

# Learning Feature Weights from Customer Return-Set Selections

L. Karl Branting

LiveWire Logic, Inc.  
2700 Gateway Centre Blvd., Suite 900  
Morrisville, NC 27560, USA  
karl.branting@livewirelogic.com

**Abstract.** This paper describes LCW, a procedure for learning customer preferences represented as feature weights by observing customers' selections from return sets. An empirical evaluation on simulated customer behavior indicated that uninformed hypotheses about customer weights lead to low ranking accuracy unless customers place some importance on almost all features or the total number of features is quite small. In contrast, LCW's estimate of the mean preferences of a customer population improved as the number of customers increased, even for larger numbers of features of widely differing importance. This improvement in the estimate of mean customer preferences led to improved prediction of individual customer's rankings, irrespective of the extent of variation among customers and whether a single or multiple retrievals were permitted. The experimental results suggest that the return set that optimizes benefit may be smaller for customer populations with little variation than for customer populations with wide variation.

## 1. Introduction

A growing proportion of sales transactions are conducted over the Internet. In business-to-customer sales, an e-commerce site must perform the *inventory selection task*, which consists of eliciting each customer's requirements, finding the item (product or service) from the business's inventory that most nearly satisfies those requirements, presenting the item to the customer, and prompting the customer to either consummate the sale or refine the requirements. An item in the inventory that satisfies the customer's requirements at least as well as any other item in the inventory is *optimal* with respect to the inventory and requirements.

Case-based reasoning (CBR) is an increasingly popular paradigm for the inventory selection task (Kohlmaier, Schmitt & Bergmann, 2001; Stahl, 2001; Wilke, 1999; Wilke, Lenz & Wess, 1998). In contrast to standard database retrieval, which is restricted to exact matches, retrieval in CBR systems can involve partial matches ordered by the degree to which each product satisfies the customer's requirements. This permits an optimal item to be presented to the customer even when nothing in the inventory completely satisfies the customer's requirements.

While discrimination nets are sometimes used in CBR for case retrieval (e.g., (Kolodner, 1984)), the most common technique in e-commerce applications of CBR is nearest-neighbor retrieval (Wettschereck & Aha, 1995). In this approach, inventory items are ordered by the similarity between customers' requirements and inventory items under a metric defined on a set of numeric or symbolic features that the requirements share with the inventory items.

The most common such metric is scaled Euclidean distance, e.g., for cases  $c1$  and  $c2$ :

$$\text{dist}(c1, c2) = \sqrt{\sum_{f=1}^n w_f \times d(c1_f, c2_f)^2}$$

where  $c1_f$  and  $c2_f$  are the values of feature  $f$  for cases  $c1$  and  $c2$ ,  $w_f$  is the weight assigned to feature  $f$ , and  $d(c1_f, c2_f)$  equals  $c1_f - c2_f$  if feature  $f$  is numeric, 1 if  $f$  is symbolic and  $c1_f = c2_f$ , and 0 otherwise.

The customer is typically presented with a *return set* consisting of the  $rs$  most similar items. The customer can select any of the items in the return set or perform another retrieval with modified requirements.

Two evaluation criteria for the product selection task can be distinguished:

- **Satisfaction**, the degree to which the customer's requirements are satisfied by the best inventory item presented to the customer. Satisfaction is maximized when the return set contains an optimal item. Satisfaction is related to the recall criterion used in Information Retrieval (IR) research in that it is a function of the case quality, but (as in (Burke, Hammond, Kulyukin, Lytinen, Tomuro & Schoenberg, 1997)) differs in that it assigns no penalty for unretrieved cases.
- **Cognitive load** (Sweller, Chandler, Tierney & Cooper, 1990) imposed on the customer to find the best inventory item that was presented. Cognitive load is a function both of the number of actions performed by the customer (e.g., "click distance") and the number of alternatives from which each action is selected. This criterion is related to the precision criterion used in IR research, in that low precision increases the number of alternatives from which a selection must be made.

*Benefit* (discussed in greater detail below) is a function of both satisfaction and cognitive load, reflecting the relative importance of both components.

In principle, there is a trade-off between satisfaction and cognitive load: satisfaction can be maximized at the expense of cognitive load by presenting the customer with all inventory items; similarly, cognitive load can be maximized at the cost of satisfaction by presenting a single choice. In practice, however, e-commerce customers tolerate only a very low cognitive load. Faced with repeated retrievals or lengthy lists of items for perusal, customers quickly abandon a site and try somewhere else (Nielson, 2000). Thus, cognitive load must be kept low, which requires minimizing (1) the number of retrievals that the customer must perform to find an optimal item, (2) the size of each return set, and (3) the rank of the best item within each return set. Of these, the first is the most important, since each retrieval entails a network delay and an increase by the return-set size in the number of alternatives that the customer must consider.

The number of retrievals can be minimized by making the probability that an optimal item is in the initial return set as high as possible. For a given return-set size, this probability depends on how accurately the weights used in the similarity metric model the customer's preferences, i.e., the relative importance the customer attaches to the case features. If the customer's preferences can be modeled by a scaled Euclidean distance metric and the system has exact values for the feature weights, then the inventory items can be ranked perfectly. Under these circumstances, even a single-item return set can be

guaranteed to contain the optimal item. In practice, however, uncertainty about feature weights means that a larger return set—and therefore larger cognitive load—is required to make it probable that an optimal item is presented.

Various approaches have been used for acquiring individual preferences. One approach is to interview the user to determine pairwise relative feature importance (Keeney & Raiffa, 1993; Branting, 1999). Less intrusive approaches, such as collaborative filtering, attempt to infer preferences from *a priori* knowledge, such as group membership (Goldberg, Nichols, Oki & Terry, 1992). A third approach strives to form user models based on observations of user decisions, either in response to system suggestions (“candidate/revision” or “learn-on-failure” (Branting & Broos, 1997; Maes, 1994)) or through passive observations (Dent, Boticario, McDermott, Mitchell & Zabowski, 1992).

This work explores the feasibility of learning customer preferences by observing customers’ selections from return sets. This approach is appropriate for customer populations that share preferences to some extent. In such customer populations, two sources of uncertainty concerning an individual customer’s preferences can be distinguished: uncertainty about the mean preferences of the entire customer population; and uncertainty about the individual’s deviation from the population’s mean on each feature. When individual variations are limited, mean customer preferences constitute a good model of most customers. Under these circumstances, a small return set is likely to contain an optimal item. The greater the variation, the larger the return set required to ensure a high-satisfaction case, and the higher the resultant cognitive load.

Mean customer preferences can be determined from a representative collection of individual customer’s preferences. Individual customers’ feature weights can, in turn, be estimated by observing individual return set selections in the following manner: whenever the customer selects an item other than the item highest ranked using the current feature weights, the system can adjust the weights to make the selected item rank first. If the customer makes repeated retrievals, the weights can be repeatedly adjusted. The mean customer preferences can be estimated from the mean of the observed individual preferences.

This paper proposes a procedure for adjusting feature weights based on return set selections and demonstrates the performance of the algorithm through experiments performed on simulated data. The experiments are intended to address the following questions:

- What is the probability of misranking an inventory item as a function of feature-weight error? Stated differently, for a given feature weight error, how large must the return set be to guarantee an acceptable probability that it will contain an optimal item?
- For a given customer variance and return-set size, how many customers must be seen to achieve an acceptable probability that the return set will contain an optimal item?
- How much difference is there in the learning rate when each customer performs only a single retrieval as opposed to repeated retrievals?
- How can the best trade-off between satisfaction and cognitive load be achieved?

In the experiments described below, several simplifying assumptions are made concerning the nature of the inventory-selection task. The inventory is assumed to consist of  $I$  items, each represented by  $n$  features having real values uniformly distributed across the interval  $[0..1]$ . The mean feature weights of the entire customer population are represented by an  $n$ -dimensional feature vector  $GW$  (meaning *global weights*), with all features weights normalized to the  $[0..1]$  interval. The preferences of the  $i^{th}$  customer

are represented by a normalized  $n$ -dimensional feature vector,  $c_i$ . Each customer's feature weights are intended to represent the relative importance of each feature to that customer. The weights of each customer's preferences are normally distributed around the mean with some standard deviation  $\sigma$ .<sup>1</sup>

The next section describes an experiment designed to gauge the importance of feature-weight accuracy. Section 3 describes LCW, a procedure for learning global weights from individual return-set selections. Section 4 then presents an experimental evaluation designed to determine the return-set size and minimum number of customers needed to achieve an acceptable probability that the return set will contain an optimal item for different values of  $\sigma$ .

## 2. The Importance of Feature-Weight Accuracy

This section investigates the degree of benefit that can be expected when a system lacks any knowledge of customers' preferences, that is, when the system does not perform any learning. An important variable in determining the importance of customer feature-weight accuracy is the actual distribution of each customer's feature weights. Intuitively, one would expect a minority of features to be of high importance, and a larger number of features to be less important (e.g., in car purchasing price may be very important, whereas seat-cover material is less important, though not inconsequential).

This distribution of feature importances, and four alternatives, can be represented by the following  $n$ -dimensional weight vectors:

1. *Two-Level*. For each feature there is 25 per cent probability that the feature weight is in the interval  $[0.75..1.0]$  and a 75 per cent probability that it is in the interval  $[0.0..0.25]$ .
2. *Exponential*. For  $0 \leq i < n$ , exactly one feature has weight  $\sqrt{0.75 * \frac{1}{2^i}}$ .
3. *Linear*. For  $0 \leq i \leq n - 1$ , exactly one feature has weight  $1 - \frac{1+i}{1+n}$ .
4. *Random*. Feature weights are uniformly distributed on the  $[0..1]$  interval.
5. *Equal*. All features are weighted equally.

The Two-Level distribution is intended to mimic human preferences by making a minority of features quite important and the majority of features less important, but not completely irrelevant. However, the assumption that the Two-Level distribution is most accurate model is mere surmise that must be tested in experiments involving human subjects. It seems possible that the distribution of feature importances is context-dependent, varies among individuals, or depends on other features not yet identified. Moreover, the division of 25% more important, 75% less important, is arbitrary.

The four alternative feature weight distributions are included because of the uncertainty about the nature of human preferences. The Exponential distribution corresponds to having one very important feature, several moderately important features, and the remaining features of little importance. Linear and Random, by contrast, correspond to situations in which there is no clear division between important and unimportant features. Finally, Equal distribution corresponds to the assumption that all weights are equally important.

<sup>1</sup> The distribution is truncated to fit within the  $[0..1]$  interval, e.g., for global weight 0.8, no values above 1.0 are permitted, to restrict values to the  $[0..1]$  interval, and no values below 0.6 are permitted, so that the samples are symmetrical around the global weight.

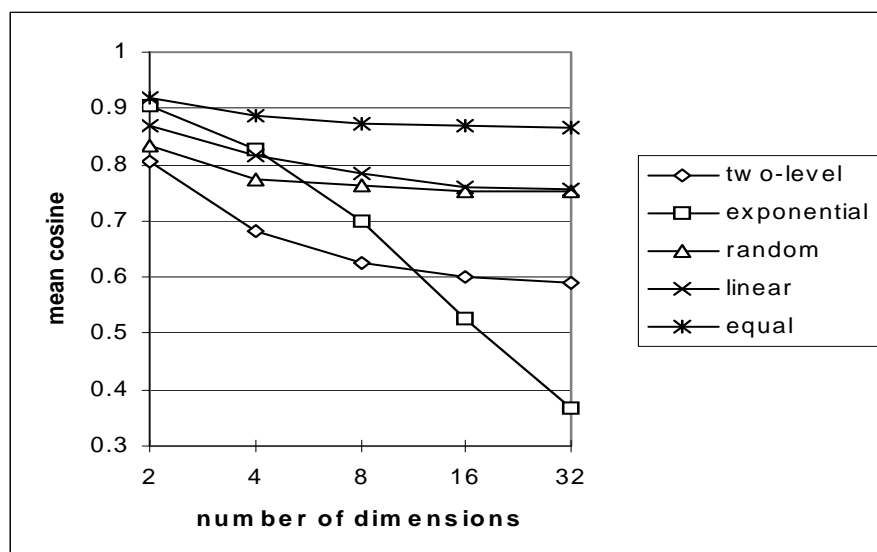


Fig. 1. The mean cosine between the actual and uninformed hypothesized weights.

In each round of this experiment, an inventory of 100 items was created, together with a weight vector of one of the five distribution types, representing the actual preferences of a simulated customer, and a weight vector of random values, the *hypothesized weights*,  $HW$ , representing an uninformed model of the customer. Finally, an additional item was created, representing the features requested by customer. This item is termed the *requirements* vector,  $r$ .

The behavior of an e-commerce site was simulated by sorting the inventory  $I$  by similarity to the requirements vector  $r$  under the hypothesized weights,  $HW$ , producing a *hypothesized ranking*. Next, the inventory item that best matched the requirements under the actual weights (i.e., the optimal item) was selected.<sup>2</sup> Finally, the position of the optimal item in the hypothesized ranking was determined and the cosine between the actual and hypothesized weight vectors was calculated.

One thousand rounds were performed for 2, 4, 8, 16, and 32 features for each distribution. In each trial three values were calculated:

- The mean rank of the optimal item in the hypothesized ranking.
- The mean cosine between the actual and hypothesized weight vectors.
- The minimum return-set size needed to include at least 90 per cent of the optimal cases.

Figures 1 and 2 display the mean cosine and minimum size needed to included at least 90 per cent of the optimal cases, respectively. These figures show that as the number of features increases, the cosine between the actual customer weights and the hypothesized weights decreases and the minimum return-set size required for a 90 per cent probability of containing the optimal case increases. These changes, which correspond to decreased benefit, are more gradual for Random and Linear weight distributions than

<sup>2</sup> With inventory items having random, real-valued feature values, there is always a unique optimal item.

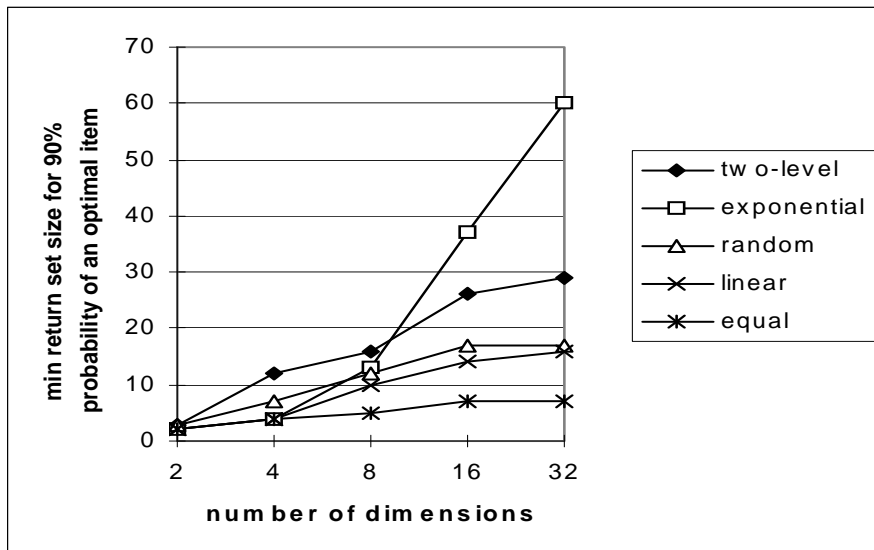


Fig. 2. The minimum return-set size needed to included at least 90 per cent of the optimal cases given uninformed weights.

for Two-Level or Exponential. Intuitively, ranking accuracy decreases as the number of unimportant features increases.<sup>3</sup>

These results indicate that unless customers' place some importance on almost all features or the total number of features is quite small, an uninformed hypothesis about customer weights can be expected to lead to low benefit. In other words, an unacceptably large return-set size is required to guarantee a high probability that the optimal case is returned. The next section presents an algorithm for modifying an initial hypothesis based on customer observations.

### 3. LCW: Learning Customer Weights

The LCW (Learning Customer Weights) algorithm rests on the assumption that customer preferences can be modeled by weights normally distributed around a population-wide mean,  $GW$ . The population-wide mean is modeled by hypothesized global weights,  $HW$ , initialized to random values.  $HW$  is used as the initial hypothesis for each individual customer's weights,  $hc_i$ . These weights are adjusted based on the customer's behavior to form revised weights  $hc'_i$ . The hypothesized global weights,  $HW$ , are then revised in light of  $hc'_i$ .

LCW consists of three connected algorithms: LearnGlobalWeights, LearnCustomerWeights, and AdjustWeights. The first algorithm, LearnGlobalWeights, takes as arguments an inventory ( $I$ ) and return-set size ( $rs$ ), and revises the global hypothesis ( $HW$ ) in light of the apparent preferences of each customer. When a small number of customers

<sup>3</sup> Very similar results are obtained using equal, rather than random, hypothesized weights, except that the equal-weight hypothesis leads to the optimal item always being ranked first and a cosine of 1.0.

have been seen,  $HW$  is simply the most recent  $hc'_i$ . When enough customers have been seen,  $HW$  is assigned the value of the mean of the last 10% of customers seen.<sup>4</sup>

```

Procedure LearnGlobalWeights(I, rs)
HW := n random values;
FOR EACH customer i DO
  r := customer i's requirements;
  hci := HW;
  hci' := LearnCustomerWeights(r, hci, I, rs);
  IF i - 1 (i.e., the number of customers previously seen) < 20
  THEN HW := hci'
  ELSE HW := the mean of the last 10% of the previous hc's;

```

The function `LearnCustomerWeights` takes hypothesized customer weights ( $hc_i$ ), the inventory ( $I$ ), and a return-set size ( $rs$ ). It elicits the customer's requirements ( $r$ ), ranks  $I$  by similarity to  $r$  under weight vector  $hc_i$ , and presents the closest  $rs$  for selection by the customer. When the customer selects an inventory item  $s$  in the return set other than the item that is most similar to  $r$  under weights  $hc_i$ ,  $hc_i$ 's weights are incrementally adjusted until  $s$  is ranked higher than any other member of the return set (or until the threshold for the maximum number of weight adjustments is exceeded). The cycle of retrieval and weight adjustment is repeated until the case ranked highest under the customer's actual weights,  $c_i$ , is ranked first under the hypothesized weights,  $hc_i$ , (corresponding to a satisfied customer making a selection) or the threshold for the maximum number of retrievals is exceeded (corresponding to a customer losing patience and terminating the sales episode). In the experiments described below, this threshold is set to either 1 or 10.

```

Function LearnCustomerWeights(hci, I, rs)
REPEAT
  r := the customer's requirements
  HypothesizedRankings := I sorted by similarity to r under hci;
  ReturnSet := the first rs members of HypothesizedRankings;
  s := the member of ReturnSet most similar to r under ci,
      customer i's actual weight vector;
  IF s is not the first element of ReturnSet
  THEN hci := AdjustWeights(s, r, ReturnSet, hci);
UNTIL s is the first element of ReturnSet OR
      Retrievals > MaximumRetrievalLimit.
return hci;

```

`AdjustWeights` iterates through each element,  $mri$  (for "misranked item"), of the return set that is incorrectly ranked higher than the customer's selection. For each such item, the weight of each feature of the hypothesized customer weights,  $hc_i$ , whose value in the requirements vector,  $r$ , is closer in value to  $mri$  than to  $s$  is multiplied by 1 minus the learning rate (to diminish the weight of a feature contributing to a misranking). Conversely, the weight of each feature of  $hc_i$  whose value in  $r$  is closer to  $s$  than to  $mri$  is multiplied by 1 plus the learning rate (to increase the weight of a feature contributing to a correct ranking).

<sup>4</sup> This approach speeds learning, because when few customers have been seen, each  $hc'_i$  is likely to be much closer to  $GW$  than any previous customer weights and, therefore, than the mean of the previous customer weights. When many customers have been seen, however, the mean of the customers may differ from  $G$  by less than the variance between individual customers.

```

Function AdjustWeights(s,r,ReturnSet,hci)
Adjustments := 0;
WHILE Adjustments < MaximumAdjustmentLimit
  SortedReturnSet := ReturnSet sorted by hci;
  MisRanked := elements of SortedReturnSet ranked above s;
  IF MisRanked = {}
  THEN return hci
  ELSE FOR EACH item mri in MisRanked
    FOR EACH feature f
      IF diff(f(mri),f(r)) < diff(f(s),f(r))
      THEN f(hci) *= (1 - LearningRate)
      ELSE IF diff(f(mri),f(r)) > diff(f(s),f(r))
      THEN f(hci) *= (1 + LearningRate);
return hci;

```

In the experiments below, LearningRate was 0.01 and MaximumAdjustmentLimit was 1000.

Intuitively, AdjustWeights “pulls” the customer’s selection  $s$  towards the customer’s requirements  $r$  while “pushing” each member of  $mri$  away from  $r$  (Bonzano, Cunningham & Smyth, 1997; Zhang & Yang, 1999). AdjustWeights is similar to RELIEF (Kira & Rendell, 1992) in rewarding weights that contribute to correct similarity assessments while penalizing weights that contribute to incorrect similarity assessments. However, AdjustWeights differs in that its weight adjustments continue until the weights are consistent with the customer’s ranking of the inventory items in the return set. Since AdjustWeights hill-climbs on a linear weight vector, it may not converge if the customer’s actual preference function is not a linear function of feature values.<sup>5</sup>

#### 4. Experimental Evaluation

LCW induces an estimate of the mean preferences of a customer population based on estimations of each individual customer derived from that customer’s behavior. Section 2 provided evidence that, absent knowledge about customers’ preferences, ranking can be expected to be inaccurate if there are wide variations in feature importances unless the number of features is very low.

LCW would constitute an improvement over an uninformed weighting scheme if it satisfied the following two conditions: (1) LCW’s estimate of the mean preferences of the customer population improved as the number of customers who have been seen increases, and (2) the improvement in the estimate of mean customer preferences resulted in an improvement in the accuracy of predictions of individual customers’ rankings.

In the following experiment, the number of dimensions was held constant at 8, on the assumption that customers are unlikely to have the patience to enter more than 8 feature values. The inventory size was held constant at 100 items. In each round of the experiment, a new inventory was created, each item of which was an element of  $[0..1]^8$ , and a set of actual global weights was created, consisting of 8 values drawn from a two-level distribution. A set of 10, 20, 50, 100, or 250 customer preferences were created, each consisting of an 8-tuple of real numbers whose values were normally distributed around the global weights with  $\sigma \in \{0, 0.1, 0.25, 0.5, 1\}$ . Each customer was assigned a requirements vector  $r$  consisting of a random 8-tuple of real num-

<sup>5</sup> See (Stahl, 2001) for an approach to feature weight optimization based on gradient ascent.



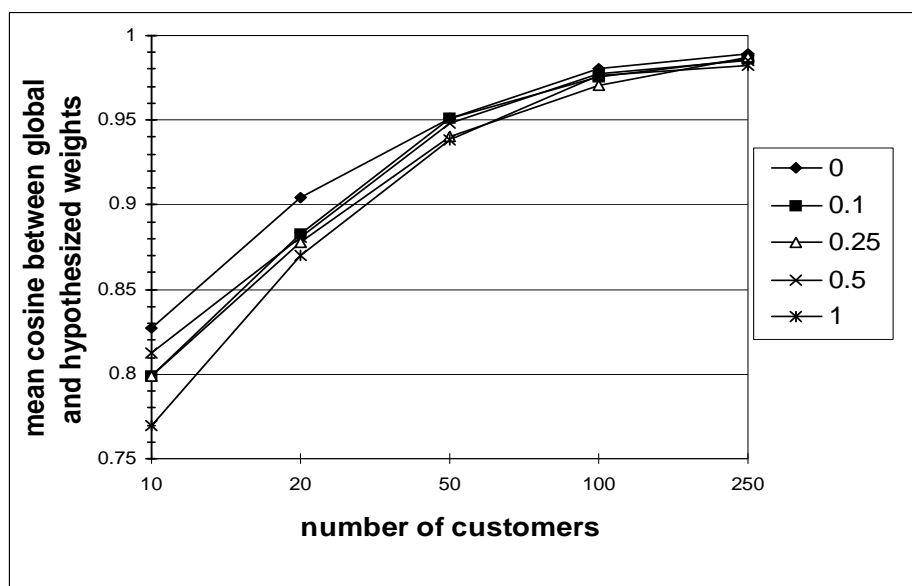


Fig. 3. The mean cosine between the actual and hypothesized global weights as a function of the number of customers seen, for  $\sigma$  from 0 to 1.0.

bers. The maximum permitted retrievals was either 1 or 10, and the return-set size was  $\in \{1, 2, 3, 5, 7, 10, 12, 15, 25, 50, 100\}$ .

In each round, the global hypothesis,  $HW$ , was initially assigned random values, and LCW was applied to each simulated customer in turn. Fifty rounds were performed for each inventory, set of customers,  $\sigma$ , maximum retrievals, and return-set size.

#### 4.1. Learning Global Weights

Figure 3 displays the mean cosine between the vectors representing the actual and hypothesized global weights as a function of the number of customers seen. Separate learning curves are shown for  $\sigma$  from 0 to 1.0. The learning curve is steeper for customer populations with lower  $\sigma$ . However, the cosine is over 0.90 even for the largest  $\sigma$  tested after 100 customers have been seen. This indicates that the first condition mentioned above—that LCW’s estimate of the mean preferences of the customer population improves as the number of customers who have been seen increases—is satisfied for a wide range of  $\sigma$ s.

#### 4.2. Improving Ranking

Evidence for the second condition—that improvement in the estimate of mean customer preferences leads to improved accuracy in predictions of individual customers’ rankings—is set forth in Figure 4, which shows that the minimum return-set size needed for a 90 per cent probability that the optimal case is in the return set decreases as the number of customers seen increases. For even the largest  $\sigma$  tested, 1.0, a return-set size of five is sufficient after 250 customers have been seen.

A different view of the experimental data is shown in Figure 5, which graphs the

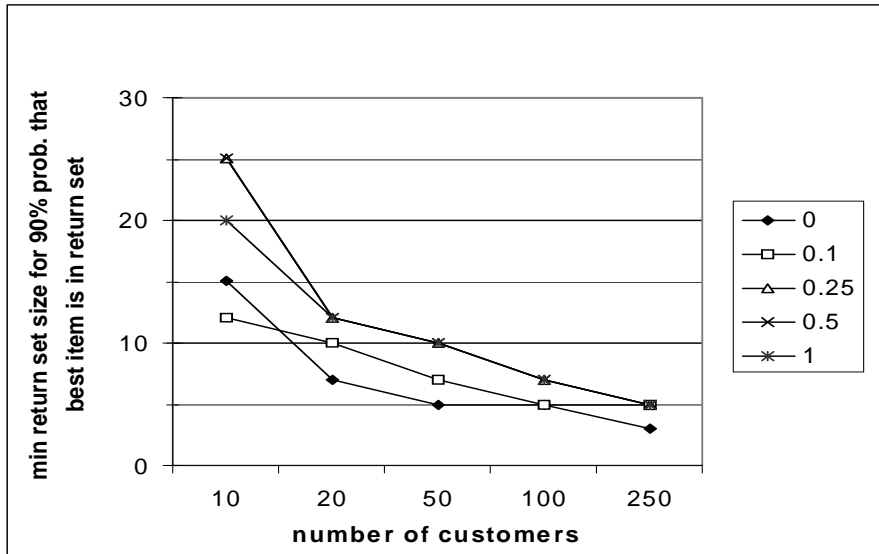


Fig. 4. The minimum return-set size needed for a 90 per cent probability that the optimal case is in the return set as a function of the number of customers, for  $\sigma$  from 0 to 1.0. The graph for  $\sigma$  of 0.5 overlaps the graph for  $\sigma$  of 0.25.

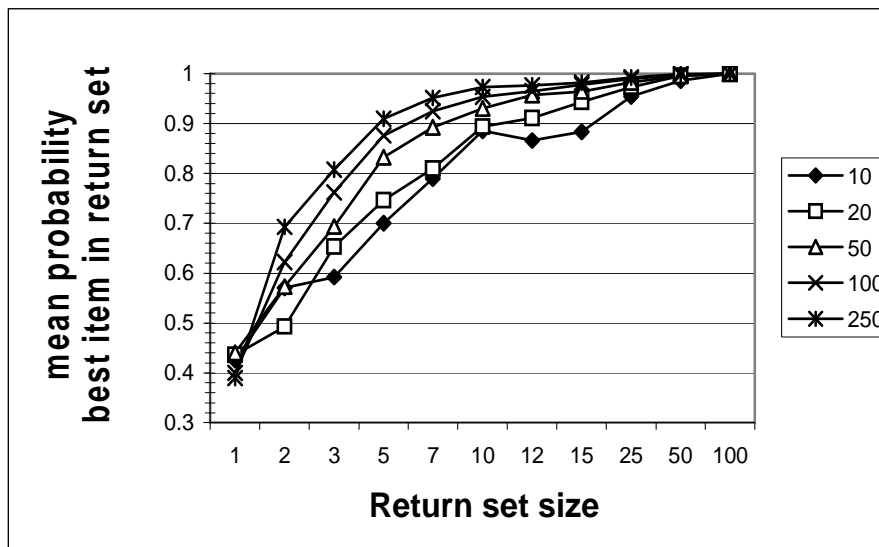


Fig. 5. The probability that the optimal item is in the return set as a function of return-set size, for 10, 20, 50, 100, and 200 customers and  $\sigma$  of 0.25.

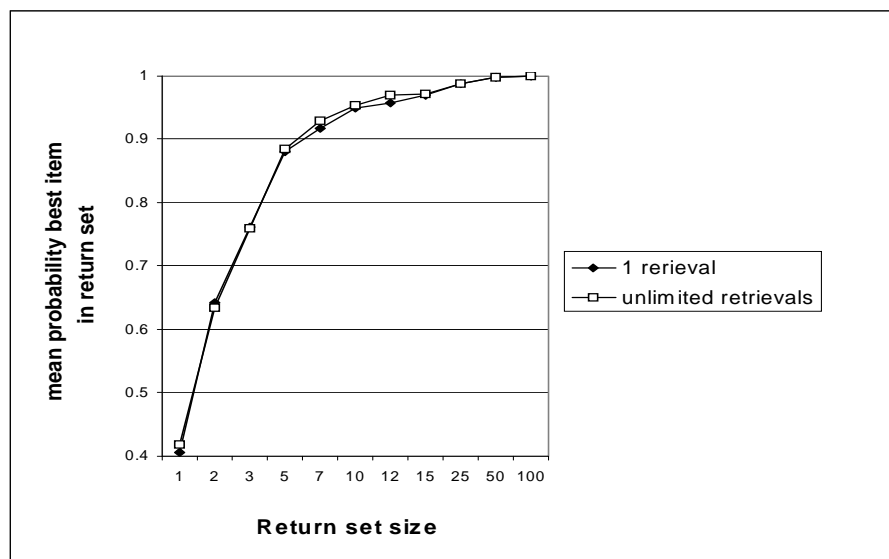


Fig. 6. The mean probability that the optimal item is in the return set as a function of return-set size, for 100 customers and  $\sigma$  of 0.25.

probability that the optimal item is in the return set as a function of return-set size, for customers seen of 10, 20, 50, 100, and 200 and  $\sigma$  of .25. The graph shows that when only 10 customers have been seen, there is a less than 90 per cent probability that the optimal case is in the return set for return-set size of less than 25. After 250 customers, however, a return-set size of 5 is sufficient for 90 per cent probability.

### 4.3. Single vs. Multiple Retrievals

The experiment included two retrieval conditions. Under the first condition, each customer made only a single retrieval. Alternatively, the customer made repeated retrievals until the best item in the return set was ranked first up to a maximum of 10 retrievals. The latter condition was intended to model relevance feedback, as in “more like this” retrievals that are permitted by some search engines. The observation above that customers seldom tolerate interactions requiring heavy cognitive load suggests that multiple retrievals, even if permitted, may be infrequent in practice.

How much opportunity for learning is lost when each customer is limited to a single retrieval? Figure 6 compares the single versus multiple retrieval conditions by the mean probability that the optimal item is in the return set as a function of return-set size. The difference in performance under the two conditions is negligible. This reflects the fact (not shown in the graph) that the mean number of retrievals never exceeded 1.5 under any of the experimental conditions. This indicates that when a customer’s estimated weights have been adjusted once, further adjustments are not likely to produce an improvement as to that customer and inventory set.

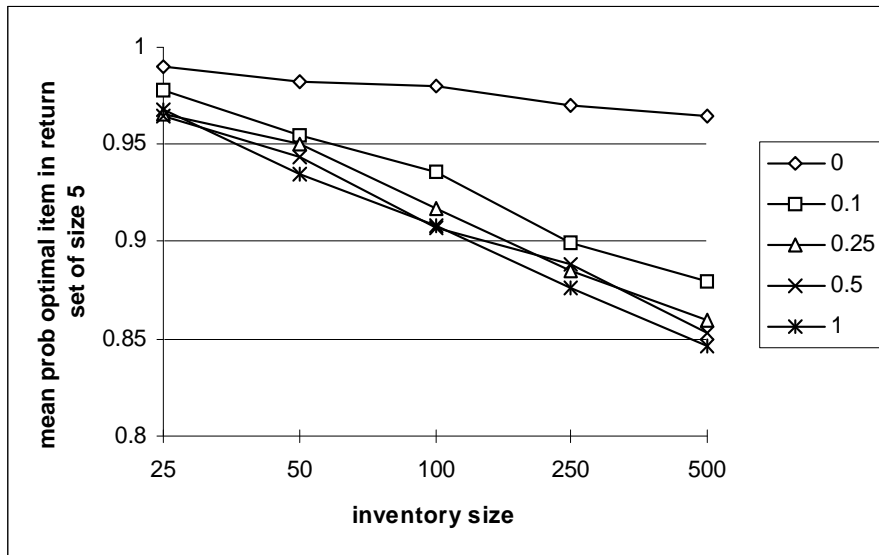


Fig. 7. The mean probability that the optimal item is in the return set as a function of inventory size, for 250 customers,  $\sigma$  of 0.25, and return-set size 5.

#### 4.4. The Effect of Inventory Size

The effect inventory size on retrieval was investigated by testing the effect of 25, 50, 100, 250, and 500 given 250 customers,  $\sigma$  of .25, and return-set size of 5. Figure 7 shows that the mean probability that an optimal item is in the return set of size 5 decreases with inventory size, falling below 90 per cent for an inventory of 250 items for nonzero  $\sigma$ . Similarly, Figure 8 shows that the mean final rank of the optimal item rises more rapidly for smaller return sets (particularly, for return-set size of 1) than for larger return sets.

#### 4.5. Maximizing Benefit

Choosing the optimal return-set size for a given customer population requires striking a balance between satisfaction and cognitive load. *Benefit* is intended to express the net trade-off between these two factors. The precise relative importance that should be attached to these two factors is impossible to determine *a priori*, as it depends on both the cognitive-load tolerance of the customer population and the importance that the vendor attaches to guaranteeing that an optimal case is returned. However, several heuristics for balancing satisfaction and cognitive load can be identified.

One heuristic is that an optimal case should almost always be returned. Under this heuristic, benefit is maximized by finding the smallest return set that insures a high probability of optimality. Another heuristic is that a return set of 1, which permits no learning, is about as bad as a return set consisting of all the cases, which maximizes cognitive load.

These ideas can be expressed in the following equation, in which  $k$  is a constant representing the relative importance of satisfaction as compared to cognitive load,  $rs$  is the return-set size, and  $sat$  represents satisfaction:

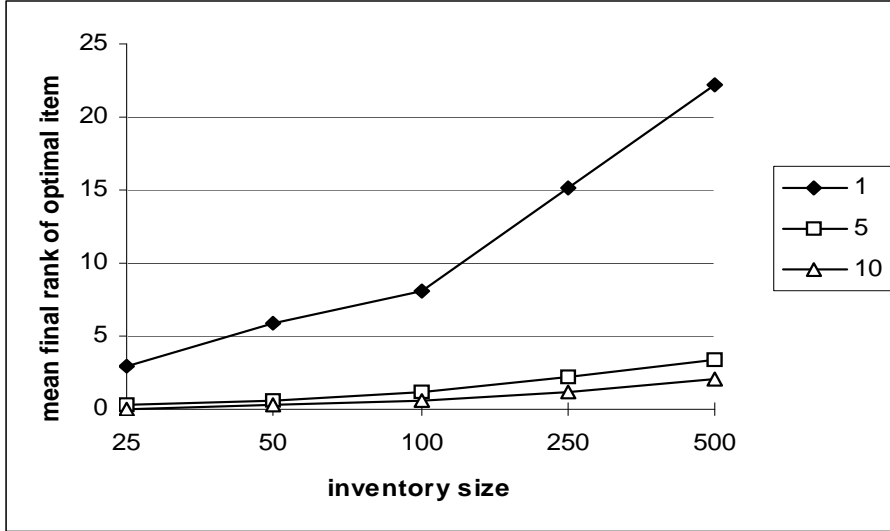


Fig. 8. The mean final rank of the optimal item as a function of inventory size, for 250 customers,  $\sigma$  of 0.25, and return-set size of 1, 5, and 10.

$$\text{benefit} = k \cdot \text{sat} - (1 - k) \cdot rs$$

Satisfaction, *sat*, can be measured by how much of an improvement the best item in the return set is over a randomly selected item. More precisely, *sat* can be measured by how much larger is (1) the mean distance between the requirements vector and every case in the inventory than (2) the distance between the requirements vector and the best item in the return set. *Sat* is normalized to the [0..1] interval by dividing it by the amount by which the optimal item is an improvement over a randomly selected item (more precisely, by the mean distance between the requirement vector and every case in the inventory minus the distance between the requirements vector and the optimal item):

$$\text{sat} = \frac{\text{MeanDistance} - \text{dist}(r, \text{BestInReturnSet})}{\text{MeanDistance} - \text{dist}(r, \text{OptimalItem})}$$

When  $k = .997$ , benefit is approximately the same for a return-set of size 1 as for a return set consisting of the entire inventory. Figure 9 displays the mean benefit as a function of return-set size for  $k = .997$ , 100 customers, and  $\sigma$  from 0 to 1.0. For small return-set sizes, the benefit is greater for customer populations with lower  $\sigma$ . For large return-set size, however,  $\sigma$  has little effect on benefit. This is because when the return set is sufficiently large, the optimal case is almost always returned regardless of  $\sigma$ . For  $\sigma = 0$ , benefit is maximized by a return-set size of 5. For larger  $\sigma$ , benefit is maximized by a return-set of 7.

## 5. Conclusion

This paper has identified the task of acquiring customer preferences from return set selections, a learning problem that can arise when case-based reasoning is applied to the inventory-selection task. The paper described LCW, a procedure for performing this task, and proposed a criterion for retrieval performance, *benefit*, which reflects both the

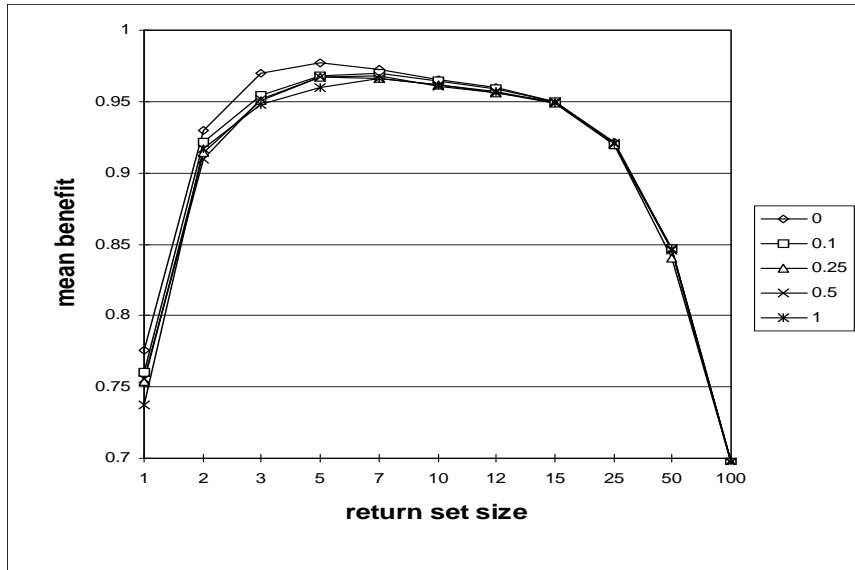


Fig. 9. Mean benefit as a function of return-set size for  $\sigma$  from 0 to 1.0.

degree to which the best item in the return set satisfies the customer's requirements and the cognitive load imposed on the customer to select that item.

An empirical evaluation simulating the behavior of a customer population indicated that unless customers place some importance on almost all features or the total number of features is quite small, an uninformed hypothesis about customer weights leads to low benefit. Simulation of LCW showed that, under the experimental conditions:

- LCW's estimate of the mean preferences of the customer population improves as the number of customers who have been seen increases.
- This improvement in the estimate of mean customer preferences leads to improved ranking performance for individual customers. Specifically, as the number of customers seen increases, the minimum return-set size needed to insure that an optimal item is returned decreases.
- Ranking performance when only one retrieval is permitted is almost identical to ranking performance when multiple retrievals are permitted.
- For small return-set size, benefit is higher for low  $\sigma$  customer populations than for high  $\sigma$  customer populations. The data suggest that the return set that optimizes benefit may be lower for low than for high  $\sigma$  customer populations.

These conclusions must be qualified by uncertainty concerning the assumptions underlying the experimental evaluation, i.e., whether:

- Customer preferences can be accurately represented by a linear weight vector.
- The mean weights of the customer population as a whole are distributed roughly into two levels.
- The weights of the individual customers for each feature are normally distributed around the mean value for that feature.

LCW illustrates how unobtrusive learning methods can improve the performance

of web-based applications by customizing their performance to customers' behavior. Techniques for learning by observation are likely to become an increasingly important aspect of e-commerce.

## References

- Bonzano, A., Cunningham, P. & Smyth, B. (1997), Using introspective learning to improve retrieval in CBR: A case study in air traffic control, in 'Proceedings of the Second International Conference on Case-Based Reasoning (ICCBR-97)', Springer, Providence, Rhode Island, pp. 291–302.
- Branting, K. (1999), Active exploration in instance-based preference modeling, in 'Proceedings of the Third International Conference on Case-Based Reasoning (ICCBR-99), Lecture Notes in Artificial Intelligence 1650', Monastery Seon, Germany.
- Branting, K. & Broos, P. (1997), 'Automated acquisition of user preferences', *International Journal of Human-Computer Studies* **46**, 55–77.
- Burke, R., Hammond, K., Kulyukin, V., Lytinen, S., Tomuro, N. & Schoenberg, S. (1997), Question answering from frequently-asked question files: Experiences with the FAQ finder system, Technical Report TR-97-05, University of Chicago, Department of Computer Science.
- Dent, L., Boticario, J., McDermott, J., Mitchell, T. & Zabowski, D. (1992), A personal learning apprentice, in 'Proceedings of Tenth National Conference on Artificial Intelligence', AAAI Press/MIT Press, San Jose, CA, pp. 96–103.
- Goldberg, D., Nichols, D., Oki, B. & Terry, D. (1992), 'Using collaborative filtering to weave an information tapestry', *Communications of the ACM* **35**(12), 61–70.
- Keeney, R. & Raiffa, H. (1993), *Decisions with Multiple Objectives: Preferences and Value Tradeoffs*, second edn, Cambridge University Press.
- Kira, K. & Rendell, L. (1992), The feature selection problem: Traditional methods and a new algorithm, in 'Proceedings of the Tenth National Conference Conference on Artificial Intelligence (AAAI-92)', MIT Press, pp. 129–134.
- Kohlmaier, A., Schmitt, S. & Bergmann, R. (2001), A similarity-based approach to attribute selection in user-adaptive sales dialogs, in D. Aha & I. Watson, eds, 'Fourth International Conference on Case-Based Reasoning, ICCBR 2001, Lecture Notes in Artificial Intelligence 2080', Springer, Vancouver, BC, Canada, pp. 306–320.
- Kolodner, J. (1984), *Retrieval and Organizational Strategies in Conceptual Memory: a Computer Model*, Lawrence Erlbaum Associates, Hillsdale, NJ.
- Maes, P. (1994), 'Agents that reduce work and information overload', *Communications of the ACM* **37**(7), 31–40.
- Nielson, J. (2000), *Designing Web Usability*, New Riders Publishing, Indianapolis, Indiana, USA.
- Stahl, A. (2001), Learning feature weights from case order feedback, in D. Aha & I. Watson, eds, 'Case-Based Reasoning Research and Development, 4th International Conference on Case-Based Reasoning, ICCBR 2001, Lecture Notes in Artificial Intelligence 2080', Springer, Vancouver, BC, Canada, pp. 502–516.
- Sweller, J., Chandler, P., Tierney, P. & Cooper, M. (1990), 'Cognitive load as a factor in the structuring of technical material', *Journal of Experimental Psychology: General* pp. 176–192.
- Wettschereck, D. & Aha, D. (1995), Weighting features, in 'Lecture Notes in Artificial Intelligence 1010', Springer, Sesimbra, Portugal, pp. 347–358.
- Wilke, W. (1999), Knowledge Management for Intelligent Sales Support in Electronic Commerce, PhD thesis, University of Kaiserslautern.
- Wilke, W., Lenz, M. & Wess, S. (1998), Intelligent sales support with CBR, in M. Lenz, B. Bartsch-Spoerl, H.-D. Burkhard & S. Wess, eds, 'Case-Based Reasoning Technology: from Foundations to Applications. LNAI 1400', Vol. 1400, Springer, pp. 91–113.
- Zhang, Z. & Yang, Q. (1999), Dynamic refinement of feature weights using quantitative introspective learning, in 'Sixteenth International Joint Conference on Artificial Intelligence (IJCAI-99)', Morgan Kaufmann, pp. 228–233.

**L. Karl Branting** is Principal Research Scientist at LiveWire Logic, Inc., where his current work focuses on design of agent architectures for customer relations management and empirical approaches to natural language processing. Dr. Branting was a professor of computer science at the University of Wyoming from 1991 to 2001. In addition, he was a Fulbright Senior Scholar at the University of Kaiserslautern, a Judicial Fellow at the United States Supreme Court, and an NSF Career grant recipient. Dr. Branting's computer science publications include a book, two edited collections, and more than 50 journal and conference papers. Dr. Branting graduated magna cum laude in philosophy from the University of Colorado and received a J.D. from Georgetown University and a Ph.D. in Computer Science from the University of Texas at Austin.