

Development of a recommender system based on navigational and behavioral patterns of customers in e-commerce sites

Yong Soo Kim, Bong-Jin Yum*, Junehwa Song, Su Myeon Kim

Korea Advanced Institute of Science and Technology, 373-1 Gusung-Dong, Yusung-Gu, Daejeon 305-701, South Korea

Abstract

In this article, a novel CF (collaborative filtering)-based recommender system is developed for e-commerce sites. Unlike the conventional approach in which only binary purchase data are used, the proposed approach analyzes the data captured from the navigational and behavioral patterns of customers, estimates the preference levels of a customer for the products which are clicked but not purchased, and CF is conducted using the preference levels for making recommendations. This also compares with the existing works on clickstream data analysis in which the navigational and behavioral patterns of customers are analyzed for simple relationships with the target variable. The effectiveness of the proposed approach is assessed using an experimental e-commerce site. It is found among other things that the proposed approach outperforms the conventional approach in almost all cases considered. The proposed approach is versatile and can be applied to a variety of e-commerce sites as long as the navigational and behavioral patterns of customers can be captured.

© 2004 Elsevier Ltd. All rights reserved.

Keywords: Recommender system; Collaborative filtering; E-commerce; Preference level

1. Introduction

Personalized services for individual customers are now popular in e-commerce sites. Properly designed and well-executed personalized services enable e-commerce companies to capture the unique needs and preferences of individual customers, help them build customer loyalty, and thereby, strengthen their competitiveness in the marketplace.

A recommender system is a typical software solution used in e-commerce for personalized services (Berson, Smith, & Thearing, 2000; Lawrence, Almasi, Korlyar, Viveros, & Duri, 2001; Sarwar, Karypis, Konstan, & Riedl, 2000; Yuan & Chang, 2001). It helps customers find the products they would like to purchase by providing recommendations based on their preferences, and is

particularly useful in e-commerce sites that offer millions of products for sale.

There are two paradigms for recommender systems, namely, collaborative filtering (CF) and content-based filtering (CBF). CF recommends products based on the similarity of the preferences of a group of customers known as a neighbor (Hill, Stead, Rosenstein, & Furnas, 1995; Resnick, Iacovou, Suchak, Bergstrom, & Riedle, 1994; Shardanand & Maes, 1995). On the other hand, CBF recommends products to a customer based on the products' similarity to the customer's past or historical preferences (Basu, Hirsh, & Cohen, 1998; Krulwich & Burkey, 1996; Lang, 1995). Therefore, CBF may not be suitable for recommending such products as music, art, movie, audio, photograph, video, etc. which are frequently sold in e-commerce sites since these products may not be easily analyzed for relevant attributive information (Balabanovic & Shoham, 1997; Shardanand & Maes, 1995). For this reason, CF is adopted in the present study which deals with recommendations in e-commerce sites.

Conventional CF is known to work well for the case where customers show their preferences for specific products in an explicit manner (e.g. rating movies). However, CF usually does not work well with binary data

* Corresponding author. Address: Department of Industrial Engineering, Korea Advanced Institute of Science and Technology, 373-1 Gusung-Dong, Yusung-Gu, Daejeon 305-701, South Korea. Tel.: +82 42869 3116; fax: +82 42869 3110.

E-mail addresses: yskim95@kaist.ac.kr (Y.S. Kim), bjyum@kaist.ac.kr (B.-J. Yum), junesong@kaist.ac.kr (J. Song), sumyeon@kaist.ac.kr (S.M. Kim).

(e.g. ‘purchase’ or ‘no purchase’ data) which are typical of e-commerce data (Hayes, Cunningham, & Smyth, 2001). To overcome this problem, recent studies proposed methods that relate the customers’ navigational and behavioral patterns with their preferences (Claypool, Le, Wased, & Brown, 2001; Kelly & Belkin, 2001; Lee, Podlaeck, Schonberg, & Hoch, 2001; Lee, Podlaeck, Schonberg, Hoch, & Gomory, 2000; Morita & Shinoda, 1994; Nichols, 1997; Rafter & Smyth, 2001). Instead of explicitly acquiring the customers’ ratings for specific products, these ‘implicit ratings’ methods passively monitor the navigational and behavioral patterns of customers (Nichols, 1997) and derive their preference levels (i.e. implicit ratings) by analyzing the clickstream data which represent the navigational and behavioral patterns of the customers (Claypool et al., 2001; Kelly & Belkin, 2001; Rafter & Smyth, 2001). In addition, several authors presented detailed case studies of the clickstream data analysis from various e-commerce sites (Lee et al., 2000, 2001). In their studies, customers’ shopping patterns (e.g. product impression, click-through, basket placement, and purchase) are analyzed, and the so-called micro-conversion rate for each adjacent pair of parameters is computed to assess the effectiveness of web merchandising. For a review and classification of various implicit measures of customer interests, the reader is referred to Kelly and Teevan (2003) or Oard and Kim (2001).

The existing works on implicit ratings mainly consider a simple correlation between a behavioral or a navigational parameter (e.g. length of reading time, number of visits, book marking variable, etc.) and the target variable (e.g. purchase/no-purchase variable). They are limited in predicting the target variable in that the observed implicit parameters are not considered in a simultaneous manner.

In this article, we extend the existing methods of implicit ratings and further develop a recommender system. The system provides a framework to analyze the inter-relationship between different behavioral and/or navigational parameters and to numerically determine customers’ preference levels from their behavioral and navigational patterns. Moreover, it can quantitatively predict the target variable from those parameters.

The proposed method consists of the following four phases. First, the data related to a customer’s purchase, navigational, and behavioral patterns are collected. Second, the customer’s preference for a certain product is numerically determined. If the product is purchased, the corresponding preference level is set to 1. If the product is clicked but not purchased, then the preference level is determined by estimating the probability of reaching the point of purchase using the data gathered from the first phase. This process is carried out using the decision tree (DT) analysis, logistic regression (LR) analysis, or artificial neural network (ANN). Third, CF is performed using the preference levels calculated in the second phase as the input values, and the preference levels of a customer for the products not clicked

are predicted. Finally, a Top-*N* list of products is generated as a recommendation to the customer.

To illustrate and assess the effectiveness of the proposed approach, an empirical study was conducted by constructing an experimental e-commerce site for compact disc (CD) albums. It was opened to the students of Korea Advance Institute of Science and Technology (KAIST) for a period of 50 days. Then, the relative performance (i.e. prediction accuracy) of the proposed recommender system is compared with that of the conventional system in which only the binary purchase data are used. The results from the above experimental study clearly show that the proposed method using the preference data is superior to the conventional method using only the binary purchase data.

In the performance study, we use the *F1* value (Sarwar et al., 2000) as the metric and come up with some additional findings: (i) constrained Pearson correlation coefficient (CPC) as a similarity measure performs consistently better than Pearson correlation coefficient and/or Jaccard coefficient for both approaches; (ii) if CPC is used, then the proposed approach outperforms the conventional approach in almost all cases considered; and (iii) the proposed approach performs best when LR is used for predicting the preference levels, CPC is used as a similarity measure, and the size of recommendation is ‘small’.

The rest of this article is organized as follows. In Section 2, details of the proposed method are presented, and the results of the experimental study are described in Section 3. Finally, Section 4 presents the conclusion and future research directions.

2. Proposed recommender system

2.1. Captured data from e-commerce sites

The proposed recommender system is developed based on the customers’ navigational and behavioral patterns in e-commerce sites. Navigational patterns include browsing, searching, product click, basket placement, and actual purchase, while behavioral patterns consist of the click ratio for a certain type of product, length of reading time spent on a specific product, number of visits to a specific product, printing, and bookmarking. Although the proposed system is developed using an experimental e-commerce site as an example, it can be applied to a variety of e-commerce sites as long as the above navigational and behavioral patterns can be captured.

The product taxonomy in an e-commerce site generally has a hierarchical structure. For instance, Fig. 1 shows such a hierarchical structure for the experimental e-commerce site used in the present study. More specifically, there are seven genres at Level 1, and each genre has 3–8 different types of CD’s at Level 2. Finally, each type at Level 2 has about 20–1000 different CD’s.

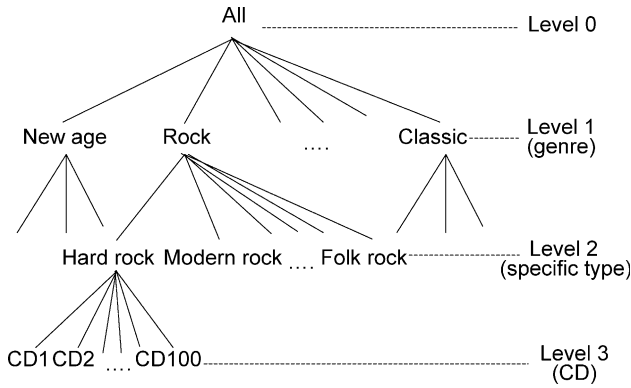


Fig. 1. Product taxonomy of experimental CD e-commerce site.

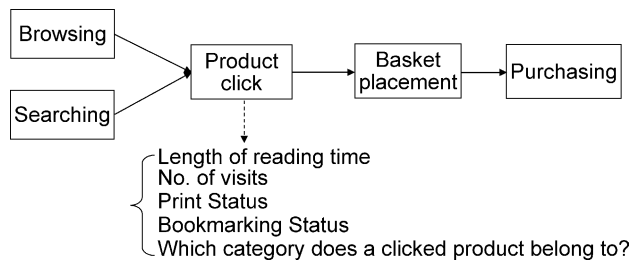


Fig. 2. Possible actions that can be taken by customers in e-commerce sites and possible data that can be obtained from such actions.

Fig. 2 illustrates possible actions and steps that customers can take in an e-commerce site, ranging from the point of logging-in to the web site to the point of actual purchase of a product. It also indicates the possible data that can be gathered from these actions.

After logging-in to the web site, a customer can either browse through the site just to check whether there are interesting products or intentionally search for a specific product to purchase. When the customer clicks a product, he or she will be provided with specific information. Then, the customer can either print or bookmark the page as a reference for a future purchase or compare the details of the product with other available goods. Other important

information that can be obtained from the customer’s actions within the site include: (i) the time it takes for the customer to read about a specific product (length of reading time); (ii) the number of visits to a specific product (number of visits); and (iii) the category to which the product belongs. A product that is frequently viewed and read can be surmised as a popular product. Furthermore, products in a certain category with a high click ratio can also be considered popular. For instance, if the click ratio for the Classic CD’s is higher than the Rock CD’s at Level 1 in Fig. 1, this could mean that the customer enjoys classic music more than rock.

Table 1 shows the parameters which describe the behavioral and navigational patterns of a customer in the experimental e-commerce site. Then, for each customer who visits the site and clicks at least one product, the corresponding parameter values are captured and summarized as shown in Table 2. In Table 2, a ‘case’ corresponds to a product clicked. Note that several cases may exist for a customer. Hereafter, the term ‘customer’ is used to represent a customer who visits the site and clicks at least one product.

2.2. Proposed methodology

The proposed methodology consists of the following four phases

- Phase I All the data related to the purchase, navigational, and behavioral patterns are gathered as shown in Tables 1 and 2. Descriptive statistics are also calculated and analyzed.
- Phase II For each customer, the preference level of a product which is clicked but not purchased is estimated (the preference level of a purchased product is set to 1).
- Phase III CF is performed using the preference levels in Phase II as input values, and the preference levels of a customer for the products not clicked are predicted.

Table 1
Data collected from the experimental e-commerce site

Parameters	Descriptions
Click type	Binary variable: searching = 1; browsing = 0
Number of visits	Discrete variable
Length of reading time	Continuous variable (s)
Print status	Binary variable: print = 1; no print = 0
Bookmarking status	Binary variable: bookmarking = 1; no bookmarking = 0
Level 1 click ratio (genre)	Continuous variable defined for each product k clicked by customer i . Let j be the category (at Level 1) to which product k belongs. Then, Level 1 click ratio for product, $k = (\text{Total number of products clicked by customer } i \text{ that belong to category } j \text{ at Level 1}) / (\text{Total number of products clicked by customer } i)$
Level 2 click ratio (specific type)	Continuous variable defined for each product k clicked by customer i . Let j be the category (at Level 2) to which product k belongs. Then, Level 2 click ratio for product, $k = (\text{Total number of products clicked by customer } i \text{ that belong to category } j \text{ at Level 2}) / (\text{Total number of products clicked by customer } i)$
Basket placement status	Binary variable: basket placement = 1; no basket placement = 0
Purchase status	Binary variable: purchase = 1; no purchase = 0

Table 2
Structure of collected data (example)

Case	Customer	CD	Click type	Length of reading time	No. of visits	Level 1 ratio	Level 2 ratio	Basket placement	Purchase
1	1	A	1	49	2	0.67	0.33	1	1
2	1	B	1	15	1	0.67	0.33	1	0
3	1	C	0	4	1	0.33	0.33	0	0
4	2	A	0	6	1	0.75	0.50	0	0
5	2	C	0	8	1	0.75	0.50	0	0
6	2	D	1	12	1	0.25	0.25	1	1
7	2	E	0	6	1	0.25	0.25	0	0
				⋮					

Phase IV After making a Top-*N* list, recommendations are made to each customer.

In Phase II, the preference level of a product which is clicked but not purchased is estimated according to the following three steps.

- (1) Estimation of the probability of purchase after basket placement (*p*):

$$p = \frac{\text{Total number of cases in which product is purchased}}{\text{Total number of cases in which product is placed in basket}}$$

- (2) Estimation of the probability of basket placement for a product which is clicked but not placed in the basket (*b*): In the case where a clicked product is not placed in the basket, the probability that the product would be purchased is difficult to estimate by simply using the parameters in Table 1. In this case, the probability that the product would be placed in the basket after being clicked (*b*) is first estimated. This is done using DT analysis, ANN, or LR analysis. In these analyses, basket placement status is considered as the target variable, while all the other variables, excluding the purchase status, are regarded as input variables. In DT analysis, the probabilities of reaching basket placement are estimated by following the paths of the constructed tree. In ANN or LR analysis, the probabilities of reaching basket placement are determined as the predicted values.

- (3) Determination of the preference level of a product which is clicked but not purchased for each customer: The preference level of a product which is placed in the basket but not purchased is set to *p*. On the other hand, the preference level of the product which is clicked but not placed in the basket is set to (*b* × *p*).

In Phase III, CF is conducted using the preference levels determined in Phase II as input values. In a conventional recommender system, only the purchase status is used for CF. In other words, only 0's (no purchase) and 1's (purchase) are used as input data (refer to Fig. 3(a)). In the proposed approach, however, the probability of reaching the point of purchase is estimated for a product clicked by a customer. Therefore, a stream of values between 0 and 1 are used as input data for the proposed CF (refer to Fig. 3(b)). In Fig. 3, blank cells indicate that the corresponding products are not clicked.

3. Experimental evaluation

3.1. Data sets

The experimental e-commerce site was opened to the students of KAIST for a period of about 50 days. Among the 2465 albums that were actually clicked by the customers (i.e. among the 2465 cases observed), 338 albums were purchased. An example data set is shown in Table 2.

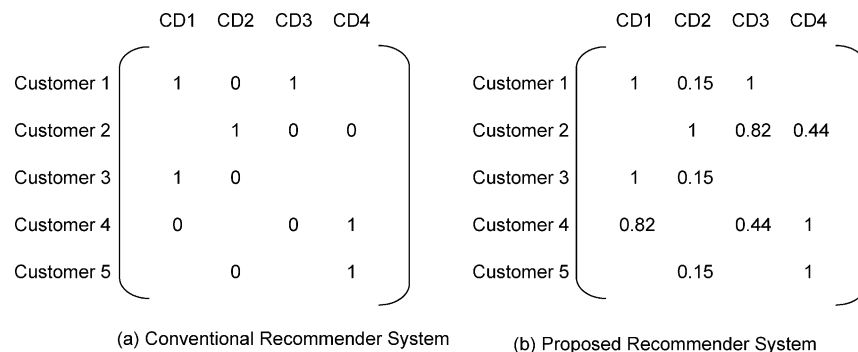


Fig. 3. 'Customer-product preference level matrix' for CF: conventional vs. proposed recommender systems.

Table 3
Basket placement vs. purchase status

	Basket placement	No basket placement	Total
Purchase	338	0	338
No purchase	74	2053	2127
Total	412	2053	2465

Since there were very few cases of printing or bookmarking, these parameters were excluded in the subsequent analyses.

3.2. Descriptive statistics: Phase I

The influence of the navigational and behavioral patterns of customers on the product purchase is first analyzed. That is, the relationship between the purchase status and each of the other parameters is evaluated as shown in Tables 3–8.

The probability of purchase after basket placement (i.e. p) is calculated as 0.82 ($=338/412$) (see Table 3). This probability is relatively high, which confirms the results of the previous studies (Lee et al., 2000, 2001). As described in Phase II of the proposed approach, the preference level of the product which is placed in the basket but not purchased is set to 0.82.

Table 4 shows the relationship between the product click type and product purchase status. When a product is clicked after being searched by a customer, the probability of its being purchased is estimated as 0.316 ($=218/689$). However, when a product is clicked after browsing through the site, it is only 0.068 ($=120/1776$). Based on these results, we may conclude that the products clicked after searching have higher preference levels than the ones clicked after browsing.

Table 5 presents the relationship between the number of visits and product purchase status. The probability of purchasing a product after the first click is 0.076 ($=136/1800$). After the second click, it becomes 0.295 ($=113/383$). For the case where the web page for a certain CD is clicked more than twice, the probability of purchase

Table 4
Click type vs. purchase status

	Product clicked through searching	Product clicked through browsing	Total
Purchase	218	120	338
No purchase	471	1656	2127
Total	689	1776	2465

Table 5
Number of visits vs. purchase status

	1 visit	2 visits	3 or more visits	Total
Purchase	136	113	89	338
No purchase	1664	270	193	2127
Total	1800	383	282	2465

Table 6
Length of reading time: results of t -test (significance level=0.05)

	N	Mean	Std dev	Std Err	$Pr > t $
Purchase	338	61.35	144.59	7.86	<0.0001
No purchase	2127	27.31	67.54	1.46	

Table 7
Level 1 (genre) click ratio: result of t -test (significance level=0.05)

	N	Mean	Std dev.	Std err.	$Pr > t $
Purchase	338	0.5750	0.3011	0.0164	0.0303
No purchase	2127	0.5378	0.2917	0.0063	

turns out to be 0.316 ($=89/282$). These results also confirm our intuition that the more frequently a product is visited, the higher becomes the probability of its being purchased.

Table 6 compares the average reading times of the purchased and not purchased products. If a product is visited more than once, the reading times for all visits are summed up for the product, and therefore, the total length of reading time for a product increases as more visits are made. This was done to verify the hypothesis that customers would take his or her time to carefully read the detailed description of a product before purchasing. A t -test with unequal variances is performed since the hypothesis of equal variances in samples ‘purchase’ and ‘no purchase’ is rejected at the 5% significance level. The result of the t -test shows that the difference between the average reading times of the purchased and not purchased products is statistically significant at the 5% significance level (see the value of the least significance probability, $Pr > |t|$), from which we may infer that a longer reading time may indicate a higher probability of purchase.

Table 7 shows the hypothesis test results on the difference between the average Level 1 (genre) click ratios for the purchased and not purchased CD’s. Similarly, Table 8 shows the hypothesis test results for the Level 2 (specific type) click ratios.

In the case of Level 1 click ratios, a t -test with equal variances is used since the hypothesis of equal variances is not rejected at the 5% significance level. However, in the case of Level 2 click ratios, a t -test with unequal variances is used since the hypothesis of equal variances is rejected at the 5% significance level.

The results in Tables 7 and 8 indicate that the means of the Level 1 or Level 2 click ratios for the purchased and not purchased products are statistically different at the 5% significance level. It is also noticed from the least

Table 8
Level 2 (specific type) click ratio: result of t -test (significant level=0.05)

	N	Mean	Std dev.	Std err.	$Pr > t $
Purchase	338	0.3666	0.2953	0.0161	<0.0001
No purchase	2127	0.2790	0.2421	0.0052	

significance probability that the significance of the mean difference in the Level 2 click ratios is statistically stronger than that in the Level 1 click ratios. These test results suggest that customers tend to click those products that belong to his or her favorite genres (Level 1) and specific types (Level 2) more often and make purchases among them.

3.3. Determination of preference levels: Phase II

In phase II, the preference levels of products that are not purchased even after being clicked are calculated based on all cases defined in Section 2.1 (also see Table 2) As shown in Fig. 3(b), a ‘Customer–Product Preference Level Matrix’ is constructed in order to conduct CF. In the matrix, the preference levels of the purchased products are set to 1, and the preference levels of the products that are not purchased even after basket placement are set to 0.82 (See Section 3.2). To determine the preference levels of the products which is clicked at least once but not placed in the basket, the probability of reaching basket placement is first predicted using DT, ANN or LR.

3.3.1. Decision tree analysis

For the intended DT analysis, the CART procedure in SPSS Answer Tree (Software/SPSS AnswerTree) is used. The basket placement status is considered as the target variable and the click type, number of visits, length of reading time, Level 1 click ratio, and Level 2 click ratio as input variables. In the CART procedure, the maximum allowable depth is set to 5, and the pruning rule is set to ‘minimum risk’. The resulting decision tree is shown in Fig. 4.

From Fig. 4, the probabilities of basket placement through different paths can be calculated. For instance, if a product is clicked through searching, the number of visits to the product is 2 or more, and the Level 2 click ratio is

greater than or equal to 0.0467, then the probability that the particular product would be placed in the basket is 0.653. This figure is then multiplied by p (i.e. the probability of purchase given basket placement) to obtain 0.535 ($=0.653 \times 0.82$), which is regarded as the preference level of the product under the above-mentioned input variable pattern. The preference levels of the products for the other seven paths in Fig. 4 can be determined in a similar manner.

A closer look at Fig. 4 reveals that the click type is one of the most significant variables that affect the chance of basket placement. As shown in Table 4, ‘searching’ has a higher probability of purchase (or equivalently, basket placement) than browsing. Furthermore, such variables as the number of visits, length of reading time, and Level 2 click ratio are also important classifiers for basket placement.

3.3.2. ANN analysis

The ANN employed in the present investigation is a multilayer feedforward network trained by backpropagation algorithm. The number of hidden layers is set to either 1 or 2, and, for each hidden layer, the number of hidden neurons changes from 3 to 12 to identify a best ANN structure. The learning rate and momentum are set to 0.1 and 0.9, respectively. It is known that a low learning rate ensures a continuous descent on the error surface and a high momentum is able to speed up training process (Sarle, 1994; Yeh, Hamey, & Westcott, 1998), and the above values are typically used in ANN training (Ting, Yunus, & Salleh, 2002).

As in the DT analysis, the basket placement status is considered as the target variable and the click type, number of visits, length of reading time, Level 1 click ratio, and Level 2 click ratio as five input variables. The whole data are randomly divided into two sets. That is, the training set consists of 70% of the data while the rest is assigned to the test set. In addition, training an ANN is stopped using the early stopping which is also called the optimal stopping rule

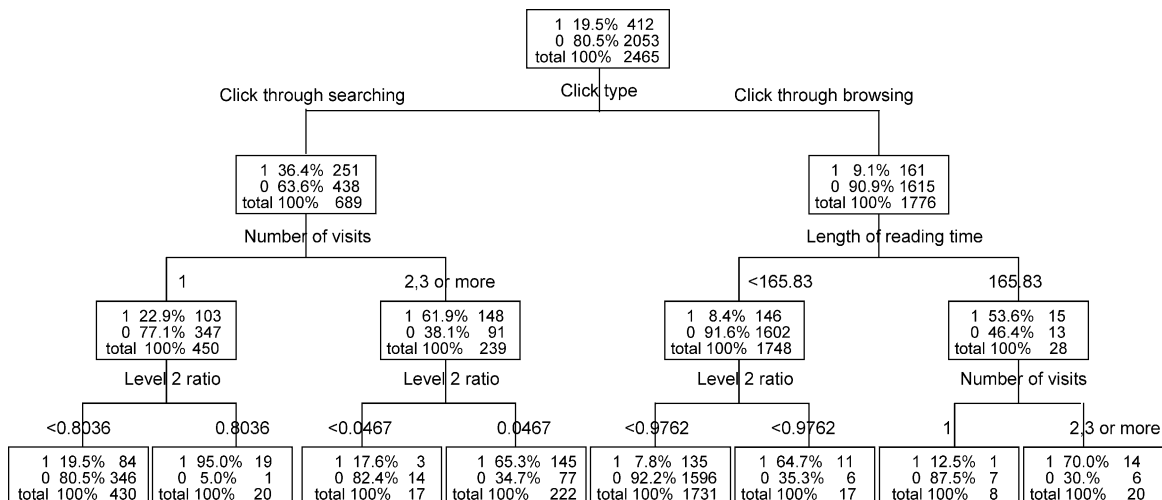


Fig. 4. Constructed decision tree: probability of reaching basket placement (1 = basket placement, 0 = no basket placement).

Table 9
Misclassification error rates for single-hidden-layer ANNs

Number of hidden neurons	Misclassification error rate
3	0.123
4	0.118
5	0.124
6	0.116
7	0.143
8	0.126
9	0.116
10	0.112
11	0.120
12	0.128

(Sarle, 1995). Under this rule, the error in the test set is also computed at the same time as the ANN is being trained, and the training is stopped once the error in the test set increases.

The results of a series of computational experiments using ‘SAS Enterprise Miner’ (Software/SAS Enterprise Miner) are shown in Tables 9 and 10. The misclassification error rate attains its minimum when a two-hidden-layer ANN is used with the numbers of neurons in the first and second hidden layers being 8 and 6, respectively. This structure is adopted in the present study to predict the probability of basket placement for a product under the navigational and behavioral patterns described by the input variables. All products that are clicked but not placed in the basket have predicted values ranging from 0 to 1. These predicted values are regarded as the probabilities of basket placement. The preference level of the product is then determined as the predicted value being multiplied by p (i.e. the probability of purchase given basket placement).

3.3.3. Logistic regression analysis

LR analysis is performed using ‘SAS Enterprise Miner’. The basket placement status is considered as a binary response variable while the click type, number of visits, length of reading time, Level 1 click ratio, and Level 2 click ratio are regarded as the predictors. In addition, the stepwise procedure is used for variable selection, and the selected variables include click type, number of visits, and Level 2 click ratio (see Table 11).

Table 10
Misclassification error rates for two-hidden-layer ANNs

Second hidden layer	First hidden layer									
	3	4	5	6	7	8	9	10	11	12
3	0.123	0.115	0.120	0.129	0.124	0.134	0.118	0.116	0.120	0.126
4	0.138	0.135	0.116	0.112	0.120	0.108	0.118	0.116	0.118	0.124
5	0.112	0.127	0.115	0.124	0.119	0.118	0.142	0.127	0.118	0.115
6	0.122	0.131	0.124	0.126	0.116	0.107	0.120	0.130	0.128	0.116
7	0.122	0.140	0.124	0.112	0.112	0.120	0.124	0.127	0.132	0.119
8	0.115	0.128	0.130	0.116	0.118	0.126	0.130	0.130	0.122	0.118
9	0.114	0.109	0.124	0.122	0.120	0.123	0.140	0.119	0.124	0.116
10	0.112	0.118	0.118	0.114	0.115	0.124	0.127	0.123	0.109	0.120
11	0.126	0.131	0.114	0.120	0.118	0.134	0.120	0.124	0.112	0.128
12	0.120	0.135	0.130	0.115	0.127	0.128	0.132	0.123	0.111	0.118

Predicted values of the response variable are regarded as the probabilities of basket placement for all products which are clicked but not placed in the basket. The preference level of the product is then determined as the predicted value being multiplied by p (i.e. the probability of purchase given basket placement).

3.4. CF and results of recommendation: Phases III and IV

3.4.1. Procedures for performance evaluation

The following procedure is used for the performance evaluation of the proposed as well as the conventional methods

- (1) In the Customer–Product Preference Level Matrix, 5 or 10% of the cells with a value of 1 are randomly selected and regarded as blank cells (refer to Fig. 3(a) and (b)).
- (2) The preference levels for these blank cells are estimated by the proposed and conventional methods using CF.
- (3) A Top- N list is generated for each customer who has blank cells in Step 1. N is varied from 5 to 30 in increments of 5 in this study.
- (4) The hidden products (i.e. the cells considered as blanks in Step 1) for each customer are checked to see if they are included in the list.

In Step 3, the Top- N list includes the products of the first N highest preference levels. Then, the performance of the proposed or conventional recommender system is evaluated by determining how effective the Top- N list is in finding the hidden products. A similar evaluation procedure was used in the previous studies (e.g. see Sarwar et al., 2000).

3.4.2. Collaborative filtering

For the proposed as well as the conventional approach, CF is performed using the Customer–Product Preference Level Matrix in which 5 or 10% of actual purchase records are intentionally hidden. Note, however, that the Customer–Product Preference Level Matrix for the conventional approach consists of binary purchase data, while that for the proposed approach consists of preference levels.

Table 11
Results of logistic regression with stepwise procedure

Parameter	Estimate	Standard Error	Chi-Square	Pr > ChiSq
Intercept	-2.6602	0.1280	431.7432	<.0001
Click Type	-0.9840	0.0623	249.4712	<.0001
Number of Visits	0.4339	0.0453	91.5328	<.0001
Level 2 Ratio	1.6351	0.2288	51.0577	<.0001

Significance level: entering=0.05; staying=0.05.

CF involves the formation of a neighborhood and making predictions. The main goal of neighborhood formation is to find similar customers for each customer. The proximity measures frequently used to determine the neighborhood of a customer are as follows.

(1) Pearson correlation coefficient (PC) (Resnick et al., 1994)

$$PC_{ab} = \frac{\sum_j(r_{aj} - \bar{r}_a)(r_{bj} - \bar{r}_b)}{\sqrt{\sum_j(r_{aj} - \bar{r}_a)^2 \sum_j(r_{bj} - \bar{r}_b)^2}}$$

PC_{ab} measures the similarity of two customers *a* and *b*, and is based on their preferences for the commonly clicked products. *r_{aj}* and *r_{bj}*, respectively, represent the preference levels of customers *a* and *b* for the commonly clicked product *j*. In addition, \bar{r}_a and \bar{r}_b , respectively, denote the average values of customer *a*'s and *b*'s preference levels for all commonly clicked products.

(2) Constrained Pearson correlation coefficient (CPC) (Shardanand & Maes, 1995)

$$CPC_{ab} = \frac{\sum_j(r_{aj} - v)(r_{bj} - v)}{\sqrt{\sum_j(r_{aj} - v)^2 \sum_j(r_{bj} - v)^2}}$$

CPC is similar to PC except that CPC is based on *v*, which is the midpoint of the scale (Shardanand & Maes, 1995). That is, PC only measures the linear tendency, while CPC measures not only the linear tendency, but also the location of the preference levels of two customers with respect to a reference value *v*. In the present study, the preference level ranges from 0 to 1, and *v* becomes 0.5 for CPC. However, the preference levels predicted using DT, LR, or ANN for the experimental data are not well scattered around 0.5, and therefore, the sample median or mean of the preference levels for all customers and products are instead used for *v*.

(3) Cosine vector

$$\cos(\vec{a}, \vec{b}) = \frac{\vec{a} \cdot \vec{b}}{\|\vec{a}\| \times \|\vec{b}\|}$$

In the case of the cosine vector, the proximity between two customer preference level vectors \vec{a} and \vec{b} is measured by computing the cosine of the angle between the two

vectors. Previous studies (e.g. Breese, Heckerman, & Kadie, 1998) have empirically shown that PC is superior to the cosine vector. Therefore, the latter is not considered in the present investigation.

(4) Jaccard coefficient for binary data (JC) (Hand, Mannila, & Smyth, 2001)

$$JC_{ab} = \frac{n_{11}}{n_{11} + n_{10} + n_{01}}$$

JC can be used to measure the similarity of two customers *a* and *b* when their preferences are represented by a binary-valued variable. *n₁₁*, *n₁₀* and *n₀₁*, respectively, denote the 'total number of products which are commonly purchased by *a* and *b*', 'total number of commonly clicked products which are purchased by *a* but not *b*', and 'total number of commonly clicked products which are purchased by *b* but not *a*'.

To measure the proximity in conducting CF, the following four measures are considered for comparison in this article: (i) Pearson correlation (PC); (ii) constrained Pearson correlation with median as the reference value (CPC_m); (iii) constrained Pearson correlation with average as the reference value (CPC_a); and (iv) Jaccard coefficient (JC). JC is used only for binary purchase data in the present study. For CPC, the sample median (*m*) and mean (*a*) of the preference levels for all customers and products are used (see Table 12).

After computing proximity measures for all pairs of customers, a neighborhood of size *k* is formed for a particular customer by selecting the *k* nearest customers based on the values of a proximity measure. To assess the effect of the size of the neighborhood on the prediction accuracy, *k* is varied to be 3, 5, and 10 in the present study.

After the neighborhood formation process, the preference levels of each customer for the products not clicked are predicted. Let *j* be a product not clicked by customer *a*, and define *N_{aj}* as the set of customers who are in the neighborhood of customer *a* and click product *j*. As mentioned earlier, for a customer in *N_{aj}*, the preference level for product *j* is either 1 if purchased or is determined using DT, LR, or ANN otherwise.

When PC is used in the neighborhood formation process, the customer *a*'s preference level for product *j* is predicted as

$$P_{aj} = \bar{r}_a + \frac{\sum_{i \in N_{aj}} PC_{ai}(r_{ij} - \bar{r}_i)}{\sum_{i \in N_{aj}} PC_{ai}}$$

Table 12
Reference values for CPC

	Binary purchase data (conventional)	Preference data (proposed)		
		Decision tree	ANN	Logistic regression
<i>M</i> (median)	0.0000	0.0780	0.0614	0.0874
<i>A</i> (average)	0.1371	0.1671	0.2514	0.1671

where \bar{r}_a and \bar{r}_i denote the average values of customer a 's and i 's preference levels, respectively, for all clicked products, and r_{ij} is the customer i 's preference level for product j (Resnick et al., 1994). When CPC is used, on the other hand, the weighted average of the preference levels for all customers in N_{aj} is computed for P_{aj} as follows (Shardanand & Maes, 1995)

$$P_{aj} = \frac{\sum_{i \in N_{aj}} \text{CPC}_{ai} \cdot r_{ij}}{\sum_{i \in N_{aj}} \text{CPC}_{ai}}$$

Finally, when JC is used, the preference level is predicted as

$$P_{aj} = \frac{\sum_{i \in N_{aj}} \text{JC}_{ai} \cdot r_{ij}}{\sum_{i \in N_{aj}} \text{JC}_{ai}}$$

Once P_{aj} 's are predicted, the products corresponding to the N highest P_{aj} 's are recommended to customer a . In the present study, N is varied from 5 to 30 in increments of 5 to assess its effect on the prediction accuracy.

3.4.3. Evaluation metrics

In order to evaluate the performance of a recommender system, the so called 'recall' and 'precision' are frequently used in the field of information retrieval (Sarwar et al., 2000) They are respectively defined as

$$\text{recall} = \frac{\sum_{i \in A} |H_i \cap \text{Top}_{N_i}|}{\sum_{i \in A} |H_i|},$$

$$\text{precision} = \frac{\sum_{i \in A} |H_i \cap \text{Top}_{N_i}|}{N \cdot |A|}$$

where

- H_i hidden products of customer i
- N total number of recommended products for each customer
- Top_{N_i} Top- N list for customer i
- A customers who has one or more hidden products.

On the other hand, these two measures are inversely related. For instance, as N increases, *recall* also increases but *precision* generally decreases. Therefore, a combined measure $F1$ is defined as the harmonic mean of recall and precision as follows (Sarwar et al., 2000).

$$F1 = \frac{2 \times \text{recall} \times \text{precision}}{\text{recall} + \text{precision}}$$

A higher value of $F1$ indicates a better performance of a recommender system. In the present study, the performance of the proposed or conventional method is evaluated based on $F1$.

3.4.4. Performance of proposed and conventional recommender systems

Computational results (i.e. $F1$ values) are summarized in Table 13 with respect to the approach taken (conventional

vs. proposed), preference prediction method for the proposed approach, similarity measure, number of recommended products (N), percentage of actually purchased products hidden, and neighborhood size (k). From Table 13, we observed the following.

Similarity measure CPC_a or CPC_m performs consistently better than PC and/or JC for both approaches. In addition, if we consider CPC_a or CPC_m only, then the proposed approach outperforms the conventional approach in almost all cases considered. It is also observed from Table 13 that $F1$ values are insensitive to the neighborhood size (k) for the k values covered in the experiment.

To assess the effects of various parameters on $F1$ in a more succinct manner, the analysis of variance (ANOVA) technique is applied to the experimental data for each approach. The experimental setting for each approach may be regarded as a full factorial design (Montgomery, 2000). For instance, factors for the proposed approach include: preference prediction method (denoted by P with three levels of DT, ANN, and LR), similarity measure (denoted by S with three levels of PC, CPC_m and CPC_a), N (with 6 levels of 5,10,...,30), percentage of actually purchased products hidden (denoted by H with 2 levels of 5 and 10%) and neighborhood size k (with three levels of 3, 5, and 10). As mentioned earlier, $F1$ values are insensitive to k , and therefore, k is not considered as a factor and is fixed at 3 in the subsequent analyses.

Table 14 shows the ANOVA table for the conventional approach. The three-way interaction effect (i.e. $N \times H \times S$) is assumed to be negligible. Note that the p -values (or the least significance probabilities) for main effects N , H and S as well as for interaction effect $H \times S$ are 'small', and therefore, those effects are considered statistically significant (Montgomery, 2000). This is also illustrated in Figs. 5 and 6, which respectively, show the three main effects and interaction effect $H \times S$. Since N is not involved in any significant interaction effect, its effect on $F1$ can be assessed independently of the other factors. That is, $F1$ achieves its highest value when $N=5$, and decreases as N increases. On the other hand, effects of H and S on $F1$ need to be assessed using their main and interaction plots together. In addition, the effect of S on $F1$ is of primary concern while that of H on $F1$ is of secondary interest since S can be regarded as a design variable for the approach while H is introduced to check the performance consistency of other factors. Note from Fig. 6 that performances of the similarity measures show different patterns depending on H . Nevertheless, CPC_a is superior to other measures regardless of the levels of H .

The ANOVA results for the proposed approach are summarized in Table 15 where the p -values indicate that all four main effects and interaction effects $P \times S$, $N \times H$, $N \times S$ and $H \times S$ are statistically significant (see also Fig. 7 for the main effects and Figs. 8–11 for the interaction effects). Main effect plot for P (see Fig. 7) and interaction plot for $P \times S$ (see Fig. 8) show that the performance of LR is better than

Table 13
Experiment results: *F1* values for conventional and proposed approaches

<i>N</i>	% of actual purchases hidden	Neighbor-hood size	Binary purchase data (conventional)				Preference data (proposed)								
							Decision tree			ANN			Logistic regression		
			JC	PC	CPC _{<i>m</i>}	CPC _{<i>a</i>}	PC	CPC _{<i>m</i>}	CPC _{<i>a</i>}	PC	CPC _{<i>m</i>}	CPC _{<i>a</i>}	PC	CPC _{<i>m</i>}	CPC _{<i>a</i>}
5	5	3	0.021	0.021	0.021	0.021	0.000	0.042	0.021	0.000	0.021	0.042	0.000	0.042	0.042
		5	0.021	0.021	0.021	0.021	0.000	0.042	0.021	0.000	0.021	0.042	0.000	0.042	0.042
		10	0.021	0.021	0.021	0.021	0.000	0.042	0.021	0.000	0.021	0.042	0.000	0.042	0.042
	10	3	0.000	0.010	0.021	0.031	0.010	0.052	0.031	0.010	0.052	0.042	0.010	0.052	0.052
		5	0.000	0.010	0.021	0.031	0.010	0.052	0.031	0.010	0.052	0.042	0.010	0.052	0.052
		10	0.000	0.010	0.021	0.031	0.010	0.052	0.031	0.010	0.052	0.042	0.010	0.052	0.052
10	5	3	0.011	0.011	0.011	0.011	0.000	0.023	0.011	0.000	0.011	0.034	0.000	0.023	0.023
		5	0.011	0.011	0.011	0.011	0.000	0.023	0.011	0.000	0.011	0.034	0.000	0.023	0.023
		10	0.011	0.011	0.011	0.011	0.000	0.023	0.011	0.000	0.011	0.034	0.000	0.023	0.023
	10	3	0.000	0.006	0.011	0.017	0.006	0.028	0.017	0.006	0.028	0.028	0.006	0.028	0.028
		5	0.000	0.006	0.011	0.017	0.006	0.028	0.017	0.006	0.028	0.028	0.006	0.028	0.028
		10	0.000	0.006	0.011	0.017	0.006	0.028	0.017	0.006	0.028	0.028	0.006	0.028	0.028
15	5	3	0.008	0.008	0.008	0.008	0.008	0.016	0.016	0.000	0.023	0.031	0.000	0.031	0.031
		5	0.008	0.008	0.008	0.008	0.008	0.016	0.016	0.000	0.023	0.031	0.000	0.031	0.031
		10	0.008	0.008	0.008	0.008	0.008	0.016	0.016	0.000	0.023	0.031	0.000	0.031	0.031
	10	3	0.000	0.004	0.008	0.012	0.004	0.020	0.016	0.004	0.027	0.020	0.004	0.027	0.031
		5	0.000	0.004	0.008	0.012	0.004	0.020	0.016	0.004	0.031	0.020	0.004	0.027	0.031
		10	0.000	0.004	0.008	0.012	0.004	0.020	0.016	0.004	0.031	0.020	0.004	0.027	0.031
20	5	3	0.006	0.006	0.006	0.012	0.006	0.018	0.018	0.000	0.018	0.024	0.000	0.024	0.024
		5	0.006	0.006	0.006	0.012	0.006	0.018	0.018	0.000	0.018	0.024	0.000	0.024	0.024
		10	0.006	0.006	0.006	0.012	0.006	0.018	0.018	0.000	0.018	0.024	0.000	0.024	0.024
	10	3	0.000	0.002	0.005	0.010	0.006	0.021	0.015	0.003	0.024	0.018	0.003	0.024	0.024
		5	0.000	0.002	0.005	0.010	0.006	0.021	0.015	0.003	0.024	0.018	0.003	0.024	0.024
		10	0.000	0.002	0.005	0.010	0.006	0.021	0.015	0.003	0.024	0.018	0.003	0.024	0.024
25	5	3	0.005	0.005	0.005	0.014	0.010	0.019	0.019	0.000	0.019	0.024	0.000	0.019	0.019
		5	0.005	0.005	0.005	0.014	0.010	0.019	0.019	0.000	0.019	0.024	0.005	0.019	0.019
		10	0.005	0.005	0.005	0.014	0.010	0.019	0.019	0.000	0.019	0.024	0.005	0.019	0.019
	10	3	0.000	0.002	0.005	0.010	0.007	0.017	0.017	0.005	0.022	0.017	0.005	0.019	0.019
		5	0.000	0.002	0.005	0.010	0.007	0.017	0.017	0.005	0.022	0.017	0.005	0.019	0.019
		10	0.000	0.002	0.005	0.010	0.010	0.017	0.017	0.005	0.022	0.017	0.005	0.019	0.019
30	5	3	0.004	0.004	0.004	0.012	0.012	0.016	0.020	0.004	0.020	0.024	0.000	0.016	0.024
		5	0.004	0.004	0.004	0.016	0.016	0.016	0.020	0.004	0.020	0.024	0.008	0.016	0.024
		10	0.004	0.004	0.004	0.016	0.012	0.016	0.020	0.004	0.016	0.028	0.004	0.016	0.024
	10	3	0.000	0.002	0.004	0.012	0.006	0.014	0.014	0.006	0.018	0.020	0.006	0.018	0.020
		5	0.000	0.002	0.004	0.012	0.006	0.014	0.014	0.006	0.018	0.020	0.004	0.018	0.018
		10	0.000	0.002	0.004	0.012	0.008	0.014	0.014	0.006	0.018	0.020	0.004	0.018	0.018

Table 14
ANOVA table for the conventional approach

Source	Degree of freedom	Sum of squares ($\times 10^4$)	Mean squares ($\times 10^4$)	<i>F</i>	<i>p</i> -value
<i>N</i>	5	9.901	1.980	16.94	0.000
<i>H</i>	1	1.050	1.050	8.99	0.009
<i>S</i>	3	6.092	2.031	17.38	0.000
<i>N</i> × <i>H</i>	5	0.196	0.039	0.34	0.884
<i>N</i> × <i>S</i>	15	0.921	0.061	0.53	0.888
<i>H</i> × <i>S</i>	3	2.336	0.779	6.66	0.004
Error	15	1.753	0.117		
Total	47	22.250			

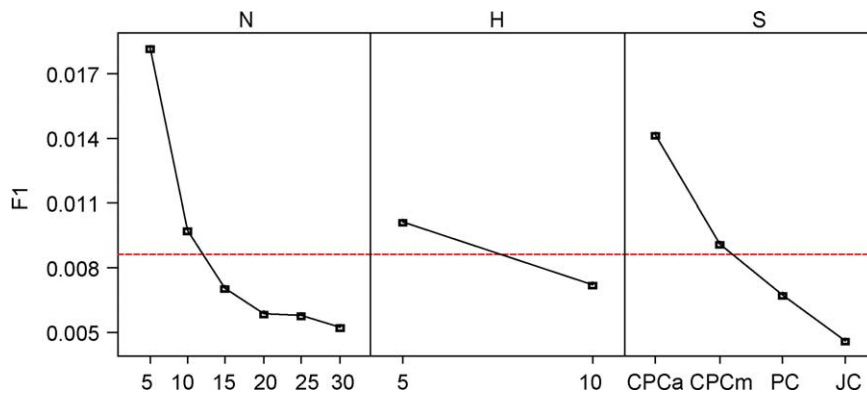


Fig. 5. Main effect plots for the conventional approach.

ANN or DT when CPC_a or CPC_m is used, while it is worse than the others when PC is used. It is noted from Figs. 7, 9 and 10 that *F1* achieves its highest value when *N*=5 except when PC is used as a similarity measure. Finally, Fig. 8 shows that CPC_a performs better than CPC_m when LR is used, and Fig. 10 shows that CPC_m performs better than CPC_a when *N*=5. Both CPC_a and CPC_m perform similarly with respect to *H* (see Fig. 11). The performance of PC is worse than CPC_a or CPC_m regardless of the levels of *P*, *N*, and *H* as shown in Figs. 8, 10 and 11, respectively.

4. Conclusion

In this article, a novel CF-based recommender system is developed for e-commerce sites. Unlike the conventional approach in which only binary purchase data are used, the proposed approach analyzes the data captured from the navigational and behavioral patterns of customers, estimates the preference levels of the products which are clicked but not purchased, and conducts CF using these preference levels for making recommendations. The proposed approach also compares with the existing works on click-stream data analysis in which customers' navigational and behavioral patterns are analyzed for simple relationships. The proposed approach is versatile and can be applied to a variety of e-commerce sites as long as the navigational and behavioral patterns of customers can be captured.

The effectiveness of the proposed approach is assessed using an experimental e-commerce site. The *F1* metric is used for performance evaluation and major findings include: (i) constrained Pearson correlation coefficients (CPC_a and CPC_m) as similarity measures perform consistently better than Pearson correlation coefficient and/or Jaccard coefficient for both approaches; (ii) if CPC_a or CPC_m is used, then the proposed approach outperforms the conventional approach in almost all cases considered; and (iii) the proposed approach performs best when logistic regression is used for predicting the preference levels, CPC_a or CPC_m is used as a similarity

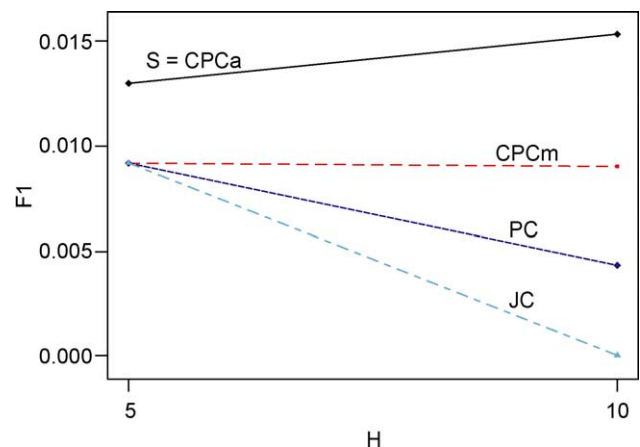


Fig. 6. Interaction plot for *S* and *H* for the conventional approach.

Table 15
ANOVA table for the proposed approach

Source	Degree of freedom	Sum of squares ($\times 10^4$)	Mean squares ($\times 10^4$)	F	p-value
P	2	1.952	0.976	6.50	0.003
N	5	28.462	5.692	37.91	0.000
H	1	1.789	1.789	11.92	0.001
S	2	100.622	50.311	335.08	0.000
P×N	10	2.176	0.218	1.45	0.180
P×H	2	0.174	0.087	0.58	0.564
P×S	4	7.208	1.802	12.00	0.000
N×H	5	5.402	1.080	7.20	0.000
N×S	10	14.500	1.450	9.66	0.000
H×S	2	1.798	0.899	5.99	0.004
Error	64	9.609	0.150		
Total	107	173.692			

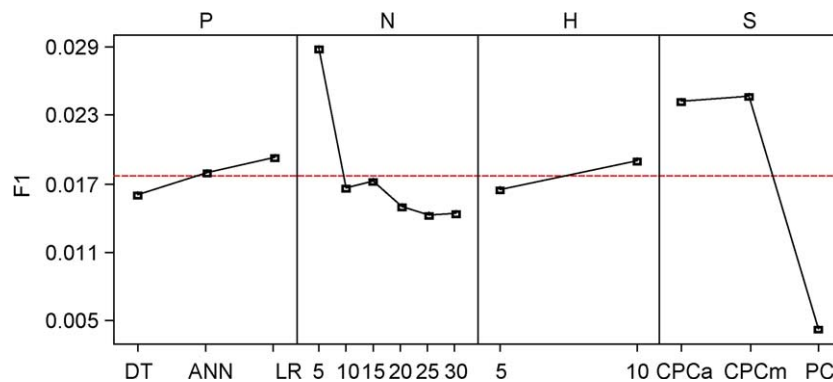


Fig. 7. Main effect plots for the proposed approach.

measure, and a ‘small’ *N* is used as a size of recommendation. In addition, Phases I (descriptive statistics) and II (preference level prediction) analyses reveal that the click type, number of visits, length of reading time, and Level 2 click ratio are important predictors or classifiers for the target variable. Especially, the importance of the click type or Level 2 click ratio has been rarely addressed in the literature as a useful measure of customer behaviors.

The above findings are based on the data captured from a relatively small experimental e-commerce site, and need further verification using a variety of actual e-commerce sites. A fruitful area of future research may also include applying the proposed approach for making recommendations in the C2C (customer to customer) environment, such as an auction site, where customers can actively participate as both sellers and buyers. In addition, customers’ sequential patterns can be used to develop better recommender systems in further study.

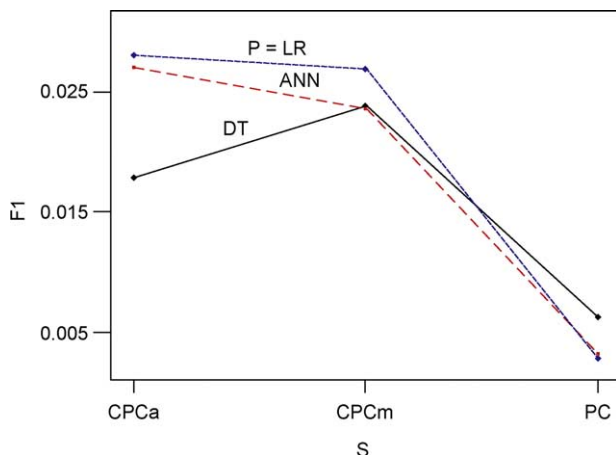


Fig. 8. Interaction plot for *P* and *S* for the proposed approach.

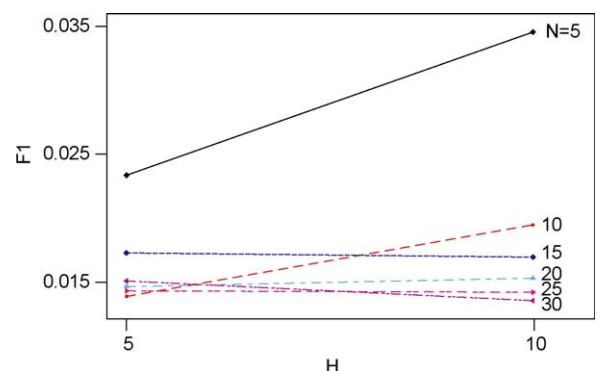
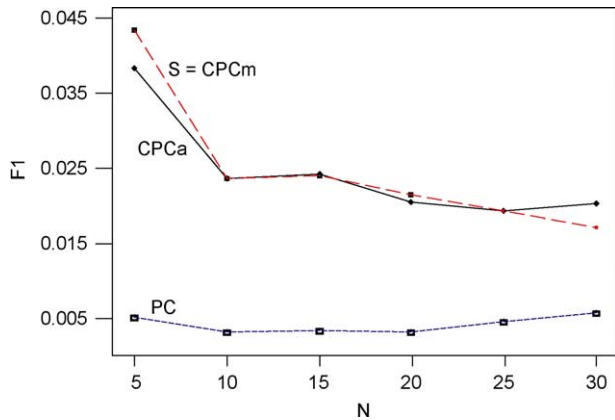
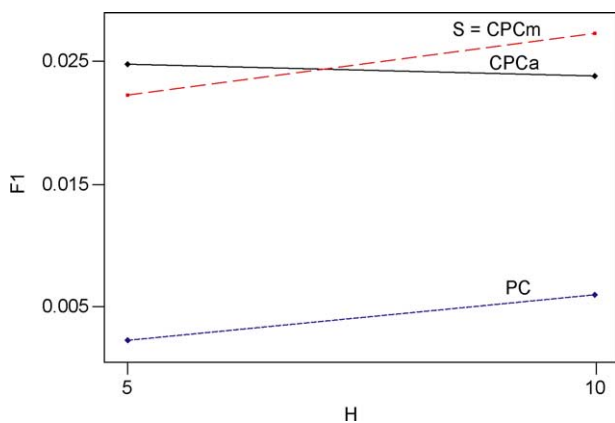


Fig. 9. Interaction plot for *N* and *H* for the proposed approach.

Fig. 10. Interaction plot for N and S for the proposed approach.Fig. 11. Interaction plot for H and S for the proposed approach.

References

- Balabanovic, M., & Shoham, Y. (1997). Content-based, collaborative recommendation. *Communications of the ACM*, 40(3), 66–72.
- Basu, C., Hirsh, H., & Cohen, W. (1998). Recommendation as classification: using social and content-based information in recommendation. *Proceedings of the 1998 workshop on recommender systems* (pp. 11–15). Menlo Park, CA.
- Berson, A., Smith, K., & Thearing, K. (2000). *Building data mining applications for CRM*. New York: McGraw-Hill.
- Breese, J.S., Heckerman, D., & Kadie, C. (1998). Empirical analysis of predictive algorithms for collaborative filtering. *Proceedings of the 14th conference on uncertainty artificial intelligence* (pp. 43–52). Madison, WI.
- Claypool, M., Le, P., Wased, M., & Brown, D. (2001). Implicit interest indicators. *Proceedings of the international conference on intelligent user interfaces* (pp. 33–40). Santa Fe, NM.
- Hand, D., Mannila, H., & Smyth, P. (2001). *Principles of data mining*. Cambridge, MA: MIT Press.
- Hayes, C., Cunningham, P., & Smyth, B. (2001). A case-based reasoning view of automated collaborative filtering. *Proceedings of the fourth international conference on case-based reasoning* (pp. 243–248). Vancouver.
- Hill, W., Stead, L., Rosenstein, M., & Furnas, G. (1995). Recommending and evaluating choices in a virtual community of use. *Proceedings of the 1995 ACM conference on factors in computing systems* (pp. 194–201). New York.
- Kelly, D., & Belkin, N.J. (2001). Reading time, scrolling, and interaction: exploring implicit sources of user preferences for relevance feedback. *Proceedings of the 24th annual ACM SIGIR conference on research and development in information retrieval* (pp. 408–409). New Orleans, LA.
- Kelly, D., & Teevan, J. (2003). Implicit feedback for inferring user preference: A bibliography. *ACM SIGIR Forum*, 37(2), 18–28.
- Krulwich, B., & Burkey, C. (1996). Learning user information interests through extraction of semantically significant phrases. *Proceedings of the AAAI spring symposium on machine learning in information access* (pp. 100–112). Stanford, CA.
- Lang, K. (1995). Newsweeder: learning to filter netnews. *Proceedings of the 12th international conference on machine learning* (pp. 331–339). Lake Tahoe, CA.
- Lawrence, R. D., Almasi, G. S., Korlyar, V., Viveros, M. S., & Duri, S. S. (2001). Personalization of supermarket product recommendations. *Data Mining and Knowledge Discovery*, 5(1), 11–32.
- Lee, J., Podlaeck, M., Schonberg, E., & Hoch, R. (2001). Visualization and analysis of click stream data of online stores for understanding web merchandising. *Data Mining and Knowledge Discovery*, 5(1/2), 59–84.
- Lee, J., Podlaeck, M., Schonberg, E., Hoch, R., & Gomory, S. (2000). Understanding merchandising effectiveness of online stores. *Electronic Markets*, 10(1), 1–9.
- Montgomery, D. C. (2000). *Design and analysis of experiments*. New York: Wiley.
- Morita, M., & Shinoda, Y. (1994). Information filtering based on user behavior analysis and best match text retrieval. *Proceedings of the 17th annual international ACM SIGIR conference on research and development in information retrieval* (pp. 272–281). Dublin.
- Nichols, D.M. (1997). Implicit rating and filtering. *Proceedings of the fifth workshop on filtering and collaborative filtering* (pp. 31–36). Budapest.
- Oard, D.W., & Kim, J. (2001). Modeling information content using observable behavior. *Proceedings of the 64th annual meeting of the American society for information science and technology* (pp. 38–45). Washington, DC.
- Rafter, R., & Smyth, B. (2001). Passive profiling from server logs in an online recruitment environment. *IJCAI's workshop on intelligent techniques for web personalisation* (pp. 35–41). Seattle, WA.
- Resnick, P., Iacovou, N., Suchak, M., Bergstrom, P., & Riedle, J. (1994). Grouplens: an open architecture for collaborative filtering of netnews. *Proceedings of the ACM 1994 conference on computer supported cooperative work* (pp. 175–186). Chapel Hill.
- Sarle, W.S. (1994). Neural network implementation in SAS software. *Proceedings of the 19th annual SAS users group international conference* (pp. 1551–1573). Cary, NC.
- Sarle, W.S. (1995). Stopped training and other remedies for overfitting. *Proceedings of the 27th symposium on the interface of computing science and statistics* (pp. 352–360). Pittsburgh, PA.
- Sarwar, B., Karypis, G., Konstan, & Riedl, J. (2000). Analysis of recommendation algorithms for e-commerce. *Proceedings of ACM e-commerce 2000 conference* (pp. 158–167). Minneapolis, MN.
- Software/SAS Enterprise Miner, <http://www.sas.com/technologies/analytics/datamining/miner>.
- Software/SPSS AnswerTree, <http://www.spss.com/answertree/>.
- Shardanand, U., & Maes, P. (1995). Social information filtering: algorithms for automating word of mouth. *Proceedings of conference on human factors in computing systems* (pp. 210–217). Denver, CO.
- Ting, H.N., Yunus, J., & Salleh, H. (2002). Speaker-independent phonation recognition for Malay Plosives using neural networks. *International joint conference on neural networks* (pp. 619–623). Honolulu, HI.
- Yeh, J.C.H., Hamey, L.G.C., & Westcott, T. (1998). Developing FENN applications using cross-validated validation training. *Proceedings of the second IEEE international conference on intelligent processing systems* (pp. 565–569). Gold Coast.
- Yuan, S., & Chang, W. (2001). Mixed-initiative synthesized learning approach for web-based CRM. *Expert Systems with Applications*, 20(2), 187–200.