance may belong to many
131 classes, their predictive power may decrease. To overcome
132 these shortcomings, we extract domain knowledge from the
133 training data and incorporate such knowledge during prediction
134 to improve prediction time and accuracy. Specifically, domain
135 knowledge is used to eliminate irrelevant classes and reduce
136 the conflict during prediction. Notice that we combine different
137 prediction models in which each model has different strengths
138 and drawbacks over other models. We strive to overcome
139 major drawbacks in each technique and improve the predictive
140 accuracy for the final hybrid model.

141 The contribution of this paper is given as follows: First, we
142 use ANN in web navigation. Second, we incorporate domain
143 knowledge in ANN prediction to eliminate irrelevant classes
144 and to improve prediction time and accuracy. Third, we fuse
145 ANN, the Markov model, and All- K th-Markov classifiers in
146 a hybrid prediction model using Dempster's rule [17] to im-
147 prove prediction accuracy and to overcome the drawbacks of
148 using each model separately. Finally, we compare our hybrid
149 model with different models, namely, Markov model, ARM,

All- K th-ARM, All- K th-Markov, and ANN using a standard 150
benchmark data set and demonstrate the superiority of our 151
method. 152

The organization of this paper is given as follows: In 153
Section II, we present the background of the N-gram concept 154
and sliding window. In Section III, we present different pre- 155
diction models that are used in web prediction. In Section IV, 156
we present the utilization of domain knowledge to improve 157
prediction. In Section V, we present a new hybrid approach 158
combining ANN, the Markov model, and the All- K th-Markov 159
model in web prediction using Dempster's rule for evidence 160
combination. In Section VI, we compare our results with 161
other methods using a standard benchmark training set. In 162
Section VII, we summarize this paper and outline some future 163
research. 164

II. N-GRAM REPRESENTATION OF PATHS 165

In web prediction, the available source of training data is 166
the users' sessions, which are the user's history of navigation 167
within a period of time. User sessions are extracted from 168
the logs of the web servers, and it contains sequences of 169
pages/clicks that the users have visited, time, data, and the pe- 170
riod of time that the user stays in each page. In web prediction, 171
the best known representation of the training session is the 172
N-gram. N-gram is tuples of the form $\langle X_1, X_2, \dots, X_n \rangle$ that 173
depict sequences of page clicks by a population of users surfing 174
a website. Each component of the N-gram takes a specific page 175
id value that identifies a web page. For example, the N-gram 176
 $\langle X_{10}, X_{21}, X_4, X_{12} \rangle$ depicts the fact that the user has visited 177
pages in the following order: page 10, page 21, page 4, and 178
finally, page 12. 179

Many models further process these N-gram sessions by 180
applying a sliding window to make training instances have 181
the same length [5], [7]. For example, if we apply a sliding 182
window of size 3 on the N-gram $\langle X_{10}, X_{21}, X_4, X_{12}, X_{11} \rangle$, 183
we will have the following 3-gram sessions: $\langle X_{10}, X_{21}, X_4 \rangle$, 184
 $\langle X_{21}, X_4, X_{12} \rangle$, and $\langle X_4, X_{12}, X_{11} \rangle$. In general, the number of 185
additional sessions using sliding window w applied on session 186
 A is $|A| - w + 1$, where $|A|$ is the length of session A . 187

In this paper, we also use the term *number of hops*, which 188
is related to the sliding window. The number of hops for a 189
session of length N is $N - 1$, i.e., the number of clicks (or 190
hops) that the user makes to reach the last page in the session. 191
When applying sliding window of size w , the number of hops 192
in the resulted subsessions is $w - 1$. In the previous example, 193
the number of hops in the resulted 3-gram sessions is 2. 194

III. PREDICTION MODELS 195

In this section, we briefly present various prediction models 196
that have been used in web prediction. First, we present the 197
Markov model; next, we present the ANN model along with 198
improvement modifications. 199

A. Markov Model 200

The basic concept of the Markov model is to predict the 201
next action, depending on the result of previous actions. In 202

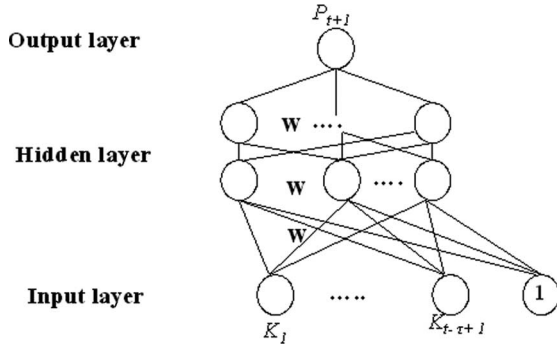


Fig. 1. Design of ANN.

203 web prediction, the next action corresponds to predicting the
 204 next page to be visited. The previous actions correspond to
 205 the previous pages that have already been visited. In web
 206 prediction, the K th-order Markov model is the probability that
 207 a user will visit the k th page, provided that he/she has visited
 208 $k - 1$ pages, i.e.,

$$\begin{aligned} & \Pr(P_k | P_{k-1}, \dots, P_{k-n}) \\ &= \Pr(S_k = P_k | S_{k-1} = P_{k-1}, \dots, S_{k-n} = P_{k-n}) \quad (1) \end{aligned}$$

209 where P_i is a web page, and S_i is the corresponding state in the
 210 Markov model state diagram. Notice that the Markov model
 211 cannot predict for a session that does occur in the training set,
 212 because such session will have zero probability. Alternatively,
 213 one can generate all orders of the Markov models and utilize
 214 them in the prediction. This model is called all- K th orders [5],
 215 [7]. The idea here is that, for a given session x of length k , the
 216 k th-order Markov model is used in the prediction. If the k th-
 217 order Markov model cannot predict for x , the $(k - 1)$ th-order
 218 Markov model is considered for prediction using a new session
 219 x' of length $k - 1$. x' is computed by ignoring the first page id in
 220 x . This process repeats until prediction is obtained. Thus, unlike
 221 the basic Markov model, the all- K th orders Markov model can
 222 predict beyond training data, and it fails only when all orders of
 223 basic Markov models fail to predict.

224 B. ANNs

225 ANN is a very powerful and robust classification technique
 226 that has been used in many applications and domains [15]. In
 227 this paper, we employ a network of two layers that uses the
 228 backpropagation algorithm for learning. The backpropagation
 229 algorithm attempts to minimize the squared-error function.

230 A typical training example in web prediction is $\langle [k_{t-\tau+1},$
 231 $\dots, k_{t-1}, k_t]^T, d \rangle$, where $[k_{t-\tau+1}, \dots, k_{t-1}, k_t]^T$ is the input
 232 to the ANN and d is the target web page. Notice that the
 233 input units of the ANN in Fig. 1 are τ previous pages that the
 234 user has recently visited, i.e., $[k_{t-\tau+1}, \dots, k_{t-1}, k_t]^T$, where
 235 k is a web page id. The output of the network is a Boolean
 236 value and not probability. We approximate the probability of the
 237 output by fitting a sigmoid function after the ANN output (see
 238 Section V-A for details). The approximated probabilistic output

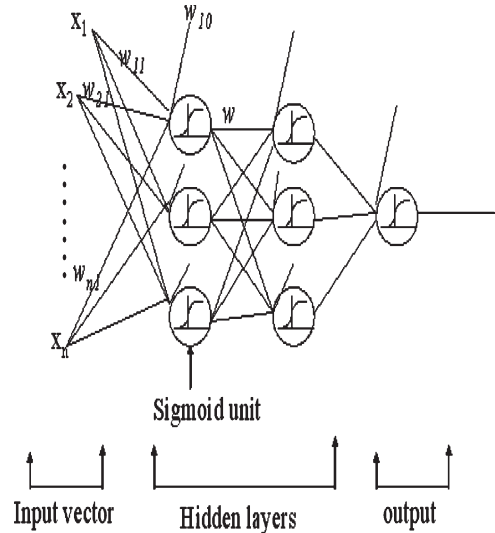


Fig. 2. ANN design in our implementation.

becomes $o' = f(o(I)) = p_{t+1}$, where I is an input session and
 239 $p_{t+1} = p(d | k_{t-\tau+1}, \dots, k_t)$. We choose the sigmoid function
 240

$$o = \sigma(w \cdot I) \quad \sigma(y) = \frac{1}{1 + e^{-y}} \quad (2)$$

as a transfer function, so that the ANN can handle nonlinearly
 241 separable data sets [15].

242 In (2), I is the input to the network, O is the output of the
 243 network, W is the matrix of weights, and σ is the sigmoid
 244 function. We implement the backpropagation algorithm for
 245 training the weights. The backpropagation algorithm employs
 246 gradient descent to attempt to minimize the squared error
 247 between network output values and the target values of these
 248 outputs. In our implementation, we set the step size, to update
 249 the ANN weights, dynamically based on the distribution of
 250 the classes in the data set. First, we set the step size to large
 251 values when updating the training examples that belong to low-
 252 distribution class and vice versa. This is because, when the
 253 distribution of the classes in the data set varies widely (for
 254 example, positive examples are equal to 10% and negative
 255 examples are equal to 90%), the network weights converge
 256 toward the examples from the class of larger distribution, which
 257 causes a slow convergence. Second, we adjust the learning
 258 rates slightly by applying a momentum constant to speed up
 259 the convergence of the network. Fig. 2 presents our multilayer
 260 ANN design that we use in our experiments. As we can see, the
 261 ANN is composed of two fully connected hidden layers. Each
 262 layer is composed of three neurons.
 263

264 IV. DOMAIN KNOWLEDGE AND 265 CLASSIFICATION REDUCTION

266 In web prediction, the number of classes/labels is large. Each
 267 page id is considered as a different label/class. For example,
 268 in our data set, we have 11 700 different page ids. Recall
 269 that, when using one-versus-one or one-versus-all, we have to
 270 consult many classifiers to resolve prediction. As a result, pre-
 271 diction time may increase, conflict can arise among classifiers,
 272

	1	2	...	N
1	0	$freq(1,2)$...	$freq(1,N)$
2	$freq(2,1)$	0	...	$freq(2,N)$
...	$freq(.,1)$	$freq(.,2)$...	$freq(.,N)$
N	0	$freq(N,2)$...	0

Fig. 3. Frequency matrix.

272 and prediction accuracy becomes low. One way to reduce/filter
273 this large number of outcomes is to use domain knowledge in
274 what we call *frequency matrix*. Frequency matrix is defined as
275 an $N \times N$ matrix, where N is the number of web pages (see
276 Fig. 3). The first row and column represent the enumeration of
277 web page ids. Each entry in the matrix represents the frequency
278 that the users have visited two pages in a sequence. For exam-
279 ple, entry (1, 2) in Fig. 3 contains the frequency of users who
280 have visited page 2 after 1. Notice that $freq(x, x)$ is always
281 zero. We can use the frequency matrix to eliminate/filter the
282 number of classifiers during prediction as follows: For a given
283 session $X = \langle x_1, x_2, \dots, x_n \rangle$ and a classifier C_i , we exclude
284 C_i in the prediction process if $freq(x_n, c_i) = 0$, where x_n is
285 the last page id that the user has visited in testing session X .

286 The frequency matrix represents the first order of Markov
287 model. One can extend that to a higher order frequency matrix.
288 In this case, an n th-order frequency matrix corresponds to the
289 n th-order Markov model. Notice that the increase of frequency
290 matrix order leads to fewer number of classes in prediction.
291 For example, given a testing session $S_3 = \langle p_1, p_2, p_3 \rangle$, the
292 following relation holds:

$$|B_1| \geq |B_2| \geq |B_3|$$

293 where

$$B_1 = \{x | \langle x \rangle < p_3, x \in T\}$$

$$B_2 = \{x | \langle x \rangle < p_2, p_3, x \in T\}$$

$$B_3 = \{x | \langle x \rangle < p_1, p_2, p_3, x \in T\}$$

294 where T is the training sessions, x is a page id, B_i is the set of
295 outcomes by applying a frequency matrix of order i , and $|B_i|$
296 is the length of set B_i . Hence, there is a tradeoff between the
297 number of classifiers during prediction (i.e., accuracy) and the
298 order of frequency matrix. Based on our observations and ex-
299 periments, we find that first-order frequency matrix is adequate
300 to balance such tradeoff and to reduce the number of classifiers
301 in prediction without affecting the accuracy. (See Section VI-D
302 for details.)

303 V. HYBRID MODEL FOR WEB PREDICTION 304 USING DEMPSTER'S RULE

305 In this section, we present our hybrid model for web predic-
306 tion, which is based on Dempster's rule for evidence combina-
307 tion, using the ANN and Markov models as bodies of evidence.
308 In our model, prediction is resolved by fusing two separate
309 classifiers models, namely: 1) ANN and 2) Markov model (see
310 Fig. 4).

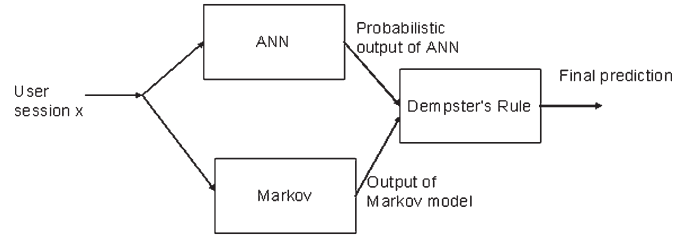


Fig. 4. Hybrid model using the Dempster's rule for evidence combination.

Dempster's rule is one part of the Dempster–Shafer Evidence
311 Combination frame for fusing independent bodies of evidence. 312
313 In Dempster's rule, the sources of evidence should be in the
314 form of basic probability. Since the ANN output value is not
315 probability [15], we need to convert/map this output into a
316 posterior probability $P(class|input)$. 316

317 In this section, we will present, first, a method to convert the
318 ANN output into a posterior probability by fitting a sigmoid
319 function after the ANN output [16]. Next, we present the
320 background of the Dempster–Shafer theory. 320

321 A. Fitting Sigmoid After ANN Output

322 We have implemented the backpropagation learning algo-
323 rithm based on minimizing the squared-error function. Hence,
324 the output of ANN cannot be considered probability. Since
325 we are using ANN as an independent body of evidence in the
326 Dempster's rule frame, we should consequently map the output
327 of ANN into probability. 327

328 One interpretation of the output of ANN, in the context of
329 the classification problem, is an estimate of the probability
330 distribution. There are several ways to interpret the ANN output
331 in terms of probability. One traditional way is to estimate
332 the probability density function (pdf) from the training data.
333 The assumption here is that we know that the training data
334 follow some distribution (typically the normal distribution).
335 The normal distribution is widely used as a model parameter
336 in which analytical techniques can be applied to estimate such
337 parameters [19], [22] as mean and standard deviation. 337

338 Another approach is to consider learning to minimize a
339 probabilistic function, instead of squared error, such as the cross
340 entropy shown in (3). Once learning is done, the output of the
341 network is an estimate of the pdf. In (3), D is the training set,
342 t_d is the target class of example d , and o_d is the output of
343 ANN, i.e., 343

$$\min - \sum_{d \in D} t_d \log(o_d) + (1 - t_d) \log(1 - o_d). \quad (3)$$

344 Since the backpropagation algorithm minimizes the squared-
345 error function, we choose to implement a parametric method to
346 fit the posterior $p(y = 1|f)$ directly, instead of estimating the
347 class-conditional densities $p(f|y)$ [16], where y is the target
348 class and f is the output function of ANN. The output of ANN
349 is computed as follows: 349

$$f(I) = \begin{cases} 1, & \text{if } \sigma(I) \geq 0.5 \\ -1, & \text{otherwise} \end{cases} \quad (4)$$

350 where I is the input to the network, and σ is the output of
 351 the sigmoid transfer function defined as in (2). It follows that
 352 class-conditional densities between the margins are apparently
 353 exponential [16]. Bayes' rule on two exponential suggests using
 354 a parametric form of sigmoid as follows:

$$P(y = 1|f) = \frac{1}{1 + \exp(Af + B)}. \quad (5)$$

355 This sigmoid model is equivalent to assuming that the output
 356 of ANN is proportional to the log odds of a positive example.
 357 Parameters A and B of (5) are fitted using maximum-likelihood
 358 estimation and can be found by minimizing the negative log
 359 likelihood of training data, which is a cross-entropy error
 360 function (3). o_d in (3) is defined as follows:

$$o_d = \frac{1}{1 + \exp(Af_d + B)}. \quad (6)$$

361 The minimization in (3) is a two-parameter minimization.
 362 Hence, it can be performed in many optimization algorithms.
 363 For robustness, we implement the model-trust minimization
 364 algorithm based on the Levenberg–Marquardt algorithm [17].

365 B. Dempster–Shafer Evidence Combination

366 The Dempster–Shafer theory is a mathematical theory
 367 of evidence [17], which is considered to be a generaliza-
 368 tion of the Bayesian theory of subjective probability. The
 369 Dempster–Shafer theory is based on two ideas. The first idea
 370 is the notion of obtaining degrees of belief for one question
 371 based on subjective probabilities for a related question, and
 372 Dempster's rule for combining such degree of belief when
 373 they are based on independent items of evidence. Since we
 374 use two independent sources of evidence, namely, ANN and
 375 Markov model, we are interested in the latter part of the
 376 Dempster–Shafer theory, namely, Dempster's rule. See [17] for
 377 more details regarding this theory.

378 Some may question why we do not use boosting and bagging
 379 rather than Dempster's rule to improve the classifier accuracy.
 380 We prefer Dempster's rule over boosting and bagging because
 381 boosting and bagging require partitioning the data set into a
 382 large number of independent bootstrap samples (> 1000) and
 383 then generating a classifier for each partition separately [12].
 384 Hence, there is a computation overhead not only in training
 385 but also in preprocessing and prediction. Furthermore, boosting
 386 and bagging cannot perform effectively if the data set does
 387 not have enough points for each class/label. In web prediction
 388 applications, many pages are sparse in the data set, because they
 389 receive very few clicks.

C. Using Dempster–Shafer Theory in Web Prediction

390

We have two sources of evidence: 1) the output of ANN 391
 and 2) the output of Markov model. These two models operate 392
 independently. Furthermore, we assume that, for any session x 393
 for which it does not appear in the Markov model, the Markov 394
 prediction probability is zero. If we use Dempster's rule for 395
 combination of evidence, we get the following: 396

$$m_{\text{ANN,Markov}}(C) = \frac{\sum_{A, B \subseteq \Theta, A \cap B = C} m_{\text{ANN}}(A) m_{\text{Markov}}(B)}{\sum_{A, B \subseteq \Theta, A \cap B \neq \phi} m_{\text{ANN}}(A) m_{\text{Markov}}(B)} \quad (7)$$

where m_{ANN} is the probabilistic output of ANN, m_{Markov} is the 397
 output of Markov model, $C \in 2^\Theta$ is a hypothesis (for example, 398
 what is the prediction of a testing session?), and Θ is the frame 399
 of discernment. A frame of discernment is an exhaustive set of 400
 mutually exclusive elements (hypothesis, propositions). 401

In web prediction, we can simplify this formulation because 402
 we have only beliefs for singleton classes (i.e., the final pre- 403
 diction is only one web page and it should not have more 404
 than one page) and the body of evidence itself ($m(\Theta)$). This 405
 means that, for any proper subset A of Θ for which we have no 406
 specific belief, $m(A) = 0$. After eliminating zero terms, we get 407
 the simplified Dempster's combination rule, for a web page P_C 408
 in (8) which is shown at the bottom of the page. Since we are 409
 interested in ranking the hypothesis, we can further simplify 410
 (8), where the denominator is independent of any particular 411
 hypothesis, as follows: 412

$$\begin{aligned} \text{rank}(P_C) &\propto m_{\text{ANN}}(P_C) m_{\text{Markov}}(P_C) \\ &+ m_{\text{ANN}}(P_C) m_{\text{Markov}}(\Theta) + m_{\text{ANN}}(\Theta) m_{\text{Markov}}(P_C). \end{aligned} \quad (9)$$

\propto is the “is proportional to” relationship. $m_{\text{ANN}}(\Theta)$ and 413
 $m_{\text{Markov}}(\Theta)$ represent the uncertainty in the bodies of evi- 414
 dence for m_{ANN} and m_{Markov} , respectively. For $m_{\text{ANN}}(\Theta)$ 415
 and $m_{\text{Markov}}(\Theta)$ in (9), we use the following. For ANN, we 416
 use the output of ANN to compute the uncertainty. We call the 417
 output of ANN for specific session x as the margin because 418
 the ANN weights correspond to the separating surface between 419
 classes and the output ANN is the distance from this surface. 420
 Uncertainty is computed based on the maximum margin of all 421
 training examples as follows: 422

$$m_{\text{ANN}}(\Theta) = \frac{1}{\ln(e + \text{ANN}_{\text{margin}})}. \quad (10)$$

$\text{ANN}_{\text{margin}}$ is the maximum distance of training examples 423
 from the margin. For Markov model uncertainty, we use the 424

$$m_{\text{ANN,Markov}}(P_C) = \frac{m_{\text{ANN}}(P_C) m_{\text{Markov}}(P_C) + m_{\text{ANN}}(P_C) m_{\text{Markov}}(\Theta) + m_{\text{ANN}}(\Theta) m_{\text{Markov}}(P_C)}{\sum_{A, B \subseteq \Theta, A \cap B \neq \phi} m_{\text{ANN}}(A) m_{\text{Markov}}(B)} \quad (8)$$

425 maximum probability of training examples as follows:

$$m_{\text{Markov}}(\Theta) = \frac{1}{\ln(e + \text{Markov}_{\text{probability}})}. \quad (11)$$

426 $\text{Markov}_{\text{probability}}$ is the maximum probability of training ex-
427 amples. Note that, in both models, the uncertainty is inversely
428 proportional to the corresponding maximum value.

429 Here, we would like to show the basic steps that are involved
430 in web prediction using Dempster's rule.

431 Step 1) Train ANN (see Section III-B).

432 Step 2) Map ANN output a probability (see Section V-A).

433 Step 3) Compute Uncertainty (ANN) (see Section V-C,
434 (10)).

435 Step 4) Construct Markov Model (see Section III-A).

436 Step 5) Compute Uncertainty of Markov model (see
437 Section V-C, (11)).

438 Step 6) **For each** testing session x , **do**

439 Step 6.1) Compute $m_{\text{ANN}}(x)$ and output ANN
440 probabilities for different pages.

441 Step 6.2) Compute $m_{\text{Markov}}(x)$ and output Markov
442 probability for different pages.

443 Step 6.3) Compute $m_{\text{ANN,Markov}}(x)$ using (9) and
444 output the final prediction.

445 Step 7) Compute prediction accuracy. // see Section VI-C.

446 VI. EVALUATION

447 In this section, we first define the prediction measurements
448 that we use in our results. Second, we present the data set that
449 we use in this paper. Third, we present the experimental setup.
450 Finally, we present our results. In all models, we use the N-gram
451 representation of paths [5], [7].

452 The following definitions will be used in the succeed-
453 ing sections to measure the performance of the prediction.
454 Pitkow and Pirolli [5] have used these parameters to mea-
455 sure the performance of the Markov model. These defin-
456 itions are given as follows: $Pr(\text{match})$ is the probability
457 that a penultimate path that was observed in a validation
458 set was matched by the same penultimate path in the train-
459 ing set. $Pr(\text{hit}|\text{match})$ is the conditional probability that
460 page x_n is correctly predicted for the testing instance $s =$
461 $\langle x_{n-1}, x_{n-2}, \dots, x_{n-k} \rangle$ and s matches a penultimate path in
462 the training set. $Pr(\text{hit})$ is defined as $pr(\text{hit}) = pr(\text{match}) \times$
463 $pr(\text{hit}|\text{match})$. $Pr(\text{miss}|\text{match})$ is the conditional proba-
464 bility that page x_n is incorrectly classified, given that its
465 penultimate path matches a penultimate path in the training
466 set. $Pr(\text{miss})$ is defined as $pr(\text{match}) \times pr(\text{miss}|\text{match})$.
467 Since we are considering the generalization accuracy and
468 the training accuracy, we add two additional measurements
469 that take into account the generalization accuracy, namely,
470 $Pr(\text{hit}|\text{mismatch})$ and overall accuracy. $Pr(\text{hit}|\text{mismatch})$
471 is the conditional probability that page x_n is correctly predicted
472 for the testing instance $s = \langle x_{n-1}, x_{n-2}, \dots, x_{n-k} \rangle$ and s does
473 not match any penultimate path in the training set. The overall
474 accuracy is defined as $pr(\text{hit}|\text{mismatch}) \times pr(\text{mismatch}) +$
475 $pr(\text{hit}|\text{match}) \times pr(\text{match})$. The overall accuracy considers
476 both matching and mismatching testing examples in computing

the accuracy. The following relations hold for the preceding
477 measurements: 478

$$Pr(\text{hit}|\text{match}) = 1 - pr(\text{miss}|\text{match}) \quad (12)$$

$$Pr(\text{hit})/Pr(\text{miss}) = Pr(\text{hit}|\text{match})/Pr(\text{miss}|\text{match}). \quad (13)$$

$Pr(\text{hit}|\text{match})$ corresponds to the training accuracy, be-
479 cause it shows the proportion of training examples that are
480 correctly classified. $Pr(\text{hit}|\text{mismatch})$ corresponds to the
481 generalization accuracy, because it shows the proportion of
482 unobserved examples that are correctly classified. The overall
483 accuracy combines both. 484

485 A. Data Set

486 For equal comparison purposes and in order to avoid dupli-
487 cating already existing work, we have used the data that were
488 collected by Pitkow and Pirolli [5] from Xerox.com for the
489 dates May 10, 1998 and May 13, 1998. Several numbers of at-
490 tributes are collected using the aforementioned method, which
491 includes the Internet Protocol address of the user, time stamp
492 with date and starting time, visiting URL address, referred URL
493 address, and the browser information or agent [20]. 493

494 B. Experimental Setup

495 We have implemented the backpropagation algorithm for
496 multilayer neural network learning. In our experiments, we
497 use a dynamic learning rate setup based on the distribution of
498 the examples from different classes. Specifically, we setup the
499 learning rate inversely to the distribution of the class, i.e., we
500 set the learning rate to a high value for low-distribution class
501 and vice versa. 501

502 In ARM, we generate the rules using the Apriori algorithm
503 proposed in [9]. We set the *minsupp* to a very low value
504 ($\text{minsupp} = 0.0001$) to capture the pages that were rarely
505 visited. We implement the recommendation engine that was
506 proposed by Mobasher *et al.* [7]. We divide the data set into
507 three partitions. Two partitions are used in training, and one
508 partition is used in testing. 508

509 C. Results

510 In this section, we present and compare the results of pre-
511 diction using four different models, namely: 1) ARM; 2) ANN;
512 3) the Markov model; and 4) the hybrid model. In addition, we
513 consider the *All-Kth-Markov model* and *All-Kth-ARM model*.
514 In the following, we will refer to the results of combining
515 the Markov model with ANN as the Dempster's rule and
516 combining ANN with the All-Kth-Markov model as the All-
517 Kth-Dempster's rule. 517

518 We consider up to seven hops in our experiments for all
519 models. Results vary based on the number of hops, because
520 different patterns are revealed for different numbers of hops.
521 Furthermore, we introduce the concept of ranking in our results.
522 Rank n means that prediction is considered to be correct if the
523 predicted page happens to be among the top n pages that has
524 the highest confidence. 524

TABLE I
PROBABILITY MEASUREMENTS USING ONE HOP AND RANK 1

	ARM	Markov	ANN	Dempster-Rule
Pr(Match)	0.590	0.590	0.590	0.590
Pr(Hit Match)	0.063	0.199	0.152	0.192
Pr(Hit)	0.037	0.118	0.09	0.114
Pr(Miss Match)	0.937	0.8	0.847	0.807
Pr(Miss)	0.555	0.474	0.502	0.478
Pr(Hit Mismatch)	0	0	0.108	0.108
Pr(Hit)/Pr(Miss)	0.067	0.249	0.179	0.238
Over all Hit/Miss	0.038	0.134	0.155	0.188
Overall accuracy	0.037	0.118	0.134	0.158

TABLE II
RESULTS OF USING THREE HOPS AND RANK 1

	ARM	All-Kth-ARM	Markov	All-Kth-Markov	ANN	Dempster's-Rule	All-Kth-Dempster's-Rule
Pr(match)	0.376	0.376	0.376	0.376	0.376	0.376	0.376
Pr(hit match)	0.043	0.103	0.231	0.230	0.156	0.232	0.232
Pr(hit)	0.016	0.038	0.087	0.086	0.059	0.087	0.087
Pr(mismatch)	0.956	0.89	0.768	0.769	0.843	0.767	0.767
Pr(miss)	0.360	0.337	0.289	0.289	0.289	0.288	0.288
Pr(hit mismatch)	0	0.044	0	0.105	0.109	0.109	0.147
Pr(hit)/Pr(miss)	0.045	.114	0.3	0.299	0.185	0.302	0.302
Over all hit/miss	0.016	0.071	0.095	0.179	0.145	0.184	0.218
Overall accuracy	0.016	0.066	0.087	0.152	0.127	0.155	0.179

525 In Table I, there are several points to note. First, the value
526 of $pr(hit|mismatch)$ is zero for both ARM and the Markov
527 model, because neither model can predict for the unobserved
528 data. Second, the Dempster's rule achieves the best scores
529 using all measurements. For example, the training accuracy
530 $pr(hit|match)$ for ARM, Markov, ANN, and Dempster's rule
531 is 6%, 19%, 15%, and 19%, respectively. The overall accu-
532 racy for ARM, Markov, ANN, and Dempsters' rule is 3%,
533 11%, 13%, and 15%, respectively. Third, even though the
534 $pr(hit|match)$ for ANN is less than that for the Markov
535 model, the overall accuracy for ANN is better. This is because
536 $pr(hit|mismatch)$ is zero in case of the Markov model, while
537 it is 10% in case of ANN. Finally, notice that ARM has the low-
538 est prediction results. The ARM uses the concept of frequent
539 item sets, instead of *item lists* (ordered item set); hence, the
540 support of one path is computed based on the frequencies of that
541 path and its combinations. In addition, ARM is very sensitive to
542 the *minsupp* values. This might cause important patterns to be
543 lost or mixed. Table II shows the results using three hops and
544 rank 1. Notice that All- K th-Markov outperforms the Markov
545 model, and the All- K th-ARM outperforms the ARM model.
546 That is because lower orders of the models are consulted in
547 case prediction is not possible for higher orders. As a result
548 of this, $pr(hit|mismatch)$ is not zero in such models. For
549 example, the $pr(hit|mismatch)$ for All- K th-Markov and All-
550 K th-ARM are 10.5% and 4.4%, respectively. In addition, com-
551 bining the All- K th-Markov model with ANN using Dempster's
552 rule has boosted the final prediction; for example, the overall
553 accuracy for ARM, All- K th-ARM, Markov, All- K th-Markov,
554 ANN, Dempster's rule, and All- K th-Dempster's rule is 1.6%,
555 6.6%, 8.7%, 15.2%, 12.7%, 15.5%, and 17.9%, respectively.

556 In Figs. 5 and 6, the accuracy approximately increases lin-
557 early with the rank. For example, in Fig. 5, the $pr(hit|match)$
558 for All- K th-Markov is 23%, 29%, 34%, 38%, 42%, 46%, 50%,
559 and 54% for ranks 1–8, respectively. In Fig. 6, the overall

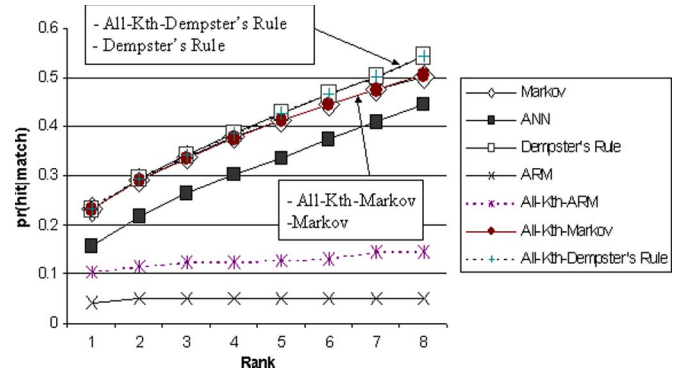


Fig. 5. $pr(hit|match)$ results using three-hop sessions and ranks from 1 to 8.

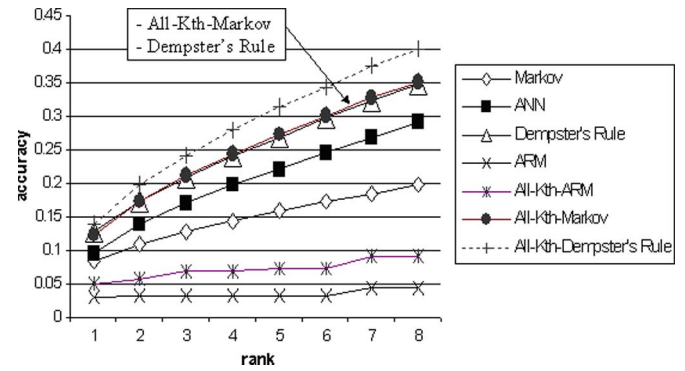


Fig. 6. Overall accuracy results using two-hop sessions and ranks from 1 to 8.

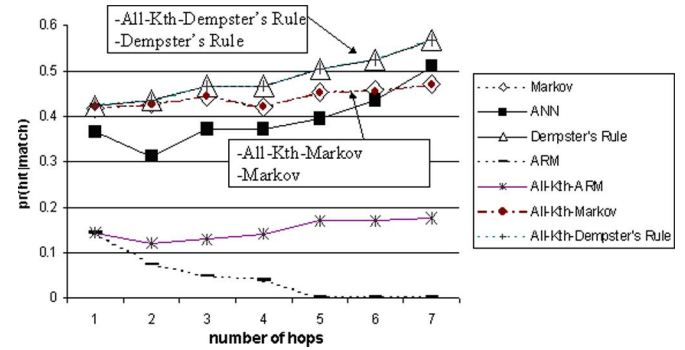


Fig. 7. Comparable results of all techniques based on $pr(hit|match)$ using rank 6.

accuracy of All- K th-Markov models is 12, 17, 21, 24, 27, 30, 560
and 35 for ranks 1–8, respectively. 561

Fig. 7 presents $pr(hit|match)$ results of using rank 6. Notice
562 that the Dempster's rule and All- K th-Dempster's rule methods
563 outperform all other techniques. 564

In Fig. 8, we notice that All- K th-Dempster's rule has
565 achieved the best overall accuracy, because it combines ANN
566 and the All- K th-Markov model. Both models have a high train-
567 ing and generalization accuracy. For example, the overall accu-
568 racy using four hops for Markov, ANN, Dempster's rule, ARM,
569 All- K th-ARM, All- K th-Markov, and All- K th-Dempster's
570 rule is 7%, 11%, 13%, 1%, 6%, 15%, and 16%, respectively. 571
In addition, notice that ANN has outperformed the Markov
572 model based on overall accuracy. This is because ANN gen-
573 eralizes better than the Markov model beyond training data. 574

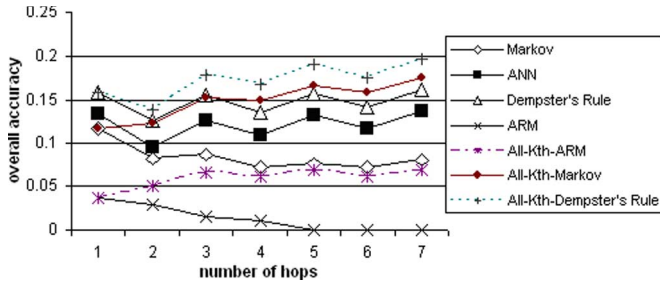


Fig. 8. Comparable results based on the overall accuracy using rank 1.

TABLE III
AVERAGE PREDICTION TIME WITH/WITHOUT DOMAIN KNOWLEDGE

Model	Average Prediction Time With Domain Knowledge (milliseconds)	Average Prediction Time Without Domain Knowledge (milliseconds)
Markov	0.567	0.544
All-Kth-Markov	1.17	0.801
ANN	6.41	556.6
Dempster's Rule	1.11	788.12

575 All- K th-Dempster's rule proves to combine the best of both
576 the ANN and All- K th-Markov models, because it has kept
577 its superiority over all techniques using all measurements and
578 using different numbers of hops.

579 D. Effect of Domain Knowledge on Prediction

580 As we mentioned previously in Section IV, we have extended
581 this model to include higher orders of domain knowledge.
582 Recall that a frequency matrix of order n corresponds to a
583 Markov model of order n .

584 Table III shows that the average prediction time using do-
585 main knowledge is 0.567, 1.17, 6.41, and 1.11 ms for the
586 Markov, All- K th-Markov, ANN, and Dempster's rule models,
587 respectively. The average prediction time without using do-
588 main knowledge is 0.544, 0.801, 556.0, and 788.0 ms for the
589 Markov, All- K th-Markov, ANN, and Dempster's rule models,
590 respectively. It is very evident that prediction time is reduced
591 dramatically for ANN and Dempster's rule. The overhead in
592 prediction without domain knowledge is a consequence of
593 loading a very large number of classifiers, i.e., 4563 classifiers,
594 and consulting them to resolve prediction. The prediction time
595 in case of the Markov model and All- K th-Markov has not been
596 affected, because such models, contrary to ANN, can handle a
597 multiclass problem without the used of an on-VS-all model.

598 In part A of Fig. 9, the overall accuracy without using
599 any domain knowledge is 18.4%, 18.4%, 0.5%, and 15.7%
600 for Markov, All- K th-Markov, ANN, and Dempster's rule, re-
601 spectively. The overall accuracy in case of using first-order
602 frequency matrix (DK) is 18.4%, 18.5%, 21.6%, and 24.5%
603 for Markov, All- K th-Markov, ANN, and Dempster's rule, re-
604 spectively. Recall that the domain knowledge is based on the
605 frequency matrix of order n , which is another representation of
606 the n th order of the Markov model; hence, the overall accuracy
607 for the basic knowledge is already included in such models. On
608 the other hand, the performance of ANN and Dempster's rule is

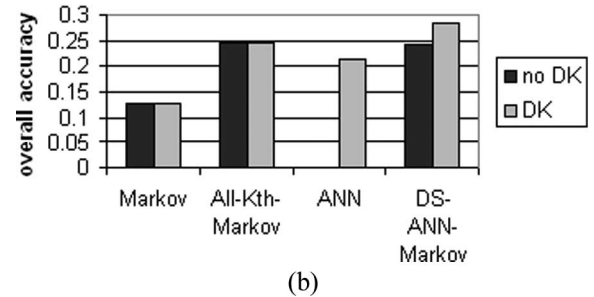
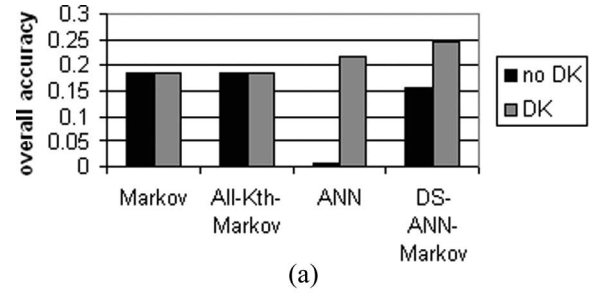


Fig. 9. Comparable results using the overall accuracy with/without domain knowledge. (a) One hop using rank 3. (b) Three hops using rank 3.

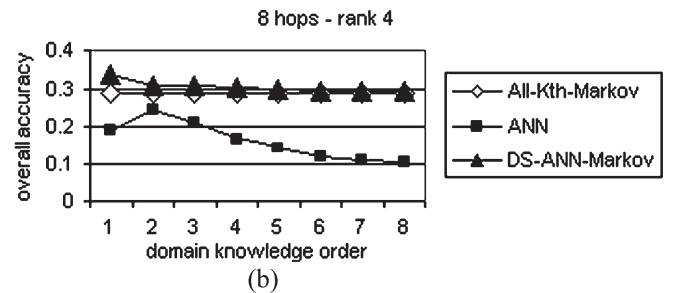
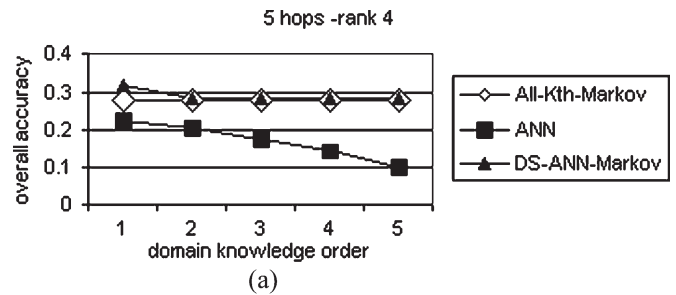


Fig. 10. Effect of domain knowledge order on the overall accuracy. (a) Five hops using rank 4. (b) Eight hops using rank 4.

affected by not using any domain knowledge, and the overall
609 accuracy has dropped largely. Similar results can be seen in
610 Fig. 9(b) when using three hops and rank 4. Fig. 10 presents
611 the effect of using different orders of domain knowledge on
612 the overall accuracy. Since we obtained similar results when
613 using different rankings and different number of hops, we
614 only show the results for five and eight hops using rank 4.
615 Recall that, in the previous experiments, we considered only
616 the first-order frequency matrix. Here, we consider a frequency
617 matrix of different orders as domain knowledge applied to the
618 All- K th-Markov model, ANN, and Dempster's rule. The three
619 curves (from top to bottom) in each subfigure represent the
620 overall accuracy of All- K th-Markov, ANN, and Dempster's
621 rule. For example, the overall accuracy for Dempster's rule
622

623 is 31%, 28%, 28%, 28%, and 28% using domain knowledge
 624 of orders 1, 2, 3, 4, and 5, respectively. Notice that the use
 625 of higher order for domain knowledge did not improve the
 626 accuracy. On the contrary, it slightly affects the overall accu-
 627 racy negatively. This can be related to the tradeoff between
 628 the number of classifiers to consult and the order of domain
 629 knowledge. Using higher order domain knowledge leads to
 630 less number of classes to consult. This may positively affect
 631 the accuracy and speed up the retrieval process. However,
 632 this might exclude correct classes, decrease the accuracy, and
 633 finally offsets the improvement of accuracy. Conversely, not
 634 using domain knowledge leads to consulting a huge number of
 635 classifiers that cause conflict. From Fig. 10, we find that using
 636 domain knowledge of order 1 or 2 can balance such tradeoffs,
 637 because accuracy is not affected dramatically.

638 VII. CONCLUSION AND FUTURE WORK

639 In this paper, we use a hybrid method in web prediction based
 640 on Dempster's rule for evidence combination to improve pre-
 641 diction accuracy. We used two sources of evidence/prediction in
 642 our hybrid method. The first body of evidence is ANNs. To im-
 643 prove the prediction of ANN further, we incorporated different
 644 orders of domain knowledge in prediction to improve prediction
 645 accuracy. The second body of evidence is the widely used
 646 Markov model, which is a probabilistic model. Furthermore,
 647 we applied the *All-Kth-Markov* model. The *All-Kth-Dempster's*
 648 *rule* proves its effectiveness by combining the best of ANN and
 649 the *All-Kth-Markov* model, as demonstrated by the fact that its
 650 predictive accuracy has outperformed all other techniques.

651 We would like to extend our research in the following direc-
 652 tions. First, we would like to study the impact/effect of other
 653 features in the session's logs by extracting statistical features
 654 from the data set to improve accuracy. Next, we would like
 655 to perform more experiments and analyses on the effect of the
 656 frequency matrix order on prediction. Finally, we would like to
 657 use boosting and bagging in the same context, and compare it
 658 with our hybrid approach.

659 REFERENCES

- 660 [1] M. Levene and G. Loizou, "Computing the entropy of user navigation in
 661 the web," *Int. J. Inf. Technol. Decis. Making*, vol. 2, no. 3, pp. 459–476,
 662 Sep. 2003.
- 663 [2] Q. Yang, H. Zhang, and T. Li, "Mining web logs for prediction models in
 664 WWW caching and prefetching," in *Proc. 7th ACM SIGKDD Int. Conf.*
 665 *KDD*, Aug. 26–29, 2001, pp. 473–478.
- 666 [3] K. Chinen and S. Yamaguchi, "An interactive prefetching proxy server for
 667 improvement of WWW latency," in *Proc. 7th Annu. Conf. INET*, Kuala
 668 Lumpur, Malaysia, Jun. 1997.
- 669 [4] V. Chung, C. H. Li, and J. Kwok, "Dissimilarity learning for nominal
 670 data," *Pattern Recognit.*, vol. 37, no. 7, pp. 1471–1477, Jul. 2004.
- 671 [5] J. Pitkow and P. Pirulli, "Mining longest repeating subsequences to predict
 672 World Wide Web surfing," in *Proc. 2nd USITS*, Boulder, CO, Oct. 1999,
 673 pp. 139–150.
- 674 [6] J. Griffioen and R. Appleton, "Reducing file system latency using a
 675 predictive approach," in *Proc. Summer USENIX Tech. Conf.*, Cambridge,
 676 MA, 1994, pp. 197–207.
- 677 [7] B. Mobasher, H. Dai, T. Luo, and M. Nakagawa, "Effective personaliza-
 678 tion based on association rule discovery from web usage data," in *Proc.*
 679 *ACM Workshop WIDM*, Atlanta, GA, Nov. 2001, pp. 9–15. Held at CIKM.
- 680 [8] R. Agrawal, T. Imielinski, and A. Swami, "Mining association rules be-
 681 tween sets of items in large databases," in *Proc. ACM SIGMOD Conf.*
 682 *Manage. Data*, Washington, DC, May 1993, pp. 207–216.

- [9] D. Duchamp, "Prefetching hyperlinks," in *Proc. 2nd USITS*, Boulder, CO, 683
 Oct. 1999, pp. 127–138. 684
- [10] A. Pandey, J. Srivastava, and S. Shekhar, "A web intelligent prefetcher for 685
 dynamic pages using association rules—A summary of results," in *Proc.* 686
SIAM Workshop Web Mining, Report No. 01-004, 2001. 687
- [11] Z. Su, Q. Yang, Y. Lu, and H. Zhang, "Whatnext: A prediction system for 688
 web requests using N-gram sequence models," in *Proc. 1st Int. Conf. Web* 689
Inf. Syst. Eng. Conf., Hong Kong, Jun. 2000, pp. 200–207. 690
- [12] L. Breiman, "Bagging predictors," *Mach. Learn.*, vol. 24, no. 2, pp. 123– 691
 140, Aug. 1996. 692
- [13] B. M. Sarwar, G. Karypis, J. Konstan, and J. Riedl, "Analysis of recom- 693
 mender algorithms for e-commerce," in *Proc. 2nd ACM EC*, Minneapolis, 694
 MN, Oct. 2000, pp. 158–167. 695
- [14] S. Brin and L. Page, "The anatomy of a large-scale hypertextual web 696
 search engine," in *Proc. 7th Int. WWW Conf.*, Brisbane, Australia, 1998, 697
 pp. 107–117. 698
- [15] E. Parzen, "On the estimation of a probability density function and mode," 699
Ann. Math. Stat., vol. 33, no. 3, pp. 1065–1076, Sep. 1962. 700
- [16] J. Platt, "Probabilistic outputs for SVMs and comparisons to regular- 701
 ized likelihood methods," in *Advances in Large Margin Classifiers*. 702
 Cambridge, MA: MIT Press, 1999. 703
- [17] G. Shafer, *A Mathematical Theory of Evidence*. Princeton, NJ: Princeton 704
 Univ. Press, 1976. 705
- [18] T. Joachims, D. Freitag, and T. Mitchell, "WebWatcher: A tour guide for 706
 the World Wide Web," in *Proc. IJCAI*, 1997, pp. 770–777. 707
- [19] O. Nasraoui and R. Krishnapuram, "An evolutionary approach to min- 708
 ing robust multi-resolution web profiles and context sensitive URL as- 709
 sociations," *Int. J. Comput. Intell. Appl.*, vol. 2, no. 3, pp. 339–348, 710
 2002. 711
- [20] R. Cooley, B. Mobasher, and J. Srivastava, "Data preparation for mining 712
 World Wide Web browsing patterns," *J. Knowl. Inf. Syst.*, vol. 1, no. 1, 713
 pp. 5–32, 1999. 714
- [21] T. Mitchell, *Machine Learning*. New York: McGraw-Hill, 1997. 715
- [22] W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery, *Numer- 716
 ical Recipes in C: The Art of Scientific Computing*, 2nd ed. Cambridge, 717
 U.K.: Cambridge Univ. Press, 1992. 718



Mamoun A. Awad received the B.Sc. degree in 719
 computer science from Baghdad University, Bagh- 720
 dad, Iraq, in June 1994, the M.S. degree in com- 721
 puter science from Wichita State University, Wichita, 722
 Kansas, in May 1999, and the Ph.D. degree in soft- 723
 ware engineering from the University of Texas at 724
 Dallas, Richardson, in December 2005. 725

He is currently an Assistant Professor in the Infor- 726
 mation Technology College, United Arab Emirates 727
 University, Al Ain, where he has taught and con- 728
 ducted research since August 2006. He has already 729

published 15 conference proceeding papers, book chapters, and journal articles. 730
 His current areas of research are data mining, intrusion detection, and Web 731
 prediction. 732



Latifur R. Khan received the B.Sc. degree in com- 733
 puter science and engineering from the Bangladesh 734
 University of Engineering and Technology, Dhaka, 735
 Bangladesh, in November 1993 and the M.S. 736
 and Ph.D. degrees in computer science from the 737
 University of Southern California, Los Angeles, in 738
 December 1996 and August 2000, respectively. 739

He is currently an Associate Professor in the De- 740
 partment of Computer Science, University of Texas 741
 at Dallas (UTD), Richardson, where he has taught 742
 and conducted research since September 2000, and 743

the Director of the UTD Database and Data Mining Laboratory. He is currently 744
 on the editorial board of North Holland's *Computer Standards and Interface* 745
Journal, published by Elsevier. He has published over 80 papers in presti- 746
 gious journals and conference proceedings. His research interests include data 747
 mining, multimedia information management, and semantic Web and database 748
 systems. 749

Dr. Khan was a Committee Member in numerous prestigious conferences, 750
 symposiums, and workshops, including the ACM SIGKDD Conference on 751
 Knowledge Discovery and Data Mining. 752

Web Navigation Prediction Using Multiple Evidence Combination and Domain Knowledge

Mamoun A. Awad and Latifur R. Khan

Abstract—Predicting users' future requests in the World Wide Web can be applied effectively in many important applications, such as web search, latency reduction, and personalization systems. Such application has traditional tradeoffs between modeling complexity and prediction accuracy. In this paper, we study several hybrid models that combine different classification techniques, namely, Markov models, artificial neural networks (ANNs), and the All- K th-Markov model, to resolve prediction using Dempster's rule. Such fusion overcomes the inability of the Markov model in predicting beyond the training data, as well as boosts the accuracy of ANN, particularly, when dealing with a large number of classes. We also employ a reduction technique, which uses domain knowledge, to reduce the number of classifiers to improve the predictive accuracy and the prediction time of ANNs. We demonstrate the effectiveness of our hybrid models by comparing our results with widely used techniques, namely, the Markov model, the All- K th-Markov model, and association rule mining, based on a benchmark data set.

Index Terms—Artificial neural networks (ANNs), association rule mining (ARM), Dempster's rule, Markov model, N-gram.

I. INTRODUCTION

WEB PREDICTION is the problem of predicting the next web page that a user might visit after surfing in a website. The importance of web prediction originates from the fact that various applications, such as latency reduction, web search, and recommendation systems, can be made more effective through the use and the improvement of web prediction.

One of the early applications of web prediction is the latency of viewing of web documents [6]. Traditional solutions are based on caching and prefetching [2], [3], [9]. Other advanced intelligent methods [10], [11] acquire knowledge from surfers' previous path history and utilize that in prediction. Pandey *et al.* [10] present an intelligent prefetching method based on a proxy server using association rule mining (ARM) to generate association rules that are later used to predict future requests.

World Wide Web (WWW) prediction can also improve search engines. The entire structure of the WWW can be pictured as a connected graph, where each node corresponds to a website, and surfers navigate from one node to another. The distribution of the visits over all WWW pages can be computed

and used in reweighting and reranking results. In such scenario, we consider the surfer path information to be more important than the keywords that were entered by the user [14].

Another application of web prediction is recommendation systems, in which we try to find the top k users having the same interests or tastes to a target user record. ARM is a well-known model that is used in recommendation systems. Mobasher *et al.* [7] propose the frequent item set graph to match an active user session with frequent item sets and predict the next page that the user is likely to visit. Prediction for the active session is based on the confidence of the corresponding association rule.

Other prediction models that are widely used in WWW predictions and its related applications include k nearest neighbors (NN), artificial neural network (ANN), fuzzy interference, and Markov model. Joachims *et al.* [18] propose the k NN-based recommender WebWatcher. The WebWatcher is a learning tour guide agent that extracts knowledge from user's previous clicks and from the hypertext structure. Nasraoui and Krishnapuram [19] propose a web recommendation system using fuzzy inference. Clustering is applied to group profiles using hierarchical unsupervised niche clustering. Context-sensitive uniform resource locator (URL) association is inferred using a fuzzy approximate-reasoning-based engine. Levene and Loizou [1] compute the information gain from the navigation trail to construct a Markov chain model to analyze user navigation pattern through the web. The main contribution of [1] is that they present a mechanism to estimate the navigation trail. Pitkow and Pirolli [5] explore pattern extraction and pattern matching based on the Markov model that predicts future surfing paths. Longest repeating subsequences (LRS) is proposed to reduce the model complexity (not predictive accuracy) by focusing on significant surfing patterns.

Our work is related to the path-based prediction model using the N-gram model [14] and the LRS model [5]. However, our approach differs from them in the following ways: First, only one path-based prediction technique is used when combining different N-gram models. Second, the main focus of LRS is to reduce the modeling complexity by reducing the data set. Third, all these models are probabilistic, i.e., it depends on the frequencies of patterns/occurrences in the training set. Therefore, our model can predict some values that Markov models cannot (i.e., our model can predict for some unobserved values). In the work of Nasraoui and Krishnapuram [19], the focus is to use a set of URL predictors by creating a neural network for each profile independently with a separate training set. Their goal is to overcome the high complexity of the architecture and training in case that one neural network is used. In our work, we not only use a set of predictors but also fuse them

Manuscript received October 17, 2007; revised February 21, 2007. This paper was recommended by Associate Editor M. Kamel.

M. A. Awad was with the University of Texas at Dallas, Richardson, TX 75083-0688 USA. He is now with the United Arab Emirates University, Al Ain, UAE (e-mail: mamoun.awad@uaeu.ac.ae).

L. R. Khan is with the Department of Computer Science, University of Texas at Dallas, Richardson, TX 75083-0688 USA (e-mail: lkhan@utdallas.edu).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TSMCA.2007.904781

92 in a hybrid model for prediction. Our goal is to improve the
93 accuracy using different prediction techniques, namely, ANN,
94 the Markov model, the All- K th model, and using different
95 N-gram models.

96 One important subtlety of web prediction is that web pre-
97 diction is a multiclass problem of a large number of classes
98 (11 700 classes in our experiments). Here, we define a class (or
99 a label) as a unique identifier that represents a web page in a
100 web site. Most multiclass techniques, such as one-versus-one
101 and one-versus-all, are based on generalizing binary classifiers,
102 and prediction is resolved by checking against all these binary
103 classifiers. As a result of that, prediction accuracy is very
104 low [4], because the prediction model has many conflicting
105 outcomes from the classifiers.

106 There are several problems with the current state-of-the-art
107 solutions. First, models such as Markov and ARM models are
108 unable to generalize beyond training data [5]. This is because
109 prediction using ARM and LRS pattern extraction is done based
110 on choosing the path of the highest probability in the training
111 set; hence, any new surfing path is misclassified, because it has
112 zero probability. Second, prediction using ARM suffers from
113 well-known limitations including scalability and efficiency [7],
114 [13]. Finally, many of the previous methods have ignored
115 domain knowledge as a means to improving prediction.

116 In this paper, we present a new approach to improving the ac-
117 curacy in web prediction. Our approach is based on generating
118 a hybrid prediction model by fusing two different classification
119 models. We use four classification models, namely: 1) ANNs;
120 2) ARM; 3) Markov model; and 4) *All- K th-model*. ARM and
121 Markov model are powerful techniques for predicting seen data,
122 i.e., already observed data; however, they cannot predict beyond
123 training data (see Section III-A). On the other hand, the All-
124 K th model and ANN are powerful techniques that can predict
125 beyond training data. In other words, the ANN and All- K th
126 models can predict some values that the Markov model and
127 ARM cannot. We combine the All- K th model with ANN by
128 fusing their outcomes using Dempster's rule.

129 Nonetheless, when dealing with a large number of classes or
130 when there is a possibility that one instance may belong to many
131 classes, their predictive power may decrease. To overcome
132 these shortcomings, we extract domain knowledge from the
133 training data and incorporate such knowledge during prediction
134 to improve prediction time and accuracy. Specifically, domain
135 knowledge is used to eliminate irrelevant classes and reduce
136 the conflict during prediction. Notice that we combine different
137 prediction models in which each model has different strengths
138 and drawbacks over other models. We strive to overcome
139 major drawbacks in each technique and improve the predictive
140 accuracy for the final hybrid model.

141 The contribution of this paper is given as follows: First, we
142 use ANN in web navigation. Second, we incorporate domain
143 knowledge in ANN prediction to eliminate irrelevant classes
144 and to improve prediction time and accuracy. Third, we fuse
145 ANN, the Markov model, and All- K th-Markov classifiers in
146 a hybrid prediction model using Dempster's rule [17] to im-
147 prove prediction accuracy and to overcome the drawbacks of
148 using each model separately. Finally, we compare our hybrid
149 model with different models, namely, Markov model, ARM,

All- K th-ARM, All- K th-Markov, and ANN using a standard 150
benchmark data set and demonstrate the superiority of our 151
method. 152

The organization of this paper is given as follows: In 153
Section II, we present the background of the N-gram concept 154
and sliding window. In Section III, we present different pre- 155
diction models that are used in web prediction. In Section IV, 156
we present the utilization of domain knowledge to improve 157
prediction. In Section V, we present a new hybrid approach 158
combining ANN, the Markov model, and the All- K th-Markov 159
model in web prediction using Dempster's rule for evidence 160
combination. In Section VI, we compare our results with 161
other methods using a standard benchmark training set. In 162
Section VII, we summarize this paper and outline some future 163
research. 164

II. N-GRAM REPRESENTATION OF PATHS 165

In web prediction, the available source of training data is 166
the users' sessions, which are the user's history of navigation 167
within a period of time. User sessions are extracted from 168
the logs of the web servers, and it contains sequences of 169
pages/clicks that the users have visited, time, data, and the pe- 170
riod of time that the user stays in each page. In web prediction, 171
the best known representation of the training session is the 172
N-gram. N-gram is tuples of the form $\langle X_1, X_2, \dots, X_n \rangle$ that 173
depict sequences of page clicks by a population of users surfing 174
a website. Each component of the N-gram takes a specific page 175
id value that identifies a web page. For example, the N-gram 176
 $\langle X_{10}, X_{21}, X_4, X_{12} \rangle$ depicts the fact that the user has visited 177
pages in the following order: page 10, page 21, page 4, and 178
finally, page 12. 179

Many models further process these N-gram sessions by 180
applying a sliding window to make training instances have 181
the same length [5], [7]. For example, if we apply a sliding 182
window of size 3 on the N-gram $\langle X_{10}, X_{21}, X_4, X_{12}, X_{11} \rangle$, 183
we will have the following 3-gram sessions: $\langle X_{10}, X_{21}, X_4 \rangle$, 184
 $\langle X_{21}, X_4, X_{12} \rangle$, and $\langle X_4, X_{12}, X_{11} \rangle$. In general, the number of 185
additional sessions using sliding window w applied on session 186
 A is $|A| - w + 1$, where $|A|$ is the length of session A . 187

In this paper, we also use the term *number of hops*, which 188
is related to the sliding window. The number of hops for a 189
session of length N is $N - 1$, i.e., the number of clicks (or 190
hops) that the user makes to reach the last page in the session. 191
When applying sliding window of size w , the number of hops 192
in the resulted subsessions is $w - 1$. In the previous example, 193
the number of hops in the resulted 3-gram sessions is 2. 194

III. PREDICTION MODELS 195

In this section, we briefly present various prediction models 196
that have been used in web prediction. First, we present the 197
Markov model; next, we present the ANN model along with 198
improvement modifications. 199

A. Markov Model 200

The basic concept of the Markov model is to predict the 201
next action, depending on the result of previous actions. In 202

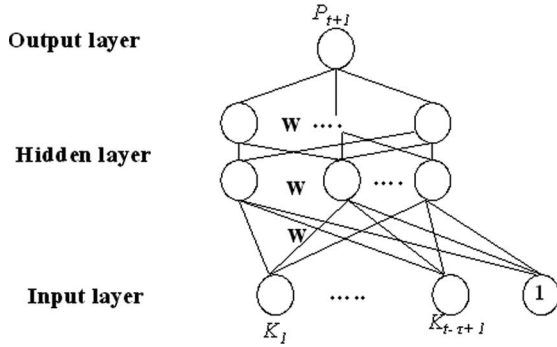


Fig. 1. Design of ANN.

203 web prediction, the next action corresponds to predicting the
 204 next page to be visited. The previous actions correspond to
 205 the previous pages that have already been visited. In web
 206 prediction, the K th-order Markov model is the probability that
 207 a user will visit the k th page, provided that he/she has visited
 208 $k - 1$ pages, i.e.,

$$\begin{aligned} & \Pr(P_k | P_{k-1}, \dots, P_{k-n}) \\ &= \Pr(S_k = P_k | S_{k-1} = P_{k-1}, \dots, S_{k-n} = P_{k-n}) \quad (1) \end{aligned}$$

209 where P_i is a web page, and S_i is the corresponding state in the
 210 Markov model state diagram. Notice that the Markov model
 211 cannot predict for a session that does occur in the training set,
 212 because such session will have zero probability. Alternatively,
 213 one can generate all orders of the Markov models and utilize
 214 them in the prediction. This model is called all- K th orders [5],
 215 [7]. The idea here is that, for a given session x of length k , the
 216 k th-order Markov model is used in the prediction. If the k th-
 217 order Markov model cannot predict for x , the $(k - 1)$ th-order
 218 Markov model is considered for prediction using a new session
 219 x' of length $k - 1$. x' is computed by ignoring the first page id in
 220 x . This process repeats until prediction is obtained. Thus, unlike
 221 the basic Markov model, the all- K th orders Markov model can
 222 predict beyond training data, and it fails only when all orders of
 223 basic Markov models fail to predict.

224 B. ANNs

225 ANN is a very powerful and robust classification technique
 226 that has been used in many applications and domains [15]. In
 227 this paper, we employ a network of two layers that uses the
 228 backpropagation algorithm for learning. The backpropagation
 229 algorithm attempts to minimize the squared-error function.

230 A typical training example in web prediction is $\langle [k_{t-\tau+1},$
 231 $\dots, k_{t-1}, k_t]^T, d \rangle$, where $[k_{t-\tau+1}, \dots, k_{t-1}, k_t]^T$ is the input
 232 to the ANN and d is the target web page. Notice that the
 233 input units of the ANN in Fig. 1 are τ previous pages that the
 234 user has recently visited, i.e., $[k_{t-\tau+1}, \dots, k_{t-1}, k_t]^T$, where
 235 k is a web page id. The output of the network is a Boolean
 236 value and not probability. We approximate the probability of the
 237 output by fitting a sigmoid function after the ANN output (see
 238 Section V-A for details). The approximated probabilistic output

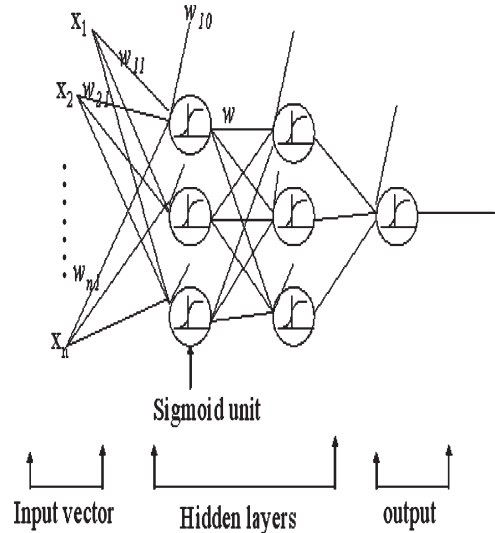


Fig. 2. ANN design in our implementation.

becomes $o' = f(o(I)) = p_{t+1}$, where I is an input session and
 239 $p_{t+1} = p(d | k_{t-\tau+1}, \dots, k_t)$. We choose the sigmoid function
 240

$$o = \sigma(w \cdot I) \quad \sigma(y) = \frac{1}{1 + e^{-y}} \quad (2)$$

as a transfer function, so that the ANN can handle nonlinearly
 241 separable data sets [15].
 242

In (2), I is the input to the network, O is the output of the
 243 network, W is the matrix of weights, and σ is the sigmoid
 244 function. We implement the backpropagation algorithm for
 245 training the weights. The backpropagation algorithm employs
 246 gradient descent to attempt to minimize the squared error
 247 between network output values and the target values of these
 248 outputs. In our implementation, we set the step size, to update
 249 the ANN weights, dynamically based on the distribution of
 250 the classes in the data set. First, we set the step size to large
 251 values when updating the training examples that belong to low-
 252 distribution class and vice versa. This is because, when the
 253 distribution of the classes in the data set varies widely (for
 254 example, positive examples are equal to 10% and negative
 255 examples are equal to 90%), the network weights converge
 256 toward the examples from the class of larger distribution, which
 257 causes a slow convergence. Second, we adjust the learning
 258 rates slightly by applying a momentum constant to speed up
 259 the convergence of the network. Fig. 2 presents our multilayer
 260 ANN design that we use in our experiments. As we can see, the
 261 ANN is composed of two fully connected hidden layers. Each
 262 layer is composed of three neurons.
 263

264 IV. DOMAIN KNOWLEDGE AND 265 CLASSIFICATION REDUCTION

In web prediction, the number of classes/labels is large. Each
 266 page id is considered as a different label/class. For example,
 267 in our data set, we have 11 700 different page ids. Recall
 268 that, when using one-versus-one or one-versus-all, we have to
 269 consult many classifiers to resolve prediction. As a result, pre-
 270 diction time may increase, conflict can arise among classifiers,
 271

	1	2	...	N
1	0	$freq(1,2)$...	$freq(1,N)$
2	$freq(2,1)$	0	...	$freq(2,N)$
...	$freq(.,1)$	$freq(.,2)$...	$freq(.,N)$
N	0	$freq(N,2)$...	0

Fig. 3. Frequency matrix.

272 and prediction accuracy becomes low. One way to reduce/filter
273 this large number of outcomes is to use domain knowledge in
274 what we call *frequency matrix*. Frequency matrix is defined as
275 an $N \times N$ matrix, where N is the number of web pages (see
276 Fig. 3). The first row and column represent the enumeration of
277 web page ids. Each entry in the matrix represents the frequency
278 that the users have visited two pages in a sequence. For exam-
279 ple, entry (1, 2) in Fig. 3 contains the frequency of users who
280 have visited page 2 after 1. Notice that $freq(x, x)$ is always
281 zero. We can use the frequency matrix to eliminate/filter the
282 number of classifiers during prediction as follows: For a given
283 session $X = \langle x_1, x_2, \dots, x_n \rangle$ and a classifier C_i , we exclude
284 C_i in the prediction process if $freq(x_n, c_i) = 0$, where x_n is
285 the last page id that the user has visited in testing session X .

286 The frequency matrix represents the first order of Markov
287 model. One can extend that to a higher order frequency matrix.
288 In this case, an n th-order frequency matrix corresponds to the
289 n th-order Markov model. Notice that the increase of frequency
290 matrix order leads to fewer number of classes in prediction.
291 For example, given a testing session $S_3 = \langle p_1, p_2, p_3 \rangle$, the
292 following relation holds:

$$|B_1| \geq |B_2| \geq |B_3|$$

293 where

$$B_1 = \{x | \langle x \rangle < p_3, x \in T\}$$

$$B_2 = \{x | \langle x \rangle < p_2, p_3, x \in T\}$$

$$B_3 = \{x | \langle x \rangle < p_1, p_2, p_3, x \in T\}$$

294 where T is the training sessions, x is a page id, B_i is the set of
295 outcomes by applying a frequency matrix of order i , and $|B_i|$
296 is the length of set B_i . Hence, there is a tradeoff between the
297 number of classifiers during prediction (i.e., accuracy) and the
298 order of frequency matrix. Based on our observations and ex-
299 periments, we find that first-order frequency matrix is adequate
300 to balance such tradeoff and to reduce the number of classifiers
301 in prediction without affecting the accuracy. (See Section VI-D
302 for details.)

303 V. HYBRID MODEL FOR WEB PREDICTION 304 USING DEMPSTER'S RULE

305 In this section, we present our hybrid model for web predic-
306 tion, which is based on Dempster's rule for evidence combina-
307 tion, using the ANN and Markov models as bodies of evidence.
308 In our model, prediction is resolved by fusing two separate
309 classifiers models, namely: 1) ANN and 2) Markov model (see
310 Fig. 4).

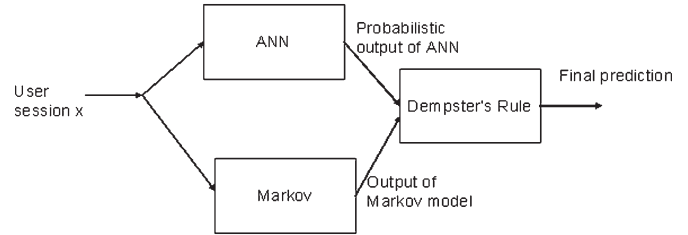


Fig. 4. Hybrid model using the Dempster's rule for evidence combination.

Dempster's rule is one part of the Dempster–Shafer Evidence
311 Combination frame for fusing independent bodies of evidence. 312
313 In Dempster's rule, the sources of evidence should be in the
314 form of basic probability. Since the ANN output value is not
315 probability [15], we need to convert/map this output into a
316 posterior probability $P(class|input)$. 316

317 In this section, we will present, first, a method to convert the
318 ANN output into a posterior probability by fitting a sigmoid
319 function after the ANN output [16]. Next, we present the
320 background of the Dempster–Shafer theory. 320

321 A. Fitting Sigmoid After ANN Output

322 We have implemented the backpropagation learning algo-
323 rithm based on minimizing the squared-error function. Hence,
324 the output of ANN cannot be considered probability. Since
325 we are using ANN as an independent body of evidence in the
326 Dempster's rule frame, we should consequently map the output
327 of ANN into probability. 327

328 One interpretation of the output of ANN, in the context of
329 the classification problem, is an estimate of the probability
330 distribution. There are several ways to interpret the ANN output
331 in terms of probability. One traditional way is to estimate
332 the probability density function (pdf) from the training data.
333 The assumption here is that we know that the training data
334 follow some distribution (typically the normal distribution).
335 The normal distribution is widely used as a model parameter
336 in which analytical techniques can be applied to estimate such
337 parameters [19], [22] as mean and standard deviation. 337

338 Another approach is to consider learning to minimize a
339 probabilistic function, instead of squared error, such as the cross
340 entropy shown in (3). Once learning is done, the output of the
341 network is an estimate of the pdf. In (3), D is the training set,
342 t_d is the target class of example d , and o_d is the output of
343 ANN, i.e., 343

$$\min - \sum_{d \in D} t_d \log(o_d) + (1 - t_d) \log(1 - o_d). \quad (3)$$

344 Since the backpropagation algorithm minimizes the squared-
345 error function, we choose to implement a parametric method to
346 fit the posterior $p(y = 1|f)$ directly, instead of estimating the
347 class-conditional densities $p(f|y)$ [16], where y is the target
348 class and f is the output function of ANN. The output of ANN
349 is computed as follows: 349

$$f(I) = \begin{cases} 1, & \text{if } \sigma(I) \geq 0.5 \\ -1, & \text{otherwise} \end{cases} \quad (4)$$

350 where I is the input to the network, and σ is the output of
 351 the sigmoid transfer function defined as in (2). It follows that
 352 class-conditional densities between the margins are apparently
 353 exponential [16]. Bayes' rule on two exponential suggests using
 354 a parametric form of sigmoid as follows:

$$P(y = 1|f) = \frac{1}{1 + \exp(Af + B)}. \quad (5)$$

355 This sigmoid model is equivalent to assuming that the output
 356 of ANN is proportional to the log odds of a positive example.
 357 Parameters A and B of (5) are fitted using maximum-likelihood
 358 estimation and can be found by minimizing the negative log
 359 likelihood of training data, which is a cross-entropy error
 360 function (3). o_d in (3) is defined as follows:

$$o_d = \frac{1}{1 + \exp(Af_d + B)}. \quad (6)$$

361 The minimization in (3) is a two-parameter minimization.
 362 Hence, it can be performed in many optimization algorithms.
 363 For robustness, we implement the model-trust minimization
 364 algorithm based on the Levenberg–Marquardt algorithm [17].

365 B. Dempster–Shafer Evidence Combination

366 The Dempster–Shafer theory is a mathematical theory
 367 of evidence [17], which is considered to be a generaliza-
 368 tion of the Bayesian theory of subjective probability. The
 369 Dempster–Shafer theory is based on two ideas. The first idea
 370 is the notion of obtaining degrees of belief for one question
 371 based on subjective probabilities for a related question, and
 372 Dempster's rule for combining such degree of belief when
 373 they are based on independent items of evidence. Since we
 374 use two independent sources of evidence, namely, ANN and
 375 Markov model, we are interested in the latter part of the
 376 Dempster–Shafer theory, namely, Dempster's rule. See [17] for
 377 more details regarding this theory.

378 Some may question why we do not use boosting and bagging
 379 rather than Dempster's rule to improve the classifier accuracy.
 380 We prefer Dempster's rule over boosting and bagging because
 381 boosting and bagging require partitioning the data set into a
 382 large number of independent bootstrap samples (> 1000) and
 383 then generating a classifier for each partition separately [12].
 384 Hence, there is a computation overhead not only in training
 385 but also in preprocessing and prediction. Furthermore, boosting
 386 and bagging cannot perform effectively if the data set does
 387 not have enough points for each class/label. In web prediction
 388 applications, many pages are sparse in the data set, because they
 389 receive very few clicks.

C. Using Dempster–Shafer Theory in Web Prediction

390

We have two sources of evidence: 1) the output of ANN 391
 and 2) the output of Markov model. These two models operate 392
 independently. Furthermore, we assume that, for any session x 393
 for which it does not appear in the Markov model, the Markov 394
 prediction probability is zero. If we use Dempster's rule for 395
 combination of evidence, we get the following: 396

$$m_{\text{ANN,Markov}}(C) = \frac{\sum_{A, B \subseteq \Theta, A \cap B = C} m_{\text{ANN}}(A) m_{\text{Markov}}(B)}{\sum_{A, B \subseteq \Theta, A \cap B \neq \phi} m_{\text{ANN}}(A) m_{\text{Markov}}(B)} \quad (7)$$

where m_{ANN} is the probabilistic output of ANN, m_{Markov} is the 397
 output of Markov model, $C \in 2^\Theta$ is a hypothesis (for example, 398
 what is the prediction of a testing session?), and Θ is the frame 399
 of discernment. A frame of discernment is an exhaustive set of 400
 mutually exclusive elements (hypothesis, propositions). 401

In web prediction, we can simplify this formulation because 402
 we have only beliefs for singleton classes (i.e., the final pre- 403
 diction is only one web page and it should not have more 404
 than one page) and the body of evidence itself ($m(\Theta)$). This 405
 means that, for any proper subset A of Θ for which we have no 406
 specific belief, $m(A) = 0$. After eliminating zero terms, we get 407
 the simplified Dempster's combination rule, for a web page P_C 408
 in (8) which is shown at the bottom of the page. Since we are 409
 interested in ranking the hypothesis, we can further simplify 410
 (8), where the denominator is independent of any particular 411
 hypothesis, as follows: 412

$$\begin{aligned} \text{rank}(P_C) &\propto m_{\text{ANN}}(P_C) m_{\text{Markov}}(P_C) \\ &+ m_{\text{ANN}}(P_C) m_{\text{Markov}}(\Theta) + m_{\text{ANN}}(\Theta) m_{\text{Markov}}(P_C). \end{aligned} \quad (9)$$

\propto is the “is proportional to” relationship. $m_{\text{ANN}}(\Theta)$ and 413
 $m_{\text{Markov}}(\Theta)$ represent the uncertainty in the bodies of evi- 414
 dence for m_{ANN} and m_{Markov} , respectively. For $m_{\text{ANN}}(\Theta)$ 415
 and $m_{\text{Markov}}(\Theta)$ in (9), we use the following. For ANN, we 416
 use the output of ANN to compute the uncertainty. We call the 417
 output of ANN for specific session x as the margin because 418
 the ANN weights correspond to the separating surface between 419
 classes and the output ANN is the distance from this surface. 420
 Uncertainty is computed based on the maximum margin of all 421
 training examples as follows: 422

$$m_{\text{ANN}}(\Theta) = \frac{1}{\ln(e + \text{ANN}_{\text{margin}})}. \quad (10)$$

$\text{ANN}_{\text{margin}}$ is the maximum distance of training examples 423
 from the margin. For Markov model uncertainty, we use the 424

$$m_{\text{ANN,Markov}}(P_C) = \frac{m_{\text{ANN}}(P_C) m_{\text{Markov}}(P_C) + m_{\text{ANN}}(P_C) m_{\text{Markov}}(\Theta) + m_{\text{ANN}}(\Theta) m_{\text{Markov}}(P_C)}{\sum_{A, B \subseteq \Theta, A \cap B \neq \phi} m_{\text{ANN}}(A) m_{\text{Markov}}(B)} \quad (8)$$

425 maximum probability of training examples as follows:

$$m_{\text{Markov}}(\Theta) = \frac{1}{\ln(e + \text{Markov}_{\text{probability}})}. \quad (11)$$

426 $\text{Markov}_{\text{probability}}$ is the maximum probability of training ex-
427 amples. Note that, in both models, the uncertainty is inversely
428 proportional to the corresponding maximum value.

429 Here, we would like to show the basic steps that are involved
430 in web prediction using Dempster's rule.

431 Step 1) Train ANN (see Section III-B).

432 Step 2) Map ANN output a probability (see Section V-A).

433 Step 3) Compute Uncertainty (ANN) (see Section V-C,
434 (10)).

435 Step 4) Construct Markov Model (see Section III-A).

436 Step 5) Compute Uncertainty of Markov model (see
437 Section V-C, (11)).

438 Step 6) **For each** testing session x , **do**

439 Step 6.1) Compute $m_{\text{ANN}}(x)$ and output ANN
440 probabilities for different pages.

441 Step 6.2) Compute $m_{\text{Markov}}(x)$ and output Markov
442 probability for different pages.

443 Step 6.3) Compute $m_{\text{ANN,Markov}}(x)$ using (9) and
444 output the final prediction.

445 Step 7) Compute prediction accuracy. // see Section VI-C.

446 VI. EVALUATION

447 In this section, we first define the prediction measurements
448 that we use in our results. Second, we present the data set that
449 we use in this paper. Third, we present the experimental setup.
450 Finally, we present our results. In all models, we use the N-gram
451 representation of paths [5], [7].

452 The following definitions will be used in the succeed-
453 ing sections to measure the performance of the prediction.
454 Pitkow and Pirolli [5] have used these parameters to mea-
455 sure the performance of the Markov model. These defin-
456 itions are given as follows: $Pr(\text{match})$ is the probability
457 that a penultimate path that was observed in a validation
458 set was matched by the same penultimate path in the train-
459 ing set. $Pr(\text{hit}|\text{match})$ is the conditional probability that
460 page x_n is correctly predicted for the testing instance $s =$
461 $\langle x_{n-1}, x_{n-2}, \dots, x_{n-k} \rangle$ and s matches a penultimate path in
462 the training set. $Pr(\text{hit})$ is defined as $pr(\text{hit}) = pr(\text{match}) \times$
463 $pr(\text{hit}|\text{match})$. $Pr(\text{miss}|\text{match})$ is the conditional proba-
464 bility that page x_n is incorrectly classified, given that its
465 penultimate path matches a penultimate path in the training
466 set. $Pr(\text{miss})$ is defined as $pr(\text{match}) \times pr(\text{miss}|\text{match})$.
467 Since we are considering the generalization accuracy and
468 the training accuracy, we add two additional measurements
469 that take into account the generalization accuracy, namely,
470 $Pr(\text{hit}|\text{mismatch})$ and overall accuracy. $Pr(\text{hit}|\text{mismatch})$
471 is the conditional probability that page x_n is correctly predicted
472 for the testing instance $s = \langle x_{n-1}, x_{n-2}, \dots, x_{n-k} \rangle$ and s does
473 not match any penultimate path in the training set. The overall
474 accuracy is defined as $pr(\text{hit}|\text{mismatch}) \times pr(\text{mismatch}) +$
475 $pr(\text{hit}|\text{match}) \times pr(\text{match})$. The overall accuracy considers
476 both matching and mismatching testing examples in computing

the accuracy. The following relations hold for the preceding 477
478 measurements: 479

$$Pr(\text{hit}|\text{match}) = 1 - pr(\text{miss}|\text{match}) \quad (12)$$

$$Pr(\text{hit})/Pr(\text{miss}) = Pr(\text{hit}|\text{match})/Pr(\text{miss}|\text{match}). \quad (13)$$

$Pr(\text{hit}|\text{match})$ corresponds to the training accuracy, be- 479
480 cause it shows the proportion of training examples that are
481 correctly classified. $Pr(\text{hit}|\text{mismatch})$ corresponds to the
482 generalization accuracy, because it shows the proportion of
483 unobserved examples that are correctly classified. The overall
484 accuracy combines both. 485

486 A. Data Set

487 For equal comparison purposes and in order to avoid dupli-
488 cating already existing work, we have used the data that were
489 collected by Pitkow and Pirolli [5] from Xerox.com for the
490 dates May 10, 1998 and May 13, 1998. Several numbers of at-
491 tributes are collected using the aforementioned method, which
492 includes the Internet Protocol address of the user, time stamp
493 with date and starting time, visiting URL address, referred URL
494 address, and the browser information or agent [20]. 495

496 B. Experimental Setup

497 We have implemented the backpropagation algorithm for
498 multilayer neural network learning. In our experiments, we
499 use a dynamic learning rate setup based on the distribution of
500 the examples from different classes. Specifically, we setup the
501 learning rate inversely to the distribution of the class, i.e., we
502 set the learning rate to a high value for low-distribution class
503 and vice versa. 504

505 In ARM, we generate the rules using the Apriori algorithm
506 proposed in [9]. We set the *minsupp* to a very low value
507 (*minsupp* = 0.0001) to capture the pages that were rarely
508 visited. We implement the recommendation engine that was
509 proposed by Mobasher *et al.* [7]. We divide the data set into
510 three partitions. Two partitions are used in training, and one
511 partition is used in testing. 512

513 C. Results

514 In this section, we present and compare the results of pre-
515 diction using four different models, namely: 1) ARM; 2) ANN;
516 3) the Markov model; and 4) the hybrid model. In addition, we
517 consider the *All-Kth-Markov model* and *All-Kth-ARM model*.
518 In the following, we will refer to the results of combining
519 the Markov model with ANN as the Dempster's rule and
520 combining ANN with the All-Kth-Markov model as the All-
521 Kth-Dempster's rule. 522

523 We consider up to seven hops in our experiments for all
524 models. Results vary based on the number of hops, because
525 different patterns are revealed for different numbers of hops.
526 Furthermore, we introduce the concept of ranking in our results.
527 Rank n means that prediction is considered to be correct if the
528 predicted page happens to be among the top n pages that has
529 the highest confidence. 530

TABLE I
PROBABILITY MEASUREMENTS USING ONE HOP AND RANK 1

	ARM	Markov	ANN	Dempster-Rule
Pr(Match)	0.590	0.590	0.590	0.590
Pr(Hit Match)	0.063	0.199	0.152	0.192
Pr(Hit)	0.037	0.118	0.09	0.114
Pr(Miss Match)	0.937	0.8	0.847	0.807
Pr(Miss)	0.555	0.474	0.502	0.478
Pr(Hit Mismatch)	0	0	0.108	0.108
Pr(Hit)/Pr(Miss)	0.067	0.249	0.179	0.238
Over all Hit/Miss	0.038	0.134	0.155	0.188
Overall accuracy	0.037	0.118	0.134	0.158

TABLE II
RESULTS OF USING THREE HOPS AND RANK 1

	ARM	All-Kth-ARM	Markov	All-Kth-Markov	ANN	Dempster's-Rule	All-Kth-Dempster's-Rule
Pr(match)	0.376	0.376	0.376	0.376	0.376	0.376	0.376
Pr(hit match)	0.043	0.103	0.231	0.230	0.156	0.232	0.232
Pr(hit)	0.016	0.038	0.087	0.086	0.059	0.087	0.087
Pr(mismatch)	0.956	0.89	0.768	0.769	0.843	0.767	0.767
Pr(miss)	0.360	0.337	0.289	0.289	0.289	0.288	0.288
Pr(hit mismatch)	0	0.044	0	0.105	0.109	0.109	0.147
Pr(hit)/Pr(miss)	0.045	0.114	0.3	0.299	0.185	0.302	0.302
Over all hit/miss	0.016	0.071	0.095	0.179	0.145	0.184	0.218
Overall accuracy	0.016	0.066	0.087	0.152	0.127	0.155	0.179

525 In Table I, there are several points to note. First, the value
526 of $pr(hit|mismatch)$ is zero for both ARM and the Markov
527 model, because neither model can predict for the unobserved
528 data. Second, the Dempster's rule achieves the best scores
529 using all measurements. For example, the training accuracy
530 $pr(hit|match)$ for ARM, Markov, ANN, and Dempster's rule
531 is 6%, 19%, 15%, and 19%, respectively. The overall accu-
532 racy for ARM, Markov, ANN, and Dempsters' rule is 3%,
533 11%, 13%, and 15%, respectively. Third, even though the
534 $pr(hit|match)$ for ANN is less than that for the Markov
535 model, the overall accuracy for ANN is better. This is because
536 $pr(hit|mismatch)$ is zero in case of the Markov model, while
537 it is 10% in case of ANN. Finally, notice that ARM has the low-
538 est prediction results. The ARM uses the concept of frequent
539 item sets, instead of *item lists* (ordered item set); hence, the
540 support of one path is computed based on the frequencies of that
541 path and its combinations. In addition, ARM is very sensitive to
542 the *minsupp* values. This might cause important patterns to be
543 lost or mixed. Table II shows the results using three hops and
544 rank 1. Notice that All- K th-Markov outperforms the Markov
545 model, and the All- K th-ARM outperforms the ARM model.
546 That is because lower orders of the models are consulted in
547 case prediction is not possible for higher orders. As a result
548 of this, $pr(hit|mismatch)$ is not zero in such models. For
549 example, the $pr(hit|mismatch)$ for All- K th-Markov and All-
550 K th-ARM are 10.5% and 4.4%, respectively. In addition, com-
551 bining the All- K th-Markov model with ANN using Dempster's
552 rule has boosted the final prediction; for example, the overall
553 accuracy for ARM, All- K th-ARM, Markov, All- K th-Markov,
554 ANN, Dempster's rule, and All- K th-Dempster's rule is 1.6%,
555 6.6%, 8.7%, 15.2%, 12.7%, 15.5%, and 17.9%, respectively.

556 In Figs. 5 and 6, the accuracy approximately increases lin-
557 early with the rank. For example, in Fig. 5, the $pr(hit|match)$
558 for All- K th-Markov is 23%, 29%, 34%, 38%, 42%, 46%, 50%,
559 and 54% for ranks 1–8, respectively. In Fig. 6, the overall

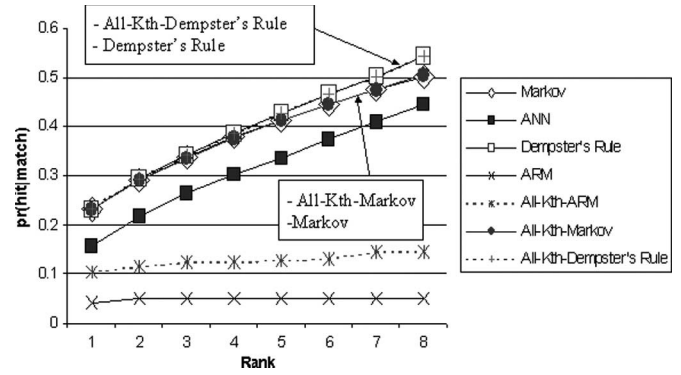


Fig. 5. $pr(hit|match)$ results using three-hop sessions and ranks from 1 to 8.

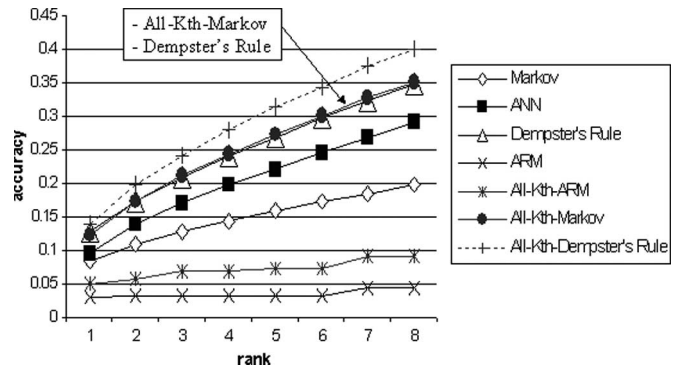


Fig. 6. Overall accuracy results using two-hop sessions and ranks from 1 to 8.

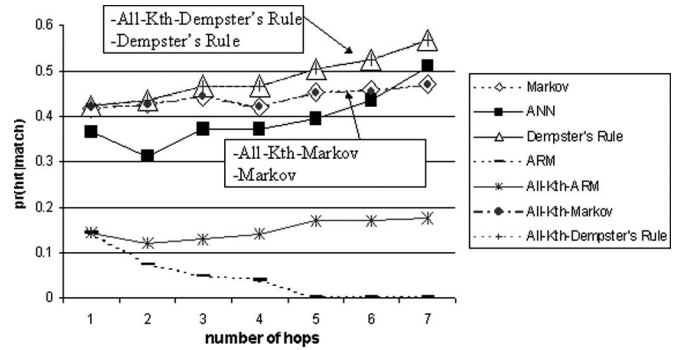


Fig. 7. Comparable results of all techniques based on $pr(hit|match)$ using rank 6.

accuracy of All- K th-Markov models is 12, 17, 21, 24, 27, 30, 560
and 35 for ranks 1–8, respectively. 561

Fig. 7 presents $pr(hit|match)$ results of using rank 6. Notice
that the Dempster's rule and All- K th-Dempster's rule methods
outperform all other techniques. 564

In Fig. 8, we notice that All- K th-Dempster's rule has
achieved the best overall accuracy, because it combines ANN
and the All- K th-Markov model. Both models have a high train-
ing and generalization accuracy. For example, the overall accu-
racy using four hops for Markov, ANN, Dempster's rule, ARM,
All- K th-ARM, All- K th-Markov, and All- K th-Dempster's
rule is 7%, 11%, 13%, 1%, 6%, 15%, and 16%, respectively.
In addition, notice that ANN has outperformed the Markov
model based on overall accuracy. This is because ANN gen-
eralizes better than the Markov model beyond training data.

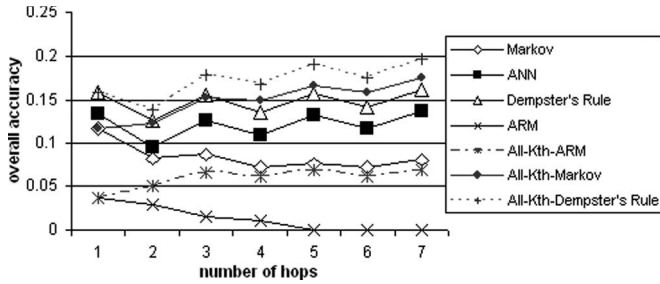


Fig. 8. Comparable results based on the overall accuracy using rank 1.

TABLE III
AVERAGE PREDICTION TIME WITH/WITHOUT DOMAIN KNOWLEDGE

Model	Average Prediction Time With Domain Knowledge (milliseconds)	Average Prediction Time Without Domain Knowledge (milliseconds)
Markov	0.567	0.544
All-Kth-Markov	1.17	0.801
ANN	6.41	556.6
Dempster's Rule	1.11	788.12

575 All- K th-Dempster's rule proves to combine the best of both
576 the ANN and All- K th-Markov models, because it has kept
577 its superiority over all techniques using all measurements and
578 using different numbers of hops.

579 D. Effect of Domain Knowledge on Prediction

580 As we mentioned previously in Section IV, we have extended
581 this model to include higher orders of domain knowledge.
582 Recall that a frequency matrix of order n corresponds to a
583 Markov model of order n .

584 Table III shows that the average prediction time using do-
585 main knowledge is 0.567, 1.17, 6.41, and 1.11 ms for the
586 Markov, All- K th-Markov, ANN, and Dempster's rule models,
587 respectively. The average prediction time without using do-
588 main knowledge is 0.544, 0.801, 556.0, and 788.0 ms for the
589 Markov, All- K th-Markov, ANN, and Dempster's rule models,
590 respectively. It is very evident that prediction time is reduced
591 dramatically for ANN and Dempster's rule. The overhead in
592 prediction without domain knowledge is a consequence of
593 loading a very large number of classifiers, i.e., 4563 classifiers,
594 and consulting them to resolve prediction. The prediction time
595 in case of the Markov model and All- K th-Markov has not been
596 affected, because such models, contrary to ANN, can handle a
597 multiclass problem without the used of an on-VS-all model.

598 In part A of Fig. 9, the overall accuracy without using
599 any domain knowledge is 18.4%, 18.4%, 0.5%, and 15.7%
600 for Markov, All- K th-Markov, ANN, and Dempster's rule, re-
601 spectively. The overall accuracy in case of using first-order
602 frequency matrix (DK) is 18.4%, 18.5%, 21.6%, and 24.5%
603 for Markov, All- K th-Markov, ANN, and Dempster's rule, re-
604 spectively. Recall that the domain knowledge is based on the
605 frequency matrix of order n , which is another representation of
606 the n th order of the Markov model; hence, the overall accuracy
607 for the basic knowledge is already included in such models. On
608 the other hand, the performance of ANN and Dempster's rule is

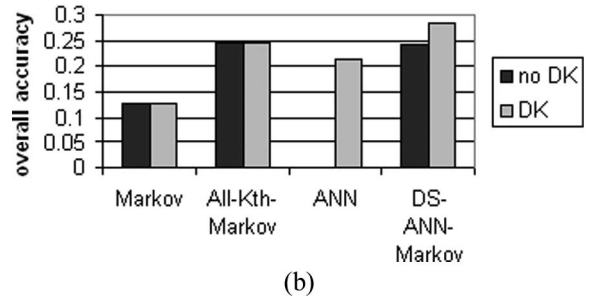
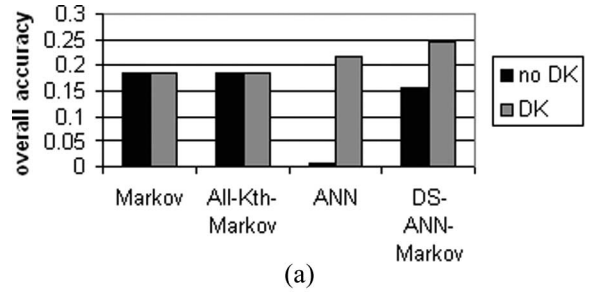


Fig. 9. Comparable results using the overall accuracy with/without domain knowledge. (a) One hop using rank 3. (b) Three hops using rank 3.

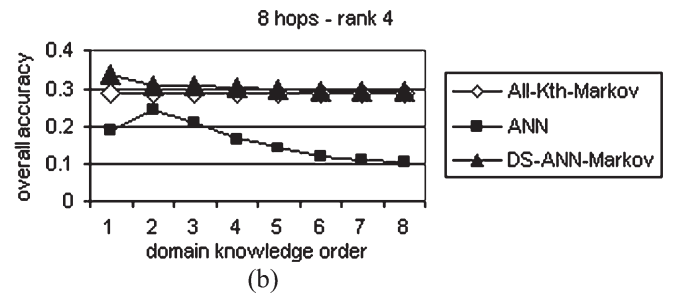
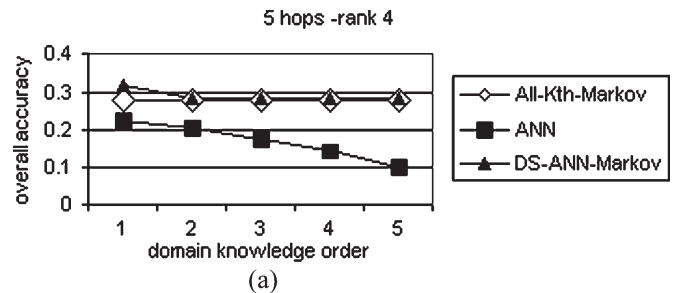


Fig. 10. Effect of domain knowledge order on the overall accuracy. (a) Five hops using rank 4. (b) Eight hops using rank 4.

affected by not using any domain knowledge, and the overall
609 accuracy has dropped largely. Similar results can be seen in
610 Fig. 9(b) when using three hops and rank 4. Fig. 10 presents
611 the effect of using different orders of domain knowledge on
612 the overall accuracy. Since we obtained similar results when
613 using different rankings and different number of hops, we
614 only show the results for five and eight hops using rank 4.
615 Recall that, in the previous experiments, we considered only
616 the first-order frequency matrix. Here, we consider a frequency
617 matrix of different orders as domain knowledge applied to the
618 All- K th-Markov model, ANN, and Dempster's rule. The three
619 curves (from top to bottom) in each subfigure represent the
620 overall accuracy of All- K th-Markov, ANN, and Dempster's
621 rule. For example, the overall accuracy for Dempster's rule
622

623 is 31%, 28%, 28%, 28%, and 28% using domain knowledge
 624 of orders 1, 2, 3, 4, and 5, respectively. Notice that the use
 625 of higher order for domain knowledge did not improve the
 626 accuracy. On the contrary, it slightly affects the overall accu-
 627 racy negatively. This can be related to the tradeoff between
 628 the number of classifiers to consult and the order of domain
 629 knowledge. Using higher order domain knowledge leads to
 630 less number of classes to consult. This may positively affect
 631 the accuracy and speed up the retrieval process. However,
 632 this might exclude correct classes, decrease the accuracy, and
 633 finally offsets the improvement of accuracy. Conversely, not
 634 using domain knowledge leads to consulting a huge number of
 635 classifiers that cause conflict. From Fig. 10, we find that using
 636 domain knowledge of order 1 or 2 can balance such tradeoffs,
 637 because accuracy is not affected dramatically.

638 VII. CONCLUSION AND FUTURE WORK

639 In this paper, we use a hybrid method in web prediction based
 640 on Dempster's rule for evidence combination to improve pre-
 641 diction accuracy. We used two sources of evidence/prediction in
 642 our hybrid method. The first body of evidence is ANNs. To im-
 643 prove the prediction of ANN further, we incorporated different
 644 orders of domain knowledge in prediction to improve prediction
 645 accuracy. The second body of evidence is the widely used
 646 Markov model, which is a probabilistic model. Furthermore,
 647 we applied the *All-Kth-Markov* model. The *All-Kth-Dempster's*
 648 *rule* proves its effectiveness by combining the best of ANN and
 649 the *All-Kth-Markov* model, as demonstrated by the fact that its
 650 predictive accuracy has outperformed all other techniques.

651 We would like to extend our research in the following direc-
 652 tions. First, we would like to study the impact/effect of other
 653 features in the session's logs by extracting statistical features
 654 from the data set to improve accuracy. Next, we would like
 655 to perform more experiments and analyses on the effect of the
 656 frequency matrix order on prediction. Finally, we would like to
 657 use boosting and bagging in the same context, and compare it
 658 with our hybrid approach.

659 REFERENCES

- 660 [1] M. Levene and G. Loizou, "Computing the entropy of user navigation in
 661 the web," *Int. J. Inf. Technol. Decis. Making*, vol. 2, no. 3, pp. 459–476,
 662 Sep. 2003.
- 663 [2] Q. Yang, H. Zhang, and T. Li, "Mining web logs for prediction models in
 664 WWW caching and prefetching," in *Proc. 7th ACM SIGKDD Int. Conf.*
 665 *KDD*, Aug. 26–29, 2001, pp. 473–478.
- 666 [3] K. Chinen and S. Yamaguchi, "An interactive prefetching proxy server for
 667 improvement of WWW latency," in *Proc. 7th Annu. Conf. INET*, Kuala
 668 Lumpur, Malaysia, Jun. 1997.
- 669 [4] V. Chung, C. H. Li, and J. Kwok, "Dissimilarity learning for nominal
 670 data," *Pattern Recognit.*, vol. 37, no. 7, pp. 1471–1477, Jul. 2004.
- 671 [5] J. Pitkow and P. Pirulli, "Mining longest repeating subsequences to predict
 672 World Wide Web surfing," in *Proc. 2nd USITS*, Boulder, CO, Oct. 1999,
 673 pp. 139–150.
- 674 [6] J. Griffioen and R. Appleton, "Reducing file system latency using a
 675 predictive approach," in *Proc. Summer USENIX Tech. Conf.*, Cambridge,
 676 MA, 1994, pp. 197–207.
- 677 [7] B. Mobasher, H. Dai, T. Luo, and M. Nakagawa, "Effective personaliza-
 678 tion based on association rule discovery from web usage data," in *Proc.*
 679 *ACM Workshop WIDM*, Atlanta, GA, Nov. 2001, pp. 9–15. Held at CIKM.
- 680 [8] R. Agrawal, T. Imielinski, and A. Swami, "Mining association rules be-
 681 tween sets of items in large databases," in *Proc. ACM SIGMOD Conf.*
 682 *Manage. Data*, Washington, DC, May 1993, pp. 207–216.

- [9] D. Duchamp, "Prefetching hyperlinks," in *Proc. 2nd USITS*, Boulder, CO, 683
 Oct. 1999, pp. 127–138. 684
- [10] A. Pandey, J. Srivastava, and S. Shekhar, "A web intelligent prefetcher for 685
 dynamic pages using association rules—A summary of results," in *Proc.* 686
SIAM Workshop Web Mining, Report No. 01-004, 2001. 687
- [11] Z. Su, Q. Yang, Y. Lu, and H. Zhang, "Whatnext: A prediction system for 688
 web requests using N-gram sequence models," in *Proc. 1st Int. Conf. Web* 689
Inf. Syst. Eng. Conf., Hong Kong, Jun. 2000, pp. 200–207. 690
- [12] L. Breiman, "Bagging predictors," *Mach. Learn.*, vol. 24, no. 2, pp. 123– 691
 140, Aug. 1996. 692
- [13] B. M. Sarwar, G. Karypis, J. Konstan, and J. Riedl, "Analysis of recom- 693
 mender algorithms for e-commerce," in *Proc. 2nd ACM EC*, Minneapolis, 694
 MN, Oct. 2000, pp. 158–167. 695
- [14] S. Brin and L. Page, "The anatomy of a large-scale hypertextual web 696
 search engine," in *Proc. 7th Int. WWW Conf.*, Brisbane, Australia, 1998, 697
 pp. 107–117. 698
- [15] E. Parzen, "On the estimation of a probability density function and mode," 699
Ann. Math. Stat., vol. 33, no. 3, pp. 1065–1076, Sep. 1962. 700
- [16] J. Platt, "Probabilistic outputs for SVMs and comparisons to regular- 701
 ized likelihood methods," in *Advances in Large Margin Classifiers*. 702
 Cambridge, MA: MIT Press, 1999. 703
- [17] G. Shafer, *A Mathematical Theory of Evidence*. Princeton, NJ: Princeton 704
 Univ. Press, 1976. 705
- [18] T. Joachims, D. Freitag, and T. Mitchell, "WebWatcher: A tour guide for 706
 the World Wide Web," in *Proc. IJCAI*, 1997, pp. 770–777. 707
- [19] O. Nasraoui and R. Krishnapuram, "An evolutionary approach to min- 708
 ing robust multi-resolution web profiles and context sensitive URL as- 709
 sociations," *Int. J. Comput. Intell. Appl.*, vol. 2, no. 3, pp. 339–348, 710
 2002. 711
- [20] R. Cooley, B. Mobasher, and J. Srivastava, "Data preparation for mining 712
 World Wide Web browsing patterns," *J. Knowl. Inf. Syst.*, vol. 1, no. 1, 713
 pp. 5–32, 1999. 714
- [21] T. Mitchell, *Machine Learning*. New York: McGraw-Hill, 1997. 715
- [22] W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery, *Numer- 716
 ical Recipes in C: The Art of Scientific Computing*, 2nd ed. Cambridge, 717
 U.K.: Cambridge Univ. Press, 1992. 718



Mamoun A. Awad received the B.Sc. degree in 719
 computer science from Baghdad University, Bagh- 720
 dad, Iraq, in June 1994, the M.S. degree in com- 721
 puter science from Wichita State University, Wichita, 722
 Kansas, in May 1999, and the Ph.D. degree in soft- 723
 ware engineering from the University of Texas at 724
 Dallas, Richardson, in December 2005. 725

He is currently an Assistant Professor in the Infor- 726
 mation Technology College, United Arab Emirates 727
 University, Al Ain, where he has taught and con- 728
 ducted research since August 2006. He has already 729

published 15 conference proceeding papers, book chapters, and journal articles. 730
 His current areas of research are data mining, intrusion detection, and Web 731
 prediction. 732



Latifur R. Khan received the B.Sc. degree in com- 733
 puter science and engineering from the Bangladesh 734
 University of Engineering and Technology, Dhaka, 735
 Bangladesh, in November 1993 and the M.S. 736
 and Ph.D. degrees in computer science from the 737
 University of Southern California, Los Angeles, in 738
 December 1996 and August 2000, respectively. 739

He is currently an Associate Professor in the De- 740
 partment of Computer Science, University of Texas 741
 at Dallas (UTD), Richardson, where he has taught 742
 and conducted research since September 2000, and 743

the Director of the UTD Database and Data Mining Laboratory. He is currently 744
 on the editorial board of North Holland's *Computer Standards and Interface* 745
Journal, published by Elsevier. He has published over 80 papers in presti- 746
 gious journals and conference proceedings. His research interests include data 747
 mining, multimedia information management, and semantic Web and database 748
 systems. 749

Dr. Khan was a Committee Member in numerous prestigious conferences, 750
 symposiums, and workshops, including the ACM SIGKDD Conference on 751
 Knowledge Discovery and Data Mining. 752