

Park Recommendation Algorithm based on User Reviews and Ratings

Chunxu Wang^a, Haiyan Wang^{a,*}, Jingwen Pi^a, and Li An^b

^a*School of Information, Beijing Forestry University, Beijing, 100083, China*

^b*Department of Geography, San Diego State University, California, 92182, USA*

Abstract

Recommendation systems are widely used in e-commerce websites as they can recommend appropriate movies, songs, books, and other items to users according to users' historical behavior. In traditional collaborative filtering algorithms, users' historical scores are usually used to predict the unknown item rating, while ignoring their textual reviews. Therefore, this paper proposes a park recommendation model based on user reviews and ratings (PRMRR). PRMRR first uses the latent Dirichlet allocation model to extract the statistical distribution of the park features. Secondly, it detects user preference distribution based on park features and user ratings. In order to measure the credibility of user ratings, user rating confidence level is considered to correct user preferences. Thirdly, it uses Kullback-Leibler divergence to calculate the similarity between different users and then predicts the unknown park rating for a specific user. Finally, the proposed algorithm is evaluated on two real park data sets, and the results on two different data sets show that the proposed approach outperforms other traditional approaches. Our recommendation algorithm thus has great potential to improve the quality of park recommendation and effectively handle the data sparsity problem.

Keywords: collaborative filtering; user preference; park recommendation; user rating confidence level

(Submitted on October 22, 2018; Revised on November 21, 2018; Accepted on December 23, 2018)

© 2019 Totem Publisher, Inc. All rights reserved.

1. Introduction

With the continuous development of computer technology, recommendation algorithms have been broadly used to deal with information overload problems in e-commerce sites. Recommendation algorithms are designed to recommend what users may like, and they have been widely applied in films, books, songs, and other fields. Content-based recommendation and collaborative filtering (CF) are two methods widely used in recommendation systems. Content-based recommendations mainly rely on additional information about items to make predictions for users. For parks, such additional information may include traffic, hardware, special activities, and so on. CF is the most commonly used recommendation algorithm [1], and it makes predictions through a user-item rating matrix. The basic idea of CF is that if users had similar preferences in the past, they will have similar preferences in the future [2]. According to different implementations of CF, CF can also be divided into memory-based methods and model-based methods. Memory-based methods include the user-based method [3] and the item-based method [2]. In the user-based method, we first search users who have the same interests as the target user (i.e., the neighbor users), and items are recommended to the target user through neighbor users' ratings on those items. The item-based method calculates the similarity of items by analyzing the users' historical behavior and then makes recommendations through similarity between items. Memory-based methods are widely used in actual situations, but there is a problem of data sparsity. In most cases, common ratings between two users are rare or absent, seriously affecting the quality of the corresponding recommendation. Model-based methods, such as the latent model [4], Singular Value Decomposition (SVD) [5], and Bayesian network [2], use the user-item rating matrix as the training data, identify the relevance between users by training machine learning models, and then make intelligent recommendations. Model-based methods alleviate the data sparsity problem, but the cost of off-line training models is often higher when the model is complex or the training data set is large in size.

As an essential part of life, many people choose to spend much leisure time in parks. In one area, there are many

* Corresponding author.

E-mail address: wanghaiyan76@126.com

different types of parks. However, it is difficult for tourists to find a park that matches their travel plans due to the wide variety of parks and lack of effective information. Therefore, research on park recommendation can help tourists better choose parks. In actual park recommendation processes, the user-item rating matrix is sparse and parks have many unique features, which makes models or methods from other fields not directly applicable. Therefore, applying CF to park recommendation directly will lead to a lower accuracy.

With the rapid development of the Internet, users usually give a corresponding comment after purchasing items or obtaining services in e-commerce websites. Compared with a single rating, user reviews contain a large amount of information. For example, the contents in Figure 1 below are from Dianping.com and describe the ratings and reviews of two users on Tao Ranting Park in the Beijing area. The two users all gave five points on the park, but their perspectives of evaluation were different. The user with ID 1224 prefers the park for its convenient access, beautiful scenery, and distinctive buildings. The user with ID 1218 thinks the park is suitable for travel and the environment is good. It can be seen that the user reviews contain abundant user preference information than a simple rating, while the traditional collaborative filtering algorithm tends to ignore this important resource. Therefore, if we can discover the park characteristics and the user preferences from user reviews, the corresponding recommendation quality must be greatly improved.

User ID: 1224 User name: Love you forever& Rating: 5 Content: Tao ranting park is located in the southern part of Beijing. You can take the subway line 4 to the Tao ranting Park station, or you can get there by bus. The scenery here is very beautiful, each season brings you different feeling. Every time I come here, there are many uncles and aunts chatting, singing, dancing and enjoying themselves in the park. In addition, the buildings of the Tao ranting Park are very distinctive.	User ID: 1218 User name: A beautiful panda Rating: 5 Content: There are memories of my childhood, and we will visit here on children's day, spring outing and autumn outing. The big snow mountain is the most famous place here, and there are still many adults playing now. In the morning, people come here to exercise and the environment in the park is pretty good.
--	---

Figure 1. User reviews example

In this paper, we present a new approach that links information hidden in user comments to user ratings. Normally, we think different users have different preferences, so we can measure their similarity through user preferences and then make recommendations for users. In our approach, we first use the latent Dirichlet allocation (LDA) topic model to discover the hidden features of different parks. Second, we calculate user preferences through user ratings and park features. In order to measure the credibility of user ratings, a user rating confidence level is introduced to magnify user preferences and make user preferences more accurate. Third, according to preference distribution data of two users, we use Kullback-Leibler (KL) divergence to calculate the similarity between the two users and then predict the unknown park scores for a specific user. Finally, we test our approach through a series of experiments on two real park data sets, and the final results on two different data sets show that the performance of our approach is better than that of the baseline methods.

2. Related Work

In recent years, review analysis has been applied to the field of recommendation system and effectively improved the quality of recommendation. The review analysis mainly includes two aspects: sentiment analysis and topic extraction.

Recent years have witnessed the increasing popularity of the LDA topic model in recommendation systems. The LDA model, also known as a three-layer Bayesian probability model, is a document topic generation model widely used in topic extraction. It was first proposed by Blei et al. [6]. Liu et al. [7] applied LDA to the user-item rating matrix by treating each user as a document and treating an item as a word. Zhao et al. [8] proposed a hybrid approach of LDA and matrix factorization. The method first predicted the probability of users' rating behaviors and then used matrix factorization to make rating predictions. These two methods only consider the rating information but ignore the user reviews. McAuley and Leskovec [9] combined the hidden factors in user ratings with the hidden topics in user reviews to create user preferences and then mapped it to SVD model for rating prediction. However, this method cannot consider both the user and the item at the same time. Zhang et al. [10] used the LDA topic model to explore topics hidden in user reviews and created portraits for users and items respectively. Then, the resultant machine learning algorithm was used for rating prediction based on the collaborative filtering algorithm. Xu et al. [11] used the LDA model to generate the topic distribution of each review to get the most important features of each user and then calculated the user similarity and predicted the rating for the corresponding user. Zhou and Wu [12] proposed a rating LDA model for collaborative filtering by adding rating information to the LDA model. Pu et al. [13] utilized user and item information to generate the user's sentiment topic distribution in an unsupervised way, which was then able to make reliable rating predictions.

There has also been a large amount of research on sentiment analysis of user reviews. Ganu et al. [14] identified topics

and sentiments of related comments from user reviews and then made rating predictions according to different levels of information in user reviews. Qu et al. [15] introduced the bag-of-opinions method, in which the all reviews consist of three parts: root, modifier, and negative word. This method predicts user's ratings by bag-of-opinions and the linear model. Leung et al. [16] proposed an approach to identify features in reviews by integrating sentiment analysis and CF, but this approach relies on natural language processing techniques to deal with unstructured, natural language texts. Ganu et al. [17] proposed an approach by aggregating similar users through topic and sentiment extraction in user reviews to predict rating for users. Although this method alleviates the data sparsity problem by using user reviews, it needs to identify the topics and sentiments of about 3,400 sentences, which is time-consuming.

On the basis of the previous studies about user review analysis, we propose a park recommendation algorithm based on user reviews and ratings. Compared with previous studies that mostly focus on movies, songs, books, or other commodities or services, we have selected a new field of study, where we mainly target parks. Second, we utilize user rating confidence level to measure the reliability of users' ratings in our approach, which can avoid or minimize users' random or false ratings and thus make the user preferences more accurate. Third, KL divergence is used to measure the similarity of user preferences instead of the traditional method of calculating similarity. Finally, the proposed model is based on user ratings and user reviews, so we can effectively find the park features that users are interested in and then find more accurate neighbors for users by analyzing these two parts. Thus, it can alleviate the data sparsity problem and enhance park recommendation accuracy.

3. The Park Recommendation Model based on User Reviews and Ratings (PRMRR)

In this section, we first briefly introduce the LDA topic model including its generation process, followed by a detailed description of the park recommendation model based on user reviews and ratings (PRMRR) model that we propose in this article.

Specifically, the PRMRR model is comprised of the following steps.

Step 1 The LDA topic model is used to find the probability distribution of the park on each attribute from the user reviews, and user preferences are calculated based on user ratings and park features. In order to make user ratings more credible, we introduce the so-called user rating confidence level to modify user preferences.

Step 2 KL divergence is used to calculate the similarity between two users.

Step 3 Find the neighbor users for the target user and then predict the score of the unknown park for a certain target user. The main symbols and definitions involved in the paper are shown in Table 1.

Table 1. Main symbols and definitions in the paper

Symbol	Meaning
$R_{m \times n}$	User-park rating matrix
$U = \{u_1, u_2, \dots, u_m\}$	User sets
$I = \{i_1, i_2, \dots, i_n\}$	Park sets
$R_{u,i}$	The rating of the user u for the park i
$C_{u,i}$	The review of the user u for the park i
X_u	The parks that have been rated by user u
Ave_u	The average rating of user u
Ave_i	The average rating of park i
P_u	User preferences

3.1. LDA Topic Model

The LDA topic model is a three-layer Bayesian probability model with three layers of documents, topics, and words. In this paper, documents, topics, and words represent users, latent preferences, and parks respectively. The model holds that every word in an article is obtained by "choosing a topic with a certain probability and choosing a word from the topic in a certain probability". Accordingly, we use it to discover user preferences hidden in user reviews.

As shown in Figure 2, N represents the total number of words in corpus, D the collection of all documents in corpus, K the number of topics, θ_d the topic probability distribution of the document d , and φ_k the probability distribution of words under the topic K . α and β are hyper-parameters of θ_d and φ_k respectively. Document processing and parameter selection are critical to the model, which will be discussed in the experimental part.

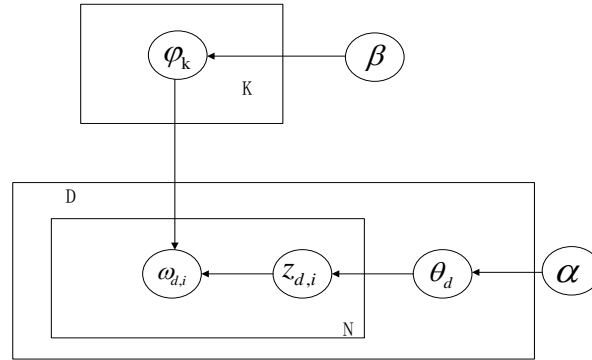


Figure 2. The model diagram of LDA

3.2. User Preference Calculation

User preferences represent a user's preference for a certain attribute of a park. For example, some users may like community parks, some users may like scenic parks, while others may prefer parks with amusement facilities. In this paper, user preferences are calculated by combining user ratings and the characteristic distribution of parks. We regard all user reviews $\{C_{1,i}, C_{2,i}, \dots, C_{n,i}\}$ included in park i as a document d_i , and corpus D is comprised of user reviews $\{d_1, d_2, \dots, d_n\}$ for all parks that are indexed from 1, 2, \dots , up to n . In order to make the corpus data accessible by computers, the data preprocessing links (e.g., words segmentation and removal of stop words) are used in the corpus. Then, we apply the LDA topic model to the corpus to calculate the probability distribution of each park on various features. The probability distribution of park i on different topics can be represented by a set of n -dimensional vectors, as shown in Equation (1):

$$\theta_i = \{\theta_{1,i}, \theta_{2,i}, \dots, \theta_{n,i}\} \quad (1)$$

Where n is the number of features and $\theta_{n,i}$ is the feature distribution for park i on topic n . For user u , each comment has a corresponding rating $R_{u,i}$, and we thus define user preferences P_u as Equation (2):

$$P_u = \frac{\sum_{i \in X_u} \theta_i \times R_{u,i}}{|X_u|} \quad (2)$$

Where X_u represents the parks that have been rated by user u . In order to identify effective user preferences, this paper adopts the user rating confidence level as an improved user preference calculation method. The user rating confidence level includes two factors: user subjective rating confidence level and user objective rating confidence level. These two factors measure the reliability of user scores from both subjective and objective aspects to reduce the impact of false comments or random comments.

Different users have different rating habits. In general, some users have a higher overall rating, while some users have a lower overall rating. In order to get more accurate user preference, the user's average rating is selected as the measurement standard to measure the user preferences for different parks, and a sigmoid function is introduced to describe the user subjective rating confidence level $Sub_{u,i}$. The formula of $Sub_{u,i}$ is defined as Equation (3):

$$Sub_{u,i} = \frac{1}{1 + \exp(-(R_{u,i} - Ave_u))} \quad (3)$$

Under this definition, $Sub_{u,i}$ varies from 0 to 1. The higher the user's rating on the park, the more interested the user is in the park. Meanwhile, when $R_{u,i} > Ave_u$, the user's preference for the park is positive, the corresponding characteristic distribution of the park is more consistent with the user's true preference distribution, and vice versa.

On the other hand, we utilize the user objective rating confidence level $Obj_{u,i}$ to measure the objective reliability of user ratings. If the user's score is close to the average score of the park, it indicates that the user is relatively objective, and vice versa. The user objective rating confidence level is defined as Equation (4):

$$Obj_{u,i} = 2 \times \left(1 - \frac{1}{1 + \exp(-|R_{u,i} - Ave_i|)} \right) \quad (4)$$

The value of $Obj_{u,i}$ varies from 0 to 1, with a higher value indicating higher objectiveness. The smaller the difference between the user's rating on the park and the average rating of the park, the more objective the user is. In other words, the user's objective rating confidence level is relatively high. To sum up, the user subjective rating confidence level $Sub_{u,i}$ and user objective rating confidence level $Obj_{u,i}$ are multiplied to get the final user rating confidence level. Then, the user preference calculation formula is calculated by Equation (5) after combining Equations (3) and (4).

$$P_u = \frac{\sum_{i \in X_u} \theta_i \times R_{u,i} \times Obj_{u,i} \times Sub_{u,i}}{|X_u|} \quad (5)$$

This definition indicates that after joining these two parts, we can effectively avoid users' false or random comments to get a more authentic distribution of user preferences.

3.3. User Similarity Calculation

After getting user preferences, the next step is to calculate the similarity between different users. The common methods of computing similarity include cosine similarity, Pearson similarity, and Jaccard correlation coefficient. In order to measure the user's preferences in a comprehensive way, this paper uses *KL* divergence to calculate the similarity between users instead of using traditional similarity calculation methods.

KL divergence, also known as relative entropy, is widely used to calculate the difference between two probability distributions [18]. Given two probability distributions P and Q , the *KL* distance of P and Q is defined as Equation (6):

$$KL(P \parallel Q) = P \log \left(\frac{P}{Q} \right) \quad (6)$$

Based on the above definition, we let p_u and p_v represent the user preference vectors for user u and user v , respectively. Accordingly, the *KL* distance between user u and user v is defined as Equation (7):

$$KL(p_u \parallel p_v) = \sum p_u \log \frac{p_u}{p_v} \quad (7)$$

The *KL* distance is asymmetrical, that is, $KL(p_u \parallel p_v) \neq KL(p_v \parallel p_u)$. Therefore, *JS* divergence is introduced to calculate the *KL* distance between users. Among them, the definition of using the *JS* divergence correction is shown as Equation (8):

$$JS(p_u \parallel p_v) = \frac{1}{2} \times KL(p_u \parallel \frac{p_u + p_v}{2}) + \frac{1}{2} \times KL(p_v \parallel \frac{p_u + p_v}{2}) \quad (8)$$

According to the *JS* divergence, the user similarity $Sim_{u,v}$ between user u and user v can be obtained as Equation (9):

$$Sim_{u,v} = \frac{1}{1 + JS(p_u || p_v)} \quad (9)$$

Using *KL* divergence to calculate the similarity between different users, all users' preference information is considered in a comprehensive way. The results are not affected by special results, i.e., they are of higher objectivity.

3.4. Rating Prediction

In the stage of rating prediction, we use the neighborhood method to predict the unknown rating. First, we find the neighbor users for the target user by the similarity between users and then predict the target user's rating on the unknown park based on the ratings of the neighbor users. The prediction formula is shown as Equation (10):

$$R(u,i) = \bar{r}_u + \frac{\sum_{v \in N_u} sim(u,v) \times (r_{v,i} - \bar{r}_v)}{\sum_{v \in N_u} (|sim(u,v)|)} \quad (10)$$

Among them, \bar{r}_u and \bar{r}_v represent the average rating of user u and user v respectively, $r_{v,i}$ the rating for user v on park i , $sim(u,v)$ the similarity between user u and user v , and N_u the set of neighbor users of the target user u . In order to judge whether the user's evaluation is close to an objective evaluation, this paper introduces the user evaluation accuracy to measure the credibility of user u . The user evaluation accuracy is defined as Equation (11):

$$Weight_u = 2 \times \left(1 - \frac{1}{1 + \exp\left(-\left| \sum_{i \in X_u} \frac{|R_{u,i} - Ave_i|}{|X_u|} \right| \right)} \right) \quad (11)$$

On the user's rating set, the average absolute error between the user's rating of the park and the actual rating of the park becomes smaller and smaller, indicating that the user's evaluation result is closer to the objective evaluation of the park and suggesting that the user is trustworthy. Otherwise, the user is less trusted.

The rating prediction formula after adding user evaluation accuracy is shown as Equation (12):

$$R(u,i) = \bar{r}_u + \frac{\sum_{v \in N_u} sim(u,v) \times (r_{v,i} - \bar{r}_v) \times Weight_u}{\sum_{v \in N_u} (|sim(u,v)|)} \quad (12)$$

Equation (12) implies that the higher the accuracy of user evaluation, the greater the weight of the user in the rating prediction formula.

4. Experiments

In this section, we describe the experiment in detail. We first introduce the data sets and evaluation metric used in experiments. Furthermore, we determine the setting of experimental parameters and compare our algorithm with baseline methods to verify the performance of the recommendation.

4.1. Data Sets

In this paper, two park data sets are used to verify the proposed method, viz. the data set from Beijing (Park-Beijing hereafter) and the data set from Shanghai (Park-Shanghai hereafter). These two data sets are crawled from Dianping.com, the world's first independent third party website that holds information from consumer reviews website. The two data sets in Park-Beijing and Park-Shanghai are from the Beijing area and Shanghai area respectively, and they include user IDs, park IDs, users' ratings (1-5), and users' reviews on parks. Park-Beijing contains more than 80000 records from 2003 to 2017, and Park-Shanghai contains more than 80000 records from 2003 to 2018. According to the experimental needs, we filter out the reviews that have no textual comments. Specifically, in order to avoid the data being too sparse, we retain users with

more than three comments and parks with more than six comments in both data sets. The detailed statistical information of the two data sets is shown in Table 2.

Among them, the sparsity of a data set is defined as Equation (13):

$$Sparsity = 1 - \frac{|Number\ of\ user\ reviews|}{|Number\ of\ users| \times |Number\ of\ items|} \quad (13)$$

The sparsity describes the sparse degree of the data set; the higher the sparsity, the sparser the data is. The user rating distribution is shown in Table 3. It can be seen that the rating is concentrated in levels 3, 4, and 5 on both data sets, and thus it is limited to utilize user ratings simply for recommendation.

Table 2. Statistical information on data sets

Data set	Park-Beijing	Park-Shanghai
Number of users	3347	2398
Number of parks	231	242
Number of reviews	28107	20850
Sparsity	96.36%	96.41%
User average comments	8.39	8.69

Table 3. Ratio of rating distribution

Rating	1	2	3	4	5
Park-Beijing	0.5%	1.3%	25.3%	41.1%	31.7%
Park-Shanghai	0.3%	2%	23%	47.9%	26.8%

4.2. Evaluation Metric

Mean absolute error (MAE) and root mean square error (RMSE) are the most widely used recommendation quality evaluation indices in recommendation systems. MAE is the average of the absolute error between the predicted value and the actual value, and RMSE is the root mean square error between the predicted value and the actual value. Let the predicted user rating be denoted as $\{p_1, p_2, \dots, p_n\}$ and the actual user rating set as $\{q_1, q_2, \dots, q_n\}$. The MAE is defined as Equation (14):

$$MAE = \sum_{i=1}^n \frac{|p_i - q_i|}{n} \quad (14)$$

Obviously, the smaller the MAE value, the higher the prediction accuracy. The RMSE is defined as Equation (15):

$$RMSE = \sqrt{\sum_{i=1}^n \frac{|p_i - q_i|^2}{n}} \quad (15)$$

Similar to MAE, the smaller the RMSE value, the higher the prediction accuracy. Compared with MAE, Netflix thinks that RMSE has increased penalties for predicting inaccurate items and is more critical for the evaluation of rating prediction [19].

4.3. Baseline

To verify the effectiveness of our algorithm, we compare our algorithm with four approaches as follows:

- SVD: The user rating matrix is decomposed into the product of three matrices, then the three matrices are used to predict the unknown rating.
- CF: A basic collaborative filtering algorithm that takes into account the mean ratings of each user.
- NMF [20]: A collaborative filtering algorithm based on non-negative matrix factorization.
- LDA: Through the user rating matrix, we use the LDA topic model to find the user potential vector and predict the unknown rating.

4.4. Experimental Results and Analysis

In this section, we design two experiments. First, we need to select the appropriate number of topics for corpus in two data sets, and then we compare our method with other methods and analyze the experimental results.

4.4.1. Topic Selection

Since the LDA topic model is applied to user comment text, the number of topics needs to be determined. Considering that users may be interested in only a few topics, other topics will have less impact on users. In order to find the most suitable number of topics on the corpus, this paper sets α at the empirical value 0.2, β at the empirical value 0.1, and the number of iterations at 100. Then, the corresponding MAE is calculated by setting different topic numbers (5-25) to select the most appropriate number of topics.

When the number of topics varies from 5 to 25, the MAE of the PRMRR on two data sets is shown in Figure 3. Overall, MAE shows a trend of initial decline and then rise on both data sets. Specifically, it can be seen from Figure 3(a) that the method generates good performance when the topic is 15 for the Park-Beijing data set (Figure 3(a)), while the method achieves the smallest MAE when the topic is 10 for the Park-Shanghai data set (Figure 3(b)). Therefore, we set the number of topics at 10 on the Park-Shanghai data set and 15 on the Park-Beijing data set.

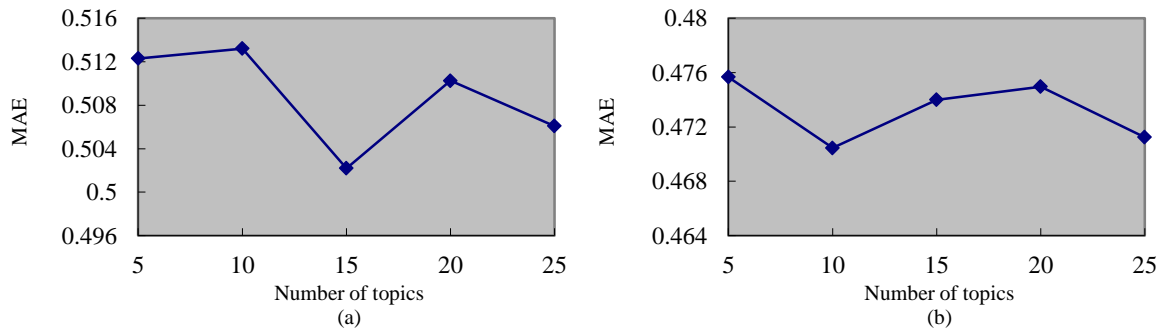


Figure 3. (a) MAE under different topics on Park-Beijing; (b) MAE under different topics on Park-Shanghai

4.4.2. Results and Analysis

In order to verify the performance of proposed algorithm, we conducted a series of experiments on the Park-Shanghai data set and Park-Beijing data set. Each of the data sets is randomly divided into 70% training data and 30% testing data.

The MAE values of different recommendation algorithms (Figure 4) are used to examine the performance of these algorithms. Regardless of what data sets (i.e., Park-Beijing and Park-Shanghai) are chosen, the MAE value of PRMRR is the smallest, suggesting that PRMRR is superior to other algorithms. On Park-Beijing, compared with NMF, SVD, LDA, and CF, the PRMRR method reduces the error by 8.6%, 6%, 3.1%, and 0.8% respectively. On Park-Shanghai, compared with NMF, SVD, LDA, and CF, the PRMRR method reduces the error by 14.5%, 9.5%, 1.9%, and 3.2% respectively. It is obvious that the performance of our method is better than that of other methods.

We also use RMSE to assess the performance of different recommendation algorithms based on two data sets (Figure 5). Our method has generated the smallest RMSE values compared with other algorithms. On Park-Beijing, the RMSE from our method is the lowest (close to 0.66), suggesting our model performs best among all the methods. In addition, CF is closest to the RMSE of the method proposed in this paper. On Park-Shanghai, the proposed method still has the lowest RMSE. Among baselines, LDA is closest to the proposed method. On both data sets, our method has generated the best recommendation performance.

Using MAE and RMSE as performance indicators, we also see that the performance of PRMRR from Park-Shanghai is better than that from Park-Beijing overall (Figures 4 and 5). The main reason is that the number of user average comments of Park-Shanghai is bigger than that of Park-Beijing, making user preferences more precise in the calculation process, although the data of Park-Shanghai is sparser than that of Park-Beijing.

All the above experimental results demonstrate that compared with the baseline approaches, our approach has

significantly reduced the error of the rating prediction, suggesting that the quality of recommendation has been significantly improved. The results also show that user preferences can be evaluated more effectively after introducing user reviews, which include more useful information than a single user rating information. Compared with traditional methods, the method proposed in this paper can effectively discover users' interests and greatly improve the quality of recommendation by fully analyzing user reviews.

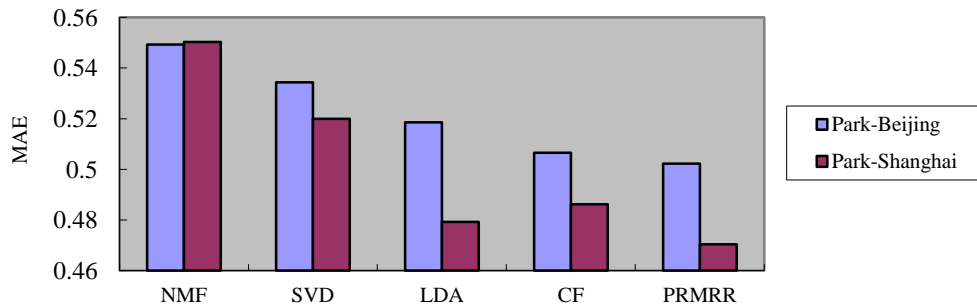


Figure 4. MAE on different algorithms

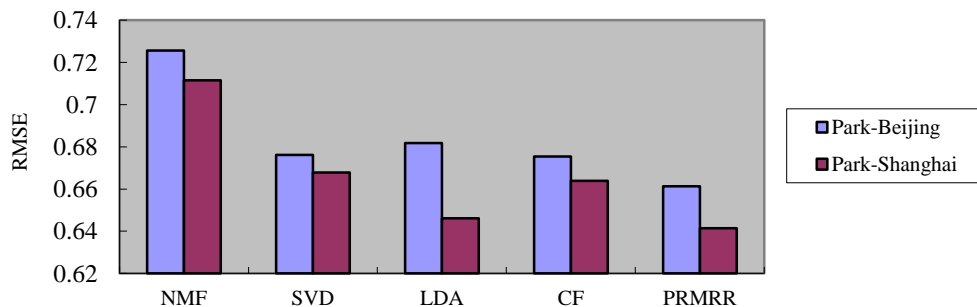


Figure 5. RMSE on different algorithms

5. Conclusions

In this paper, we proposed a park recommendation algorithm based on user reviews and ratings. The algorithm extracts potential topic characteristics of the corresponding park through the LDA topic model, which is capable of fully extracting user preferences based on user ratings and park features. Compared with traditional collaborative filtering algorithms, our algorithm takes advantage of a large amount of information included in user reviews, outperforms all other traditional methods, and effectively alleviates the data sparsity problem. Although our algorithm was tested on two exemplar data sets, we believe it is applicable to other parks and can greatly improve the quality of park recommendation.

However, our park recommendation algorithm still has many aspects worth additional efforts. In the near future, we will collect more park data from other regions to further verify the quality of the algorithm. Additionally, the influence of the time when users make comments as well as the impact of user social relationships on user preferences will be considered to further improve the quality of recommendation.

Acknowledgements

This work is supported by the National Natural Science Foundation of China (No. 61772078).

References

1. J. S. Breese, D. Heckerman, and C. Kadie, "Empirical Analysis of Predictive Algorithms for Collaborative Filtering," in *Proceedings of 14th Conference on Uncertainty in Artificial Intelligence*, pp. 43-52, Madison, Wisconsin, July 1998
2. B. Sarwar, G. Karypis, J. Konstan, and J. Riedl, "Item-based Collaborative Filtering Recommendation Algorithms," in *Proceedings of the 10th International Conference on World Wide Web*, pp. 285-295, Hong Kong, China, May 2001
3. P. Resnick, N. Iacovou, M. Suchak, P. Bergstrom, and J. Riedl, "GroupLens: An Open Architecture for Collaborative Filtering of Netnews," in *Proceedings of the 1994 ACM Conference on Computer Supported Cooperative Work*, pp. 175-186, Chapel Hill,

North Carolina, USA, October 1994

4. T. Hofmann, "Latent Semantic Models for Collaborative Filtering," *Acm Transactions on Information Systems*, Vol. 22, No.1, pp. 89-115, 2004
5. B. Sarwar, G. Karypis, J. Konstan, and J. Riedl, "Application of Dimensionality Reduction in Recommender System-A Case Study," in *Proceedings of the ACM WebKDD Workshop on Web Mining for E-Commerce*, pp. 82-90, Boston, MA, USA, 2000
6. D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent Dirichlet Allocation," *Journal of Machine Learning Research*, Vol. 3, pp. 993-1022, 2003
7. Q. Liu, E. H. Chen, H. Xiong, C. H. Ding, and J. Chen, "Enhancing Collaborative Filtering by User Interest Expansion via Personalized Ranking," *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, Vol. 42, No. 1, pp. 218-233, 2012
8. X. Y. Zhao, Z. D. Niu, W. Chen, C. Y. Shi, K. Niu, and D. L. Liu, "A Hybrid Approach of Topic Model and Matrix Factorization based on Two-Step Recommendation Framework," *Journal of Intelligent Information Systems*, Vol. 44, No. 3, pp. 335-353, 2014
9. J. McAuley and J. Leskovec, "Hidden Factors and Hidden Topics: Understanding Rating Dimensions with Review Text," in *Proceedings of the 7th ACM Conference on Recommender Systems*, pp. 165-172, Hong Kong, China, October 2013
10. R. Zhang, Y. F. Gao, W. Z. Yu, P. F. Chao, X. Y. Yang, M. Gao, et al., "Review Comment Analysis for Predicting Ratings," in *Proceedings of the 16th International Conference on Web-Age Information Management*, pp. 247-259, Qingdao, China, 2015
11. J. N. Xu, X. L. Zheng, and W. F. Ding, "Personalized Recommendation based on Reviews and Ratings Alleviating the Sparsity Problem of Collaborative Filtering," in *Proceedings of IEEE Ninth International Conference on E-Business Engineering*, pp. 9-16, Washington, DC, USA, September 2012
12. X. Z. Zhou and S. X. Wu, "Rating LDA Model for Collaborative Filtering," *Knowledge-based Systems*, Vol. 110, pp. 135-143, 2016
13. X. J. Pu, G. S. Wu, and C. F. Yuan, "User-Aware Topic Modeling of Online Reviews," *Multimedia Systems*, pp. 1-11, 2017
14. G. Ganu, N. Elhadad, and A. Marian, "Beyond the Stars: Improving Rating Predictions using Review Text Content," in *Proceedings of International Workshop on the Web and Databases*, pp. 1-6, Providence, Rhode Island, USA, June 2009
15. L. Qu, G. Ifrim, and G. Weikum, "The Bag-of-Opinions Method for Review Rating Prediction from Sparse Text Patterns," in *Proceedings of the 23rd International Conference on Computational Linguistics*, pp. 913-921, Beijing, China, August 2010
16. W. K. Leung, C. F. Chan, and F. L. Chung, "Integrating Collaborative Filtering and Sentiment Analysis: A Rating Inference Approach," in *Proceedings of the ECAI 2006 Workshop on Recommender Systems*, pp. 62-66, Riva del Garda, Italy, August 2006
17. G. Ganu, "Improving the Quality of Predictions using Textual Information in Online User Reviews," *Information Systems*, Vol. 38, No. 1, pp. 1-15, 2013
18. A. Huang, "Similarity Measures for Text Document Clustering," in *Proceedings of the Sixth New Zealand Computer Science Research Student Conference*, pp. 49-56, Christchurch, New Zealand, 2008
19. L. Xiang, "Recommendation System Practice," The Posts and Telecommunications Press, Beijing, China, 2012
20. S. Zhang, W. H. Wang, J. Ford, and F. Makedon, "Learning from Incomplete Ratings using Non-Negative Matrix Factorization," in *Proceedings of the Sixth SIAM International Conference on Data Mining*, pp. 549-553, Bethesda, Md, USA, April 2006

Chunxu Wang is currently pursuing her Master's degree in management science and engineering at Beijing Forestry University. Her research interests include recommender systems.

Haiyan Wang is an associate professor in the School of Information at Beijing Forestry University. Her main research interests include data mining and data analysis.

Jingwen Pi is currently a Master's student in the School of Information at Beijing Forestry University. Her main research interests include data analysis.

Li An is currently a professor at San Diego State University. His research interests include reciprocal human-environment relationships and space time data analysis and modeling.