

# Recommendations on grocery shopping: customer or product similarities?

K.Y.M. de Kruiff

444677

Supervisor: F.J.L. van Maasakkers

Second assessor: S.I. Birbil

Bachelor Thesis: Econometrics and Operational Research  
Erasmus School of Economics, ERASMUS UNIVERSITY ROTTERDAM

5th July 2019

## **Abstract**

Online grocery shopping becomes more popular every day and benefits both customer and grocery store. To increase revenue, grocery stores recommend products to customers to add to their online baskets. Such recommendations can be either item- or consumer-based. In this research I investigate recommender systems based on the Collaborative Filtering algorithm. First, I replicate part of the paper of Li, Dias, Jarman, El-Deredy and Lisboa (2009), which investigates different standard item-based Collaborative Filtering models and introduces a new personalised recommender system, which outperforms the standard methods. Second, I investigate a user-based Collaborative Filtering model and compare these results to results of the item-based models. The results of the replication do not entirely agree with the results found by Li et al. (2009), which is mainly due to the higher number of items used in the replication. Despite the challenges of the user-based Collaborative Filtering model, the performance of the user-based model is similar to the performance of some of the item-based models, but does not outperform the item-based models. To generate all the results I have used a real-world dataset from the grocery store Ta-Feng.

The views stated in this thesis are those of the author and not necessarily those of Erasmus School of Economics or Erasmus University Rotterdam.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Literature review</b>	<b>4</b>
<b>3</b>	<b>Data</b>	<b>6</b>
<b>4</b>	<b>Methodology</b>	<b>7</b>
4.1	Random Walk model . . . . .	7
4.2	Basket-Sensitive Random Walk on bipartite network . . . . .	8
4.2.1	First-Order Transition Probability Matrix . . . . .	8
4.2.2	Basket-Sensitive Random Walk . . . . .	9
4.3	Performance metrics and evaluation protocols . . . . .	10
4.4	User-based Collaborative Algorithm . . . . .	12
<b>5</b>	<b>Results</b>	<b>13</b>
5.1	Results replication . . . . .	14
5.2	Results user-based CF algorithm . . . . .	17
<b>6</b>	<b>Conclusion</b>	<b>18</b>
<b>7</b>	<b>Discussion</b>	<b>19</b>
<b>A</b>	<b>Matlab code</b>	<b>24</b>

# 1 Introduction

Grocery shopping becomes easier every day. Nowadays a lot of brick-and-mortar grocery stores, i.e. Albert Heijn, Jumbo and Coop, have an online shopping website where you can order your groceries and have them delivered on your doorstep. While ordering groceries online, a grocery shop recommends products which are not yet in your shopping basket, by highlighting other products: ‘You might also like...’ or ‘Did you forget...’. For example, if you have cake mix in your shopping basket but no butter or eggs, you probably get a recommendation for butter and eggs as these products are often sold with the cake mix. These recommendations are generated by recommender systems.

A good recommender system is beneficial for both customer and grocery store. A customer benefits from accurate recommendations, thereby diminishing the chance that crucial groceries do not get ordered. At the same time, accurate recommendations lead to more revenue for the grocery store as they will sell more products<sup>1</sup>. Therefore, a grocery store wants to optimise the recommender system used by its online webshop.

The main part of a recommender system is the personalisation algorithm that models consumer<sup>2</sup> shopping behaviour. This algorithm is used to search for items that are likely of interest to the individual customer. In general, the algorithm generates a rank-ordered list of items, which are not already present in the individual’s basket. It is expected that a customer is more likely to accept recommendations on top of the list than items on the bottom of that list.

Several frameworks are used for item recommendations. For example, the Collaborative Filtering (CF) framework and the Random Walk model. These algorithms are used in a many papers (e.g. Breese, Heckerman & Kadie, 1998; Brin & Page, 1998; Sarwar, Karypis, Konstan, Riedl et al., 2001) that focus on recommendations for leisure products (e.g. movies, books and music). In contrast to products in the grocery shopping domain, leisure products are mostly purchased just once. In the domain of grocery shopping, customers buy products repeatedly. In addition, movies, books and music are often rated by their users, which is again not applicable to grocery shopping products. Many existing recommendation techniques are based on these user ratings. Hence, these techniques cannot simply be applied to grocery shopping data. Therefore, this research focuses on

---

<sup>1</sup>In this research, the terms products and items are used interchangeably when talking about recommender systems.

<sup>2</sup>In this research, the terms consumer and customer are used interchangeably.

the application of recommender systems for online grocery shopping. I investigate an algorithm introduced by Li et al. (2009) that appears to suit the characteristics of grocery shopping data better than the currently applied algorithms developed for leisure products. My research question is: How does the item-based algorithm introduced by Li et al. (2009) perform on grocery shopping data compared to a user-based CF algorithm? To answer my research question I first replicate part of the research of Li et al. (2009). Second, I investigate the user-based algorithm as described in Breese, Heckerman and Kadie (1998) for the same dataset as used in the replication and compare these results to the results of the replication.

The remainder of this research paper is organised as follows. Section 2 gives an overview of the empirical literature. A description of the retrospective data is given in Section 3. Section 4 describes the recommender systems used for the replication and the extension and the performance measures used to evaluate the recommender systems. The results of the replication and the extension are presented in Section 5. Finally, Section 6 provides the overall conclusion based on the results and Section 7 lists suggestions for improvement to this research and possibilities for future research.

## 2 Literature review

Previous research has demonstrated that CF is an effective framework to generate item recommendations (Breese et al., 1998; Konstan et al., 1997; Liu & Yang, 2008; Sarwar, Karypis, Konstan, Riedl et al., 2001; Yildirim & Krishnamoorthy, 2008). Researchers invented a number of CF algorithms that can be split into two categories: user-based and item-based algorithms (Breese et al., 1998).

User-based algorithms use the entire user-item database to generate a recommendation. These algorithms find a set of users, also called neighbours, that have the same preference (i.e. they either tend to buy a similar set of items or rate different products the same way) as the target user. The preferences of the neighbours are then used to generate one or more recommendations for the target user. This technique, known as nearest-neighbour or user-based CF, is very popular and widely used in practice. However, sparsity of the data and scalability of the process are challenges in user-based CF algorithms (Sarwar et al., 2001).

Item-based CF is a very simple method and has a relatively good performance. Therefore, it is one of the most popular recommendation algorithms. The item-based CF is a similarity-based algorithm, which assumes that consumers are likely to accept an item recommendation that is

similar to an item bought previously. In this research, similarity is defined as how well a product supplements the products already present in the basket. Hence, the success of the item-based CF depends on the quality of the similarity measure. There are two types of similarity measures that are widely used in empirical literature (Li et al., 2009). First, the cosine-based measure, which is symmetric and determines the similarity between two products as the normalised inner product of two feature vectors of a user-item matrix. The feature vectors contain the number of times a user bought a certain product. To reduce the problem of overestimating the relevance of items that have an extremely high purchase frequency, the raw purchase frequency is converted to a log scale with  $\log(\text{frequency} + 1)$ . Afterwards the  $n \times m$  user-item matrix is normalised, where  $n$  is the number of users and  $m$  is the number of items. The cosine similarity is computed as:

$$\text{sim}(i, j) = \frac{M_{*,i} \cdot M_{*,j}}{|M_{*,i}| |M_{*,j}|}, \quad (1)$$

where  $M$  reflects the  $n \times m$  normalised user-item matrix and  $M_{*,i}$  is the  $i$ th column of  $M$ .

Another similarity measure is the conditional probability-based similarity measure, which reflects the ratio between the number of users who have bought both items and the number of users who bought just one of the two items. This measure is asymmetric. A challenge of an asymmetric measure is that items often have high conditional probabilities with respect to very popular products caused by customers who purchase many items (Deshpande & Karypis, 2004). Therefore, Deshpande and Karypis (2004) propose a similarity measure that puts more emphasis on customers who have bought fewer items as these customers are more reliable indicators when determining the similarity. The measure is defined as:

$$\text{sim}(i, j) = \frac{\sum_{\forall q: M_{q,i} > 0} M_{q,j}}{\text{Freq}(i) \times (\text{Freq}(j))^\alpha}, \quad (2)$$

where  $\alpha \in [0, 1]$  and is used to penalise very popular products,  $\text{Freq}(i)$  the number of users that have purchased item  $i$ , and  $M$  the same normalised  $n \times m$  user-item matrix as used in (1). As (2) uses the non-zero entries of the normalised user-item matrix  $M$ , customers contribute less to the total similarity if they purchased more items. Hence, more emphasis is put on the purchasing decisions of the customers who bought fewer items (Deshpande & Karypis, 2004).

After determining the similarity matrices, the recommendations are determined based on the products already in the customer’s basket. The rows of the similarity matrix belonging to the products already in the customer’s basket are added together, which results in a single row matrix containing the total similarity of each possible item. Then, the items with the highest values in the resulting row matrix form a list of recommendations for the user.

Previous research has shown that item-based CF is able to achieve recommendation accuracies similar to, or better than, the user-based approach (Sarwar et al., 2001). However, a challenge of the item-based CF is that it is not able to recommend an item that has never been co-purchased with the other items before. A solution to this deficiency is to use a Random Walk based recommendation algorithm, which randomly jumps from product to product.

Besides the different CF algorithms, recommender systems can be based on a Random Walk which will be described in Section 4. Further, there are several machine learning techniques that can also be applied to generate recommendations for customers, i.e. the nearest-neighbour technique and the Random-Forest framework. The machine learning techniques are beyond the scope of this research.

### 3 Data

In this research I use a basket dataset from the grocery store Ta-Feng.<sup>3</sup> Ta-Feng is a Chinese membership retailer warehouse that sells a wide variety of goods (Hsu, Chung & Huang, 2004). The dataset consists of shopping records of different customers from November 2000 until February 2001 and contains, among other data, the shopping date, customer ID, type of product and the amount of the product bought. If several shopping records have the same shopping date and the same customer ID, they are seen as one transaction. To determine the number of transactions, I order the data on shopping date and afterwards on customer ID. For evaluation purposes I exclude the records from customers that only ordered products from Ta-Feng once in the four month period. The resulting dataset consists of 20,388 customers, 107,700 transactions and 2,000 product sub-classes. The sparsity, which represents the ratio of empty entries in the co-occurrence matrix, is 0.625. The other descriptive statistics of the resulting dataset are given in Table 1.

After deleting all customers with only one transaction during the four month period, I split the data into two parts: a training and a test set. The test set consists of the last transaction of each customer and the training set consists of all other transactions aggregated per customer into one basket to maximise the information content. Some of the product sub-classes are never purchased in the training set and are therefore excluded, which results in 1,973 remaining sub-classes. From the test set, I delete the transactions containing less than four items for evaluation purposes. This

---

<sup>3</sup><https://www.kaggle.com/chiranjivdas09/ta-feng-grocery-dataset>

results in a training set of 87,312 transactions and a test set of 11,813 transactions.

Table 1: Descriptive statistics adjusted Ta-Feng dataset

	Min.	Median	Mean	Max.
# products per basket	1	6	9.24	1,200
# products per subclass	1	57	497.81	21,744
# baskets per customer	2	4	5.28	86

## 4 Methodology

In this section, I describe the methods that I apply to answer my research question stated in the introduction. Section 4.1 describes the basic Random Walk model and Section 4.2 describes the personalised Basket-Sensitive Random Walk introduced by Li et al. (2009). The performance measures used for evaluation of the recommender systems are explained in Section 4.3 and Section 4.4 describes the user-based CF model.

### 4.1 Random Walk model

Many researchers have used a Random Walk model to generate recommendations for movies (Liu & Yang, 2008; Yildirim & Krishnamoorthy, 2008; Wijaya & Bressan, 2008; Huang, Zeng & Chen, 2004). The Random Walk model is very similar to Google’s PageRank algorithm, which is a link analysis algorithm that assigns an importance rank to each page. The importance rank of a page reflects the probability of visiting that page in a Random Walk (Brin & Page, 1998). Yildirim and Krishnamoorthy (2008) use this idea to develop an item graph where the nodes are items and the edges reflect the similarities between the items. The algorithm introduced by Yildirim and Krishnamoorthy (2008) generates a user-item matrix,  $R_{basket}$ , which contains the ratings of the items for each user and is computed as follows:

$$R_{basket} = MdP(I - dP)^{-1}, \quad (3)$$

where  $M$  is the same user-item matrix as in (1) and (2),  $P$  is the transition probability matrix and  $d \in (0, 1)$  is the damping factor used to define the stopping probability of the Random Walk. There are two possibilities for the transition matrix  $P$ , namely the normalised cosine and the normalised conditional probability-based similarity matrix. The normalisation of the matrices is needed to

derive a transition probability matrix from the similarity matrices. The products with the highest rating in  $R_{basket}$  that are not yet contained in a customer’s basket, are recommended to the customer. Thus, the Random Walk model is personalised but does not compute the product ratings based on the products already in the basket.

## 4.2 Basket-Sensitive Random Walk on bipartite network

As described in Section 4.1, the Random Walk model does not take the current basket into account when recommending new products. Therefore, Li et al. (2009) introduce the Basket-Sensitive Random Walk algorithm. This approach computes the transaction probabilities between products (not users) from a bipartite network as described in Section 4.2.1. Further, the model computes the ratings of the products while taking into account the products already in the basket. In this way, relevant products are given a higher rating on the list of recommendations, with less bias to the most popular products.

### 4.2.1 First-Order Transition Probability Matrix

The shopping basket data can be seen as a bipartite network with two types of nodes: consumers and products. Each edge, being between a consumer and a product, reflects the consumer’s purchase frequency of a product as shown in Figure 1.

The set of consumers is denoted by  $C = c_1, c_2, \dots, c_{|C|}$ , the set of products by  $P = p_1, p_2, \dots, p_{|P|}$  and the set of purchase frequencies by  $F = f(1, 1), f(1, 2), \dots, f(|C|, |P|)$ , which results in the bipartite network  $BN = \{C, P, F\}$ . An example of a bipartite network is given in Figure 1.

In a Random Walk, the first-order transition probability  $P(p_j|p_i)$  is the probability that a random surfer jumps from the item node  $i$  to  $j$  via all connected consumer nodes  $c_k$ . As the Random Walk only jumps via consumers who bought *both* products, the transition matrix consists of merely *first-order* transition probabilities which can be written as:

$$P(p_j|p_i) = \sum_{k=1}^{|C|} P(p_j|c_k)P(c_k|p_i), \quad (4)$$

where  $P(p_j|c_k)$  is the probability that a random surfer jumps from consumer node  $c_k$  to product node  $p_j$ ,  $P(c_k|p_i)$  the probability that the random surfer jumps to consumer node  $c_k$  from the product node  $p_i$  and  $P(p_j|p_i)$  denotes the marginal probability distribution over all consumers which reflects



the similarity between products  $i$  and  $j$ . The intuition behind probability  $P(p_j|p_i)$  is that it can be seen as the preference for product  $p_j$  from all customers  $C$  who have already bought product  $p_i$ . Every preference  $P(p_j|c_k)$  from the  $k$ th customer is weighed proportionally to the number of times the customer bought product  $p_i$ , namely  $P(c_k|p_i)$ . The conditional probabilities used in (4) are calculated as:

$$P(p_j|c_k) = \frac{f(c_k, p_j)}{(\sum f(c_k, \cdot))^{\alpha_1}}, \quad (5)$$

$$P(c_k|p_i) = \frac{f(c_k, p_i)}{(\sum f(\cdot, p_i))^{\alpha_2}}. \quad (6)$$

where  $\alpha_1, \alpha_2 \in [0, 1]$  correct for products that have high purchases frequencies. For convenience reasons this research applies  $\alpha_1 = \alpha_2$ . As in (2), the parameter  $\alpha$  corrects for the fact that people who buy a lot of different products are less informative than people who bought just a few products. Also, more popular products are less informative for the personal preferences of the shoppers than unpopular products. However, the first-order similarity measure does not capture the possible similarity between products which are never bought together, like newly launched products (Li et al., 2009). Hence, when data is sparse, the first-order similarity measure suffers and results in less effective recommendation.

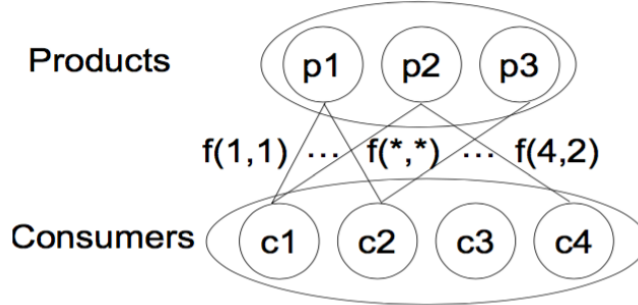


Figure 1: Bipartite network for shopping basket data where  $f(i, j)$  is the amount of product  $j$  bought by customer  $i$  (Li et al., 2009)

#### 4.2.2 Basket-Sensitive Random Walk

As described in Section 4.2.1, the transition probability matrix has a disadvantage when data is sparse. This limitation can be solved by combining the similarities of all orders as:

$$P^* = \sum_{t=1}^{\infty} \frac{(dP)^t}{|(dP)^t|}, \quad (7)$$

where  $P$  is either the cosine or conditional probability matrix described in ?? or the first-order transition probability matrix explained in 4.2.1. Again, the parameter  $d \in (0, 1)$  is the damping

factor which determines the length of the Random Walk. When  $t$ , the number of purchases of a product, goes to infinity,  $P^*$  converges to  $dP(I - dP)^{-1}$  (Yildirim & Krishnamoorthy, 2008).

As the Random Walk  $P^*$  is not based on the current shopping behaviour, the model has to be biased towards the items already in the shopping basket in order to get more accurate results. To do so, an additional rating can be assigned to items currently in the shopper’s basket during each iteration of the Random Walk computation. Therefore, Li et al. (2009) introduce a non-uniform personalisation vector  $U_{basket}$  to compute the importance score during each iteration as follows:

$$R_{basket} = d \cdot P \cdot R_{basket} + (1 - d) \cdot U_{basket}, \quad (8)$$

where  $R_{basket}$  contains the ratings for all products derived from the basket-based importance scores and  $U_{basket}^i = \frac{1}{m}$  if the  $i^{th}$  product is in the basket and zero otherwise, where  $m$  is the number of products in the current basket. The matrix  $P$  can be the normalised cosine or conditional probability-based similarity matrix or the first-order transition probability matrix as described in Section 4.2.1. The first term of (8),  $d \cdot P \cdot R_{basket}$ , can be interpreted as the probability that the Random Walk stops and that the customer is done shopping. The second term of the equation,  $(1 - d) \cdot U_{basket}$ , can be seen as the probability that the Random Walk proceeds, meaning that the customer is still adding items to the basket. However, applying (8) in practice would be too time consuming due to the re-computations of the ratings in each iteration of the Random Walk. Therefore, Li et al. (2009) also proposed an approximation,  $\hat{R}_{basket}$  defined as:

$$\hat{R}_{basket} = \sum_{p_i \in basket} R_{item_i}, \quad (9) \quad R_{item_i} = d \cdot P \cdot R_{item_i} + (1 - d) \cdot U_{item_i} \quad (10)$$

where  $R_{item}$  is the item-based importance score and  $U_{item_i}$  is the personalisation vector with the  $i^{th}$  entry set to one, and the rest set to zero. The computation of  $R_{item_i}$  can be simplified to

$$R_{item_i} = (1 - d)(I - dP)^{-1}U_{item_i}, \quad (11)$$

which also simplifies the computation of (9). The products with the highest rating in  $\hat{R}_{basket}$ , which are not already present in the basket, are recommended to the consumer.

### 4.3 Performance metrics and evaluation protocols

There are several ways to evaluate the accuracy of a recommender system. It is important that the performance evaluation results are representative of live, interactive behaviour. Sordo-Garcia, Dias,

Li, El-Deredy and Lisboa (2007) looked into three evaluation strategies which differed in how to split the retrospective basket into evidence and target components, where the evidence products are used to predict the target products. The three approaches were to split the data randomly, based on popularity and via the leave- $n$ -out approach (Breese et al., 1998). According to the results the popularity-based approach was the only one that ranked the recommender systems consistently with their live performance (Sordo-Garcia et al., 2007).

To evaluate the accuracy of the recommender systems, this research uses the popularity-based binary hit rate,  $\mathbf{bHR}(\mathbf{pop})$ , and the weighed hit rate,  $\mathbf{wHR}(\mathbf{100})$ , introduced by Li et al. (2009). The popularity based binary hit rate is based on a leave-three-out principle; the least three popular products of the test basket, based on the training set, are the targets, the rest of the test basket is used as evidence. When testing the recommender systems, the evidence items are used to predict the target items. The binary hit rate is computed as the proportion of test baskets having at least one out of three correctly predicted target items (Li et al., 2009). Besides leaving out the three least popular items, I also evaluate the performance of the recommender system when using a random leave-three-out approach ( $\mathbf{bHR}(\mathbf{rnd})$ ), where three target items are randomly selected from the test basket. However,  $\mathbf{bHR}(\mathbf{rnd})$  favors items which occur with a high frequency and therefore tends to overestimate the model performance. Moreover, previous research has shown that when using a leave-one-out cross validation the hit rates also over-emphasize the performance of popular products (Sordo-Garcia et al., 2007; Li, Dias, El-Deredy & Lisboa, 2007).

To solve the issue of over-estimation for a leave-one-out cross validation, Li et al. (2009) propose a new performance measure, which weighs the hit of items inversely to their popularity. The weighed hit rate based on a leave-one-out approach,  $\mathbf{wHR}(\mathbf{100})$ , is computed as:

$$\mathbf{wHR}(\mathbf{100}) = \frac{\sum_i (1 - p(x_i)) \dot{HIT}(x_i)}{\sum_i (1 - p(x_i))} \quad \text{with} \quad HIT(x_i) = \begin{cases} 1, & \text{if } x_i \text{ is predicted correctly} \\ 0, & \text{otherwise} \end{cases}, \quad (12)$$

where  $x_i$  is the target item and  $p(x_i)$  is its prior probability based on the popularity of the product. The popularity and hence the prior probability  $p(x_i)$  is based on the converted user-item matrix and is computed as the ratio between the number of times an item was bought and the total number of items bought by all customers. When all items are predicted correctly, the basket will achieve a hit rate of one. The final hit rate is computed by averaging over all test baskets. Another measure based on a leave-one-out cross validation that biases the results towards the performance of small classes, is the macro-average hit rate,  $\mathbf{macroHR}(\mathbf{100})$ . The  $\mathbf{macroHR}(\mathbf{100})$  measure computes the hit

rate for each product in a basket, which results in a vector of zeros, reflecting no hit, and ones, reflecting a hit. The final hit rate is averaged over all products, which is the length of the vector, instead of all baskets as for the `wHR(100)`.

In total I evaluate three types of standard item-based CF models, the cosine (1), the conditional probability (2) and the bipartite network (4) based similarities, referred to as  $cf(cos)$ ,  $cf(cp)$  and  $cf(bn)$  respectively. I also evaluate the standard Random Walk item-based CF model ( $rw$ ) described in Section 4.1 and the Basket-Sensitive Random Walk model ( $bsrw$ ) described in Section 4.2. All these methods are item-based models and are compared to a baseline method,  $pop$ , which merely recommends the most popular items not contained in the basket.

#### 4.4 User-based Collaborative Algorithm

Besides the item-based CF algorithm, I investigate the user-based CF algorithm for the same Ta-Feng dataset, which is not yet done by Li et al. (2009). The idea of the user-based algorithm is to predict the ratings of the products for a particular user, also referred to as active user, based on the user-ratings of similar users from the training set, where ratings are based on the number of purchases of a product. The user database for the Ta-Feng dataset is not too big, namely 20,388 customers. Hence, the scalability problem, which was mentioned in Section 2, of the user-based CF algorithm might be avoided. Therefore, it is interesting to investigate whether for the size of this user database a user-based CF approach could also work, despite its ‘cold start’ problem and scalability issues (Konstan et al., 1997; Shardanand & Maes, 1995).

The user-based CF algorithm makes use of the cosine similarity matrix containing the similarities between users instead of items. Again, the normalised  $n \times m$  user-item matrix  $M$  is constructed as described in Section 2 and can be used to compute the cosine similarity matrix for the users as:

$$sim(i, j) = \frac{M_{i,*} \cdot M_{j,*}}{|M_{i,*}| |M_{j,*}|}, \quad (13)$$

where  $M_{i,*}$  is the  $i$ th row of the  $n \times m$  user-item matrix  $M$  representing the quantity of the items bought by customer  $i$  (Deshpande & Karypis, 2004).

Before predicting ratings for the active user for the test set, the  $k$  most similar users are determined from the cosine similarity matrix. Then, the predicted ratings of the active user for item  $j$ ,

$r_{a,j}$ , are computed as a weighed sum of the ratings of the  $k$  most similar users as:

$$r_{a,j} = \bar{v}_a + \sum_{i=1}^k w(a,i)(v_{i,j} - \bar{v}_i), \quad (14)$$

where  $\bar{v}_i$  is the mean of all the ratings of the products in the basket of customer  $i$  for the training set,  $v_{i,j}$  is the rating from user  $i$  for item  $j$  and  $w(a,i)$  is the normalised cosine similarity between user  $a$  and  $i$  (Breese et al., 1998). The mean of all ratings is computed as:

$$\bar{v}_i = \frac{1}{|I_i|} \sum_{j \in I_i} v_{i,j}, \quad (15)$$

where  $I_i$  is the set of items of user  $i$  for which there exists a rating in the training set. As for the item-based algorithm, the user-based algorithm is evaluated with the test set which is split into targets and evidence. With the targets and evidence products, the performance metrics **bHR(pop)**, **bHR(rnd)** and **wHR(100)** can be computed based on the exact same training and test set as used to generate the replication results. In this way an accurate comparison can be made to the results of the item-based CF models. However, this also makes it hard to compare the results of the user-based CF model to the results of Li et al. (2009) as this research makes use of a different composition of the Ta-Feng dataset.

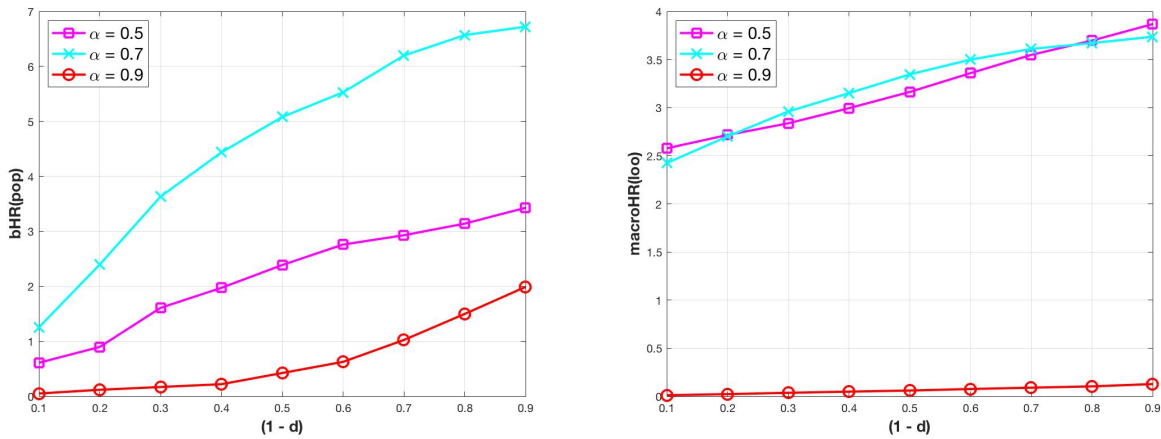
Besides evaluating the performance of the user-based CF algorithm compared to the item-based CF algorithm and the baseline method *pop*, I intend to find an optimal level for  $k$ . To determine the optimal level for  $k$  I investigate a trade-off between the values of the performance metrics and the computation time of the user-based CF algorithm. For the computation time I focus on the time needed to compute (14) as this will differ most between the different values for  $k$ . Suggestions for an optimal value for  $k$  are widely varying in empirical literature (Mobasher, Dai, Luo & Nakagawa, 2001; Sarwar et al., 2001; Al Mamunur Rashid, Karypis & Riedl, 2006). Therefore, I consider values going from three to 1000.

## 5 Results

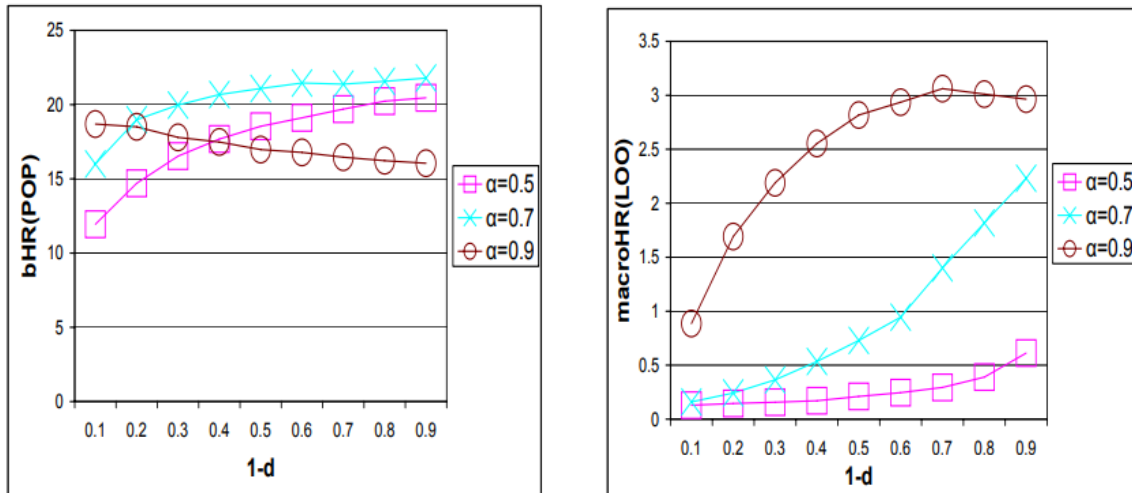
This section gives an overview of the results of the replication of Li et al. (2009) in paragraph 5.1 and paragraph 5.2 gives the results of the user-based CF algorithm and a comparison to the results of the replication.

## 5.1 Results replication

The results of the replication in this paper are similar to the results obtained by Li et al. (2009), but not exactly the same. Figure 2 shows the performance of the metrics  $\text{bHR}(\text{pop})$  and  $\text{macroHR}(\text{loo})$  with respect to the parameters  $\alpha$  and  $d$  for the results of the replication (Table 2a) and for the results of Li et al. (2009) (Table 2). Table 2 reports the best performance of the models with respect to the penalty factor  $\alpha$ , and the damping factor  $d$  for the measures  $\text{bHR}(\text{pop})$ ,  $\text{bHR}(\text{rnd})$  and  $\text{wHR}(\text{loo})$  described in Section 4.3.



(a) Replication results for  $cf(bn)+bsrw$



(b) Plot result Li, Dias, Jarman, El-Deredy and Lisboa (2009) for  $cf(bn)+bsrw$

Figure 2: Predictive performance for different settings of  $\alpha$  and  $d$  for the Basket-Sensitive Random Walk model

It is expected that the performance measures increase when the parameter values increase. Namely, if  $\alpha$  increases, the similarity between popular items decreases and if  $d$  decreases, so  $(1 - d)$  increases, then the ratings increase for those similar but less frequent items (Li et al., 2009). However, this does not hold for the results of the replication in this paper, which can be seen from Figure 2a. Moreover, this expectation does not hold for the results of Li et al. (2009) either. Figure 2a shows that the performance measures do increase when  $(1 - d)$  increases, but they decrease when  $\alpha$  increases. The difference with the result of Li et al. (2009) can be explained by the fact that when  $\alpha$  and  $(1 - d)$  increase, less frequent items appear more often on top of the list of recommendations. However, Li et al. (2009) use 1,093 items and the replication 1,973 items. Hence, when the parameters increase, the probability of correctly predicting the least frequent items is smaller for the replication than for Li et al. (2009), which causes the inconsistencies seen in Figure 2a. The inconsistencies of the results of Li et al. (2009) can be explained by the fact that the macro-averaged measure (`macroHR(100)`) is biased to performance in small classes compared to the basket-averaged measure (`bHR(pop)`). Hence, `macroHR(100)` increases when more infrequent items are predicted correctly and `bHR(pop)` decreases as popular items are probably missed (Li et al., 2009).

Based on Table 2, some remarks can be made on the results of the replication with respect to the results generated by Li et al. (2009). First of all, Table 2a shows that almost all methods outperform the baseline method, *pop*, which is in line with the results of Li et al. (2009) presented in Table 2b. Second, it is remarkable to see that for almost every method the values of the binary hit rate based on popularity of the replication are significantly smaller than the results generated by Li et al. (2009) except for the Random Walk method. This difference can again be explained by the fact that this research considers 1,973 items for the customers to choose from instead of 1,093 as in the study of Li et al. (2009). With more items, the probability of predicting the least popular items correctly becomes smaller and therefore the binary hit rate based on popularity decreases.

Third, from Li et al. (2009) it was not entirely clear whether they performed a personalised or a non-personalised Random Walk algorithm. The Random Walk model applied in this research is personalised as it includes a multiplication with the user-item matrix as described in Section 4.1. The personalisation of the Random Walk model could explain the relatively high results for the performance metrics for the Random Walk model. Hence, Li et al. (2009) probably applied a non-personalised Random Walk model which explains the difference in results. However, it should be taken into consideration that the number of items used differs and could also be of influence.

Further, as mentioned in Section 4.3,  $\text{bHR}(\text{rnd})$  favors very popular items and therefore tends to overestimate the model performance. This is in line with the results of the replication and the results of Li et al. (2009) as the values for the  $\text{bHR}(\text{rnd})$  metric are higher than for the  $\text{bHR}(\text{pop})$  metric.

Finally, the results for the Basket-Sensitive Random Walk of the replication are lower than the results generated by Li et al. (2009) for the  $\text{bHR}(\text{pop})$  measure, which is already explained above. The rest of the evaluation measures generate more or less the same results as Li et al. (2009) do. However, the results for the Basket-Sensitive Random Walk model based on the bipartite network transition probability matrix for the replication are not as outstanding as the results of Li et al. (2009). This could again be due to the number of items used in the replication or a misinterpretation of the method as described in Li et al. (2009). When only comparing the results of the Basket-Sensitive Random Walk model to the results of the standard CF model, the  $cf(\text{bn}) + \text{bsrw}$  method is always among the top two performers. Also, the bipartite network-based Basket-Sensitive Random Walk model has the highest hit rates compared with the cosine and conditional probability-based Basket-Sensitive Random Walk model, which is in line with the results of Li et al. (2009). This can be explained by the fact that the transition probability of the bipartite network is directly structured from the graph and hence is more effective than the transition probability indirectly derived by normalising the similarity matrices (Li et al., 2009).

Table 2: Predictive performance comparison of different models

(a) Results replication (%)				(b) Results Li et al. (2009) (%)			
Method	L-3-O		L-1-O CV	Method	L-3-O		L-1-O CV
	bHR(pop)	bHR(rnd)	wHR(100)		bHR(pop)	bHR(rnd)	wHR(100)
cf(cos)	6.2728	20.7483	<b>4.3438</b>	cf(cos)	18.1945	22.0730	3.9281
cf(cp)	4.7067	21.0023	4.1406	cf(cp)	14.6064	19.9129	3.3541
cf(bn)	5.4770	<b>24.1260</b>	2.4651	cf(bn)	21.7341	25.2556	4.6646
cf(cos) + bsrw	6.7891	18.6405	3.2600	cf(cos) + bsrw	18.2429	22.6599	3.9652
cf(cp) + bsrw	4.5374	20.7822	4.1205	cf(cp) + bsrw	16.7605	20.6571	3.4367
cf(bn) + bsrw	6.7214	24.0159	4.3078	cf(bn) + bsrw	<b>21.7886</b>	<b>26.1814</b>	<b>4.8151</b>
cf(cos) + rw	11.1149	23.8551	3.9529	cf(cos) + rw	7.7994	5.4880	0.7289
cf(cp) + rw	<b>12.6132</b>	19.9611	2.9326	cf(cp) + rw	14.0679	19.1807	3.1493
pop	0.5418	18.3357	2.4944	pop	7.9900	16.5700	2.2800

Note: For the results of the replication the metric based on random selection of targets is evaluated just once as it did not differ much per run and spared computation time.



## 5.2 Results user-based CF algorithm

Table 3 shows the results for the user-based CF model for different values of  $k$ , which reflects the number of users used to determine the ratings of the active user. Further, Table 3 shows that the user-based CF outperforms the baseline method, *pop*, denoted in Table 2a for all values of  $k$ . The user-based CF also outperforms the standard item-based CF model based on the leave-three-out metrics for almost every  $k$ . However, the best methods for the  $\text{bHR}(\text{pop})$  metric of the Random Walk model and Basket-Sensitive Random Walk model shown in Table 2a outperform the user-based CF model, but for the  $\text{bHR}(\text{rnd})$  metric the user-based CF model is better for  $k \geq 20$ . The above results are based on the results for different values of  $k$  for the test set. However, for a more accurate comparison between the results of the item- and user-based CF model, a validation set should be used to determine the optimal level for  $k$  first and afterwards the results of the performance metrics should be determined for the test set for the optimal value of  $k$ . These latter results should then be used for comparison with the results of the item-based CF model.

Besides comparing the user-based CF model to the item-based CF models, it is interesting to see what the optimal value for  $k$  is. The evaluation of the performance metrics with respect to  $k$  is shown in Figure 3. The exact values of the performance metrics displayed in Figure 3 correspond to the results denoted in Table 3. Table 3 and Figure 3 show that the macro-averaged metric does not differ much amongst different values of  $k$ . This could be explained by the fact that  $\text{macroHR}(\text{loo})$  is based on a leave-one-out cross validation and hence that the recommended item with the highest rating does not differ much when taking more customers into consideration. The basket-averaged binary hit rate based on popularity decreases once  $k$  increases. This can be due to the number of items taken into consideration determining the ratings. When  $k$  increases, more and more users are used who are not very similar to the active user. This leads to less accurate rating prediction as a lot of irrelevant information is used and more products are probably considered, which causes the  $\text{bHR}(\text{pop})$  to decrease. Furthermore, the random basket-averaged binary hit rate and the elapsed time both increase when  $k$  increases. As seen before, the  $\text{bHR}(\text{rnd})$  tends to overestimate the model performance as it favors very popular products. When using more information as  $k$  increases, the very popular products will probably appear more often, which causes overestimation of the metric.

Depending on what metric is important for a research, one could determine the optimal value for  $k$ . However, I believe that  $k = 10$  can be seen as the most appropriate value for  $k$  for this dataset as

Table 3: Predictive performance user-based CF algorithm for different sizes of neighbourhood

Method	L-3-O (%)		L-1-O CV (%)	Elapsed time (sec)
	bHR(pop)	bHR(rnd)	macroHR(lo)	
k = 3	<b>10.5223</b>	22.6191	3.5400	<b>32.4475</b>
k = 5	10.1840	23.3472	3.6800	38.8772
k = 10	9.7858	23.9228	3.7500	47.8670
k = 20	9.5149	24.3799	3.8700	53.1561
k = 50	9.1425	25.1926	3.9800	87.7133
k = 100	8.6346	24.5916	<b>4.0200</b>	146.1390
k = 500	7.4410	25.6328	3.9700	728.6756
k = 1000	6.8315	<b>26.4708</b>	3.9300	1083.4911

Note: The metric based on random selection of targets is evaluated just once as it did not differ much per run and spared computation time.

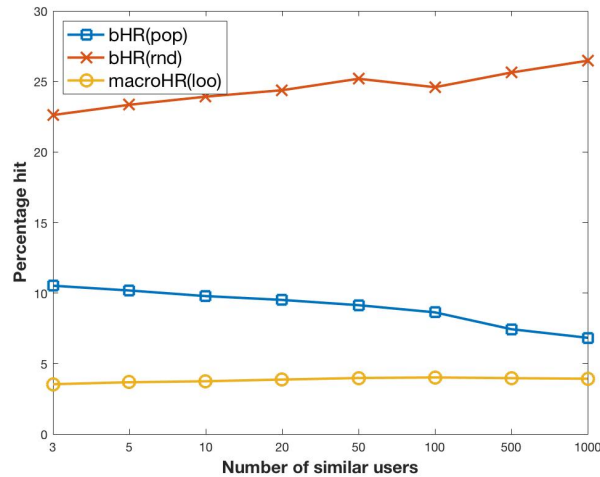


Figure 3: Evaluation of different metrics for different sizes of neighbourhood

the computation time for the ratings is less than one minute, the metric **bHR(pop)** is still relatively high and the measure **macroHR(lo)** does not increase much when increasing the value of  $k$ . As the **bHR(rnd)** is probably overestimated even more for a higher value of  $k$ , a lower value of  $k$  is also preferable for **bHR(rnd)**.

## 6 Conclusion

Grocery stores implement a variety of recommender systems to recommend products that are not already present in a customer’s basket. As it benefits both customer and grocery store to have the

best possible working recommender system, this research investigates a basket-sensitive recommender system introduced by Li et al. (2009). I replicate the paper from Li et al. (2009) for the Ta-Feng dataset and apply a user-based CF model to the same dataset. The results of the user-based CF model are compared to the results of the item-based models. Furthermore, I assess into the optimal number of users used to determine the ratings for the active user.

The results of the replication of Li et al. (2009) in this paper approach the results of Li et al. (2009) to a large extent. However, with a different set of items used in this paper for the replication, the replication does not yield the exact same conclusion as Li et al. (2009). Li et al. (2009) conclude that the Basket-Sensitive Random Walk model performs better than all other models they investigated. However, this paper finds that the Basket-Sensitive Random Walk model can be outperformed by some of the other models. The results of the replication show that every performance metric has a different optimal method, but that overall  $cf(bn) + bsrw$  performs well on all measures.

The user-based CF model outperforms the standard item-based CF model, but in turn outperformed by some of the Random Walk and Basket-Sensitive Random Walk models. This can be explained by the fact that the user-based CF model does not take into account the items already in the customer’s basket, whereas the Random Walk and Basket-Sensitive Random Walk model do. Overall, the Basket-Sensitive Random Walk is the most appropriate recommender system in the domain of grocery shopping as the difference with the best performing models per metric is very small and the model is among the top two performers for every metric, which is not the case for the other models.

Finally, the optimal number of neighbours used to determine the ratings of the active user is determined by examining the three different hit rates,  $bHR(pop)$ ,  $bHR(rnd)$  and  $macroHR(100)$ , and by elapsed time for computing the ratings of the active user. For this research the optimal number of neighbours used is ten, based on a trade-off between the decrease of  $bHR(pop)$  and an increase in  $bHR(rnd)$  and the computation time.

## 7 Discussion

As some of the results of the replication do not agree with the results generated by Li et al. (2009) there is room for improvement to this research. One of the big differences between the replication and Li et al. (2009) is the number of items used as it was not clear from Li et al. (2009) how

they determined the number of items used in their research. Therefore, it would be interesting to determine whether the results are more alike when taking the 1,093 most popular items and base the replication on only those 1,093 items. From Li et al. (2009) it is not clear if they took the 1,093 most popular items. However, with further investigation of the results, this might become clear. Further, the results for the popularity based binary hit rate based on 1,093 items will show whether the low hit rates for  $\text{bHR}(\text{pop})$  are due to the number of items or whether there is a different underlying reason.

There are several other limitations to this research. First of all, in this research the items are all organised in sub-classes. This is probably done to overcome the sparsity and scalability problem, but it also causes a loss of information. However, a sub-class item recommendation is not as precise as a specific product recommendation and could therefore not be optimal. Second, people are price-sensitive and this research does not take the price of products into consideration while generating recommendations. Not everyone has the same income, and therefore the same budget, to do groceries. The price of a product is a really important factor and by including the price factor in the recommender system, recommendations can be more accurate. Furthermore, it is also of interest for the grocery stores to recommend products with a high profit margin. Hence, including this in the recommender system benefits the grocery store. Third, there are a lot of other factors that may also be of interest when generating recommendations for grocery shoppers, which are not yet included in this paper, i.e. age of the customer, location of product in the store, and whether it is a seasonal product yes or no. Hence, future research needs to point out how including more details in the recommender system benefits the performance of recommender systems.

Besides the ideas for future research proposed by Li et al. (2009), I think it is interesting to investigate the user-based approach in more detail. There are several other methods that can be applied to determine similarity between customers. For example, it can be interesting to cluster customers by age, average price of the basket and average size of the basket and provide recommendations based on the cluster to which a customer belongs. Furthermore, it is also interesting to examine how machine learning techniques like Random-Forest and nearest-neighbour perform on grocery recommendations. These machine learning techniques are, just like the CF algorithm, widely explored in the field of movie recommendation, but have not yet been applied on the domain of grocery shopping. Besides, it might even be possible to apply an item-based machine learning approach to the dataset. One of the main advantages of machine learning techniques is that they

learn from their 'mistakes'. By applying the machine learning technique on a greater amount of data, the machine learning technique is able to recognise patterns and provide even more accurate recommendations.

## References

- Al Mamunur Rashid, S. K. L., Karypis, G. & Riedl, J. (2006). Clustknn: A highly scalable hybrid model-& memory-based cf algorithm. *Proceeding of webKDD, 2006*.
- Breese, J. S., Heckerman, D. & Kadie, C. (1998). Empirical analysis of predictive algorithms for collaborative filtering. In *Proceedings of the fourteenth conference on uncertainty in artificial intelligence* (pp. 43–52). Morgan Kaufmann Publishers Inc.
- Brin, S. & Page, L. (1998). The anatomy of a large-scale hypertextual web search engine. *Computer networks and ISDN systems*, 30(1-7), 107–117.
- Deshpande, M. & Karypis, G. (2004). Item-based top-n recommendation algorithms. *ACM Transactions on Information Systems (TOIS)*, 22(1), 143–177.
- Hsu, C.-N., Chung, H.-H. & Huang, H.-S. (2004). Mining skewed and sparse transaction data for personalized shopping recommendation. *Machine Learning*, 57(1-2), 35–59.
- Huang, Z., Zeng, D. & Chen, H. (2004). A link analysis approach to recommendation under sparse data. *AMCIS 2004 Proceedings*, 239.
- Konstan, J. A., Miller, B. N., Maltz, D., Herlocker, J. L., Gordon, L. R. & Riedl, J. (1997). GroupLens: Applying collaborative filtering to usenet news. *Communications of the ACM*, 40(3), 77–87.
- Li, M., Dias, B. M., Jarman, I., El-Deredy, W. & Lisboa, P. J. (2009). Grocery shopping recommendations based on basket-sensitive random walk. In *Proceedings of the 15th acm sigkdd international conference on knowledge discovery and data mining* (pp. 1215–1224). ACM.
- Li, M., Dias, B., El-Deredy, W. & Lisboa, P. J. (2007). A probabilistic model for item-based recommender systems. In *Proceedings of the 2007 acm conference on recommender systems* (pp. 129–132). ACM.
- Liu, N. N. & Yang, Q. (2008). Eigenrank: A ranking-oriented approach to collaborative filtering. In *Proceedings of the 31st annual international acm sigir conference on research and development in information retrieval* (pp. 83–90). ACM.
- Mobasher, B., Dai, H., Luo, T. & Nakagawa, M. (2001). Effective personalization based on association rule discovery from web usage data. In *Proceedings of the 3rd international workshop on web information and data management* (pp. 9–15). ACM.
- Sarwar, B. M., Karypis, G., Konstan, J. A., Riedl, J. et al. (2001). Item-based collaborative filtering recommendation algorithms. *Www*, 1, 285–295.
- Shardanand, U. & Maes, P. (1995). Social information filtering: Algorithms for automating

- . In *Chi* (Vol. 95, pp. 210–217). Citeseer.
- Sordo-Garcia, C. M., Dias, M. B., Li, M., El-Dereby, W. & Lisboa, P. J. (2007). Evaluating retail recommender systems via retrospective data: Lessons learnt from a live-intervention study. In *Dmin* (Vol. 95, pp. 197–206). Citeseer.
- Wijaya, D. T. & Bressan, S. (2008). A random walk on the red carpet: Rating movies with user reviews and pagerank. In *Proceedings of the 17th acm conference on information and knowledge management* (Vol. 95, pp. 951–960). ACM.
- Yildirim, H. & Krishnamoorthy, M. S. (2008). A random walk method for alleviating the sparsity problem in collaborative filtering. In *Proceedings of the 2008 acm conference on recommender systems* (Vol. 95, pp. 131–138). ACM.

## A Matlab code

Listing 1: Code to generate the descriptive statistics for the Ta-Feng dataset

```
1 transactions2 = 1;
2 baskets2 = [];
3 transDate2 = [];
4 customerTransaction2 = [];
5 subclassPerBasket = zeros(107700,2000);
6 indexNewBasket = [];
7 %amount of baskets
8 for i = 1:720164
9     if i == 1
10        baskets2(transactions2) = newAmount(i);
11        transDate2(transactions2) = newTransactionDate(i);
12        customerTransaction2(transactions2) = newCustomerID(i);
13        indexNewBasket = [indexNewBasket; i];
14    else if newTransactionDate(i+1) == newTransactionDate(i) && newCustomerID(i+1)
        == newCustomerID(i)
15        baskets2(transactions2) = baskets2(transactions2) + newAmount(i+1);
16    else
17        indexNewBasket = [indexNewBasket; i+1];
18        transactions2 = transactions2 + 1;
19        baskets2(transactions2) = newAmount(i+1);
20        transDate2(transactions2) = newTransactionDate(i+1);
21        customerTransaction2(transactions2) = newCustomerID(i+1);
22    end
23 end
24 end
25 meanBasket = mean(baskets2);
26 medianBasket = median(baskets2);
27 minBasket = min(baskets2);
28 maxBasket = max(baskets2);
29 %amount of customers
30 [sortedCustomers2, ia, indexSortedCustomers2] = unique(customerTransaction2);
31 %amount of subclasses
32 [sortedSubClasses2, is, indexSortedSubclasses2] = unique(newSubClasses);
33 %amount of products per subclass
34 amountPerSubclass = zeros(2000, 1);
35 for j=1:2000
36     for i = 1:size(newAmount,2)
37         if indexSortedSubclasses2(i) == j
38             amountPerSubclass(j) = amountPerSubclass(j) +newAmount(i);
39         end
40     end
41 end
42 meanAmountSC = mean(amountPerSubclass);
43 medianAmountSC = median(amountPerSubclass);
44 maxAmountSC = max(amountPerSubclass);
45 minAmountSC = min(amountPerSubclass);
46 %amount of baskets per customer
47 basketsPerCustomer = zeros(20388,1);
```



```

48 for j=1:20388
49     for i = 1:transactions2
50         if indexSortedCustomers2(i) == j
51             basketsPerCustomer(j) = basketsPerCustomer(j) + 1;
52         end
53     end
54 end
55 minbasketsPC = min(basketsPerCustomer);
56 maxbasketsPC = max(basketsPerCustomer);
57 meanbasketsPC = mean(basketsPerCustomer);
58 medianbasketsPC = median(basketsPerCustomer);
59 %aggregated baskets (only necessary for trainingset)
60 aggregatedBaskets = zeros(20388, 1);
61
62 for j = 1:20388
63     for i= 1:size(baskets2,1)
64         if sortedCustomers2(j) == customerTransaction2(i)
65             aggregatedBaskets(j) = aggregatedBaskets(j) + baskets2(i);
66         end
67     end
68 end
69 subclassesInBasket = zeros(20388, 2000);
70 for i = 1:20388
71     i
72     newIndexSubClass = find(newCustomerID == sortedCustomers2(i));
73     for j = 1:size(newIndexSubClass)
74         subclassesInBasket(i,j) = newSubClasses(newIndexSubClass(j));
75     end
76 end
77 %co occurence matrix
78 coOccurenceMatrix = zeros(2000, 2000);
79 for i = 1:20388
80     uniqueSubClassesPerBasket = unique(subclassesInBasket(i,:));
81     for j = 2:size(uniqueSubClassesPerBasket, 2)-1
82         for k = j+1:size(uniqueSubClassesPerBasket,2)
83             coOccurenceMatrix(indexMatrix, indexMatrix2) = 1;
84         end
85     end
86 end
87 %generate test set
88 [sortedCustomersLast, iaLast, indexSortedCustomersLast] = unique(
    customerTransaction2, 'last');
89 newCustomerIDtest = customerTransaction2(iaLast)';
90 newIndexNewBasketTest = indexNewBasket(iaLast);
91 newBasketsTest = baskets2(iaLast)';
92 newTransDateTest = transDate2(iaLast)';
93 subclassPerBasketTest = subclassPerBasket(iaLast,:);
94 newAmountTest = newAmount(iaLast,:);
95 %remove removed subclasses from baskets
96 removedSubclasses = uniqueSubClasses(indexZero);
97 adjustments = 0;
98 for i = 1:20388
99     for j = 1:2000

```

```

100     for k = 1:27
101         if removedSubclasses(k) == subclassPerBasketTest(i,j)
102             subclassPerBasketTest(i,j) = 0;
103             adjustments = adjustments + 1;
104         end
105     end
106 end
107 end
108 end

```

Listing 2: Code to generate the training and test set and the similarity matrices

```

1  indexNewCustomerTrain = [];
2  for i = 1:20388
3      if i == 1
4          indexNewCustomerTrain = [indexNewCustomerTrain; i];
5      else if sortedCustomerIDtrain(i+1) == sortedCustomerIDtrain(i)
6          %do nothing
7      else
8          indexNewCustomerTrain = [indexNewCustomerTrain; i+1];
9      end
10 end
11 end
12 %determine baskets
13 userItemMatrix = zeros(20388, 2000);
14 for i=1:size(indexNewCustomer,1)
15     if i == 1
16         for j = 1:indexNewCustomer(i+1)-1
17             for k = 1:nnz(sortedSubclassPerBasket(indexIDtrain(j),:))
18                 indexSubClass = find(subclassPerBasketTrain(indexIDtrain(j),k) ==
19                     uniqueSubClasses);
20                 userItemMatrix(i,indexSubClass) = userItemMatrix(i,indexSubClass) +
21                     size(indexSubClass,2);
22             end
23         end
24     else if i == size(indexNewCustomer,1)
25         for j = indexNewCustomer(i):87312
26             for k = 1:nnz(sortedSubclassPerBasket(indexIDtrain(j),:))
27                 indexSubClass = find(subclassPerBasketTrain(indexIDtrain(j),k)
28                     == uniqueSubClasses);
29                 userItemMatrix(i,indexSubClass) = userItemMatrix(i,indexSubClass
30                     ) + size(indexSubClass,2);
31             end
32         end
33     else
34         for j = indexNewCustomer(i):indexNewCustomer(i+1) - 1
35             for k = 1:nnz(sortedSubclassPerBasket(indexIDtrain(j),:))
36                 indexSubClass = find(subclassPerBasketTrain(indexIDtrain(j),k)

```

```

37     end
38 end
39 %delete subclasses which are never bought
40 SUM = sum(userItemMatrix,1);
41 indexZero = [];
42 for i=1:2000
43     if SUM(i) == 0
44         indexZero = [indexZero;i];
45     end
46 end
47 userItemMatrix2 = userItemMatrix;
48 userItemMatrix2(:,indexZero)=[];
49 %delete subclasses which are never bought in the training set
50 uniqueSubClasses2 = uniqueSubClasses;
51 uniqueSubClasses2(indexZero) = [];
52 %define training baskets
53 trainingBaskets = zeros(20388, 1973);
54 for i = 1:20388
55     for j = 1:1973
56         if userItemMatrix2(i,j) > 0
57             trainingBaskets(i,j) = uniqueSubClasses2(j);
58         end
59     end
60 end
61 trainingBaskets = sort(trainingBaskets,2, 'descend');
62 %normalize user-item matrix
63 logUserItemMatrix = log(userItemMatrix2 + 1);
64 sumUser = sum(logUserItemMatrix,2);
65 for i = 1:20388
66     for j = 1:1973
67         normUserItemMatrix(i,j) = logUserItemMatrix(i,j)/sumUser(i);
68     end
69 end
70 %cosine based similarity
71 cosineSim = eye(1973);
72 for i = 1:1972
73     for j=i+1:1973
74         cosineSim(i,j) = dot(normUserItemMatrix(:,i), normUserItemMatrix(:,j))/( (
75             sqrt(sum(abs(normUserItemMatrix(:,i)).^2))) * (sqrt(sum(abs(
76                 normUserItemMatrix(:,j)).^2))) );
77         cosineSim(j,i) = cosineSim(i,j);
78     end
79 end
80 %conditional prob. based similarity
81 for i=1:1973
82     freq(i) = nnz(userItemMatrix2(:,i));
83 end
84 totAmountOfBoughtProducts = 0;
85 for i = 1:20388
86     totAmountOfBoughtProducts = totAmountOfBoughtProducts + sum(logUserItemMatrix(i
87         ,:));
88 end
89 for i = 1:1973

```

```

87     popularity(i) = sum(logUserItemMatrix(:,i))/totAmountOfBoughtProducts;
88 end
89 condProbSimAlpha5 = eye(1973);
90 condProbSimAlpha7 = eye(1973);
91 condProbSimAlpha9 = eye(1973);
92 for i = 1:1972
93     for j = i+1:1973
94         sumPosEntries = 0;
95         for c = 1:20388
96             if normUserItemMatrix(c,i)>0
97                 sumPosEntries = sumPosEntries + normUserItemMatrix(c,j);
98             end
99         end
100        sumPosEntries2 = 0;
101        for c = 1:20388
102            if normUserItemMatrix(c,j)>0
103                sumPosEntries2 = sumPosEntries2 + normUserItemMatrix(c,i);
104            end
105        end
106        condProbSimAlpha5(i,j) = sumPosEntries / (freq(i) * freq(j)^0.5);
107        condProbSimAlpha5(j,i) = sumPosEntries2 / (freq(j) * freq(i)^0.5);
108        condProbSimAlpha7(i,j) = sumPosEntries / (freq(i) * freq(j)^0.7);
109        condProbSimAlpha7(j,i) = sumPosEntries2 / (freq(j) * freq(i)^0.7);
110        condProbSimAlpha9(i,j) = sumPosEntries / (freq(i) * freq(j)^0.9);
111        condProbSimAlpha9(j,i) = sumPosEntries2 / (freq(j) * freq(i)^0.9);
112    end
113 end

```

Listing 3: Code to determine the transition probability matrix based on a bipartite network for three values of  $\alpha$

```

1  totProductsPerUser = sum(logUserItemMatrix,2);
2  totAmountProduct = sum(logUserItemMatrix, 1);
3  for i = 1:20388
4      for j = 1:1973
5          probCustToProd5(i,j) = logUserItemMatrix(i,j)/( (totProductsPerUser(i))^0.5
6              );
7          probProdToCust5(i,j) = logUserItemMatrix(i,j)/( (totAmountProduct(j))^0.5 );
8          probCustToProd7(i,j) = logUserItemMatrix(i,j)/( (totProductsPerUser(i))^0.7
9              );
10         probProdToCust7(i,j) = logUserItemMatrix(i,j)/( (totAmountProduct(j))^0.7 );
11         probCustToProd9(i,j) = logUserItemMatrix(i,j)/( (totProductsPerUser(i))^0.9
12             );
13         probProdToCust9(i,j) = logUserItemMatrix(i,j)/( (totAmountProduct(j))^0.9 );
14     end
15 end
16 %determine transition probability matrix of bipartite network
17 bnMatrix5 = zeros(1973, 1973);
18 bnMatrix7 = zeros(1973, 1973);
19 bnMatrix9 = zeros(1973, 1973);
20 for j = 1:1973
21     for k=1:1973
22         for i = 1:20388

```

```

20     bnMatrix5(j,k) = bnMatrix5(j,k) + probCustToProd5(i,j)*probProdToCust5(i
      ,k);
21     bnMatrix5(k,j) = bnMatrix5(k,j) + probCustToProd5(i,k)*probProdToCust5(i
      ,j);
22     bnMatrix7(j,k) = bnMatrix7(j,k) + probCustToProd7(i,j)*probProdToCust7(i
      ,k);
23     bnMatrix7(k,j) = bnMatrix7(k,j) + probCustToProd7(i,k)*probProdToCust7(i
      ,j);
24     bnMatrix9(j,k) = bnMatrix9(j,k) + probCustToProd9(i,j)*probProdToCust9(i
      ,k);
25     bnMatrix9(k,j) = bnMatrix9(k,j) + probCustToProd9(i,k)*probProdToCust9(i
      ,j);
26     end
27 end
28 end

```

Listing 4: Code to determine the ratings for products based on the Random Walk model

```

1  %determine personalization vector U
2  Uitem = eye(1973);
3  %rating CS
4  RitemCS = inv(eye(1973) - .1*normCS) * (1 - .1) * Uitem;
5  RitemCS = inv(eye(1973) - .2*normCS) * (1 - .2) * Uitem;
6  RitemCS = inv(eye(1973) - .3*normCS) * (1 - .3) * Uitem;
7  RitemCS = inv(eye(1973) - .4*normCS) * (1 - .4) * Uitem;
8  RitemCS = inv(eye(1973) - .5*normCS) * (1 - .5) * Uitem;
9  RitemCS = inv(eye(1973) - .6*normCS) * (1 - .6) * Uitem;
10 RitemCS = inv(eye(1973) - .7*normCS) * (1 - .7) * Uitem;
11 RitemCS = inv(eye(1973) - .8*normCS) * (1 - .8) * Uitem;
12 RitemCS = inv(eye(1973) - .9*normCS) * (1 - .9) * Uitem;
13 %rating CP alpha 0.5
14 RitemCP5_1 = inv(eye(1973) - .1*normCP5) * (1 - .1) * Uitem;
15 RitemCP5_2 = inv(eye(1973) - .2*normCP5) * (1 - .2) * Uitem;
16 RitemCP5_3 = inv(eye(1973) - .3*normCP5) * (1 - .3) * Uitem;
17 RitemCP5_4 = inv(eye(1973) - .4*normCP5) * (1 - .4) * Uitem;
18 RitemCP5_5 = inv(eye(1973) - .5*normCP5) * (1 - .5) * Uitem;
19 RitemCP5_6 = inv(eye(1973) - .6*normCP5) * (1 - .6) * Uitem;
20 RitemCP5_7 = inv(eye(1973) - .7*normCP5) * (1 - .7) * Uitem;
21 RitemCP5_8 = inv(eye(1973) - .8*normCP5) * (1 - .8) * Uitem;
22 RitemCP5_9 = inv(eye(1973) - .9*normCP5) * (1 - .9) * Uitem;
23 %rating CP alpha 0.7
24 RitemCP7_1 = inv(eye(1973) - .1*normCP7) * (1 - .1) * Uitem;
25 RitemCP7_2 = inv(eye(1973) - .2*normCP7) * (1 - .2) * Uitem;
26 RitemCP7_3 = inv(eye(1973) - .3*normCP7) * (1 - .3) * Uitem;
27 RitemCP7_4 = inv(eye(1973) - .4*normCP7) * (1 - .4) * Uitem;
28 RitemCP7_5 = inv(eye(1973) - .5*normCP7) * (1 - .5) * Uitem;
29 RitemCP7_6 = inv(eye(1973) - .6*normCP7) * (1 - .6) * Uitem;
30 RitemCP7_7 = inv(eye(1973) - .7*normCP7) * (1 - .7) * Uitem;
31 RitemCP7_8 = inv(eye(1973) - .8*normCP7) * (1 - .8) * Uitem;
32 RitemCP7_9 = inv(eye(1973) - .9*normCP7) * (1 - .9) * Uitem;
33 %rating CP alpha 0.9
34 RitemCP9_1 = inv(eye(1973) - .1*normCP9) * (1 - .1) * Uitem;
35 RitemCP9_2 = inv(eye(1973) - .2*normCP9) * (1 - .2) * Uitem;
36 RitemCP9_3 = inv(eye(1973) - .3*normCP9) * (1 - .3) * Uitem;

```

```

37 RitemCP9_4 = inv(eye(1973) - .4*normCP9) * (1 - .4) * Uitem;
38 RitemCP9_5 = inv(eye(1973) - .5*normCP9) * (1 - .5) * Uitem;
39 RitemCP9_6 = inv(eye(1973) - .6*normCP9) * (1 - .6) * Uitem;
40 RitemCP9_7 = inv(eye(1973) - .7*normCP9) * (1 - .7) * Uitem;
41 RitemCP9_8 = inv(eye(1973) - .8*normCP9) * (1 - .8) * Uitem;
42 RitemCP9_9 = inv(eye(1973) - .9*normCP9) * (1 - .9) * Uitem;
43 %rating BN alpha 0.5
44 RitemBN5_1 = inv(eye(1973) - .1*normBnMatrix5) * (1 - .1) * Uitem;
45 RitemBN5_2 = inv(eye(1973) - .2*normBnMatrix5) * (1 - .2) * Uitem;
46 RitemBN5_3 = inv(eye(1973) - .3*normBnMatrix5) * (1 - .3) * Uitem;
47 RitemBN5_4 = inv(eye(1973) - .4*normBnMatrix5) * (1 - .4) * Uitem;
48 RitemBN5_5 = inv(eye(1973) - .5*normBnMatrix5) * (1 - .5) * Uitem;
49 RitemBN5_6 = inv(eye(1973) - .6*normBnMatrix5) * (1 - .6) * Uitem;
50 RitemBN5_7 = inv(eye(1973) - .7*normBnMatrix5) * (1 - .7) * Uitem;
51 RitemBN5_8 = inv(eye(1973) - .8*normBnMatrix5) * (1 - .8) * Uitem;
52 RitemBN5_9 = inv(eye(1973) - .9*normBnMatrix5) * (1 - .9) * Uitem;
53 %rating BN alpha 0.7
54 RitemBN7_1 = inv(eye(1973) - .1*normBnMatrix7) * (1 - .1) * Uitem;
55 RitemBN7_2 = inv(eye(1973) - .2*normBnMatrix7) * (1 - .2) * Uitem;
56 RitemBN7_3 = inv(eye(1973) - .3*normBnMatrix7) * (1 - .3) * Uitem;
57 RitemBN7_4 = inv(eye(1973) - .4*normBnMatrix7) * (1 - .4) * Uitem;
58 RitemBN7_5 = inv(eye(1973) - .5*normBnMatrix7) * (1 - .5) * Uitem;
59 RitemBN7_6 = inv(eye(1973) - .6*normBnMatrix7) * (1 - .6) * Uitem;
60 RitemBN7_7 = inv(eye(1973) - .7*normBnMatrix7) * (1 - .7) * Uitem;
61 RitemBN7_8 = inv(eye(1973) - .8*normBnMatrix7) * (1 - .8) * Uitem;
62 RitemBN7_9 = inv(eye(1973) - .9*normBnMatrix7) * (1 - .9) * Uitem;
63 %rating BN alpha 0.9
64 RitemBN9_1 = inv(eye(1973) - .1*normBnMatrix9) * (1 - .1) * Uitem;
65 RitemBN9_2 = inv(eye(1973) - .2*normBnMatrix9) * (1 - .2) * Uitem;
66 RitemBN9_3 = inv(eye(1973) - .3*normBnMatrix9) * (1 - .3) * Uitem;
67 RitemBN9_4 = inv(eye(1973) - .4*normBnMatrix9) * (1 - .4) * Uitem;
68 RitemBN9_5 = inv(eye(1973) - .5*normBnMatrix9) * (1 - .5) * Uitem;
69 RitemBN9_6 = inv(eye(1973) - .6*normBnMatrix9) * (1 - .6) * Uitem;
70 RitemBN9_7 = inv(eye(1973) - .7*normBnMatrix9) * (1 - .7) * Uitem;
71 RitemBN9_8 = inv(eye(1973) - .8*normBnMatrix9) * (1 - .8) * Uitem;
72 RitemBN9_9 = inv(eye(1973) - .9*normBnMatrix9) * (1 - .9) * Uitem;

```

Listing 5: Code to determine the ratings of the products based on the evidence products with help of the similarity matrices. The code also contains the computation of the binary hit rate based on a popularity selection of targets. The code for the hit rate based on a random selection of targets is similar to the code based on a popularity selection.

```

1 %determine targets test basket based on popularity
2 for i = 1:11813
3     for j = 1:nnz(testBaskets(i,:))
4         indexSubClassTest = find(testBaskets(i,j) == uniqueSubClasses2);
5         freqSubclass(i,j) = popularity(indexSubClassTest);
6     end
7 end
8 freqSubclass( freqSubclass == 0 ) = Inf;
9 [sortedFreqSubclass, indexFreqSubclass] = sort(freqSubclass,2,'ascend');
10 freqSubclass( isinf(freqSubclass) ) = 0;
11 sortedFreqSubclass( isinf(sortedFreqSubclass) ) = 0;

```

```

12 for i = 1:11813
13     for j = 1:3
14         targetSubclass(i,j) = testBaskets(i,indexFreqSubclass(i,j));
15     end
16     for j = 4:nnz(testBaskets(i,:))
17         evidenceSubclass(i,j) = testBaskets(i,indexFreqSubclass(i,j));
18     end
19 end
20 %determine recommendations
21 csRows1=zeros(11813,1973);csRows2=zeros(11813,1973);csRows3=zeros(11813,1973);
    csRows4=zeros(11813,1973);csRows5=zeros(11813,1973);csRows6=zeros(11813,1973);
    csRows7=zeros(11813,1973);csRows8=zeros(11813,1973);csRows9=zeros(11813,1973);
22 cp5Rows1=zeros(11813,1973);cp5Rows2=zeros(11813,1973);cp5Rows3=zeros(11813,1973);
    cp5Rows4=zeros(11813,1973);cp5Rows5=zeros(11813,1973);cp5Rows6=zeros(11813,1973)
    ;cp5Rows7=zeros(11813,1973);cp5Rows8=zeros(11813,1973);cp5Rows9=zeros
    (11813,1973);
23 cp7Rows1=zeros(11813,1973);cp7Rows2=zeros(11813,1973);cp7Rows3=zeros(11813,1973);
    cp7Rows4=zeros(11813,1973);cp7Rows5=zeros(11813,1973);cp7Rows6=zeros(11813,1973)
    ;cp7Rows7=zeros(11813,1973);cp7Rows8=zeros(11813,1973);cp7Rows9=zeros
    (11813,1973);
24 cp9Rows1=zeros(11813,1973);cp9Rows2=zeros(11813,1973);cp9Rows3=zeros(11813,1973);
    cp9Rows4=zeros(11813,1973);cp9Rows5=zeros(11813,1973);cp9Rows6=zeros(11813,1973)
    ;cp9Rows7=zeros(11813,1973);cp9Rows8=zeros(11813,1973);cp9Rows9=zeros
    (11813,1973);
25 bn5Rows1=zeros(11813,1973);bn5Rows2=zeros(11813,1973);bn5Rows3=zeros(11813,1973);
    bn5Rows4=zeros(11813,1973);bn5Rows5=zeros(11813,1973);bn5Rows6=zeros(11813,1973)
    ;bn5Rows7=zeros(11813,1973);bn5Rows8=zeros(11813,1973);bn5Rows9=zeros
    (11813,1973);
26 bn7Rows1=zeros(11813,1973);bn7Rows2=zeros(11813,1973);bn7Rows3=zeros(11813,1973);
    bn7Rows4=zeros(11813,1973);bn7Rows5=zeros(11813,1973);bn7Rows6=zeros(11813,1973)
    ;bn7Rows7=zeros(11813,1973);bn7Rows8=zeros(11813,1973);bn7Rows9=zeros
    (11813,1973);
27 bn9Rows1=zeros(11813,1973);bn9Rows2=zeros(11813,1973);bn9Rows3=zeros(11813,1973);
    bn9Rows4=zeros(11813,1973);bn9Rows5=zeros(11813,1973);bn9Rows6=zeros(11813,1973)
    ;bn9Rows7=zeros(11813,1973);bn9Rows8=zeros(11813,1973);bn9Rows9=zeros
    (11813,1973);
28 for i = 1:11813
29     for j = 4:nnz(testBaskets(i,:))
30         indexSubClassTest = find(testBaskets(i,indexFreqSubclass(i,j)) ==
            uniqueSubClasses2);
31         csRows1(i,:) = csRows1(i,:) + RitemCS1(indexSubClassTest,:);
32         csRows2(i,:) = csRows2(i,:) + RitemCS2(indexSubClassTest,:);
33         csRows3(i,:) = csRows3(i,:) + RitemCS3(indexSubClassTest,:);
34         csRows4(i,:) = csRows4(i,:) + RitemCS4(indexSubClassTest,:);
35         csRows5(i,:) = csRows5(i,:) + RitemCS5(indexSubClassTest,:);
36         csRows6(i,:) = csRows6(i,:) + RitemCS6(indexSubClassTest,:);
37         csRows7(i,:) = csRows7(i,:) + RitemCS7(indexSubClassTest,:);
38         csRows8(i,:) = csRows8(i,:) + RitemCS8(indexSubClassTest,:);
39         csRows9(i,:) = csRows9(i,:) + RitemCS9(indexSubClassTest,:);
40         cp5Rows1(i,:) = cp5Rows1(i,:) + RitemCP5_1(indexSubClassTest,:);
41         cp5Rows2(i,:) = cp5Rows2(i,:) + RitemCP5_2(indexSubClassTest,:);
42         cp5Rows3(i,:) = cp5Rows3(i,:) + RitemCP5_3(indexSubClassTest,:);
43         cp5Rows4(i,:) = cp5Rows4(i,:) + RitemCP5_4(indexSubClassTest,:);

```

```

44     cp5Rows5(i,:) = cp5Rows5(i,:) + RitemCP5_5(indexSubClassTest,:);
45     cp5Rows6(i,:) = cp5Rows6(i,:) + RitemCP5_6(indexSubClassTest,:);
46     cp5Rows7(i,:) = cp5Rows7(i,:) + RitemCP5_7(indexSubClassTest,:);
47     cp5Rows8(i,:) = cp5Rows8(i,:) + RitemCP5_8(indexSubClassTest,:);
48     cp5Rows9(i,:) = cp5Rows9(i,:) + RitemCP5_9(indexSubClassTest,:);
49     cp7Rows1(i,:) = cp7Rows1(i,:) + RitemCP7_1(indexSubClassTest,:);
50     cp7Rows2(i,:) = cp7Rows2(i,:) + RitemCP7_2(indexSubClassTest,:);
51     cp7Rows3(i,:) = cp7Rows3(i,:) + RitemCP7_3(indexSubClassTest,:);
52     cp7Rows4(i,:) = cp7Rows4(i,:) + RitemCP7_4(indexSubClassTest,:);
53     cp7Rows5(i,:) = cp7Rows5(i,:) + RitemCP7_5(indexSubClassTest,:);
54     cp7Rows6(i,:) = cp7Rows6(i,:) + RitemCP7_6(indexSubClassTest,:);
55     cp7Rows7(i,:) = cp7Rows7(i,:) + RitemCP7_7(indexSubClassTest,:);
56     cp7Rows8(i,:) = cp7Rows8(i,:) + RitemCP7_8(indexSubClassTest,:);
57     cp7Rows9(i,:) = cp7Rows9(i,:) + RitemCP7_9(indexSubClassTest,:);
58     cp9Rows1(i,:) = cp9Rows1(i,:) + RitemCP9_1(indexSubClassTest,:);
59     cp9Rows2(i,:) = cp9Rows2(i,:) + RitemCP9_2(indexSubClassTest,:);
60     cp9Rows3(i,:) = cp9Rows3(i,:) + RitemCP9_3(indexSubClassTest,:);
61     cp9Rows4(i,:) = cp9Rows4(i,:) + RitemCP9_4(indexSubClassTest,:);
62     cp9Rows5(i,:) = cp9Rows5(i,:) + RitemCP9_5(indexSubClassTest,:);
63     cp9Rows6(i,:) = cp9Rows6(i,:) + RitemCP9_6(indexSubClassTest,:);
64     cp9Rows7(i,:) = cp9Rows7(i,:) + RitemCP9_7(indexSubClassTest,:);
65     cp9Rows8(i,:) = cp9Rows8(i,:) + RitemCP9_8(indexSubClassTest,:);
66     cp9Rows9(i,:) = cp9Rows9(i,:) + RitemCP9_9(indexSubClassTest,:);
67     bn5Rows1(i,:) = bn5Rows1(i,:) + RitemBN5_1(indexSubClassTest,:);
68     bn5Rows2(i,:) = bn5Rows2(i,:) + RitemBN5_2(indexSubClassTest,:);
69     bn5Rows3(i,:) = bn5Rows3(i,:) + RitemBN5_3(indexSubClassTest,:);
70     bn5Rows4(i,:) = bn5Rows4(i,:) + RitemBN5_4(indexSubClassTest,:);
71     bn5Rows5(i,:) = bn5Rows5(i,:) + RitemBN5_5(indexSubClassTest,:);
72     bn5Rows6(i,:) = bn5Rows6(i,:) + RitemBN5_6(indexSubClassTest,:);
73     bn5Rows7(i,:) = bn5Rows7(i,:) + RitemBN5_7(indexSubClassTest,:);
74     bn5Rows8(i,:) = bn5Rows8(i,:) + RitemBN5_8(indexSubClassTest,:);
75     bn5Rows9(i,:) = bn5Rows9(i,:) + RitemBN5_9(indexSubClassTest,:);
76     bn7Rows1(i,:) = bn7Rows1(i,:) + RitemBN7_1(indexSubClassTest,:);
77     bn7Rows2(i,:) = bn7Rows2(i,:) + RitemBN7_2(indexSubClassTest,:);
78     bn7Rows3(i,:) = bn7Rows3(i,:) + RitemBN7_3(indexSubClassTest,:);
79     bn7Rows4(i,:) = bn7Rows4(i,:) + RitemBN7_4(indexSubClassTest,:);
80     bn7Rows5(i,:) = bn7Rows5(i,:) + RitemBN7_5(indexSubClassTest,:);
81     bn7Rows6(i,:) = bn7Rows6(i,:) + RitemBN7_6(indexSubClassTest,:);
82     bn7Rows7(i,:) = bn7Rows7(i,:) + RitemBN7_7(indexSubClassTest,:);
83     bn7Rows8(i,:) = bn7Rows8(i,:) + RitemBN7_8(indexSubClassTest,:);
84     bn7Rows9(i,:) = bn7Rows9(i,:) + RitemBN7_9(indexSubClassTest,:);
85     bn9Rows1(i,:) = bn9Rows1(i,:) + RitemBN9_1(indexSubClassTest,:);
86     bn9Rows2(i,:) = bn9Rows2(i,:) + RitemBN9_2(indexSubClassTest,:);
87     bn9Rows3(i,:) = bn9Rows3(i,:) + RitemBN9_3(indexSubClassTest,:);
88     bn9Rows4(i,:) = bn9Rows4(i,:) + RitemBN9_4(indexSubClassTest,:);
89     bn9Rows5(i,:) = bn9Rows5(i,:) + RitemBN9_5(indexSubClassTest,:);
90     bn9Rows6(i,:) = bn9Rows6(i,:) + RitemBN9_6(indexSubClassTest,:);
91     bn9Rows7(i,:) = bn9Rows7(i,:) + RitemBN9_7(indexSubClassTest,:);
92     bn9Rows8(i,:) = bn9Rows8(i,:) + RitemBN9_8(indexSubClassTest,:);
93     bn9Rows9(i,:) = bn9Rows9(i,:) + RitemBN9_9(indexSubClassTest,:);
94     end
95 end
96 %determine place ratings

```



```

97 [maxCS1, indexMaxCS1] = sort(csRows1, 2, 'descend');
98 [maxCS2, indexMaxCS2] = sort(csRows2, 2, 'descend');
99 [maxCS3, indexMaxCS3] = sort(csRows3, 2, 'descend');
100 [maxCS4, indexMaxCS4] = sort(csRows4, 2, 'descend');
101 [maxCS5, indexMaxCS5] = sort(csRows5, 2, 'descend');
102 [maxCS6, indexMaxCS6] = sort(csRows6, 2, 'descend');
103 [maxCS7, indexMaxCS7] = sort(csRows7, 2, 'descend');
104 [maxCS8, indexMaxCS8] = sort(csRows8, 2, 'descend');
105 [maxCS9, indexMaxCS9] = sort(csRows9, 2, 'descend');
106 [indexCP5_1, indexMaxCP5_1] = sort(cp5Rows1, 2, 'descend');
107 [indexCP5_2, indexMaxCP5_2] = sort(cp5Rows2, 2, 'descend');
108 [indexCP5_3, indexMaxCP5_3] = sort(cp5Rows3, 2, 'descend');
109 [indexCP5_4, indexMaxCP5_4] = sort(cp5Rows4, 2, 'descend');
110 [indexCP5_5, indexMaxCP5_5] = sort(cp5Rows5, 2, 'descend');
111 [indexCP5_6, indexMaxCP5_6] = sort(cp5Rows6, 2, 'descend');
112 [indexCP5_7, indexMaxCP5_7] = sort(cp5Rows7, 2, 'descend');
113 [indexCP5_8, indexMaxCP5_8] = sort(cp5Rows8, 2, 'descend');
114 [indexCP5_9, indexMaxCP5_9] = sort(cp5Rows9, 2, 'descend');
115 [indexCP7_1, indexMaxCP7_1] = sort(cp7Rows1, 2, 'descend');
116 [indexCP7_2, indexMaxCP7_2] = sort(cp7Rows2, 2, 'descend');
117 [indexCP7_3, indexMaxCP7_3] = sort(cp7Rows3, 2, 'descend');
118 [indexCP7_4, indexMaxCP7_4] = sort(cp7Rows4, 2, 'descend');
119 [indexCP7_5, indexMaxCP7_5] = sort(cp7Rows5, 2, 'descend');
120 [indexCP7_6, indexMaxCP7_6] = sort(cp7Rows6, 2, 'descend');
121 [indexCP7_7, indexMaxCP7_7] = sort(cp7Rows7, 2, 'descend');
122 [indexCP7_8, indexMaxCP7_8] = sort(cp7Rows8, 2, 'descend');
123 [indexCP7_9, indexMaxCP7_9] = sort(cp7Rows9, 2, 'descend');
124 [indexCP9_1, indexMaxCP9_1] = sort(cp9Rows1, 2, 'descend');
125 [indexCP9_2, indexMaxCP9_2] = sort(cp9Rows2, 2, 'descend');
126 [indexCP9_3, indexMaxCP9_3] = sort(cp9Rows3, 2, 'descend');
127 [indexCP9_4, indexMaxCP9_4] = sort(cp9Rows4, 2, 'descend');
128 [indexCP9_5, indexMaxCP9_5] = sort(cp9Rows5, 2, 'descend');
129 [indexCP9_6, indexMaxCP9_6] = sort(cp9Rows6, 2, 'descend');
130 [indexCP9_7, indexMaxCP9_7] = sort(cp9Rows7, 2, 'descend');
131 [indexCP9_8, indexMaxCP9_8] = sort(cp9Rows8, 2, 'descend');
132 [indexCP9_9, indexMaxCP9_9] = sort(cp9Rows9, 2, 'descend');
133 [indexBN5_1, indexMaxBN5_1] = sort(bn5Rows1, 2, 'descend');
134 [indexBN5_2, indexMaxBN5_2] = sort(bn5Rows2, 2, 'descend');
135 [indexBN5_3, indexMaxBN5_3] = sort(bn5Rows3, 2, 'descend');
136 [indexBN5_4, indexMaxBN5_4] = sort(bn5Rows4, 2, 'descend');
137 [indexBN5_5, indexMaxBN5_5] = sort(bn5Rows5, 2, 'descend');
138 [indexBN5_6, indexMaxBN5_6] = sort(bn5Rows6, 2, 'descend');
139 [indexBN5_7, indexMaxBN5_7] = sort(bn5Rows7, 2, 'descend');
140 [indexBN5_8, indexMaxBN5_8] = sort(bn5Rows8, 2, 'descend');
141 [indexBN5_9, indexMaxBN5_9] = sort(bn5Rows9, 2, 'descend');
142 [indexBN7_1, indexMaxBN7_1] = sort(bn7Rows1, 2, 'descend');
143 [indexBN7_2, indexMaxBN7_2] = sort(bn7Rows2, 2, 'descend');
144 [indexBN7_3, indexMaxBN7_3] = sort(bn7Rows3, 2, 'descend');
145 [indexBN7_4, indexMaxBN7_4] = sort(bn7Rows4, 2, 'descend');
146 [indexBN7_5, indexMaxBN7_5] = sort(bn7Rows5, 2, 'descend');
147 [indexBN7_6, indexMaxBN7_6] = sort(bn7Rows6, 2, 'descend');
148 [indexBN7_7, indexMaxBN7_7] = sort(bn7Rows7, 2, 'descend');
149 [indexBN7_8, indexMaxBN7_8] = sort(bn7Rows8, 2, 'descend');

```

```

150 [indexBN7_9, indexMaxBN7_9] = sort(bn7Rows9, 2, 'descend');
151 [indexBN9_1, indexMaxBN9_1] = sort(bn9Rows1, 2, 'descend');
152 [indexBN9_2, indexMaxBN9_2] = sort(bn9Rows2, 2, 'descend');
153 [indexBN9_3, indexMaxBN9_3] = sort(bn9Rows3, 2, 'descend');
154 [indexBN9_4, indexMaxBN9_4] = sort(bn9Rows4, 2, 'descend');
155 [indexBN9_5, indexMaxBN9_5] = sort(bn9Rows5, 2, 'descend');
156 [indexBN9_6, indexMaxBN9_6] = sort(bn9Rows6, 2, 'descend');
157 [indexBN9_7, indexMaxBN9_7] = sort(bn9Rows7, 2, 'descend');
158 [indexBN9_8, indexMaxBN9_8] = sort(bn9Rows8, 2, 'descend');
159 [indexBN9_9, indexMaxBN9_9] = sort(bn9Rows9, 2, 'descend');
160 %determine recommendations and hitrates
161 for i = 1:11813
162     a = 0;
163     b = 0;
164     c = 0;
165     d = 0;
166     e = 0;
167     f = 0;
168     g = 0;
169     h = 0;
170     z = 0;
171     for j = 4:1973
172         if a < 3
173             if ismember(uniqueSubClasses2(indexMaxCS1(i,j-3)),evidenceSubclass(i,:))
174                 > 0
175                 %do nothing
176             else
177                 a = a+1;
178                 recCS1(i,a) = uniqueSubClasses2(indexMaxCS1(i,j-3));
179             end
180         if b < 3
181             if ismember(uniqueSubClasses2(indexMaxCS2(i,j-3)),evidenceSubclass(i,:))
182                 > 0
183                 %do nothing
184             else
185                 b = b+1;
186                 recCS2(i,b) = uniqueSubClasses2(indexMaxCS2(i,j-3));
187             end
188         if c < 3
189             if ismember(uniqueSubClasses2(indexMaxCS3(i,j-3)),evidenceSubclass(i,:))
190                 > 0
191                 %do nothing
192             else
193                 c = c+1;
194                 recCS3(i,c) = uniqueSubClasses2(indexMaxCS3(i,j-3));
195             end
196         if d < 3
197             if ismember(uniqueSubClasses2(indexMaxCS4(i,j-3)),evidenceSubclass(i,:))
198                 > 0
199                 %do nothing

```

```

199         else
200             d = d+1;
201             recCS4(i,d) = uniqueSubClasses2(indexMaxCS4(i,j-3));
202         end
203     end
204     if e < 3
205         if ismember(uniqueSubClasses2(indexMaxCS5(i,j-3)),evidenceSubclass(i,:))
206             > 0
207             %do nothing
208         else
209             e = e+1;
210             recCS5(i,e) = uniqueSubClasses2(indexMaxCS5(i,j-3));
211         end
212     end
213     if f < 3
214         if ismember(uniqueSubClasses2(indexMaxCS6(i,j-3)),evidenceSubclass(i,:))
215             > 0
216             %do nothing
217         else
218             f = f+1;
219             recCS6(i,f) = uniqueSubClasses2(indexMaxCS6(i,j-3));
220         end
221     end
222     if g < 3
223         if ismember(uniqueSubClasses2(indexMaxCS7(i,j-3)),evidenceSubclass(i,:))
224             > 0
225             %do nothing
226         else
227             g = g+1;
228             recCS7(i,g) = uniqueSubClasses2(indexMaxCS7(i,j-3));
229         end
230     end
231     if h < 3
232         if ismember(uniqueSubClasses2(indexMaxCS8(i,j-3)),evidenceSubclass(i,:))
233             > 0
234             %do nothing
235         else
236             h = h+1;
237             recCS8(i,h) = uniqueSubClasses2(indexMaxCS8(i,j-3));
238         end
239     end
240     if z < 3
241         if ismember(uniqueSubClasses2(indexMaxCS9(i,j-3)),evidenceSubclass(i,:))
242             > 0
243             %do nothing
244         else
245             z = z+1;
246             recCS9(i,z) = uniqueSubClasses2(indexMaxCS9(i,j-3));
247         end
248     end
249     end
250     end
251     %determine hit rate similarities

```

```

247 bHR_CS1 = zeros(11813,3);bHR_CS2 = zeros(11813,3);bHR_CS3 = zeros(11813,3);bHR_CS4 =
      zeros(11813,3);bHR_CS5 = zeros(11813,3);bHR_CS6 = zeros(11813,3);bHR_CS7 =
      zeros(11813,3);bHR_CS8 = zeros(11813,3);bHR_CS9 = zeros(11813,3);
248 for i = 1:11813
249     for j = 1:3
250         for z= 1:3
251             if targetSubclass(i,j) == recCS1(i,z)
252                 bHR_CS1(i,j) = 1;
253             end
254             if targetSubclass(i,j) == recCS2(i,z)
255                 bHR_CS2(i,j) = 1;
256             end
257             if targetSubclass(i,j) == recCS3(i,z)
258                 bHR_CS3(i,j) = 1;
259             end
260             if targetSubclass(i,j) == recCS4(i,z)
261                 bHR_CS4(i,j) = 1;
262             end
263             if targetSubclass(i,j) == recCS5(i,z)
264                 bHR_CS5(i,j) = 1;
265             end
266             if targetSubclass(i,j) == recCS6(i,z)
267                 bHR_CS6(i,j) = 1;
268             end
269             if targetSubclass(i,j) == recCS7(i,z)
270                 bHR_CS7(i,j) = 1;
271             end
272             if targetSubclass(i,j) == recCS8(i,z)
273                 bHR_CS8(i,j) = 1;
274             end
275             if targetSubclass(i,j) == recCS9(i,z)
276                 bHR_CS9(i,j) = 1;
277             end
278         end
279     end
280 end
281 hitCS1=0;hitCS2= 0;hitCS3=0;hitCS4=0;hitCS5=0;hitCS6=0;hitCS7=0;hitCS8=0;hitCS9=0;
282 for i = 1:11813
283     if nnz(bHR_CS1(i,:))>0
284         hitCS1 = hitCS1 + 1;
285     end
286     if nnz(bHR_CS2(i,:))>0
287         hitCS2 = hitCS2 + 1;
288     end
289     if nnz(bHR_CS3(i,:))>0
290         hitCS3 = hitCS3 + 1;
291     end
292     if nnz(bHR_CS4(i,:))>0
293         hitCS4 = hitCS4 + 1;
294     end
295     if nnz(bHR_CS5(i,:))>0
296         hitCS5 = hitCS5 + 1;
297     end

```

```

298     if nnz(bHR_CS6(i,:))>0
299         hitCS6 = hitCS6 + 1;
300     end
301     if nnz(bHR_CS7(i,:))>0
302         hitCS7 = hitCS7 + 1;
303     end
304     if nnz(bHR_CS8(i,:))>0
305         hitCS8 = hitCS8 + 1;
306     end
307     if nnz(bHR_CS9(i,:))>0
308         hitCS9 = hitCS9 + 1;
309     end
310 end
311 percentageHitCS(1) = (hitCS1/11813);
312 percentageHitCS(2) = (hitCS2/11813);
313 percentageHitCS(3) = (hitCS3/11813);
314 percentageHitCS(4) = (hitCS4/11813);
315 percentageHitCS(5) = (hitCS5/11813);
316 percentageHitCS(6) = (hitCS6/11813);
317 percentageHitCS(7) = (hitCS7/11813);
318 percentageHitCS(8) = (hitCS8/11813);
319 percentageHitCS(9) = (hitCS9/11813);

```

Listing 6: Code for weighed hit rate based on conditional probability similarity matrix for  $\alpha = 0.5$ . Same approach can be applied for  $\alpha = 0.7$  and  $\alpha = 0.9$

```

1  sumPreviousProb = zeros(11813,1);
2  sumHitCP5_1 = zeros(11813,1);sumHitCP5_2 = zeros(11813,1);
3  sumHitCP5_3 = zeros(11813,1);sumHitCP5_4 = zeros(11813,1);
4  sumHitCP5_5 = zeros(11813,1);sumHitCP5_6 = zeros(11813,1);
5  sumHitCP5_7 = zeros(11813,1);sumHitCP5_8 = zeros(11813,1);
6  sumHitCP5_9 = zeros(11813,1);
7  for i=1:11813
8      for j = 1:nnz(testBaskets(i,:))
9          %determine target and evidence
10         target = testBaskets(i,j);
11         evidence = testBaskets(i,:);
12         evidence(j) = 0;
13         evidence = sort(evidence, 'descend');
14         % %determine similarity evidence cosine based
15         cp5RowsL10_1 = zeros(1,1973);cp5RowsL10_2 = zeros(1,1973);
16         cp5RowsL10_3 = zeros(1,1973);cp5RowsL10_4 = zeros(1,1973);
17         cp5RowsL10_5 = zeros(1,1973);cp5RowsL10_6 = zeros(1,1973);
18         cp5RowsL10_7 = zeros(1,1973);cp5RowsL10_8 = zeros(1,1973);
19         cp5RowsL10_9 = zeros(1,1973);
20         for m = 1:nnz(evidence)
21             indexSubClassTest = find(evidence(m) == uniqueSubClasses2);
22             cp5RowsL10_1(1,:) = cp5RowsL10_1(1,:) + RitemCP5_1(indexSubClassTest,:);
23             cp5RowsL10_2(1,:) = cp5RowsL10_2(1,:) + RitemCP5_2(indexSubClassTest,:);
24             cp5RowsL10_3(1,:) = cp5RowsL10_3(1,:) + RitemCP5_3(indexSubClassTest,:);
25             cp5RowsL10_4(1,:) = cp5RowsL10_4(1,:) + RitemCP5_4(indexSubClassTest,:);
26             cp5RowsL10_5(1,:) = cp5RowsL10_5(1,:) + RitemCP5_5(indexSubClassTest,:);
27             cp5RowsL10_6(1,:) = cp5RowsL10_6(1,:) + RitemCP5_6(indexSubClassTest,:);

```

```

28     cp5RowsL10_7(1,:) = cp5RowsL10_7(1,:) + RitemCP5_7(indexSubClassTest,:);
29     cp5RowsL10_8(1,:) = cp5RowsL10_8(1,:) + RitemCP5_8(indexSubClassTest,:);
30     cp5RowsL10_9(1,:) = cp5RowsL10_9(1,:) + RitemCP5_9(indexSubClassTest,:);
31 end
32 [maxCP5_1, indexMaxCP5_L10_1] = sort(cp5RowsL10_1, 2, 'descend');
33 [maxCP5_2, indexMaxCP5_L10_2] = sort(cp5RowsL10_2, 2, 'descend');
34 [maxCP5_3, indexMaxCP5_L10_3] = sort(cp5RowsL10_3, 2, 'descend');
35 [maxCP5_4, indexMaxCP5_L10_4] = sort(cp5RowsL10_4, 2, 'descend');
36 [maxCP5_5, indexMaxCP5_L10_5] = sort(cp5RowsL10_5, 2, 'descend');
37 [maxCP5_6, indexMaxCP5_L10_6] = sort(cp5RowsL10_6, 2, 'descend');
38 [maxCP5_7, indexMaxCP5_L10_7] = sort(cp5RowsL10_7, 2, 'descend');
39 [maxCP5_8, indexMaxCP5_L10_8] = sort(cp5RowsL10_8, 2, 'descend');
40 [maxCP5_9, indexMaxCP5_L10_9] = sort(cp5RowsL10_9, 2, 'descend');
41 hitCP5_L10_1 = 0;hitCP5_L10_2 = 0;hitCP5_L10_3 = 0;hitCP5_L10_4 = 0;
42 hitCP5_L10_5 = 0;hitCP5_L10_6 = 0;hitCP5_L10_7 = 0;hitCP5_L10_8 = 0;
43 hitCP5_L10_9 = 0;
44 %determine recommendations cosine sim
45 k = 0;
46 a = 0;
47 b = 0;
48 c = 0;
49 d = 0;
50 e = 0;
51 f = 0;
52 g = 0;
53 h = 0;
54 for z = 1:1973
55     if k<1
56         if ismember(uniqueSubClasses2(indexMaxCP5_L10_1(1,z)),evidence)>0
57             %do nothing
58         else
59             k = k+1;
60             recCP5_L10_1 = uniqueSubClasses2(indexMaxCP5_L10_1(1,z));
61             if target == recCP5_L10_1
62                 hitCP5_L10_1= 1;
63             end
64         end
65     end
66 end
67 for z = 1:1973
68     if a < 1
69         if ismember(uniqueSubClasses2(indexMaxCP5_L10_2(1,z)),evidence)>0
70             %do nothing
71         else
72             a = a+1;
73             recCP5_L10_2 = uniqueSubClasses2(indexMaxCP5_L10_2(1,z));
74             if target == recCP5_L10_2
75                 hitCP5_L10_2 = 1;
76             end
77         end
78     end
79 end
80 for z = 1:1973

```

```

81     if b < 1
82         if ismember(uniqueSubClasses2(indexMaxCP5_L10_3(1,z)),evidence)>0
83             %do nothing
84         else
85             b = b+1;
86             recCP5_L10_3 = uniqueSubClasses2(indexMaxCP5_L10_3(1,z));
87             if target == recCP5_L10_3
88                 hitCP5_L10_3= 1;
89             end
90         end
91     end
92 end
93 for z = 1:1973
94     if c < 1
95         if ismember(uniqueSubClasses2(indexMaxCP5_L10_4(1,z)),evidence)>0
96             %do nothing
97         else
98             c = c+1;
99             recCP5_L10_4 = uniqueSubClasses2(indexMaxCP5_L10_4(1,z));
100            if target == recCP5_L10_4
101                hitCP5_L10_4= 1;
102            end
103        end
104    end
105 end
106 for z = 1:1973
107     if d < 1
108         if ismember(uniqueSubClasses2(indexMaxCP5_L10_5(1,z)),evidence)>0
109             %do nothing
110         else
111             d = d+1;
112             recCP5_L10_5 = uniqueSubClasses2(indexMaxCP5_L10_5(1,z));
113             if target == recCP5_L10_5
114                 hitCP5_L10_5= 1;
115             end
116         end
117     end
118 end
119 for z = 1:1973
120     if e < 1
121         if ismember(uniqueSubClasses2(indexMaxCP5_L10_6(1,z)),evidence)>0
122             %do nothing
123         else
124             e = e+1;
125             recCP5_L10_6 = uniqueSubClasses2(indexMaxCP5_L10_6(1,z));
126             if target == recCP5_L10_6
127                 hitCP5_L10_6= 1;
128             end
129         end
130     end
131 end
132 for z = 1:1973
133     if f < 1

```

```

134         if ismember(uniqueSubClasses2(indexMaxCP5_L10_7(1,z)),evidence)>0
135             %do nothing
136         else
137             f = f+1;
138             recCP5_L10_7 = uniqueSubClasses2(indexMaxCP5_L10_7(1,z));
139             if target == recCP5_L10_7
140                 hitCP5_L10_7= 1;
141             end
142         end
143     end
144 end
145 for z = 1:1973
146     if g < 1
147         if ismember(uniqueSubClasses2(indexMaxCP5_L10_8(1,z)),evidence)>0
148             %do nothing
149         else
150             g = g+1;
151             recCP5_L10_8 = uniqueSubClasses2(indexMaxCP5_L10_8(1,z));
152             if target == recCP5_L10_8
153                 hitCP5_L10_8= 1;
154             end
155         end
156     end
157 end
158 for z = 1:1973
159     if h < 1
160         if ismember(uniqueSubClasses2(indexMaxCP5_L10_9(1,z)),evidence)>0
161             %do nothing
162         else
163             h = h+1;
164             recCP5_L10_9 = uniqueSubClasses2(indexMaxCP5_L10_9(1,z));
165             if target == recCP5_L10_9
166                 hitCP5_L10_9= 1;
167             end
168         end
169     end
170 end
171 %previous probability target
172 indexTarget = find(target == uniqueSubClasses2);
173 sumPreviousProb(i) = sumPreviousProb(i) + (1 - popularity(indexTarget));
174 sumHitCP5_1(i) = sumHitCP5_1(i) + (1 - popularity(indexTarget))*hitCP5_L10_1
175     ;
176 sumHitCP5_2(i) = sumHitCP5_2(i) + (1 - popularity(indexTarget))*hitCP5_L10_2
177     ;
178 sumHitCP5_3(i) = sumHitCP5_3(i) + (1 - popularity(indexTarget))*hitCP5_L10_3
179     ;
180 sumHitCP5_4(i) = sumHitCP5_4(i) + (1 - popularity(indexTarget))*hitCP5_L10_4
181     ;
182 sumHitCP5_5(i) = sumHitCP5_5(i) + (1 - popularity(indexTarget))*hitCP5_L10_5
183     ;
184 sumHitCP5_6(i) = sumHitCP5_6(i) + (1 - popularity(indexTarget))*hitCP5_L10_6
185     ;
186 sumHitCP5_7(i) = sumHitCP5_7(i) + (1 - popularity(indexTarget))*hitCP5_L10_7

```



```

181         ;
182         sumHitCP5_8(i) = sumHitCP5_8(i) + (1 - popularity(indexTarget))*hitCP5_L10_8
183         ;
184         sumHitCP5_9(i) = sumHitCP5_9(i) + (1 - popularity(indexTarget))*hitCP5_L10_9
185         ;
186     end
187     %determine weighted hit rate basket
188     wHR_L10_CP5_1(i) = sumHitCP5_1(i)/sumPreviousProb(i);
189     wHR_L10_CP5_2(i) = sumHitCP5_2(i)/sumPreviousProb(i);
190     wHR_L10_CP5_3(i) = sumHitCP5_3(i)/sumPreviousProb(i);
191     wHR_L10_CP5_4(i) = sumHitCP5_4(i)/sumPreviousProb(i);
192     wHR_L10_CP5_5(i) = sumHitCP5_5(i)/sumPreviousProb(i);
193     wHR_L10_CP5_6(i) = sumHitCP5_6(i)/sumPreviousProb(i);
194     wHR_L10_CP5_7(i) = sumHitCP5_7(i)/sumPreviousProb(i);
195     wHR_L10_CP5_8(i) = sumHitCP5_8(i)/sumPreviousProb(i);
196     wHR_L10_CP5_9(i) = sumHitCP5_9(i)/sumPreviousProb(i);
197 end
198 perCP5(1) = sum(wHR_L10_CP5_1)/11813;
199 perCP5(2) = sum(wHR_L10_CP5_2)/11813;
200 perCP5(3) = sum(wHR_L10_CP5_3)/11813;
201 perCP5(4) = sum(wHR_L10_CP5_4)/11813;
202 perCP5(5) = sum(wHR_L10_CP5_5)/11813;
203 perCP5(6) = sum(wHR_L10_CP5_6)/11813;
204 perCP5(7) = sum(wHR_L10_CP5_7)/11813;
205 perCP5(8) = sum(wHR_L10_CP5_8)/11813;
206 perCP5(9) = sum(wHR_L10_CP5_9)/11813;

```

Listing 7: Code for macro-averaged hit rate based on the bipartite network transition probability matrix for  $\alpha = 0.5$ . The same approach can be applied for  $\alpha = 0.7$  and  $\alpha = 0.9$

```

1 hitBN5_L10_1 = [];hitBN5_L10_2 = [];hitBN5_L10_3 = [];hitBN5_L10_4 = [];
2 hitBN5_L10_5 = [];hitBN5_L10_6 = [];hitBN5_L10_7 = [];hitBN5_L10_8 = [];
3 hitBN5_L10_9 = [];
4 for i=1:11813
5     for j = 1:nnz(testBaskets(i,:))
6         %determine target and evidence
7         target = testBaskets(i,j);
8         evidence = testBaskets(i,:);
9         evidence(j) = 0;
10        evidence = sort(evidence, 'descend');
11        %determine similarity evidence cosine based
12        bn5RowsL10_1 = zeros(1,1973);bn5RowsL10_2 = zeros(1,1973);
13        bn5RowsL10_3 = zeros(1,1973);bn5RowsL10_4 = zeros(1,1973);
14        bn5RowsL10_5 = zeros(1,1973);bn5RowsL10_6 = zeros(1,1973);
15        bn5RowsL10_7 = zeros(1,1973);bn5RowsL10_8 = zeros(1,1973);
16        bn5RowsL10_9 = zeros(1,1973);
17        for m = 1:nnz(evidence)
18            indexSubClassTest = find(evidence(m) == uniqueSubClasses2);
19            bn5RowsL10_1(1,:) = bn5RowsL10_1(1,:) + RitemBN5_1(indexSubClassTest,:);
20            bn5RowsL10_2(1,:) = bn5RowsL10_2(1,:) + RitemBN5_2(indexSubClassTest,:);
21            bn5RowsL10_3(1,:) = bn5RowsL10_3(1,:) + RitemBN5_3(indexSubClassTest,:);
22            bn5RowsL10_4(1,:) = bn5RowsL10_4(1,:) + RitemBN5_4(indexSubClassTest,:);
23            bn5RowsL10_5(1,:) = bn5RowsL10_5(1,:) + RitemBN5_5(indexSubClassTest,:);

```

```

24     bn5RowsL10_6(1,:) = bn5RowsL10_6(1,:) + RitemBN5_6(indexSubClassTest,:);
25     bn5RowsL10_7(1,:) = bn5RowsL10_7(1,:) + RitemBN5_7(indexSubClassTest,:);
26     bn5RowsL10_8(1,:) = bn5RowsL10_8(1,:) + RitemBN5_8(indexSubClassTest,:);
27     bn5RowsL10_9(1,:) = bn5RowsL10_9(1,:) + RitemBN5_9(indexSubClassTest,:);
28 end
29 [maxBN5_1, indexMaxBN5_L10_1] = sort(bn5RowsL10_1, 2, 'descend');
30 [maxBN5_2, indexMaxBN5_L10_2] = sort(bn5RowsL10_2, 2, 'descend');
31 [maxBN5_3, indexMaxBN5_L10_3] = sort(bn5RowsL10_3, 2, 'descend');
32 [maxBN5_4, indexMaxBN5_L10_4] = sort(bn5RowsL10_4, 2, 'descend');
33 [maxBN5_5, indexMaxBN5_L10_5] = sort(bn5RowsL10_5, 2, 'descend');
34 [maxBN5_6, indexMaxBN5_L10_6] = sort(bn5RowsL10_6, 2, 'descend');
35 [maxBN5_7, indexMaxBN5_L10_7] = sort(bn5RowsL10_7, 2, 'descend');
36 [maxBN5_8, indexMaxBN5_L10_8] = sort(bn5RowsL10_8, 2, 'descend');
37 [maxBN5_9, indexMaxBN5_L10_9] = sort(bn5RowsL10_9, 2, 'descend');
38 %determine recommendations cosine sim
39 k = 0;
40 a = 0;
41 b = 0;
42 c = 0;
43 d = 0;
44 e = 0;
45 f = 0;
46 g = 0;
47 h = 0;
48 for z = 1:1973
49     if k<1
50         if ismember(uniqueSubClasses2(indexMaxBN5_L10_1(1,z)),evidence)>0
51             %do nothing
52         else
53             k = k+1;
54             recBN5_L10_1 = uniqueSubClasses2(indexMaxBN5_L10_1(1,z));
55             if target == recBN5_L10_1
56                 hitBN5_L10_1 = [hitBN5_L10_1; 1];
57             else
58                 hitBN5_L10_1 = [hitBN5_L10_1; 0];
59             end
60         end
61     end
62 end
63 for z = 1:1973
64     if a < 1
65         if ismember(uniqueSubClasses2(indexMaxBN5_L10_2(1,z)),evidence)>0
66             %do nothing
67         else
68             a = a+1;
69             recBN5_L10_2 = uniqueSubClasses2(indexMaxBN5_L10_2(1,z));
70             if target == recBN5_L10_2
71                 hitBN5_L10_2 = [hitBN5_L10_2; 1];
72             else
73                 hitBN5_L10_2 = [hitBN5_L10_2; 0];
74             end
75         end
76     end
end

```

```

77     end
78     for z = 1:1973
79         if b < 1
80             if ismember(uniqueSubClasses2(indexMaxBN5_L10_3(1,z)),evidence)>0
81                 %do nothing
82             else
83                 b = b+1;
84                 recBN5_L10_3 = uniqueSubClasses2(indexMaxBN5_L10_3(1,z));
85                 if target == recBN5_L10_3
86                     hitBN5_L10_3 = [hitBN5_L10_3; 1];
87                 else
88                     hitBN5_L10_3 = [hitBN5_L10_3; 0];
89                 end
90             end
91         end
92     end
93     for z = 1:1973
94         if c < 1
95             if ismember(uniqueSubClasses2(indexMaxBN5_L10_4(1,z)),evidence)>0
96                 %do nothing
97             else
98                 c = c+1;
99                 recBN5_L10_4 = uniqueSubClasses2(indexMaxBN5_L10_4(1,z));
100                if target == recBN5_L10_4
101                    hitBN5_L10_4 = [hitBN5_L10_4; 1];
102                else
103                    hitBN5_L10_4 = [hitBN5_L10_4; 0];
104                end
105            end
106        end
107    end
108    for z = 1:1973
109        if d < 1
110            if ismember(uniqueSubClasses2(indexMaxBN5_L10_5(1,z)),evidence)>0
111                %do nothing
112            else
113                d = d+1;
114                recBN5_L10_5 = uniqueSubClasses2(indexMaxBN5_L10_5(1,z));
115                if target == recBN5_L10_5
116                    hitBN5_L10_5 = [hitBN5_L10_5; 1];
117                else
118                    hitBN5_L10_5 = [hitBN5_L10_5; 0];
119                end
120            end
121        end
122    end
123    for z = 1:1973
124        if e < 1
125            if ismember(uniqueSubClasses2(indexMaxBN5_L10_6(1,z)),evidence)>0
126                %do nothing
127            else
128                e = e+1;
129                recBN5_L10_6 = uniqueSubClasses2(indexMaxBN5_L10_6(1,z));

```

```

130         if target == recBN5_L10_6
131             hitBN5_L10_6 = [hitBN5_L10_6; 1];
132         else
133             hitBN5_L10_6 = [hitBN5_L10_6; 0];
134         end
135     end
136 end
137 end
138 for z = 1:1973
139     if f < 1
140         if ismember(uniqueSubClasses2(indexMaxBN5_L10_7(1,z)),evidence)>0
141             %do nothing
142         else
143             f = f+1;
144             recBN5_L10_7 = uniqueSubClasses2(indexMaxBN5_L10_7(1,z));
145             if target == recBN5_L10_7
146                 hitBN5_L10_7 = [hitBN5_L10_7; 1];
147             else
148                 hitBN5_L10_7 = [hitBN5_L10_7; 0];
149             end
150         end
151     end
152 end
153 for z = 1:1973
154     if g < 1
155         if ismember(uniqueSubClasses2(indexMaxBN5_L10_8(1,z)),evidence)>0
156             %do nothing
157         else
158             g = g+1;
159             recBN5_L10_8 = uniqueSubClasses2(indexMaxBN5_L10_8(1,z));
160             if target == recBN5_L10_8
161                 hitBN5_L10_8 = [hitBN5_L10_8; 1];
162             else
163                 hitBN5_L10_8 = [hitBN5_L10_8; 0];
164             end
165         end
166     end
167 end
168 for z = 1:1973
169     if h < 1
170         if ismember(uniqueSubClasses2(indexMaxBN5_L10_9(1,z)),evidence)>0
171             %do nothing
172         else
173             h = h+1;
174             recBN5_L10_9 = uniqueSubClasses2(indexMaxBN5_L10_9(1,z));
175             if target == recBN5_L10_9
176                 hitBN5_L10_9 = [hitBN5_L10_9; 1];
177             else
178                 hitBN5_L10_9 = [hitBN5_L10_9; 0];
179             end
180         end
181     end
182 end

```

```

183     end
184 end
185 perBN5(1) = sum(hitBN5_L10_1)/size(hitBN5_L10_1,1);
186 perBN5(2) = sum(hitBN5_L10_2)/size(hitBN5_L10_2,1);
187 perBN5(3) = sum(hitBN5_L10_3)/size(hitBN5_L10_3,1);
188 perBN5(4) = sum(hitBN5_L10_4)/size(hitBN5_L10_4,1);
189 perBN5(5) = sum(hitBN5_L10_5)/size(hitBN5_L10_5,1);
190 perBN5(6) = sum(hitBN5_L10_6)/size(hitBN5_L10_6,1);
191 perBN5(7) = sum(hitBN5_L10_7)/size(hitBN5_L10_7,1);
192 perBN5(8) = sum(hitBN5_L10_8)/size(hitBN5_L10_8,1);
193 perBN5(9) = sum(hitBN5_L10_9)/size(hitBN5_L10_9,1);

```

Listing 8: Code for user-based Collaborative Filtering model

```

1  %normalize user user sim matrix
2  sumUser = sum(userCosineSim, 2);
3  for i=1:20388
4      for j = 1:20388
5          normUserCosineSim(i,j) = userCosineSim(i,j)/sumUser(i);
6      end
7  end
8
9  %determine mean of the ratings
10 mean_ratings = zeros(20388,1);
11 for i = 1:20388
12     mean_ratings(i,1) = sum(normUserItemMatrix(i,:))/nnz(normUserItemMatrix(i,:));
13 end
14
15 %determine mean ratings test customers/active user
16 mean_ratings_test = mean_ratings(indexCustomersTestSet);
17 tic
18 %determine k most similar customers
19 ratings = zeros(11813, 1973);
20 for i = 1:11813
21     indexCustomer = indexCustomersTestSet(i);
22
23     [maxUCS, indexMaxUserCS] = sort(userCosineSim(indexCustomer,:), 'descend');
24
25     for j = 1:1973
26         sumRatings = 0;
27         for k = 1:3 %nog bepalen tm hoe veel
28             sumRatings = sumRatings + normUserCosineSim(indexCustomer,
29                 indexMaxUserCS(k)) * (normUserItemMatrix(indexMaxUserCS(k),j) -
30                 mean_ratings(indexMaxUserCS(k)));
29         end
30
31         ratings(i,j) = mean_ratings_test(i) + sumRatings;
32     end
33 end
34 toc
35 %POPULARITY BASED HIT RATE
36 %determine targets test basket
37 for i = 1:11813
38     for j = 1:nnz(testBaskets(i,:))

```

```

39     indexSubClassTest = find(testBaskets(i,j) == uniqueSubClasses2);
40     freqSubclass(i,j) = popularity(indexSubClassTest);
41 end
42 end
43 freqSubclass( freqSubclass == 0 ) = Inf;
44 [sortedFreqSubclass, indexFreqSubclass] = sort(freqSubclass,2,'ascend');
45 freqSubclass( isinf(freqSubclass) ) = 0;
46 sortedFreqSubclass( isinf(sortedFreqSubclass) ) = 0;
47 for i = 1:11813
48     for j = 1:3
49         targetSubclass(i,j) = testBaskets(i,indexFreqSubclass(i,j));
50     end
51     for j = 4:nnz(testBaskets(i,:))
52         evidenceSubclass(i,j) = testBaskets(i,indexFreqSubclass(i,j));
53     end
54 end
55 %determine recommendations
56 [maxRating, indexMaxRating] = sort(ratings, 2, 'descend');
57 for i = 1:11813
58     k = 0;
59     for j = 4:1973
60         if k<3
61             if ismember(uniqueSubClasses2(indexMaxRating(i,j-3)),evidenceSubclass(i
62                 ,:)) > 0
63                 %do nothing
64             else
65                 k = k+1;
66                 recExtension_POP(i,k) = uniqueSubClasses2(indexMaxRating(i,j-3));
67             end
68         end
69     end
70 %determine hit rate similarities
71 binaryHitExtension = zeros(11813,3);
72
73 for i = 1:11813
74     for j = 1:3
75         for k= 1:3
76             if targetSubclass(i,j) == recExtension_POP(i,k)
77                 binaryHitExtension(i,j) = 1;
78             end
79         end
80     end
81 end
82 hitPopExtension = 0;
83 for i = 1:11813
84     if nnz(binaryHitExtension(i,:))>0
85         hitPopExtension = hitPopExtension + 1;
86     end
87 end
88 percentageHitPop = hitPopExtension/11813;
89 %RANDOM BASED HIT RATE
90 %determine targets test basket

```

```

91 for i = 1:11813
92     randomIndices = randperm(nnz(testBaskets(i,:)),3);
93     for j = 1:3
94         targetSubclassRandom(i,j) = testBaskets(i,randomIndices(j));
95     end
96     evidenceSubclassRandom(i,:) = testBaskets(i,:);
97     evidenceSubclassRandom(i,randomIndices) = 0;
98 end
99 evidenceSubclassRandom = sort(evidenceSubclassRandom,2, 'descend');
100 %determine recommendations
101 recExtension_RND = [];
102 [maxRatingRND, indexMaxRating_RND] = sort(ratings, 2, 'descend');
103 for i = 1:11813
104     k = 0;
105     for j = 1:1973
106         if k<3
107             if ismember(uniqueSubClasses2(indexMaxRating_RND(i,j)),
108                 evidenceSubclassRandom(i,:)) > 0
109                 %do nothing
110             else
111                 k = k+1;
112                 recExtension_RND(i,k) = uniqueSubClasses2(indexMaxRating_RND(i,j));
113             end
114         end
115     end
116 %determine hit rate similarities
117 binaryHitRandom_EX = zeros(11813,3);
118
119 for i = 1:11813
120     for j = 1:3
121         for k= 1:3
122             if targetSubclassRandom(i,j) == recExtension_RND(i,k)
123                 binaryHitRandom_EX(i,j) = 1;
124             end
125         end
126     end
127 end
128 hitRandom_EX = 0;
129 for i = 1:11813
130     if nnz(binaryHitRandom_EX(i,:))>0
131         hitRandom_EX = hitRandom_EX + 1;
132     end
133 end
134 percentageHitRnd = hitRandom_EX/11813;
135 %WEIGHED HIT RATE
136 hitExtension = [];
137 [maxRating_loo, indexMaxRating_loo] = sort(ratings, 2, 'descend');
138 for i = 1:11813
139     for j = 1:nnz(testBaskets(i,:))
140         %determine target and evidence
141         target = testBaskets(i,j);
142         evidence = testBaskets(i,:);

```

```

143     evidence(j) = 0;
144     evidence = sort(evidence, 'descend');
145     %determine recommendation
146     k = 0;
147     for m = 1:1973
148         if k<1
149             if ismember(uniqueSubClasses2(indexMaxRating_loo(i,m)),evidence)>0
150                 %do nothing
151             else
152                 k = k+1;
153                 recEx_loo = uniqueSubClasses2(indexMaxRating_loo(i,m));
154                 if target == recEx_loo
155                     hitExtension = [hitExtension; 1];
156                 else
157                     hitExtension = [hitExtension; 0];
158                 end
159             end
160         end
161     end
162 end
163 %determine total amount of hits
164 percentageHitEx = sum(hitExtension)/size(hitExtension,1);
165

```

Listing 9: Code to generate plots of the performance of the metrics

```

1 d = [0.1;0.2;0.3;0.4;0.5;0.6;0.7;0.8;0.9];
2 %plot basket-sensitive random walk based on bipartite network
3 plot(d,bsrw_bn_5,'-s','Color','magenta','Markersize',10,'LineWidth',2)
4 hold on
5 plot(d,bsrw_bn_7,'-x','Color','cyan','Markersize',10,'LineWidth',2)
6 hold on
7 plot(d,bsrw_bn_9,'-o','Color','red','Markersize',10,'LineWidth',2)
8 hold off
9 grid on
10 x = xlabel('(1 - d)','fontweight','bold');
11 x.FontSize = 14;
12 y = ylabel('bHR(pop)','fontweight','bold');
13 y.FontSize = 14;
14 lgd = legend({'\alpha = 0.5','\alpha = 0.7', '\alpha = 0.9'},'Location','northwest')
15 ;
16 lgd.FontSize = 15;
17 %plot macro-averaged hit rate
18 plot(d,macroBN5,'-s','Color','magenta','Markersize',10,'LineWidth',2)
19 hold on
20 plot(d,macroBN7,'-X','Color','cyan','Markersize',10,'LineWidth',2)
21 hold on
22 plot(d,macroBN9,'-o','Color','red','Markersize',10,'LineWidth',2)
23 hold off
24 grid on
25 x = xlabel('(1 - d)','fontweight','bold');
26 x.FontSize = 14;
27 y = ylabel('macroHR(loo)','fontweight','bold');
28 y.FontSize = 14;

```



```
28 lgd = legend({'\alpha = 0.5', '\alpha = 0.7', '\alpha = 0.9'}, 'Location', 'northwest')  
    ;  
29 lgd.FontSize = 15;
```