## Research Article
## Improved Fuzzy C-Means Clustering for Personalized Product Recommendation

[1,2]Juebo Wu and [3]Zongling Wu
[1]Shenzhen Angelshine Co., Ltd., Shenzhen, China
[2]Department of Geography, National University of Singapore, 1 Arts Link, 117570, Singapore
[3]International School of Software, Wuhan University, Wuhan 430079, China

**Abstract:** With rapid development of e-commerce, how to better understand users' needs to provide more satisfying personalized services has become a crucial issue. To overcome the problem, this study presents a novel approach for personalized product recommendation based on Fuzzy C-Means (FCM) clustering. Firstly, the traditional FCM clustering algorithm is improved by membership adjustment and density function, in order to address the issues that the number of clusters is difficult to determine and the convergence of objective function is slow. Then, the personal preferences are divided into different groups, one of which the users have the similar tendencies in. The association rules of user preferences are mined for each group and the personalized knowledge base is established. After that, the recommendation can be generated by knowledge base and historical logs. A case study is illustrated by the proposed approach and the results show that the method of personalized product recommendation is reasonable and efficient with high performance.

**Keywords:** E-commerce, fuzzy clustering, fuzzy c-means, personalized product recommendation

## INTRODUCTION

In recent years, e-commerce has made continuous progress and improvement, resulting in commodity diversification and intense competition. Enhancing the personalization of e-commerce system is the guarantee of business to win (Qian, 2011). It is able to prevent the loss of customers and improve the sales of e-commerce systems (Huang and Duan, 2012). Product recommendation system is an effective means to improve the personalized preferences in e-commerce, which can make every customer feel that the website is defined for his own. At the same time, it also has the ability to assist customers to decide what products to buy according to the purchasing logs. The personalized service is increasingly important (Shen, 2011).

There are a great number of technologies used for recommendation systems in e-commerce, including collaborative filtering (Liu *et al*., 2009), content-based filtering (Chu and Park, 2009), data mining and knowledge discovery (Wang, 2012), multiple technologies combination (Ngai *et al*., 2009), interactive recommendation (Alon *et al*., 2009) and other technologies based on statistics (Shishehchi *et al*., 2011), Bayesian networks (Yi and Deng, 2009) and neural networks (Chou *et al*., 2010) etc. However, these technologies have still some shortcomings of their own. For example, collaborative filtering has the problem that it takes a long time for computing when there are a large number of products and users in the system. Although the dimensionality reduction method can be adopted, the recommendation quality may be decreased. For the recommendation system of content filtering, it's easy for the provider to change the product characteristics in order to change the products recommendation on purpose. Besides, the attributes of products are hard to be extracted. Interactive recommendation requires that the knowledge has to be provided by experts in different fields and the user requirements are inconsistent.

In order to improve the accuracy and performance of recommendation system, we present a novel approach for personalized product recommendation by using improved FCM clustering according to users' preferences in this study. Our idea is firstly to divide all the users into different classes and then generate recommendations for users by mining user preferences from their similar groups. Since FCM clustering algorithm cannot be applied into recommendation system directly, we modify the membership function so as to avoid the results deviation caused by membership normalization. Meanwhile, we create the weighted objective function by setting a quantitative value for each sample point, with the aim of meeting the demands of huge data in e-commerce system. The process of the improved FCM for personalized product recommendation is given and a case study demonstrates the feasibility and effectiveness of the presented approach at last.
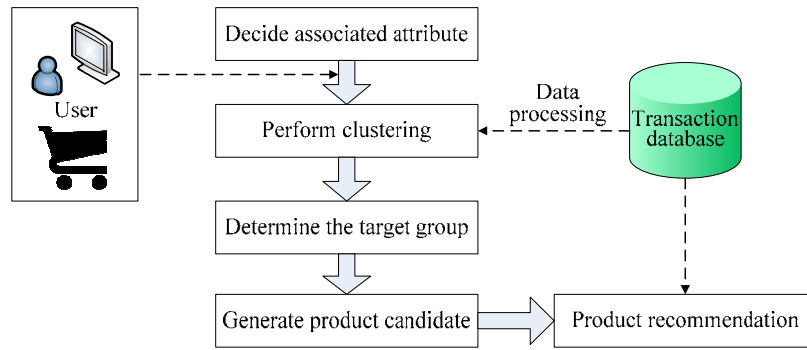
Fig. 1: Model of personalized product recommendation based on clustering algorithm

## MOTIVATION AND OUR APPROACH

Personalized product recommendation is a decision support system for assisting the user on purchasing products according to some strategy based on the characteristics and needs of individual consumer. It is to provide a personalized shopping experience for consumers, in order to solve the problem of information overload. It can help consumers find out their needs in a broad array of products.

The goals of personalized product recommendation in e-commerce are as follows:

- To turn browsers into buyers on e-commerce websites
- To improve the cross-selling capabilities
- To reduce consumer costs (time, money, etc.) to meet customer needs and increase their satisfactions
- To enhance the amount of sales and thus to improve the income of the seller

Cluster analysis is an effective way to personalized product recommendation. By clustering, the similar browsing or purchasing behaviors of customers can be identified. Common characteristics of customers can be analyzed and then the e-commerce companies are able to understand their customers better. In this way, companies can provide customers with a more appropriate and more comprehensive services to facilitate the development and implementation of the future market strategy.

The cluster analysis can be divided into the clustering for customer groups and the clustering for web pages. The clustering for customer groups plays an important role in the application of personalized service. Its aim is to extract common characteristics for each clustering group based on the historical behaviors of customers. By means of this model, the system can recommend the interest products for customers and track the recommendation effect for future use. Through clustering, customers who are close to each other within a certain range are partitioned into a group. A number of different classes are produced after clustering.

The customers with 1 group have high similarity with each other while the customers in different groups have low similarity. In many applications, the data objects in a cluster can be treated as a whole. By clustering, dense and sparse areas can be identified. Therefore, the global distribution patterns can be found including some interesting correlations among data attributes.

According to the analysis above, we present a model of personalized product recommendation based on clustering algorithm as shown in Fig. 1, which has 3 core steps:

- Decide the associated attributes between customer and product. The information can be customer buying history or reviews
- Perform cluster analysis. Two ways can be used in this step that is, clustering by customers and clustering by product items
- Do collaborative filtering and product recommendation

**FCM clustering algorithm:** Fuzzy clustering (Xie and Beni, 1991), as 1 of the unsupervised machine learning technologies, is an important method in data analysis and modeling in fuzzy theory. It can create uncertainty description for samples and reflect the real world objectively. On the basis of general classification, Bezdek (1984) improved the concept of fuzzy classification, where samples in 1 class are partly belonged to other classes by membership. A class is considered to be a fuzzy subset of sample collection. Each classification result is corresponding to a matrix, that is, fuzzy matrix (Bezdek, 1984). The FCM clustering algorithm is described as follows:

Assume that the sample collection is $X = \{x_1, x_2, \dots, x_n\}$, the goal is to divide it into c group with cluster center $c_j$ ($j = 1, 2, \dots, c$) and minimize the objective function. The objective function is defined as:

$$J_C = \sum_{j=1}^{C} \sum_{i=1}^{n} \mu_{ij}^{\alpha} \left\| x_i - c_j \right\|^2, \ 1 \leq \alpha \leq \infty \qquad (1)$$

And the objective function meets:

$$\sum_{j=1}^{C} \mu_{ij} = 1, \ \forall \ i = 1,2,...,n \qquad (2)$$

where, $\mu_{ij} \in [0,1]$ is the degree of membership of the *ith* data belonging to the *jth* cluster center. $c_j$ is the *jth* cluster center and its initial value is selected randomly. $\alpha$ is any real number greater than 1, which controls the fuzzy degree. Fuzzy clustering is executed by an iterative optimization of the objective function given above, with the update of membership $\mu_{ij}$ and the cluster centers $c_j$ by:

$$\mu_{ij} = \frac{1}{\sum_{k=1}^{C} \left( \frac{\|x_i - c_i\|^2}{\|x_i - c_k\|^2} \right)^{2(\alpha-1)}} \qquad (3)$$

$$c_j = \frac{\sum_{i=1}^{n} \mu_{ij}^{\alpha} x_i}{\sum_{i=1}^{n} \mu_{ij}^{\alpha}} \qquad (4)$$

This procedure starts from a random cluster center and then adjusts the cluster center and the degree of membership for each sample. The goal is to be converged to a saddle point of $J_C$ or a local minimum.

**Improved FCM clustering algorithm:**
**Cluster number determination:** Firstly, density function algorithm with identified cluster number is adopted to select initial cluster centers in the area of high density and the farthest away from sample points. Then, the definitions of degrees of coupling and separation are given by means of distance metric and fuzzy iteration is done for cluster number c. Finally, the cluster result evaluation for different c values is carried out by effectiveness criterion function. The sample points with highest similarity and lowest similarity between classes are chosen as the cluster centers (Dongbo *et al.*, 2002).

**Definition 1**: The density of sample point $x_i$ is defined as:

$$D_i = \sum_{k=1}^{n} \frac{1}{1 + f_d \|x_i - x_k\|^2} \quad i = 1,2,...,n \qquad (5)$$

In Eq. (5), $f_d$ is defined as:

$$f_d = \frac{4}{r_d^2} \qquad (6)$$

where, $r_d$ is the effective radius of the field density, defined as:

$$r_d = \frac{1}{2} \times \sqrt{\frac{1}{(n-1)} \sum_{k=1}^{n} \sum_{i=1}^{n} \|x_i - x_k\|^2} \qquad (7)$$

The greater $D_i$ is, the bigger the density of $x_i$ will be.

After generating the density of each sample point, the sample set $D$ in high density area is obtained by deleting the points in low density area. The point $v_1$ with highest density is selected as the 1st cluster center. And the 2nd point $v_2$ that is farthest away from point $v_1$ is chosen as the second cluster center. Then, choose a point from the rest sample points in $D$ and compute the distances from $v_1$ and $v_2$ respectively. The point with biggest value is selected as the 3rd cluster center. According to this rule, all the cluster centers can be produced.

The clustering results with different granularity are corresponding to the different divisions of the sample point set. We utilize the coupling degree and the separation degree to measure the quality of the results. The coupling degree stands for the tightness of the class while separation degree indicates the non-similarity among classes. The results will be good if the coupling degree is high and the separation degree is low.

**Definition 2:** The coupling degree of sample points is defined as:

$$Cd(c) = \frac{1}{n} \sum_{i=1}^{c} \sum_{j=1}^{n} \mu_{ij}^{m} d_{ij}^2 \quad i = 1,2,...,c, \ j = 1,2,...,n \qquad (8)$$

**Definition 3:** The separation degree of sample points is defined as:

$$Sd(c) = \frac{\sum_{i,k=1;i \neq k}^{c} d_{ik}^2}{[c(c-1)]/2} \quad i,k = 1,2,...,c \qquad (9)$$

**Definition 4:** The effectiveness criterion function of results is defined as:

$$GD(c) = \alpha Cd(b) + \beta - \frac{1}{Sd(b)} \quad \alpha + \beta = 1 \qquad (10)$$

where, $\alpha$ and $\beta$ stand for the weights of the coupling degree and the separation degree. If the values of the coupling degree and the separation degree fluctuate in a big range, the bigger one should be set a smaller weight for better results. Usually, the value of coupling degree is bigger than the value of separation degree. In this study, we set $\alpha = 0.6$ and $\beta = 0.4$.

Therefore, the smaller value of the GD(c) represents a better clustering result and the minimum value of GD(c) corresponding to the value of $C$ is the optimal number of categories. A small GD(c) denotes a

good cluster results and the minimum value is corresponding to the optimal number of clusters.

**The improved fuzzy membership:** The normalization condition defined in Eq. (2) may lead to bad results when the sample set is not an ideal situation. For example, if a sample point is far away from the center of various types of clustering, the memberships it is strictly belonged to are very small. However, due to the requirement of normalization, the memberships it is strictly belonged to are very high, which have impact on the results. In order to reduce restrictions, the sum of the membership in various clusters is defined as n, that is:

$$\sum_{j=1}^{C}\sum_{i=1}^{n}\mu_{ij} = n, \ \forall \ i = 1,2,...,n \tag{11}$$

Under this circumstance, the function of membership degree in Eq. (3) is redefined as:

$$\mu_{ij} = \frac{n\left(\frac{1}{\|x_i - c_k\|^2}\right)^{\frac{1}{\alpha-1}}}{\sum_{k=1}^{C}\sum_{l=1}^{n}\left(\frac{1}{\|x_i - c_k\|^2}\right)^{\frac{1}{\alpha-1}}}, i = 1,2,...,n, \ j = 1,2,...,C \tag{12}$$

The value of membership is not restricted within (0, 1). It can be normalized according to the actual situation as follows:

$$\mu_{ij} = \frac{\mu_{ij}}{\sum_{j=1}^{C}\mu_{ij}} \tag{13}$$

The improved degree of membership can improve the clustering results of FCM algorithm in the case of the presence of isolated points, so that the final clustering results are not very sensitive to the pre-determined number of clusters.

**Weighted sample points:** In e-commerce, the importance of each product for different customers is not the same. By assigning different weights, it can reflect the degree of importance of the products to different customers. Some sample points are more important to the classification while other sample points are less important. For example, the noise point is not important in the sample points. Therefore, we can give a weight to different sample points to distinguish the different importance, which defined as:

$$w_i = \frac{\sum_{i=1}^{n}w_i'}{w_i'} \tag{14}$$

$$w_i' = \frac{1}{n}\sum_{j=1}^{n}s(x_i, x_j) \tag{15}$$

where, $s(x_i, x_j)$ is the dissimilarity between $x_i$ and $x_j$. The value of dissimilarity is close to 0 if the distance or similarity between $x_i$ and $x_j$ is in close proximity and vice verse.

For more intensive data points, their weights are relatively close when their distances are similar away from the center. The influence can be eliminated by setting a smaller weight to the noise and the outlier points.

**Steps of the improved FCM clustering algorithm:** The main idea of the improved algorithm is as below. Firstly, determine the initial cluster number. Secondly, use the improved membership function for cluster analysis. Thirdly, compute the minimum value of weighted objective function by combining membership function of joint update. The objective function of the improved algorithm is defined as:

$$J_C = \sum_{j=1}^{C} \ \sum_{i=1}^{n} w_i \mu_{ij}^{\alpha} \|x_i - c_j\|^2, \ 1 \le \alpha \le \infty \tag{16}$$

And the cluster centers $c_j$ is computed by:

$$c_j = \frac{\sum_{i=1}^{n}w_i\mu_{ij}^{\alpha}x_i}{\sum_{i=1}^{n}w_i\mu_{ij}^{\alpha}} \tag{17}$$

The key steps of the improved FCM clustering algorithm are described as follows:

- Compute the cluster number of the sample points by Eq. (5)
- Compute the fuzzy degree of the membership for $x_i$ in cluster *j* by Eq. (12)
- Update the fuzzy degree of the membership $\mu_{ij}(t)$ as $\mu_{ij}(t+1)$
- Compute objective function by Eq. 16. If the change value of this time over last time is smaller than a threshold ε, continue next step. Else, go to step 3
- Normalize the degree of membership by Eq. (13)

In this algorithm, the threshold ε is set in advance, which should meet the required accuracy.

**Personalized product recommendation with improved FCM clustering:**
**Association rules mining:** Association rules mining is used to find out the potential relation and useful knowledge in the data. After clustering, we use

association rules mining to extract the recommendation information from the user historical logs. In order to eliminate the difference between the products, diverse weights are set to different kinds of products, defined as:

$$user = \{(product1, w1), (product2, w2), \dots\} \quad (18)$$

The weight of two products is defined as:

$$sim(i, j) = \frac{\sum_{k=1}^{m} R_{ik} \times R_{jk}}{\sqrt{\sum_{k=1}^{m} R_{ik}^2 \times \sum_{k=1}^{m} X_{jk}^2}} \quad (19)$$

where, $R_{ik}$ and $R_{jk}$ are the weights of the user i and user j for product k, respectively.

Then, the cluster analysis is carried out by the improved FCM described in order to divide the products for different groups of users.

Assume a cluster as $I = \{i_1, i_2, \dots, i_m\}$ and $i_j (0 \leq j \leq m)$ means the *jth* item in *I*. Given the associated data *D* is a set of the data and the transaction *T* is a group of cluster where $T \subseteq I$. Each transaction has a unique identifier, that is, *TID*. Given *A* is an itemset and $A \subseteq T$ if and only if *T* contains *A*. An association rule is a form of implication $A \Rightarrow B$, where $A \subseteq I$, $B \subseteq I$ and $B = \Phi$.

The condition of the association rule $A \Rightarrow B$ is as follows:

- It has the support with S%, which is the percentage of $A \cup B$ contained in transaction set *D*. It means there is S% of records with $A \cup B$ at least in data sets, defined as:

$$S\% = \sup port(A \Rightarrow B) = P(A \cup B) \quad (20)$$

- It has confidence with C%, which is the percentage of transaction set *D* including both transaction *A* and transaction *B*. It means there is C% of records with A∩B at least in data sets, defined as:

$$C\% = confidence \ (A \Rightarrow B) = P(B \mid A) = P(A \cup B)/P(A) \quad (21)$$

The support defines the item in the proportion of entire data and the confidence defines the strength of the rule. The set that meets the min support degree is called frequent itemset. The relation that meets min support and confidence is regarded as strong association rule, defined as:

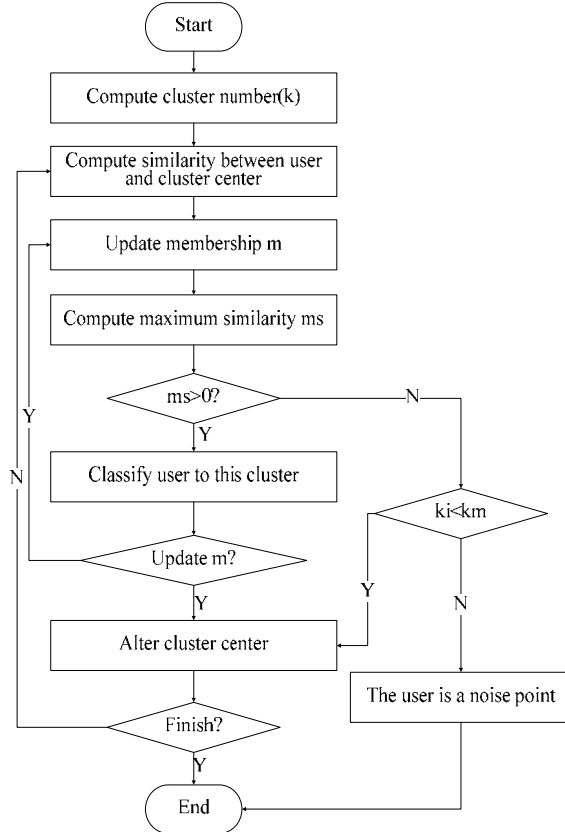$$A \Rightarrow B(S\%, C\%) \quad (22)$$



Fig. 2: Process of the presented product recommendation

**The procedure:** In personalized product recommendation, the application of association rules is to find valuable rules between users in the purchase of products. If 2 products or more products meet a certain degree of confidence and support, then they have association rules among them. The system can recommend some products to the users in order to assist them for decision making when they browsing the products. The products list that the user is interested can also be generated by establishing the knowledge base of association rules.

However, the association rules cannot be obtained only by the support and credibility from data directly, which has a high degree of support and credibility of the rules in 1 cluster. Therefore, as described in Fig. 1, we 1st classify the data into different groups of users by the improved FCM Clustering Algorithm. After that, we compute the association rules for users. And the personalized knowledge base is able to be created to provide a better set of rules for the recommendation system.

According to the characteristics of product data in e-commerce, Fig. 2 gives the detail of the algorithm for personalized product recommendation by the improved FCM.

The approach has the ability of simplicity and high rapidity and the impact of the noise data can be reduced by using automatic cluster number determination.

Table 1: Records of users' historical purchase

| User name | Book number |
|---|---|
| User1 | book1, book2, book9 |
| User2 | book3, book4, book5 |
| User3 | book2, book5, book8 |
| User4 | book1, book7, book9 |
| User5 | book1, book4, book6, book9 |
| User6 | book2, book5, book6, book8 |
| User7 | book4, book8, book9 |
| User8 | book2, book5, book6, book8, book9 |
| User9 | book1, book2, book8 |
| User10 | book3, book6, book7 |

Table 2: Clusters of users' historical purchase

| Cluster number | User name |
|---|---|
| 1 | User2, user3, user6, user8 |
| 2 | User1, user4, user9 |
| 3 | User5, user7 |
| 4 | User10 |

Table 3: Parts of association rules

| Items | Support | Confidence |
|---|---|---|
| Book2, book5 | 0.5 | 1 |
| Book5, book6 | 0.5 | 0.5 |
| Book6, book8 | 0.5 | 1 |
| Book5, book8 | 0.75 | 0.75 |
| Book2, book8 | 0.5 | 1 |

Table 4: Details of purchase records for user2

| Book number | Book name |
|---|---|
| Book3 | Java database best practices |
| Book4 | Building java enterprise applications: architecture v.1 |
| Book5 | First, catch your weka: A story of new Zealand cooking |

**Case study and analysis:** In this section we carry out 2 experiments for discussing the feasibility and effectiveness of the improved approach -- 1 for books recommendation in an online bookstore and the other for comparison between the proposed FCM with the traditional FCM.

**Books recommendation:** We select the purchase history data of users from an online bookstore as the demonstration. 10, 000 records are chosen, including 17, 980 books. 10 users are selected from such data set for discussion. In order to facilitate the analysis, the users' names are described from user1 to user10 while the book number from book1 to book10. The users' historical purchase data are shown in Table 1.

According to the presented determination method for cluster number, four clusters are generated as shown in Table 2.

Take cluster 1 for example, we elaborate the establishment of personalized knowledge base. The association rules can be obtained by analyzing the characteristics of users in cluster 1(Table 3).

Table 4 gives the details of user2 for his purchase records.

By the generated association rules for cluster 1, book6 and book8 are recommended to user 2 with the names "Data Mining: Practical Machine Learning Tools and Techniques, Third Edition" and "Data Structures and Algorithm Analysis in Java (3rd Edition)". The 2 recommended books belong to the scope of user2's personalized interest, which has high possibility that user 2 would like to buy.

**Comparison with traditional FCM:** To compare the 2 algorithms, we regard the recommendation as information retrieval problem. In this case, the evaluation standard in the field of information retrieval is used to measure the results, that is, precision rate and recall rate:

$$precision = \frac{right\_recommendation\_items}{all\_items\_of\_recommendation} \quad (23)$$

$$recall = \frac{right\_recommendation\_items}{all\_useful\_items\_of\_recommendation} \quad (24)$$

Since the precision rate and the recall rate are conflicting indicators in a certain degree, the high precision means that the recall rate is low. In order to balance them, we adopt F-measure to comprehensive evaluation (Hripcsak and Rothschild, 2005). The bigger the value of F-measure is, the higher the quality will be:

$$F - measure = \frac{2 \times precision \times recall}{precision + recall} = \frac{2}{1/precision + 1/recall} \quad (25)$$

We randomly select 50 books and their related purchase records as example. We divide the samples to 2 parts, one part for generating the recommendation while the other part for evaluating the accuracy of the recommended results. The results of F-measure for 50 books by the improved FCM (IFCM) and the traditional FCM (TFCM) are shown in Fig. 3, respectively.

The average of F-measure by the improved FCM is 76.46% and the average of F-measure by the traditional FCM is 63.3%.

The experimental results show that the personalized product recommendation with improved FCM clustering is feasible and effective. The results of recommendation are better than the traditional FCM. The system can recommend the personalized products to target users.

## CONCLUSION AND RECOMMENDATIONS

The recommendation system is an important technical means to help consumers obtain useful knowledge when searching information on the internet, which is a hot issue in e-business applications. By means of FCM clustering, we presented a novel approach for personalized product recommendation to assist consumers to make the decisions. The traditional FCM clustering was improved by modifying the
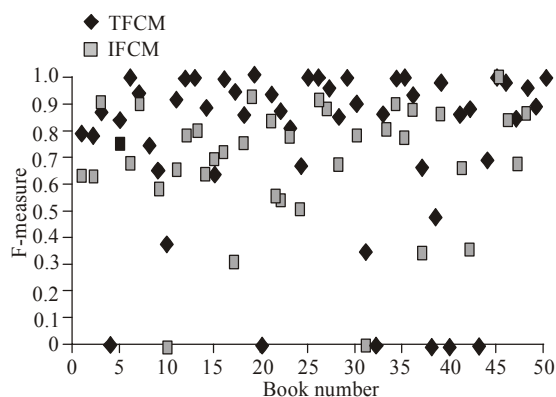
Fig. 3: Results of F-measure for 50 books by the improved
FCM and the traditional FCM

membership function and using the weighted objective
function, which can not only greatly reduce the
sensitivity to outliers and noise but also ensure the
clustering consistency. The workflow of personalized
product recommendation was given by the improved
FCM clustering. A case study of electronic bookshop
was carried out and the results showed that the
proposed approach is feasible and effective, which can
recommend interesting books to readers.

As a future work, we will focus on leveraging the
advantages of other clustering algorithms into our
model, such as hierarchical clustering and Self
Organizing Maps (SOM) clustering. Moreover, other
attributes should be investigated and incorporated in
modeling users' interests, for example, click behavior,
user location and login time etc.

## ACKNOWLEDGMENT

## REFERENCES

Alon, N., B. Awerbuch, Y. Azar and B. Patt-Shamir,
2009. Tell me who I am: An interactive
recommendation system. Theory Comp. Syst.,
45(2): 261-279.
Bezdek, J., 1984. FCM: The fuzzy c-means clustering
algorithm. Comp. Geosci., 10(2-3): 191-203.
Chou, P.H., P.H. Li, K.K. Chen and M.J. Wu, 2010.
Integrating web mining and neural network for
personalized e-commerce automatic service. Exp.
Syst. Appl., 37(4): 2898-2910.

Chu, W. and S.T. Park, 2009. Personalized
recommendation on dynamic content using
predictive bilinear models. Proceedings of the 18th
International Conference on World Wide Web, pp:
691-700.
Dongbo, B., B. Shuo and L. Guojie, 2002. Principle of
granularity in clustering and classification. Chinese
J. Comp., 8: 810-826.
Hripcsak, G. and A.S. Rothschild, 2005. Agreement,
the f-measure and reliability in information
retrieval. J. Am. Med. Inform. Assoc., 12: 296-298.
Huang, S.Y. and L.Z. Duan, 2012. E-Commerce
Recommendation Algorithm based on Multi-Level
Association Rules. In: Jin, D. and S. Lin (Eds.),
Advances in Electronic Commerce Web Applied
Communication (ECWAC). Springer-Verlag,
Berlin, Heidelberg, 148: 479-485.
Liu, D.R., C.H. Lai and W.J. Lee, 2009. A hybrid of
sequential rules and collaborative filtering for
product recommendation. Inform. Sci., 179(20):
3505-3519.
Ngai, E.W.T., L. Xiu and D.C.K. Chau, 2009.
Application of data mining techniques in customer
relationship management: A literature review and
classification. Exp. Syst. Appl., 36(2): 2592-2602.
Qian, Z.S., 2011. Research on web navigations. Res. J.
Appl. Sci. Eng. Technol., 3(10): 1171-1176.
Shen, H., 2011. A personalized E-commerce
recommendation method based on case-based
reasoning. Adv. Comp. Sci. Env. Ecoinform.
Educ., 215: 490-495.
Shishehchi, S., S.Y. Banihashem, N.A.M. Zin and
S.A.M. Noah, 2011. Review of personalized
recommendation techniques for learners in e-
learning systems. Proceeding of International
Conference on Semantic Technology and
Information Retrieval (STAIR), pp: 277-281.
Wang, T.Z., 2012. The ontology recommendation
system in E-commerce based on data mining and
web mining technology. Adv. Elec. Comm. Web
Appl. Commun., 149: 533-536.
Xie, X.L. and G. Beni, 1991. A validity measure for
fuzzy clustering. IEEE T. Pattern Anal., 13(8):
841-847.
Yi, M. and W. Deng, 2009. A utility-based
recommendation approach for e-commerce
websites based on bayesian networks. Proceeding
of International Conference on Business
Intelligence and Financial Engineering (BIFE'09),
pp: 571-574.