



Personalization of Supermarket Product Recommendations

R.D. LAWRENCE

G.S. ALMASI

V. KOTLYAR

M.S. VIVEROS

S.S. DURI

IBM T.J. Watson Research Center, P.O. Box 218, Yorktown Heights, New York 10598, USA

ricklawr@us.ibm.com

almasi@us.ibm.com

vkotlyar@us.ibm.com

viveros@us.ibm.com

sastry@us.ibm.com

Editors: Ron Kohavi and Foster Provost

Abstract. We describe a personalized recommender system designed to suggest new products to supermarket shoppers. The recommender functions in a pervasive computing environment, namely, a remote shopping system in which supermarket customers use Personal Digital Assistants (PDAs) to compose and transmit their orders to the store, which assembles them for subsequent pickup. The recommender is meant to provide an alternative source of new ideas for customers who now visit the store less frequently. Recommendations are generated by matching products to customers based on the expected appeal of the product and the previous spending of the customer. Associations mining in the product domain is used to determine relationships among product classes for use in characterizing the appeal of individual products. Clustering in the customer domain is used to identify groups of shoppers with similar spending histories. Cluster-specific lists of popular products are then used as input to the matching process.

The recommender is currently being used in a pilot program with several hundred customers. Analysis of results to date have shown a 1.8% boost in program revenue as a result of purchases made directly from the list of recommended products. A substantial fraction of the accepted recommendations are from product classes new to the customer, indicating a degree of willingness to expand beyond present purchase patterns in response to reasonable suggestions.

Keywords: recommender systems, personalization, collaborative filtering, data mining, clustering, associations, pervasive computing

1. Introduction

We describe a personalized recommender system designed to suggest new products to supermarket shoppers based upon their previous purchase behavior. The recommender system has been implemented as part of the “SmartPad” remote shopping system (Kotlyar et al., 1999) developed by IBM and Safeway Stores plc, a major supermarket retailer in the UK. This remote shopping system allows customers to prepare their shopping lists on a personal digital assistant (PDA) device such as a PalmPilot and transmit their order for subsequent pickup at the store without having to walk the aisles of the store. Although this latter feature is viewed as a convenience by a number of shoppers, it does remove the opportunity to suggest new or previously unpurchased products via special displays

in the store and so forth. The personalized recommendation system was developed as a substitute “spontaneous purchase” mechanism for this remote shopping system in a weakly-connected “pervasive computing” environment. The recommendations are computed on the server, and delivered to an individual customer’s PDA; obviously, the recommendations could also be delivered via more conventional mechanisms such as a web browser, electronic mail, or postal mail.

A number of web-based personalized recommender systems have been proposed recently (Resnick and Varian, 1997; Shardanand and Maes, 1995; Konstan et al., 1997; Borchers et al., 1998; Aggarwal et al., 1999; Personalization Summit, 1999). Personalization works by *filtering* a candidate set of *items* (such as products or web pages) through some representation of a *personal profile*. Two main paradigms for the filtering have emerged: *content-based* and *collaborative*.

A *content-based filtering* system recommends items based on their similarity to what a given person has liked in the past. Typically, both items and profiles are represented as vectors in the space of *features* and their similarity is computed via a standard distance metric, such as cosine coefficient. This approach has its roots in the vector-based model of Information Retrieval (IR), where text documents and user queries (or preferences) are both represented as vectors in the space of keywords or phrases, often referred to as *terms*. The coordinates of the vectors depend on the discriminating value of the respective terms. For example, if we were to recommend web pages from a site that reports news in the computer industry, then the term “computer” has low discriminating value, since it likely occurs in most of the pages. The term “product recommender” is likely to have high discriminating value, since we expect only a subset of pages to refer to product recommenders. A standard metric, called TFIDF (“term frequency/inverse document frequency”) is used in the IR literature to quantify discriminating value of document features (see Salton, 1989; Salton and McGill, 1983 textbooks for details).

In order to use content-based filtering in recommending grocery products, we must define the space of features. We started with product taxonomy that was available in the Safeway database. As discussed in Section 3.2 below, the taxonomy divides products in coarse-grain classes, such as “Pet foods”, and, further, into sub-classes, such as “Canned Dog Food”. We can directly use the classes and sub-classes as features of products and personal profiles. A person indicates interest in a particular feature by buying products within the corresponding class or sub-class. Products to be recommended can then be determined by computing a measure of distance between vectors representing personal preferences and vectors representing products.

The above strategy suffers from the problem of *overspecialization* (Balabanovic and Shoham, 1997): it provides us with no rigorous basis for introducing shoppers to new kinds of products beyond those classes and sub-classes that they already buy. For example, a person who buys dog food might also be interested in carpet cleaners, but we have no way of gauging this interest. We can solve this problem by assigning an implicit feature such as “appeals to dog owners” to both products. This demonstrates both the flexibility and weakness of content-based filtering systems. As new information about relationships between products becomes available it can be incorporated by defining new features. But, for this same reason, it might not be possible to use a content-based filtering system without

human intervention, since we have to recognize the need to introduce new features into content.

*Collaborative filtering*¹ aims to sidestep the problem of feature design by recommending items that other people, who are similar to the person in question, have liked. A collaborative filtering system, such as Ringo (Shardanand and Maes, 1995) and GroupLens (Konstan et al., 1997), works by collecting explicit user ratings of items in question (e.g., movies, CDs or USENET postings). Users are then compared based on how similar their ratings are, and they are recommended items favored by other people with similar interests. To compute these “word-of-mouth” recommendations, it has been suggested in the literature to use clusters in the space of user profiles in order to define prototypical profiles (Ungar and Foster, 1998).

Viewed broadly, collaborative filtering suggests using the information about a group, which can be the whole population of users or a cluster, in order to produce individual recommendations. In our recommender system we use two sources of such information. First, we apply *associations mining* (Agrawal and Srikant, 1994) to customer purchase data in order to derive relationships between product classes and sub-classes. Since these relationships are based on actual purchases, we expect to identify additional product-class relationships (e.g., purchasers of dog-food products also buy carpet cleaners) that are not captured by the product taxonomy (or even by more sophisticated keyword matching). Second, we use *clustering* (Everitt, 1993) to assign customers into groups with similar interests, based on prior purchase patterns. By itself, content-based filtering does not incorporate information about the relative popularity of products among other customers. For this reason, we use the cluster analysis to build ranked lists of the most popular products among customers assigned to each cluster; recommendations for a specific customer are then drawn from products popular among other members of the customer’s cluster.

In summary, our recommender system uses content-based filtering at its core, with the ideas from collaborative filtering utilized both to refine the content model and to make recommendations dependent on shared interests within customer clusters.

The organization of this paper is as follows. Section 2 provides a brief overview of the overall SmartPad system, with an emphasis on those aspects relevant to recommender systems. Details of the recommender system are provided in Section 3, while Section 4 describes the associations mining and clustering analyses. Section 5 describes early results and user feedback obtained to date from a field trial of the system, and Section 6 provides an overall summary.

2. The SmartPad remote shopping system

Figure 1 provides an overview of the SmartPad environment (Kotlyar et al., 1999). The overall system consists of a network of mobile devices (PDAs) connected to the SmartPad server through a dial-in service. Each customer participating in the program is issued a PDA which runs a consumer application enabling the user to build a shopping list and send it to the server. Products in the order are chosen from three personal databases (PDBs) stored on the PDAs: personal catalog, recommendations, and promotions. (Promotions refer to discounted products offered by the participating Safeway store, while recommendations

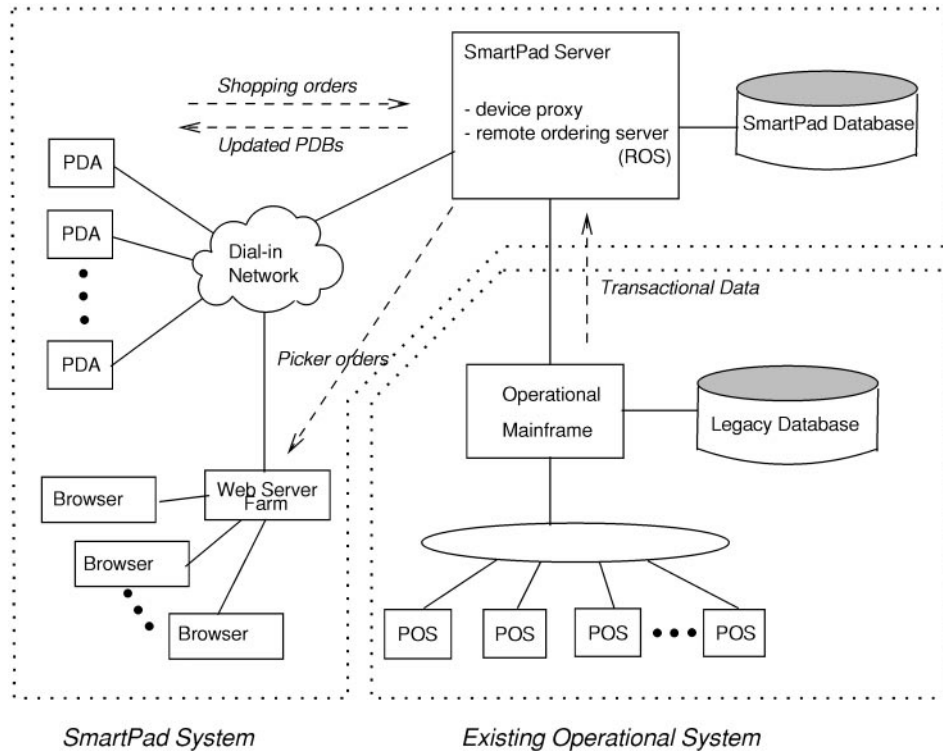


Figure 1. Overview of the SmartPad system.

are products suggested by the system discussed in this paper.) Together, the PDBs serve as a surrogate for the store catalog, which is too large to be stored or shown to the user on the device. Products not appearing on any of the PDBs can be added to the order list manually, either as free-form text or as product numbers from a printed catalog.

As shown in figure 1, transaction data from the existing Safeway operational system is made available to the SmartPad server. The SmartPad database contains detailed product information as well as customer spending histories for a large number of Safeway customers, including all participants in the actual SmartPad program. As discussed in the following section, this data provides the raw input data for the recommender system including the required data mining analysis.

All computations associated with the personalized recommender system are performed on the SmartPad server. In the current trial program, new sets of recommendations are generated weekly for all SmartPad participants, and stored on the server in the SmartPad database. These recommendations, along with the new promotions and the updated personal catalog, are transmitted to the targeted customer by synchronizing the PDBs on the customer's device at the next server connection initiated by this customer.

Figure 2 shows a view of the PDA main screen for the SmartPad consumer application. The various tabs shown here are used to access lists such as the personal catalog, store



Figure 2. SmartPad main screen.

promotions, and recommendations, as well as other useful information such as bonus points accumulated by the customer.

3. The personalized recommender system

3.1. Overview of the recommender system

Figure 3 provides an overview of the analysis involved in the personalized recommender system. The input data from the SmartPad database consists of descriptions of approximately 30,000 products and a database containing summarized customer purchase data for roughly 20,000 Safeway customers, including the customers enrolled in the SmartPad pilot program. Not all products are eligible for recommendation; for example, we avoid

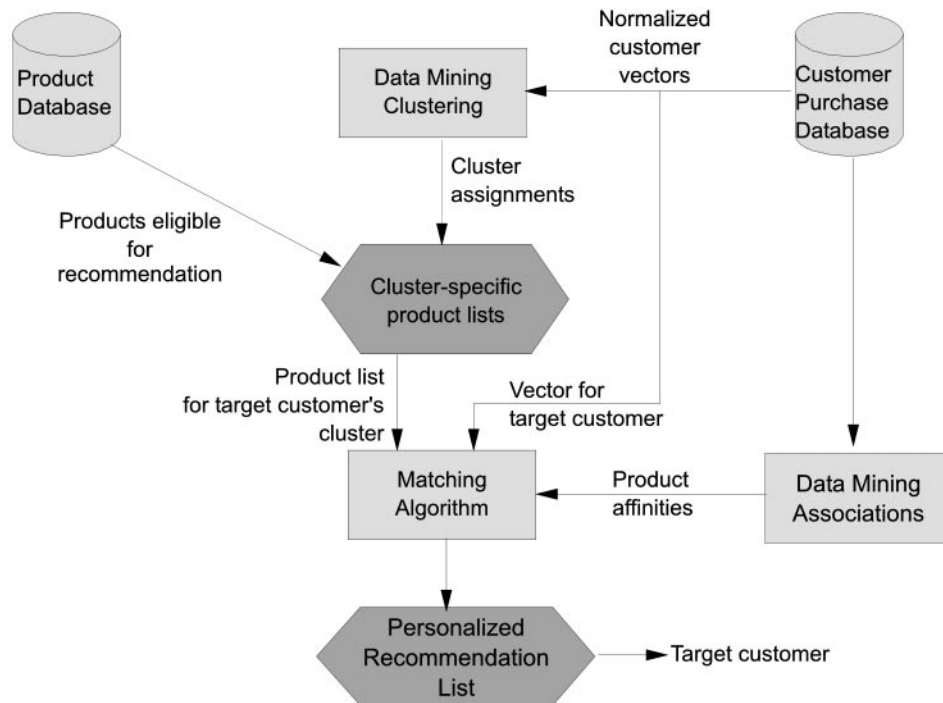


Figure 3. Overview of the personalized recommender system.

recommending tobacco, health products, and other inappropriate product classes. As discussed in Sections 3.3 and 4.2, all customers in the database are clustered based on normalized spending vectors constructed using their previous purchase behavior. Each customer is assigned to a single cluster, allowing us to build separate lists for each cluster containing the most frequently purchased products among customers in the cluster. Recommendations for a specific target customer enrolled in the SmartPad program are then drawn from the list of popular products in the assigned cluster for this customer. This cluster-specific list is passed to a “matching engine”, which ranks this list of products according to the expected appeal to the target customer. The matching engine also utilizes product relationships or “affinities” computed using associations mining, as described in Sections 3.4 and 4.1. The personalized recommendation list, comprising the 10 to 20 products with the highest scores, are returned to the customer’s device during the next server-connection session. It is important to note that the recommendation list, by design, will contain no products previously purchased by this customer. Unlike promotions, no price discounts or other incentives are offered on recommended products.

3.2. Product taxonomy

Figure 4 illustrates the three-level hierarchical product taxonomy used by Safeway. Products are divided across $G = 99$ product classes. Each product class is subdivided into fewer

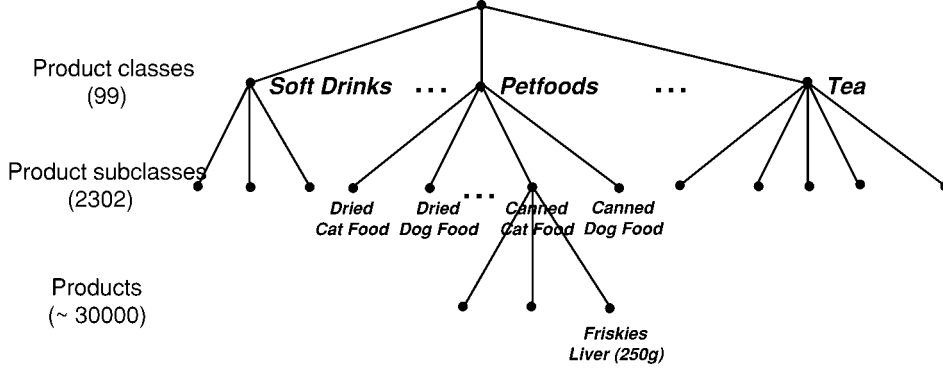


Figure 4. Safeway product taxonomy.

than 100 subclasses, generating a total of $S = 2302$ product subclasses. Absolute customer spending is available at any level in this hierarchy. Although specific to Safeway, the taxonomy used here is typical of the multi-level hierarchy often used to represent retail product catalogs.

Given a specific target customer, the recommendation system seeks to determine products which are best matched to this customer's spending profile. For this reason, we need to be able to construct a profile or vector for each customer representing the customer's "interest" across a range of attributes. This representation is designed to be similar to a vector of explicit ratings typically used in collaborative filtering applications (Shardanand and Maes, 1995). Similarly, we need to be able to represent the "appeal" of each product measured in this same attribute space. We use the subclass level of the Safeway product taxonomy to construct this attribute space. With reference to figure 4, this defines an $S = 2302$ dimensional space in which we will match products to customers. In the following subsections, we describe the construction of the customer and product vectors.

3.3. Customer model

We represent the absolute spending for customer m as

$$\mathbf{C}^{(m)} = [C_{m1}, \dots, C_{ms}, \dots, C_{mS}]^T, m = 1, \dots, M, \quad (1)$$

where C_{ms} denotes the absolute spending of customer m across all products contained in subclass s , M is the total number of customers, and S is the number of product subclasses described above. C_{ms} is computed from the raw transaction data as the sum over the past several months of spending data, and hence reflects this customer's grocery purchases over multiple shopping visits.

We then apply two separate normalizations to this result to obtain the final customer vector. First, we convert absolute spending for each subclass to fractional spending in this

subclass by normalizing to this customer’s total spending over the period:

$$\hat{C}_{ms} = \frac{C_{ms}}{\sum_{s'=1, \dots, S} C_{ms'}} \quad (2)$$

The resulting *fractional spending* vector characterizes this customer’s interest in each subclass relative to other subclasses. Since commonly purchased subclasses such as fresh vegetables will tend to dominate the fractional spending, it is also useful to take the ratio of the individual customer’s fractional spending in a subclass to the mean value for this subclass taken over all other customers:

$$\hat{\hat{C}}_{ms} = \frac{\hat{C}_{ms}}{\frac{1}{M} \sum_{m'=1, \dots, M} \hat{C}_{m's}} \quad (3)$$

This *normalized fractional spending* thus quantifies the customer’s interest in this subclass relative to the customer database as a whole, normalized such that entries equal to 1 imply an average level of interest in a subclass relative to all other customers.

3.4. Product model

We use the following notation to represent the product classifications illustrated in figure 4:

$$\begin{aligned} \mathcal{S}(n) &\equiv \text{the product subclass number for the product } n \\ \mathcal{C}(s) &\equiv \text{the product class for the product subclass } s \end{aligned} \quad (4)$$

For brevity, we will use “ $\mathcal{C}(n)$ ” instead of “ $\mathcal{C}(\mathcal{S}(n))$ ” to denote the class that product n belongs to.

Each product $n = 1, \dots, N$ is represented by an S -dimensional vector $\mathbf{P}^{(n)}$, and hence has the same dimensionality as the customer vectors discussed in the previous section. The individual entries $P_s^{(n)}$, $s = 1, \dots, S$, reflect the “affinity” the product has to the subclass s , i.e., the extent to which this product will appeal to a customer with an interest in this subclass. The simplest solution is to set the entries of the vector $\mathbf{P}^{(n)}$ as:

$$P_s^{(n)} = \begin{cases} 1 & \text{if } s = \mathcal{S}(n) \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

This captures the affinity between different products in a subclass: a purchase of a product in a subclass implies an interest in other products in the same subclass. This simple form leads to straightforward content-based filtering: if you purchased from this product subclass, you will receive recommendations only for other products within this subclass. Clearly, this eliminates the possibility of “cross-selling” potentially interesting products outside of this immediate subclass.

More specifically, Eq. (5) ignores other possible kinds of affinity between products:

1. affinities between products that belong to the same class (but to different subclasses), and

2. affinities that are derived from data-mining analysis such as associations mining (Agrawal and Srikant, 1994) performed independently at both the class and subclass level.

Case (1) is also content-based filtering since only product-taxonomy information is used. However, case (2) incorporates associations derived from “global” customer behavior, and hence reflects aspects of collaborative filtering. In order to account for the first kind of affinity, we have heuristically chosen to set $P_s^{(n)} = 0.5$ for every subclass s that is in the same class as the product n . To account for the second kind, we set $P_s^{(n)} = 0.5$ for every association $\mathcal{C}(n) \implies \mathcal{C}(n')$ between the class of the product n and the class of product n' . We also compute associations at the subclass level, and use these as described below to set affinities between associated subclasses.

Some examples of the use of associations to specify the product vector are as follows. If the use of bacon strongly implies the use of eggs (i.e., bacon and eggs are associated classes), then the vector for a particular bacon product will receive a 0.5 in the entries that correspond to all subclasses of the class “Eggs” (e.g., “Eggs/Large” or “Eggs/Medium”). Also, if there is a strong association between the subclasses “Bacon/LowFat” and “Eggs/Whites”, then the vector for the particular low-fat bacon product will receive a 1.0 in the entry for the subclass “Egg/Whites”. Given the class-level association between “Bacon” and “Eggs”, the vector for a particular low-fat bacon product will also receive contributions (0.25) in the entries associated with *all* subclasses with “Eggs”.

Overall, the formula for the entries in the product vector is:

$$P_s^{(n)} = \begin{cases} 1.0 & \text{if } s = \mathcal{S}(n) & (\text{within same subclass}) \\ 1.0 & \text{if } \mathcal{S}(n) \implies s & (\text{within associated subclass}) \\ 0.5 & \text{if } \mathcal{C}(s) = \mathcal{C}(n) & (\text{subclass within same class}) \\ 0.25 & \text{if } \mathcal{C}(n) \implies \mathcal{C}(s) & (\text{within subclass of associated class}) \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

The multipliers for subclass associations were set higher than those at the product class level in order to encourage customers to try popular new products outside their current shopping pattern, but not too far outside. The results in Section 5 seem to indicate that this intended effect was indeed achieved.

3.5. Matching algorithm

In the preceeding sections, we have developed descriptions of customers and products in the vector space defined by product subclasses. The final step in the recommendation process [see figure 3] is to score each candidate product for a specific customer and select the best matches. This score should reflect the degree of similarity between the customer vector and each product vector. Possible metrics include Pearson correlation (Shardanand and Maes, 1995), Euclidian distance (Shardanand and Maes, 1995), and cosine projection (Salton, 1989). As noted in Salton (1989), the choice among these methods is heuristic, and we chose cosine projection based on inspection of the generated recommendations. Hence, the score σ_{mn} between customer m and product n is computed using a cosine coefficient

between the corresponding vectors, $\mathbf{C}^{(m)}$ and $\mathbf{P}^{(n)}$:

$$\sigma_{mn} = \rho_n \frac{\mathbf{C}^{(m)} \cdot \mathbf{P}^{(n)}}{\|\mathbf{C}^{(m)}\| \|\mathbf{P}^{(n)}\|}, \quad (7)$$

where the \cdot operator denotes inner (dot) product:

$$\mathbf{x} \cdot \mathbf{y} = \sum_i x_i y_i, \quad (8)$$

$$\|\mathbf{x}\| = \sum_i |x_i|, \quad (9)$$

and ρ_n is a modulation factor discussed below. Both customer and product vectors are sparse, and we use this sparse structure to minimize the number of arithmetic operations necessary to evaluate Eq. (7).

If $\rho_n \equiv 1$ for $n = 1, \dots, N$ in Eq. (7), all products in a product subclass will have identical scores for a specific customer since the product vectors are constructed at the subclass level. Hence, ρ_n can be used as a “tie-breaker” to differentiate products. In general, ρ_n can reflect the desire to “push” one product over another, for example, due to higher inventories, profit margins, or other market conditions. Since we have product profit margins available, we initially computed this factor as

$$\rho_n = \left[\frac{PM_n}{\overline{PM}} \right]^\alpha, \quad (10)$$

where PM_n is the profit margin associated with product n , \overline{PM} is the mean value over all products, and α is an empirical factor designed to control the overall influence. However, we discovered that more profitable items were not necessarily more popular, and we chose not to implement this feature in the production system. Instead, we break ties by random sampling from up to 5 products with the top scores within a specific product subclass. In order to distribute the recommendations across product classes, we limit the number of recommendations for each customer to 1 product per product subclass, and 2 products per product class. In the end, we return a list of approximately 10 to 20 products to the customer’s PDA as described in Section 2; this number is determined primarily by the size of the display on the PDA device.

Previously bought products are excluded from the recommend list, since the recommender is meant to broaden each customer’s purchase patterns. It is also worth noting that the list of the most popular products in each product subclass is a dynamic list, created by periodic querying of the transaction database. This prevents recommending seasonal items at inappropriate times of the year. Since recommendations are generated on a weekly basis, the list of candidate products could also be updated weekly; in practice, we do it monthly, with an optional weekly edit to comb out special situations such as Easter candy the week after Easter. It should also be noted that several of the 99 product classes, such as Petrol and Tobacco, are excluded from the recommender by design.

4. Data mining analysis

Intelligent Miner for Data (Intelligent Miner, 1998) was used to perform the associations mining and clustering analysis shown in figure 3. Customer transaction data were used as input to both calculations, but, naturally, the data were formatted and normalized differently.

4.1. Associations mining

As discussed in Section 3.5, associations mining is used to compute the product affinities necessary to construct the product vectors in Eq. (6). We used Intelligent Miner’s “apriori” associations algorithm (Agrawal and Srikant, 1994) to extract associations (independently) among the 99 product classes and also among the 2302 product subclasses shown in figure 4. The raw input data for this analysis was 8 weeks of product-level transaction data for 8000 customers with above-average spending. In order to compute product-class associations, the transactions were first binned according to product class, and then aggregated over time. The final input table consisted of (customer ID, product-class ID) tuples, each representing a product class in which the customer has purchased over this 8-week period. An analogous procedure was used to extract associations at the product-subclass level. Note that our objective here differs from conventional market-basket analysis in which associations are computed among products purchased as part of the same transaction.

Consistent with the form of Eq. (6), we computed only simple associations containing a single item in both the body and the head of the rule, e.g., $A \Rightarrow B$. For several reasons, we needed to limit the number of rules we used to about 100 from each of the two levels. These reasons included clarity of understanding as well as computational complexity, especially in the case of the product-class rules, since the recommender algorithm assumes that all the subclasses within two related classes are also related via affinities. We used a combination of thresholds on the support, confidence, and lift of a rule to achieve this filtering.² After experimenting with several datasets and inspecting the resulting rules, we chose a combination of minimum support in the 1% – 4% range, minimum confidence of 30% to 40%, and minimum lift of 2 to 3. Although heuristic, these choices produced a reasonable number of associations for the various Safeway datasets we examined. Applications to other purchase data may require different limits on these parameters.

The associations analysis is performed periodically as a “batch” process. Rules meeting the filtering criterion are stored in the SmartPad database (see figure 1) for subsequent access during the computation of the personalized recommendations.

Table 1 shows a subset of affinities computed at the product-class and subclass levels. The rule head is on the right, so the textual format of the first rule is “5.9% of all customers buy both baby products and canned pasta. When a customer buys baby products, the customer also buys canned pasta in 41% of the cases, which is 2.4 times the rate one would expect if the sales of these two were statistically independent”. In tags of the form “Baby:Disposable Nappies”, the items before and after the colon refer to product class and product subclass, respectively. The numbers preceding these classes are actual class and subclass numbers; we will refer to these in the discussion of specific recommendations in Section 5.1.

Table 1. Sample associations computed at the product-class and product-subclass levels.

Sup	Conf	Lift	Class or subclass	Relevant affinities
0.059	0.41	2.4	20(Baby products)	\Rightarrow 41(Canned pasta)
0.082	0.47	2.2	66(Table wines)	\Rightarrow 68(Beer/Lager/Spirits)
0.125	0.50	2.0	90(Fresh beef)	\Rightarrow 91(Pork/Lamb)
0.025	0.38	9.0	2010(Baby:Disposable nappies)	\Rightarrow 2007(Baby:Wipes)
0.016	0.33	4.9	2010(Baby:Disposable nappies)	\Rightarrow 1012(Dairy:Childrens' yogurt)
0.01	0.33	4.9	2010(Baby:Disposable nappies)	\Rightarrow 3115(Instore:Babysitting center)
0.012	0.37	3.4	1020(Dairy:Childrens' fromage)	\Rightarrow 3115(Instore:Babysitting center)
0.016	0.52	5.2	2306(Biscuits:Kids biscuits)	\Rightarrow 3115(Instore:Babysitting center)
0.022	0.30	4.9	9015(Fresh beef:Beef joints)	\Rightarrow 9120(Pork/Lamb:Pork joints)

4.2. Clustering analysis

Safeway assigns each customer to a pre-defined customer class (or segment) based on purely demographic information derived from questionnaires. Examples of these classes are “empty nesters”, “young parents”, and so on. While generally useful, this demographic-based segmentation does not provide the necessary detail on each customer’s actual spending preferences. For this reason, we did not use these classifications, but instead clustered customers on the basis of their spending in the 99 product classes. We produced lists of the most popular products per subclass for each cluster, and then used the cluster-specific list of popular products as input to the generation of recommendations for a customer in a particular cluster. The five most popular products in a subclass can be quite different in different clusters, as shown in more detail further below.

Given a large number of database records, clustering (Michaud, 1999) can be useful for identifying a small number of prototypes (i.e., cluster centers) which represent dominant characteristics or features present in the input data set. Key issues in any clustering analysis are the dimensionality and the normalization of the input data. If the dimensionality of the input space is too large, and a typical record has only a few nonzero entries, then it is difficult to find records with many nonzero attributes in common. For this reason, we did not cluster at the product spending level. We also found that clustering at the subclass spending level did not produce useful results, and thus all clustering analysis was done using customer spending records at the 99-dimensional product-class level (see figure 3). The input data for this analysis were normalized fractional customer spending vectors, computed as in Eqs. (1)–(3), but evaluated at the product-class level rather than at the product-subclass level. These values were limited to a maximum value of 5 in order to limit the influence of very high product-class spending.

We applied both the neural clustering algorithm (Lawrence et al., 1999) and the so-called “demographic clustering” algorithm³ (Michaud, 1999) that are available in Intelligent Miner. Identical input data were provided to both algorithms. Using some of the classical measures of cluster validity, these results seemed to suggest that the resulting clusters were relatively diffuse, i.e., were not particularly compact. For example, the

demographic clustering method, which seeks to maximize the Condorcet criterion⁴ for cluster quality, produced an overall Condorcet criterion of 0.56 for 9 clusters, but did so by putting 97% of the records in one cluster. Neural clustering produced clusters that were much more uniform in size, but the Condorcet criterion was in the range of 0.1 to 0.2, and values of the modified Dunn's index⁵ used in Bezdek (1998) were in the range 0.3 to 0.8 compared with the values between 1 and 2 that they obtain for "good" clusters.

However, these cluster validity measures tend to favor categorical data with a relatively small number of attributes, and may be less relevant for our continuous data and relatively large number of attributes. There is a good deal of overlap in people's spending, for understandable reasons. One group of people may have extraordinary spending in baby products and none in tobacco, and another may be just the reverse, but both groups presumably need paper towels, milk, bread, lettuce, and many other everyday necessities. For this reason, it seems plausible to evaluate the clusters in terms of their dominant attributes, rather than strictly by their compactness and distinctiveness relative to other clusters.

For example, figure 5 shows a 3 by 3 self-organizing feature map (Kohonen, 1995) produced by the neural clustering method. Each cell represents a cluster, with the largest cluster containing 17% of the customers, the smallest 9%. The product classes shown represent the 3 attributes which most distinguish members of the cluster from the database background. For example, the cluster in the center of this map is likely to represent "families with young children", based on significantly above-average spending in baby products, clothing, and dairy products. The cluster at the lower right appears to represent "serious bakers", given dominant spending in sugar, home baker products, and canned fruits.

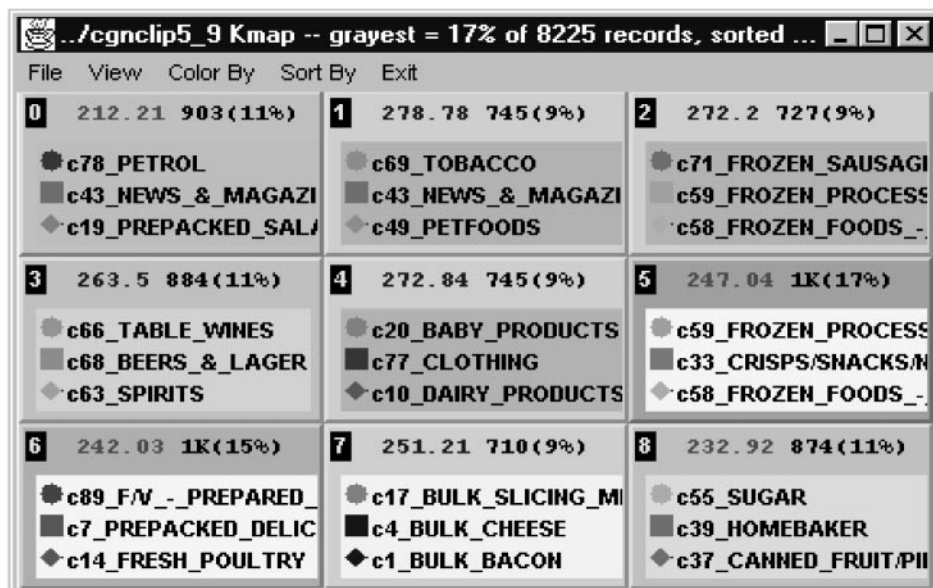


Figure 5. Cluster analysis at the product-class level.

Table 2. Summary of cluster characteristics.

Cluster number	Fraction of records	Most significant product class	Cluster share	Cluster-share enrichment
0	0.11	Petrol	0.24	2.2
1	0.09	Tobacco	0.50	5.6
2	0.09	Frozen sausage	0.67	7.4
3	0.11	Wines	0.39	3.5
4	0.09	Baby products	0.56	6.2
5	0.17	Frozen foods	0.24	1.4
6	0.15	Fresh vegetables	0.25	1.7
7	0.09	Bulk sliced meat	0.34	3.8
8	0.11	Sugar	0.22	2.0

The clusters in the upper left focus on specific non-food items, e.g., petrol,⁶ tobacco, and wines/beer/spirits. Feature maps generated with more clusters show increasing refinements in these customer purchase patterns.

Table 2 summarizes these same 9 clusters, showing the most significant product class which defines the cluster, as well as the fraction of total customers assigned to this cluster. We define the “cluster share” of a product class as the spending in that product class by a cluster of customers divided by the total spending in that product class:

$$\text{Cluster share in class } i = \frac{\sum_{m \in \text{cluster}} C_{mi}}{\sum_{m=1, \dots, M} C_{mi}}, \quad (11)$$

where C_{mi} is the absolute spending for customer m in product class i . The enrichment in cluster share is defined as the ratio of cluster share divided by the fraction of the customer population which is in the cluster. For example, cluster number 4 is characterized by spending in baby products: the 9% of the customers in this cluster generated 56% of the total database spending for baby products, an enrichment by a factor 6.2. Popular products among members of this cluster are likely to be quite different than the most frequently purchased products among members of the “tobacco” cluster shown as cluster 1 in Table 2.

Figure 6 shows a specific example of preferred products in a cluster, compared to popular products across the database as a whole. For the background population, three of the five most popular chocolate bars are made by Mars and none by Nestle, but for cluster 4 (baby products) three of the top five are made by Nestle and only one by Mars, suggesting a preference for Nestle among young children, or at least among those who shop for them.

5. Results

The results that we report here were obtained during the first 8 months of the SmartPad program at Safeway (Kotlyar et al., 1999). Phase 1 of the program lasted approximately 7 months, and included 200 customers from one store. As of this writing, Phase 2 has been in

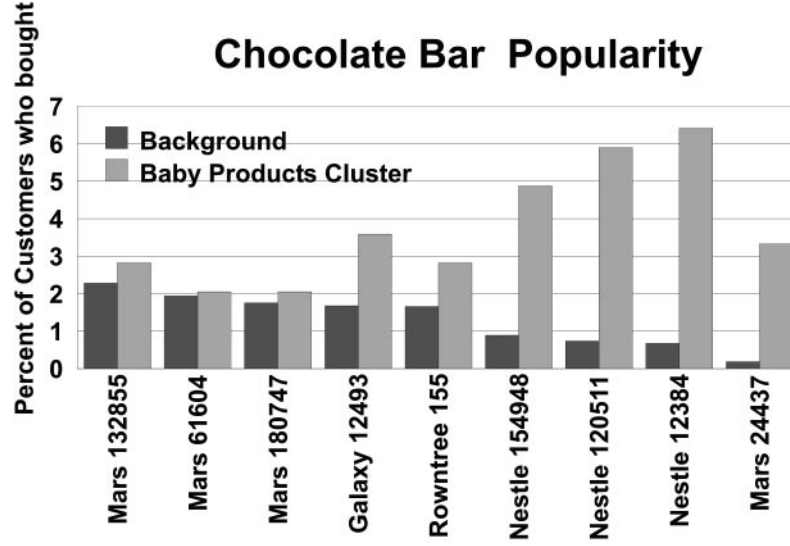


Figure 6. Cluster-specific product preferences.

progress for approximately two months, and includes a second participating Safeway store. Different versions of the recommender were used in the two phases. In this section, we first show a sample recommendation list and describe how the use of data mining influences its content. We then describe our experience with the first version of the recommender, the lessons we learned and consequent changes we made, and the performance of the current version. Full implementation of the cluster-specific input lists was not completed for these field tests: the input to the matching process was constructed using popular products within a single cluster consisting of all customers in the database.

5.1. Analysis of a sample recommendation list

Table 3 shows the fractional spending at the product subclass level for a customer who has spending in the “baby-products” subclasses. Figure 7(d) shows the recommendation list generated on the basis of this spending profile. The figure also traces the development of this recommendation list through four stages:

- As it would appear without the use of either data-mining associations or clustering, i.e., on the basis of content-based filtering alone.
- As it would appear if we only used association rules from the product class level. Note the appearance of an item from subclass 6812 due to a boost in its score resulting from the association rule in Table 1 between class 68 and 66, several of whose subclasses are present in the spending list in Table 3.
- As it would appear using association rules from both the class and subclass levels, but without clustering. Four new products have appeared in the top ten as a result of product subclass association rules found in Table 1. Note that the product that appeared

Table 3. Sample customer fractional spending at the product subclass level.

Subclass number	Spending fraction	Subclass name
7801	0.077	PETROL:PETROL
2010	0.058	DISPOSABLE NAPPIES:BABY PRODUCTS
735	0.046	INTERNATIONAL:PREPACKED DELICATESS
6652	0.035	WHITE ENGLISH:TABLE WINES
3005	0.032	CHOCOLATE:CONFECTIONERY
7734	0.027	GIRLS PAJAMAS:CLOTHING
4005	0.026	LIQUIDS:HOME LAUNDRY/LIQUID
6203	0.021	LUXURY ICE CREAM:ICE CREAM
9015	0.021	BEEF JOINTS:FRESH BEEF
4001	0.017	SOAP/DETERGENTS POWD:HOME LAUNDRY/LIQUID
2001	0.016	BABY FOOD DRY:BABY PRODUCTS
901	0.015	S/W FULL CREAM MILK:MILK
8004	0.015	ENT UK VIDEO:HOME ENTERTAINMENT
7741	0.015	BOYS BRIEFS:CLOTHING
7733	0.015	BOYS PAJAMAS:CLOTHING
6651	0.015	WHITE SOUTH AFRICA:TABLE WINES
4601	0.014	TOILET TISSUE:PAPER PRODUCTS
905	0.013	S/W SEMI SKIMMED MIL:MILK
6670	0.013	WHITE—NORTH AMERIC:TABLE WINES
2007	0.012	BABY WIPES:BABY PRODUCTS

in (b) has disappeared. However, although not shown in this figure, the score for the product from subclass 9120 was influenced by association rules from both the product class and subclass levels.

- (d) As it appears when the products for the top ten product subclasses are chosen from the list of popular products for the cluster that this customer falls into, namely, cluster 4 in figure 5. Note in particular the appearance of the Nestle candy bar instead of the Mars, consistent with the cluster preferences shown in figure 6.

The much greater effect of subclass association rules (figure 7(c) versus (b)), which results from the choice of values in Eq. (6), is reflected in the statistics for the full set of recommendation lists as well: on average, 33% of the items on a recommendation list are in product subclasses that are new to the customer (no spending in the subclass within the past 3 months) and 16% are in product classes new to the customer.

5.2. Results for phase 1

During Phase 1 of the trial, a total of 1957 complete orders were processed by the Smart-Pad system. Of these, 120 orders (6.1%) contained at least 1 product chosen from the

No Associations or Clustering...			+Class Associations		
SUBCLASS	SCORE	PRODUCT	SUBCLASS	SCORE	PRODUCT
2010	0.31	Safeway Nappies Maxi Plus Un	2010	0.31	Safeway Nappies Maxi Plus Un
6652	0.21	Safeway Mendoza Chardonnay	6652	0.21	Safeway Mendoza Chardonnay
735	0.20	Safeway Indian Selection 25	735	0.2	Safeway Indian Selection 25
4005	0.17	Ariel Future Washing Lqd Po	4005	0.17	Ariel Future Washing Lqd Po
3005	0.16	Mars Bar 5Pack	3005	0.16	Mars Bar 5Pack
9015	0.11	Safeway Beef Mini Joint 510g	9015	0.11	Safeway Beef Mini Joint 510g
6203	0.09	Ben and Jerry's Ice Cream Ph	6203	0.09	Ben and Jerry's Ice Cream P
901	0.09	Safeway Pasteurised Milk 1.136	901	0.09	Safeway Pasteurised Milk 2.
4601	0.08	Safeway Savers Toilet Tissue	4601	0.08	Safeway Savers Toilet Tissue
5621	0.07	Safeway Cordelia Luxury Bath	6812	0.07	Bacardi Breezer Tropical Li

(a) (b)

+Subclass Associations			+Clustering		
SUBCLASS	SCORE	PRODUCT	SUBCLASS	SCORE	PRODUCT
2010	0.38	Safeway Nappies Maxi Plus Un	2010	0.38	Safeway Supersoft Baby Wipes
1012	0.34	Safeway Monster Pots Set Yo	1012	0.34	Safeway Monster Pots Set Yo
3115	0.33	Safeway Creche Facility 0.5 HR	3115	0.33	Safeway Creche Facility 1 HR
6652	0.21	Safeway Mendoza Chardonnay	6652	0.21	Safeway Mendoza Chardonnay
735	0.21	Safeway Indian Selection 25	735	0.21	Safeway Ready Meal Indian Meal
4005	0.17	Ariel Future Washing Lqd Po	4005	0.17	Ariel Future Washing Lqd Po
3005	0.16	Mars Bar 5Pack	3005	0.16	Nestle Smarties 4Pack
4715	0.12	Heinz Spaghetti Bolognese 4	4715	0.12	Heinz Spaghetti Bolognese 4
9120	0.12	Safeway Pork Leg Boneless O	9120	0.12	Safeway Pork Leg Boneless R
9015	0.11	Safeway Beef Brisket Small	9015	0.11	Safeway Beef Brisket Small

(c) (d)

Figure 7. Progressive development of a sample recommendation list: (a) uninfluenced by data mining, (b) with product-class association rules, (c) with product subclass association rules added, (d) with clustering.

recommendation list. (It is important to recall that the recommendation list, by design, will contain no products previously purchased by this customer.) An objective of the product recommender is to provide a boost in revenue comparable to the spontaneous purchases a shopper might make while walking through the store or after receiving a flyer in the mail. By this measure, the results for the initial recommender were somewhat disappointing: the corresponding boost in revenue was 0.3% over and above the revenue generated by products bought from the main “personal catalog” shopping list.

As trial program progressed, we noticed with interest that the distribution of spending in the SmartPad product categories⁷ was different for items bought from the recommendation list versus the personal catalog (see figure 8), even though the distribution of items *available* from each list were quite similar. For example, wines accounted for only 3.5% of the revenue from the main shopping list, but 8.7% of the revenue from the recommendation list. By contrast, products in the household care category accounted for 12.1% of the revenue from the main shopping list but only 4.6% from the recommendation list. We interpreted these results to mean that there is a set of categories in which recommendations are more welcome than others, and interviews with participating customers confirmed that interpretation. They wanted more “interesting” recommendations, and wines fit that description but household care products did not.

Armed with this insight, we proceeded to trim the list of subclasses from which recommendable products were drawn, emphasizing those product classes in which the spending percentage from the recommendation list exceeded that on the main shopping list, and de-emphasizing the others, with the aim of creating a more “fun” or welcome set of

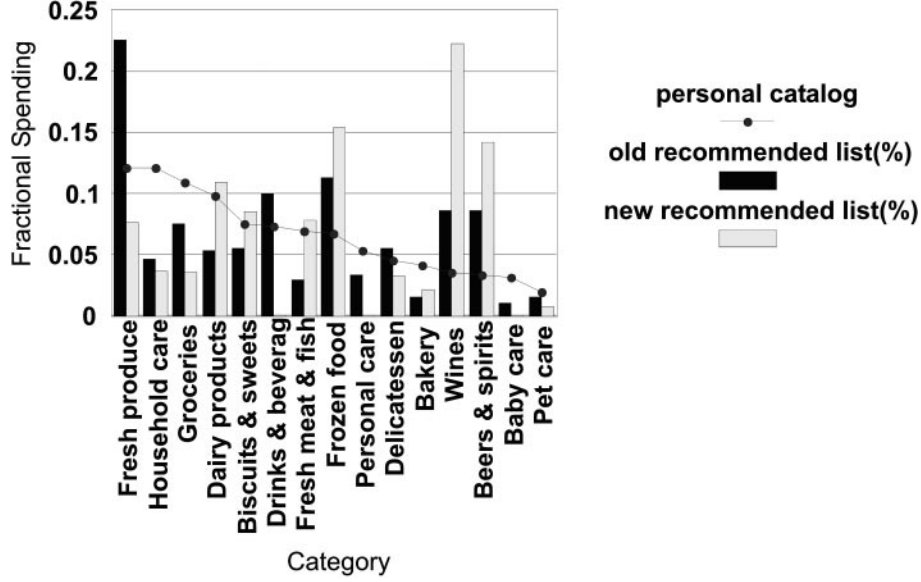


Figure 8. Product mix comparison of items bought from the main shopping list, from the old recommendation list, and from the new recommendation list. The fraction spent in 16 product categories is shown for each list.

recommendable products. We also added a second source of items eligible for recommendation, namely, new products introduced within the last month. We allowed new products to come from any category, on the rationale that their novelty made them interesting per se. This new recommender went into operation two months before this writing, and is the subject of the next section.

5.3. Results with the current version

As mentioned above, a second store and a new set of customers were added to the SmartPad program at the same time that the new recommender was introduced with the changes described above. The number of recommendations sent to the customer's device was increased from 10 products in the initial phase to 20 products in the current version.

For both stores, the results with the new recommender were better than those obtained with the old recommender at the original store, but this was much more pronounced for the new store, bearing out the saying that one gets only one chance to make a first impression.⁸ For the old store, the fraction of orders containing at least one recommended product increased from 6.1 to 7.7%, and the revenue boost rose to 0.5% from its previous value of 0.3%, a modest increase. But for the new store, the returns were much greater: 25% of the orders included at least one recommendation, with a revenue boost of 1.8%. Safeway considered this boost to be quite respectable given their experience with other promotional methods.

We were reasonably successful in meeting our goal of encouraging shoppers to try new things but not drastically new things: 51% of the acceptances from the recommendation lists corresponded to subclasses in which the shopper had spent no money in the previous three months, but only 4% corresponded to new product classes. For the recommendation lists themselves, 33% of the recommended products on average were from subclasses the shopper had not spent in before, and 16% were from product classes that were new to the customer. Outside the environment of the recommender, the rate of trying new subclasses is substantially lower, and the rate of trying new product classes is practically zero.

5.4. *Distributions of computed scores*

In order to quantify the impact of the recommender system, it would have been useful to have a control group of customers who received “placebo” recommendations, such as a list of randomly chosen products. This approach was not feasible, however, since we were dealing with a live system with real customers doing real shopping. Another approach would be to compute the ratio of accepted recommendations to total recommendations offered. However, looking at the recommendation list is voluntary on the part of the customer (it involves clicking the “light bulb” tab in figure 2), and we have no way of knowing whether the customer actually looked at the list of recommended products.⁹

A related issue is the extent to which recommendations with higher scores are accepted preferentially over recommendations with lower scores. We address this issue by comparing the distribution of scores computed from Eq. (7) for accepted recommendations with the analogous distribution for offered recommendations. The results are shown in figure 9. The scores for the accepted recommendations are based on 243 products accepted from 183 distinct recommendation lists. The distribution for the offered recommendations is taken from approximately 20,000 recommendations made to the customers who accepted at least one recommendation during the pilot program.

Figure 9 shows that the scores of the accepted recommendations are higher than the scores of a large number of offered recommendations. For example, 80% of the products placed onto the recommendation lists have scores below 0.1, but only 25% of the accepted recommendations fall in this lower bin. The mean and median scores for the offered recommendations are 0.063 and 0.034, respectively, while the mean and median scores for the accepted recommendations are 0.16 and 0.13. The difference between the two means, 0.10, falls well within the 95% confidence interval (0.090, 0.105) computed using Student’s *t*-test for the difference between means (Robbins and Van Ryzin, 1975). These results suggest that the score computed using Eq. (7) is indeed a useful indicator of a previously unbought product’s appeal to the targeted customer.

6. **Summary and future work**

We have described a product recommendation system developed as part of an overall pervasive computing solution for grocery shopping. The recommendation algorithm combines aspects of content and collaborative filtering to rate new products for a customer based on their prior purchase behavior. Analysis of results obtained during a field test of the system have shown revenue boosts of roughly 1 to 2%, with the interesting observation that people

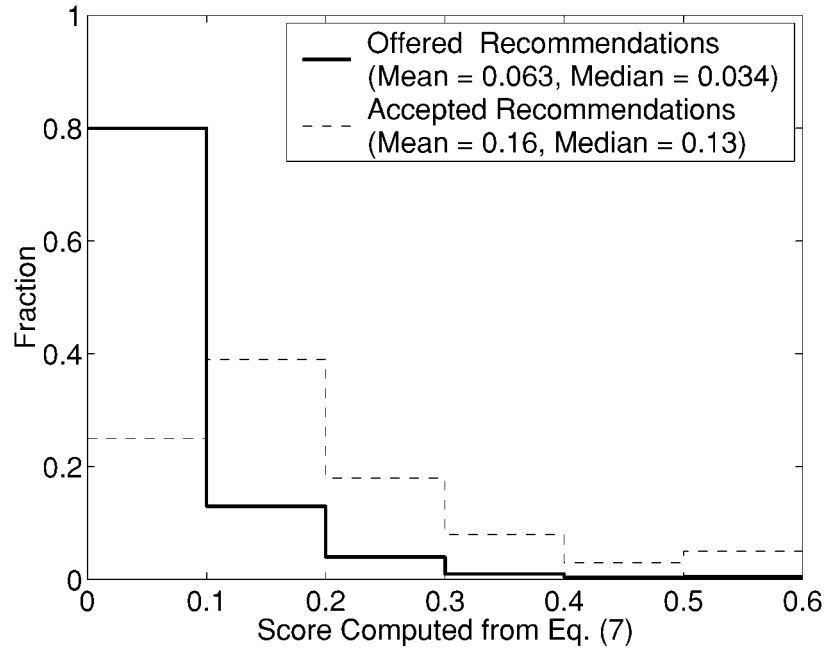


Figure 9. Distributions of scores for offered and accepted recommendations.

often choose recommendations from product classes in which they have not purchased previously.

The overall SmartPad project is an example of enabling access to server-based data and computation from mobile devices such as PDAs. In our case, the recommendation analysis runs completely on the server, with the PDA providing a mechanism to access this analysis remotely. Given the increasing computational power of these devices, it will be possible to move more of the analysis to the device itself, including compute-intensive portions which heretofore have been possible only on conventional servers. We are exploring this issue in the context of a PDA-based application designed to recommend wine selections based on specific meal choices (Almasi and Lee, 2000). It is likely that many other applications, such as financial analysis, will exploit this opportunity to do increasingly complex analysis on hand-held devices operating in a weakly connected environment.

Acknowledgments

We are especially grateful to our colleague Harry Stavropoulos and to Jeremy Wyman of Safeway UK for all their help.

Notes

1. Also referred to in literature as *social filtering* (Sharadanand and Maes, 1995).

2. Given the rule $A \Rightarrow B$, where A and B are itemsets and T is the total number of customers, the *support* for the rule, $S(A \Rightarrow B)$, is the percentage of customers who have spent in both A and B ; the *confidence* in the rule is $S(A \Rightarrow B)/S(A)$, and the *lift* is $S(A \Rightarrow B)/(S(A)*S(B))$. The lift is the ratio of actual confidence to the expected confidence, where the latter is computed assuming that A and B are statistically independent. Support and lift are symmetric in A and B ; confidence is not.
3. The word “demographic” here refers to a particular clustering *algorithm* that was applied to customers’ spending data, whereas its use above in “demographic-based segmentation” refers to the demographic *data* (age etc.) that were used as the basis for creating those segments.
4. The Condorcet criterion (Michaud, 1999) is the difference of two factors, one of which measures how similar the records within a cluster are, while the other measures how different the records within a cluster are from all records not in the cluster. For perfect clustering, the first factor is 1.0 and the second is 0.0.
5. Dunn’s index is the ratio of intercluster distance to cluster diameter. We computed it using the approximation called v_{53} in Bezdek (1998).
6. The clustering is based on customer spending in all 99 product classes, even though a few classes such as petrol and tobacco are excluded from our lists of recommendable products. We found, for example, that tobacco purchasers also buy more and stronger deodorizers and a different set of wines than non-tobacco purchasers.
7. The 16 SmartPad product categories are a superset of the 99 Safeway product classes, minus those classes like petrol and tobacco that are not available via SmartPad.
8. It is also possible that SmartPad participants in the new store were inherently more inclined to accept recommendations than those from the original store. We note that the demographics of the two stores are similar.
9. Late in the pilot program, a capability was added to the PDA which tagged whether a customer submitting an order had actually looked at the recommendation list.

References

- Agrawal, R. and Srikant, R. 1994. Fast algorithms for mining association rules. In Proc. of the 20th Int’l Conference on Very Large Databases, Santiago, Chile, Sept. 1994.
- Aggarwal, C.C., Wolf, J.L., Wu, K.-L., and Yu, P.S. 1999. Horting hatches an egg: A new graph-theoretic approach to collaborative filtering. In KDD-99 Proceedings, pp. 201–212.
- Almasi, G.S. and Lee, A.J. 2000. A PDA-based personalized recommender agent. In Proc. Fifth International Conf. on the Practical Application of Intelligent Agent and Multi Agent Technology, Manchester, England, pp. 299–309, April 2000.
- Balabanovic, M. and Shoham, Y. 1997. Fab: Content-based, collaborative recommendation. Communications of the ACM, 40(3):66–72.
- Bezdek, J.C. and Pal, N.R. 1998. Some new indexes of cluster validity. IEEE Transactions on Systems, Man, and Cybernetics—Part B: Cybernetics, 28(3):301–315.
- Borchers, A., Herlocker, J., Konstan, J., and Riedl, J. 1998. Ganging up on information overload. Computer, 31(4):106–108.
- Everitt, B.S. 1993. Cluster Analysis, London: Edward Arnold.
- First Annual Personalization Summit. 1999. www.personalization.com, San Francisco, Nov. 15–16, 1999.
- Intelligent Miner for Data. www.ibm.com/software/data/iminer/fordata.
- Kohonen, T. 1995. Self-Organizing Maps. Springer-Verlag.
- Konstan, J., Miller, B., Maltz, D., Herlocker, J., Gordon, L.R., and Riedl, J. 1997. GroupLens: Applying collaborative filtering to usenet news. Communications of ACM, 40(3).
- Kotlyar V., Viveros, M.S., Duri, S.S., Lawrence, R.D., and Almasi, G.S. 1999. A case study in information delivery to mass retail markets. In Proceedings of the 10th International Conference on Database and Expert Systems Applications (DEXA), Florence, Italy, Aug./Sept. 1999, Springer-Verlag. Lecture Notes in Computer Science, vol. 1677.
- Lawrence, R.D., Almasi, G.S., and Rushmeier, H.E. 1999. A scalable parallel algorithm for self-organizing maps with applications to sparse data mining problems. Data Mining and Knowledge Discovery, 3:171–195.
- Michaud, P. 1999. Clustering techniques. Future Generation Computer Systems, 13(2).

- Resnick, P. and Varian, H.R. 1997. Recommender systems. *Communications of the ACM*, 40(3):56–58. Also see other articles in this special issue.
- Robbins, H. and Van Ryzin, J. 1975. *Introduction to Statistics*. Scientific Research Associates, Inc..
- Salton, J. 1989. *Automatic Text Processing: The Transformation, Analysis and Retrieval of Information By Computer*. Reading, MA: Addison-Wesley.
- Salton, J. and McGill, M.J. 1983. *Introduction to Modern Information Retrieval*. New York: McGraw-Hill.
- Shardanand, U. and Maes, P. 1995. Social information filtering: Algorithms for automating word of mouth. In *Proc. CHI 95*, ACM Press, pp. 202–209.
- Ungar, L.H. and Foster, D.P. 1998. Clustering methods for collaborative filtering. In *Proceedings of 1998 AAAI Workshop on Recommender Systems*. Available as AAAI Technical Report WS-98-08.

Richard D. Lawrence is a Research Staff Member and Manager, Deep Computing Applications, at the IBM T.J. Watson Research Center. He received the B.S. degree from Stanford University in Chemical Engineering, and the Ph.D. degree from University of Illinois in Nuclear Engineering. Prior to joining IBM Research in 1990, he held research positions in the Applied Physics Division at Argonne National Laboratory and at Schlumberger-Doll Research. His current work is in the development of high-performance data mining applications in the areas of financial analysis and product recommendation systems.

George S. Almasi is a Research Staff Member at the IBM T.J. Watson Research Center. He received his Ph.D. in electrical engineering from the Massachusetts Institute of Technology. He has a variety of technical and management experiences in memory technology, display systems, and parallel computing, and is co-author of the book “Highly Parallel Computing”. His recent interests have been in high-performance parallel applications, including datamining, the visualization of its results, and its application in personal recommender systems. He is currently at work on a visualizer for Blue Gene, a supercomputer designed to perform protein folding.

Vladimir Kotlyar is currently with CrossGain Corporation of Redmond, WA. Prior to joining CrossGain, Dr. Kotlyar was a Research Staff Member at the IBM T.J. Watson Research Center. He received his Ph.D. in Computer Science from Cornell University in 1998. His main interest is performance optimization of data intensive applications.

Marisa Viveros is a senior manager of the Pervasive Computing Solutions group at the IBM T.J. Watson Research Center. She is responsible for the creation of emerging applications in the areas of wireless technology, pervasive devices, and their seamless integration in business environments. Her research areas include data management for mobile computing, business applications, data mining, and parallel databases.

Sastry S. Duri is an Advisory Software Engineer at the IBM T. J. Watson Research Center, where his research interests include e-business, mobile commerce, and pervasive computing applications. Duri received his B.Tech. in electronics and communications from Regional Engineering College, Warangal, India, the M.S. in computer science from Indian Institute of Technology, Madras, India, and his Ph.D. in electrical engineering and computer sciences from the University of Illinois at Chicago.