

Improving Collaborative Filtering with Long-Short Interest Model

No Author Given

No Institute Given

Abstract. Collaborative filtering (CF) has been widely employed within recommender systems in many real-world situations. The basic assumption of CF is that items liked by the same user would be similar and users like same items would share similar interest. But it is not always true since the user's interest changes over time. It should be more reasonable to assume that if these items are liked by the same user in the same time period, there is a strong possibility that they are similar, but the possibility will shrink if the user likes them in different time period. In this paper, we propose a long-short interest model (LSIM) based on the new assumption to improve collaborative filtering. In special, we introduce a neural network based language model to extract the sequential features on user's preference over time. Then, we integrate the sequential features to solve the rating prediction task in a feature based collaborative filtering framework. Experimental results on three MovieLens datasets demonstrate that our approach can achieve the state-of-the-art performance.

Keywords: Recommender System, Collaborative Filtering, Long-Short Interest Model

1 Introduction

In the modern era of information overload, recommender system (RS) has become more and more popular in many real-world situations. Recommender system aims to help users find the items, they are more likely to be interested in, from huge amounts of candidates. Lots of websites (e.g. Amazon, Netflix, Alibaba and Hulu) use recommender system to target customers and provide them with useful information. An excellent recommendation system can effectively increase the amount of sales. For instance, 80% of movies watched on Netflix come from their recommender system [6].

A widely used setting of recommender system [22] is to predict the rating a user will evaluate on a new item (such as a movie) given the past rating history of the users. Lots of classical recommendation methods have been proposed during the last decade, and they can be categorized into two classes: content based methods and collaborative filtering based methods. Content based methods [18] take advantage of user profiles and item properties for recommendation. While collaborative filtering based approaches [27] utilize the past interactions or preferences, such as users' ratings on items, without using user or product content

information for recommendation. Collaborative filtering based approaches have attracted more attention due to their impressive performance, and developed for many years and keep to be a hot area in both academia and industry.

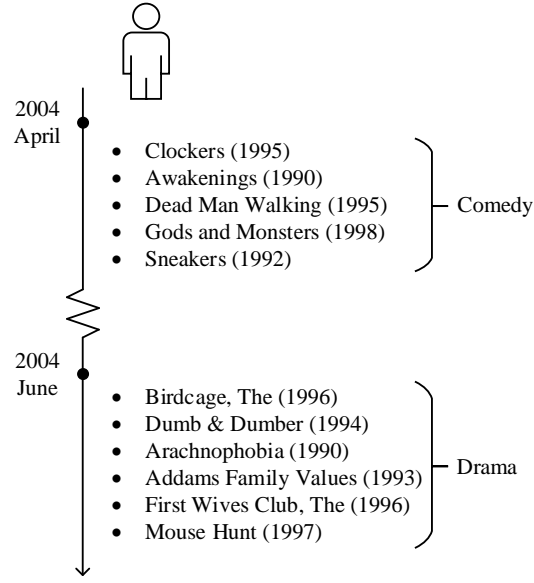


Fig. 1. The preference records of user whose id is 5988 in MovieLens-1M dataset, which are sorted by their rated time.

Collaborative filtering assumes that items liked by the same user would be similar and users like same items would share similar interest. However, it is not always true because the user's interest changes over time. For example, given a user in MovieLens-1M dataset whose id is 5988, Figure shows the movies he watched sorted by the rating time. We can find that this user liked watching comedy movies in April 2004 and changed to love watching drama movies in June 2004. These movies are going to be treated similar in conventional collaborative filtering, but they are not in actual. A more reasonable assumption, aka long-short interest assumption, should be that items liked by the same user in the same time period have a higher possibility to be similar than items liked by the same user in different time period.

Inspired by paragraph2vec algorithm [11] for learning vector representations of words which take advantage of a word order observed in a sentence, we introduce a long-short interest model (LSIM) to extract sequential features of users and items based on the new assumption. As illustrated in Figure 2, user is similar with the sentence, both of them contains a sequence follow some or-

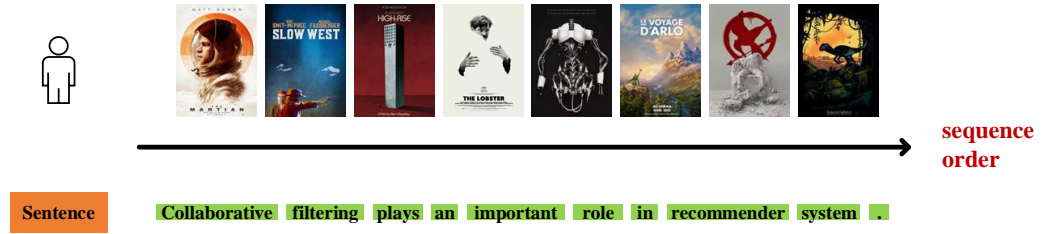


Fig. 2. Paragraph2vec learns vector representations of sentences and words based on the word order while LSIM extracts sequential features of users and items based on the rating order.

der, and items are similar with words because both of them follow the law that the more close they are, the more similar they are. To verify the effectiveness of the learned sequential features of users and items, we integrate them as side information to solve the rating prediction task in a feature based collaborative filtering framework.

The main contributions of this paper include: (1) We introduce a long-short interest model (LSIM) to extract sequential features of users and items based on the long-short interest assumption. (2) We demonstrate the effectiveness of the sequential features via integrating them as side information to solve the rating prediction task. (3) Experiments on three public MovieLens shows LSIM can achieve the state-of-the-art performance.

The rest of the paper is organized as follows. Section 2 gives an overview of the related work. Then, we describe our long-short interest model and the feature based collaborative filtering framework in Section 3. The experimental results as well as the comparisons with baseline system are shown in Section 4. Finally, we conclude the paper and outline our future work in Section 5.

2 Related Work

Our work is closely related to collaborative filtering and neural network language model. We will discuss them in the following subsections.

2.1 Collaborative Filtering

Collaborative filtering based methods can mainly divided into three categories: user-based collaborative filtering, item-based collaborative filtering and model-based collaborative filtering. User-based collaborative filtering [21] recommends items liked by users who are similar with you while item-based collaborative filtering [25] aims to recommend items similar with ones you liked in the past. Matrix factorization (MF) is the most popular model-based collaborative filtering methods, their success at the Netflix competition [1, 9] have demonstrated

their amazing strength, and lots of variants of it have been proposed in the following works.

Basically, the given ratings matrix $\mathbf{R} \in \mathbb{R}^{N \times M}$ consisting of the item preferences of the users can be decomposed as a product of two low dimensional matrices $\mathbf{U} \in \mathbb{R}^{N \times K}$ and $\mathbf{V} \in \mathbb{R}^{K \times M}$. \mathbf{U} could be treated as a user-interest matrix while \mathbf{V} could be treated as a item-interest matrix. K is the amount of interest. The decomposition can be carried out by a variety of methods such as singular value decomposition (SVD) based approaches [15], non-negative matrix factorization approach [12] and regularized alternative least square (ALS) algorithm [32]. Meanwhile, non-linear algorithms are proposed to catch subtle factors, such as Non Linear Probabilistic Matrix Factorization [10], Factorization Machines [19] and Local Low Rank Matrix Approximation [13]. However, these methods group users and treat items they rated equally, which will lose the sequential features to describe the long-short interest.

Matrix factorization methods suffer from the cold start problem, i.e. what recommendations to make when a new user/item arrives in the system. Another problem often presented in many real world applications is data sparsity. Incorporating side information has shown promising performance in collaborative filtering in such scenarios. In recent years, deep learning [7, 8] has attracted a lot of attention due to its amazing performance to learn representations on various tasks, especially in computer vision and natural language processing. Hence, some works make use of deep learning to learn effective features from side information to improve the performance of recommender system, such as [14, 24, 28, 31].

Neural Networks have attracted little attention in the collaborative filtering community. Salakhutdinov *et al.* [24] were the first work to tackle the Netflix challenge using Restricted Boltzmann Machines (RBM). They modified the RBM as a two-layer undirected graphical model consisting of binary hidden units and softmax visible units, and tested their model on the Netflix dataset and showed a comparable result with the start-of-the-art. On music recommendation, Van *et al.* [28] directly use Convolutional Neural Network (CNN) to learn effective representations of songs and use them in content-based collaborative filtering framework. Wang *et al.* [31] directly coupled matrix factorization with deep learning models, and proposed a hierarchical Bayesian model called Collaborative Deep Learning (CDL) which tightly couples Stacked Denoising AutoEncoders (SDA) [29] and Collaborative Topic Regression (CTR) [30] to solve the cold start problem. Li *et al.* [14] proposed a generate learning framework to combine rating matrix and side information, they used Probabilistic matrix factorization (PMF) [23] and marginalized Stacked Denoising AutoEncoders (mSDA) [2] in their approach, which is close to CDL but more efficient and scalable. Meanwhile, their model can learn deep features for both items and users while CDL only extracts deep features for items. In our work, although sequential features learned by LSIM are used as side information in the feature based collaborative framework, but actually they aims to describe the long-short interest, and can't be utilized to solve the cold start problem.

2.2 Neural Network Language Model

Traditional language model uses a one-hot representation to represent each word as a feature vector, where these feature vectors have the same length as the size of vocabulary, and the position that corresponds to the observed word is equal to 1, and 0 otherwise. However, this approach often exhibits significant limitations in practical tasks, suffering from high dimensionality and severe data sparsity.

Mikolov *et al.* [16, 17] proposed the word2vec algorithm to address these issues. They take advantage of the word order in text documents, explicitly modeling the assumption that closer words in the word sequence are statistically more dependent, and have generalized the classic n-gram language models by using continuous variables to represent words in a vector space. The continuous bag-of-words (CBOW) and skip-gram (SG) language models are highly scalable for learning word representations from large-scale corpora. The word2vec algorithm breaks the semantic gap between words. For example, “trade” and “deal” are totally different words in the one-hot representation, but they are similar in word2vec distribution representation. Le *et al.* [11] followed the above work and proposed the paragraph2vec algorithm to simultaneously learn vector representations of sentence and words by considering the sentence as a “global context”. Our long-short interest model shares similar idea with paragraph2vec algorithm, but we aim to simultaneously learn vector representations of user and items correspondingly by considering the user as a “global context”.

3 Our Approach

In this section, we first describe the definition of the rating prediction task and the notation we are going to use in this paper. Then we introduce our long-short interest model to extract the sequential features of users and items based on the long-short interest assumption. In the last, we utilize the sequential features as side information in the feature based collaborative filtering framework to make the final prediction.

3.1 Problem Definition

Given N users and M items, the rating r_{ij} is the rating given by the i^{th} user for the j^{th} item. In the common real-world situations, users usually rate on a fraction of items, not on the whole items. Therefore, those ratings entail a big and sparse matrix $\mathbf{R} \in \mathbb{R}^{N \times M}$. The goal of recommender system is to make a prediction on the missing ratings. Based on that, we will know the preference of a user on the items he never rates, and recommend high score items to him.

Matrix Factorization is a classic method to solve this problem. It aims to find a K dimensional low rank matrix $\hat{\mathbf{R}} \in \mathbb{R}^{N \times M}$ where $\hat{\mathbf{R}} = \mathbf{U}\mathbf{V}^T$ with $\mathbf{U} \in \mathbb{R}^{N \times K}$ and $\mathbf{V} \in \mathbb{R}^{M \times K}$ are two matrices of rank K encoding a dense representation of the users and items with

$$\underset{\mathbf{U}, \mathbf{V}}{\operatorname{argmin}} \sum_{(i,j) \in \mathcal{K}(\mathbf{R})} (r_{ij} - \mathbf{u}_i^T \mathbf{v}_j)^2 + \lambda (\|\mathbf{u}_i\|_{Fro}^2 + \|\mathbf{v}_j\|_{Fro}^2) \quad (1)$$

where $\mathcal{K}(\mathbf{R})$ is the set of indices of known ratings, \mathbf{u}_i and \mathbf{v}_j are the corresponding line vectors of \mathbf{U} and \mathbf{V} , λ is the coefficient that controls the influence of L2 regularization, and $\|\cdot\|_{Fro}$ is the Frobenius norm.

Table 1 summarizes the symbols used in our approach. In the next section, we will propose a long-short interest model to extract sequential features of users and items based on the long-short interest assumption.

Table 1. Summary of notations.

Notation	Description
N	Number of users
M	Number of items
K	Dimension of latent factors
D	Dimension of sequential features
$\mathbf{R} \in \mathbb{R}^{N \times M}$	Rating matrix
$\mathbf{U} \in \mathbb{R}^{N \times K}$	Latent factors of users
$\mathbf{V} \in \mathbb{R}^{M \times K}$	Latent factors of items
$\mathbf{X} \in \mathbb{R}^{N \times D}$	sequential features of users
$\mathbf{Y} \in \mathbb{R}^{M \times D}$	sequential features of items

3.2 Long-short Interest Model

Collaborative filtering aims at estimating the ratings a user is going to give to all other items he never interact with by using the ratings of all the other users. The basic assumption of collaborative filtering is that items liked by the same user would be similar or users like same items would share similar interest. However, in real-world situations, it is not always true because users' interest may change over a long time period. Meanwhile, the interest distribution of a user in a fixed time period are stable and don't change too much.

To describe this phenomenon that interest changes over a long time period but keep stable in a short time period, we propose the definition of **long interest** and **short interest**.

- **long interest** reflects the interest distribution of a user in a long time period, and it is reflected in the whole items list of the user's preference.
- **short interest** reflects the interest distribution of a user in a short time period, and it is reflected in a fraction of the whole items list of the user's preference in a fixed length sliding window.

Under this definition, each user will have a long interest, and each long interest will correspond several short interest. At the same time, we propose two assumption based on long interest and short interest.

1. items liked in the same short interest of the same long interest have a higher possibility to be similar than ones liked in different short interest of the same long interest.
2. the more times items show in the same short interest of different long interest, the higher possibility they are similar.

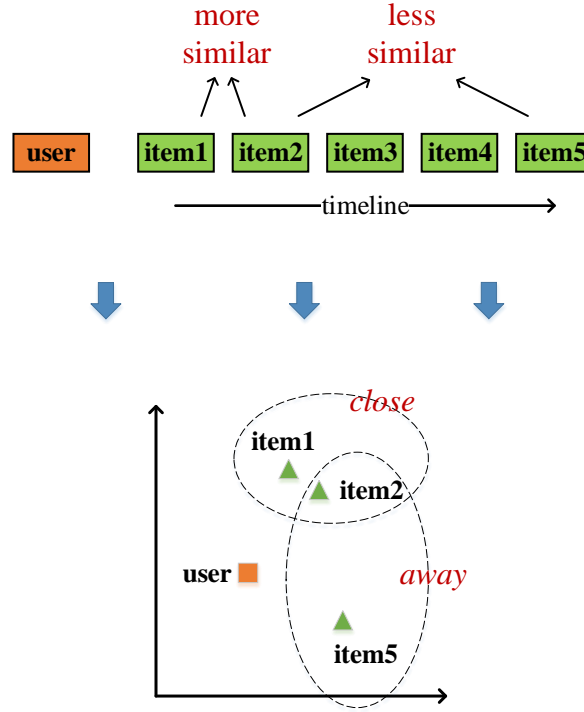


Fig. 3. A Example for Short Interest and Long Interest

For example, as illustrated in the top of Figure 3, given a certain **user**, we can see five items **item1**, **item2**, ..., **item5**, he liked in the last, are list on the timeline according to the liked time order. Those items compose the long interest of this **user**. Let's assume that the time window size of short interest is set to 2, That means **item1** and **item2** are in the same short interest, **item2** and **item3** are in the same short interest, and so on. Meanwhile **item2** and **item5** are in

different short interest. Therefore **item2** has a higher possibility to be similar with **item1** than **item5** based on our long-short interest assumption.

In special, we embed users and items into a low dimensional space to characterize these similarity in mathematical sense. We use European space distance between items in this low dimensional space to represent the similarity between them. Hence, **item2** has a higher possibility to be similar with **item1** than **item5** means **item2** is close to **item1** and far away from **item5** in this low dimensional space, as shown in the bottom of Figure 3.

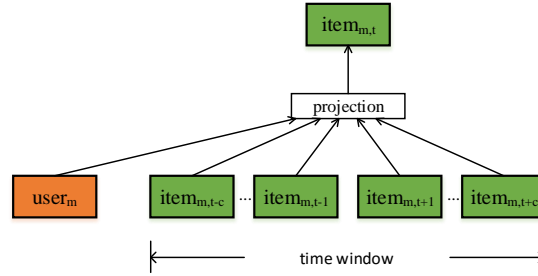


Fig. 4. Embedding Model for Extracting Interest Similarity from Users and Items

Inspired by paragraph2vec algorithm [11] for learning vector representations of words which take advantage of a word order observed in a sentence, we introduce a neural network based language model to carry out the embedding of sequential features. The embedding model simultaneously learns vector representations of users and items by considering the user as a global context, and the architecture of the embedding model is illustrated in Figure 4.

The training data set was derived from users interaction timeline T , which comprises users $x_i (i = 1, 2, \dots, N)$ and their interacted items ordered by the interacted time, $y_{i_1}, y_{i_2}, \dots, y_{i_{L_i}}$ ¹, where L_i denotes number of items interacted by user x_i , which is much less than the amount of items M . To characterize the *long interest*, we consider the whole items list as the context and generate the long interest of the current user. To characterize the *short interest*, we consider items in the same local interest as the context and generate items in it one by one with the help of long interest. More formally, objective of the embedding model is to maximize the log-likelihood over the set of T of all the interaction timeline,

¹ we use symbol x and y instead of classic u and v to avoid confusion between v and vector symbol \mathbf{v} in neural network language model.

$$\sum_{i=1}^N \left(p(x_i | y_{i_1}, y_{i_2}, \dots, y_{i_{L_i}}) + \sum_{j=1}^{L_i} p(y_{i_j} | y_{i_{j-c}} : y_{i_{j+c}}, x_i) \right) \quad (2)$$

where c is the time window size, $y_{i_{j-c}} : y_{i_{j+c}}$ denotes the sequence $y_{i_{j-c}}, y_{i_{j-c+1}}, \dots, y_{i_{j+c}}$ excluding y_{i_j} .

$p(x_i | y_{i_1}, y_{i_2}, \dots, y_{i_{L_i}})$ is the probability to generate the long interest of u_i based on all items he interacted. The prediction task is typically done via a multiclass classifier, such as softmax. There, we have

$$p(x_i | y_{i_1}, y_{i_2}, \dots, y_{i_{L_i}}) = \frac{\exp(\bar{\mathbf{v}}_1^T \mathbf{v}'_{x_i})}{\sum_{x'} \exp(\bar{\mathbf{v}}_1^T \mathbf{v}'_{x'})} \quad (3)$$

where \mathbf{v}'_{x_i} is the output vector representation of x_i , and $\bar{\mathbf{v}}_1$ is averaged input vector representation of all the items interacted by user x_i , i.e.

$$\bar{\mathbf{v}}_1 = \frac{\sum_{j=1}^{T_i} \mathbf{v}_{y_{i_j}}}{T_i} \quad (4)$$

$p(y_{i_j} | y_{i_{j-c}} : y_{i_{j+c}}, x_i)$ is the probability to generate y_{i_j} based on items in the same short interest and the user's long interest. Similarly, using softmax multiclass classifier we have

$$p(y_{i_j} | y_{i_{j-c}} : y_{i_{j+c}}, x_i) = \frac{\exp(\bar{\mathbf{v}}_2^T \mathbf{v}'_{y_{i_j}})}{\sum_{y'} \exp(\bar{\mathbf{v}}_2^T \mathbf{v}'_{y'})} \quad (5)$$

where $\mathbf{v}'_{y_{i_j}}$ is the output vector representation of y_{i_j} , and $\bar{\mathbf{v}}_2$ is averaged input vector representation of items in the same short interest and corresponding long interest x_i .

$$\bar{\mathbf{v}}_2 = \frac{\mathbf{v}_{x_i} + \sum_{-c \leq k \leq c, k \neq 0} \mathbf{v}_{y_{i_{j+k}}}}{2c + 1} \quad (6)$$

Stochastic Gradient Descent (SGD) are used as the training method, hierarchical softmax and negative sampling are two main approaches to accelerate the computation, and we use negative sampling approach in this paper.

3.3 Feature based Collaborative Filtering

Feature based collaborative filtering [3, 4] is a variety of collaborative filtering, it allows us to build factorization models incorporating side information such as temporal dynamics, neighborhood relationship, and hierarchical information compare to conventional collaborative filtering, it can also be capable of both rate prediction and collaborative ranking.

There are two kinds of side information in collaborative filtering: user side information and item side information. User side information could be the profiles of users, such as gender, age and occupation. Item side information is usually the properties of items, it mainly depends on the recommendation sceneries. In movie recommendation, item side information could be actor, director and genre of the movie. In product recommendation, item side information could be the price and category of the product. Feature based collaborative filtering summarizes the two factors as feature vectors (denoted by $\mathbf{u}_i \in \mathbb{R}^n$ and $\mathbf{v}_j \in \mathbb{R}^m$) and predicts the preference score \hat{r} as

$$\hat{r}_{ij} = \sum_{k=1}^n \alpha_k \mathbf{u}_{ik} + \sum_{k=1}^m \beta_k \mathbf{v}_{jk} + \left(\sum_{k=1}^n \mathbf{u}_{ik} \mathbf{p}_k \right)^T \left(\sum_{k=1}^m \mathbf{v}_{jk} \mathbf{q}_k \right) \quad (7)$$

where α and β controls the influence of each feature, $\mathbf{p}_k \in \mathbb{R}^K$ and $\mathbf{q}_k \in \mathbb{R}^K$ are K dimensional latent factors associated with each feature. If we represent one-hot representation of users and items like

$$\mathbf{u}_{ik} = \begin{cases} 1, k = i \\ 0, k \neq i \end{cases} \quad (8)$$

$$\mathbf{v}_{jk} = \begin{cases} 1, k = j \\ 0, k \neq j \end{cases} \quad (9)$$

the equation for rating prediction will reduce to the basic matrix factorization

$$\hat{r}_{ij} = \alpha_i + \beta_j + \mathbf{p}_i^T \mathbf{q}_j \quad (10)$$

where α , β are the biases and \mathbf{p}_i , \mathbf{q}_j are the latent factors for user \mathbf{u}_i and item \mathbf{v}_j .

Various kinds of side information can be utilized to enhance these factors. In some sense, the representation of users and items learned from long-short interest model can be treated as a kind of side information, because it distinguish the similarity and difference between users and items.

First, we define $\tilde{\mathbf{u}}_i$ as the learned vector representation of users and $\tilde{\mathbf{v}}_j$ as the learned vector representation of items. Then, we get new features of users and items by add the learned vector representation to origin one-hot representation.

$$\mathbf{u}_i = \{\mathbf{u}_i, \tilde{\mathbf{u}}_i\} \quad (11)$$

$$\mathbf{v}_j = \{\mathbf{v}_j, \tilde{\mathbf{v}}_j\} \quad (12)$$

There, the equation for rating prediction will change to

$$\begin{aligned}
\hat{r}_{ij} &= \sum_{k=1}^{N+D} \alpha_k \{\mathbf{u}_i, \tilde{\mathbf{u}}_i\}_k + \sum_{k=1}^{M+D} \beta_k \{\mathbf{v}_j, \tilde{\mathbf{v}}_j\}_k + \\
&\quad \left(\sum_{k=1}^{N+D} \{\mathbf{u}_i, \tilde{\mathbf{u}}_i\}_k \mathbf{p}_k \right)^T \left(\sum_{k=1}^{M+D} \{\mathbf{v}_j, \tilde{\mathbf{v}}_j\}_k \mathbf{q}_k \right) \\
&= \sum_{k=N+1}^{N+D} \alpha_k \tilde{\mathbf{u}}_{ik} + \sum_{k=M+1}^{M+D} \beta_k \tilde{\mathbf{v}}_{jk} + \\
&\quad \left(\sum_{k=N+1}^{N+D} \tilde{\mathbf{u}}_{ik} \mathbf{p}_k \right)^T \left(\sum_{k=M+1}^{M+D} \tilde{\mathbf{v}}_{jk} \mathbf{q}_k \right) + \\
&\quad \alpha_i + \beta_j + \mathbf{p}_i^T \mathbf{q}_j
\end{aligned} \tag{13}$$

where N is the number of users, M is the number of items, D is the dimension of sequential features of users and items learned from our long-short interest model, $\mathbf{p}_k \in \mathbb{R}^K$ and $\mathbf{q}_k \in \mathbb{R}^K$ are K dimensional latent factors associated with each feature.

4 Experiment

In this section, we conduct several experiments to evaluate the effectiveness of our proposed long-short interest model on three public MovieLens² datasets. In these experiments, we also conduct corresponding analysis to investigate: (1) the rating prediction performance of our long-short interest model compare to other benchmark models; (2) the effect of our long-short interest model on users own different interaction number.

4.1 Experimental Setup

Dataset We conduct experiments on three MovieLens datasets, i.e. MovieLens-1M, MovieLens-10M and MovieLens-20M, which are commonly used for evaluating collaborative filtering algorithms.

The MovieLens-1M dataset consists of about 1 million anonymous ratings of 3706 movies made by 6040 MovieLens users who joined MovieLens in 2000, and each rating is an integer between 1 (worst) and 5 (best). Timestamp is represented in seconds and each user has at least 20 ratings.

The MovieLens-10M dataset contains about 10 million ratings and 95580 tags applied to 10677 movies by 69878 users of the online movie recommender service MovieLens. Users were selected at random for inclusion. All users selected had rated at least 20 movies. Unlike previous MovieLens data sets, no demographic information is included. Ratings are made on a 5-star scale, with half-star increments, i.e., 0.5 is the worst score and 5.0 is the best score.

² <http://grouplens.org/datasets/movielens/>

the MovieLens-20M dataset describes 5-star rating and free-text tagging activity. It contains about 20 million ratings and 465564 tag applications across 26744 movies. These data were created by 138493 users between January 09, 1995 and March 31, 2015. This dataset was generated on March 31, 2015. Users were selected at random for inclusion. All selected users had rated at least 20 movies. No demographic information is included. Each user is represented by an id, and no other information is provided.

These datasets also provide some side information about users and items, such as gender, age and occupation of users and movies' genres, but we don't use the side information in all our experiments. Table 2 summarizes the statistics of three datasets.

Table 2. Statistics of three MovieLens datasets and Netflix dataset.

Dataset	#Users	#Items	#Ratings	Sparsity
ML-1M	6,040	3,706	1,000,209	95.53%
ML-10M	69,878	10,677	10,000,054	98.66%
ML-20M	138,493	26,744	20,000,263	99.46%

Metrics We employ two metrics, the Root Mean Square Error (RMSE) and the Mean Absolute Error (MAE), to measure the rating prediction quality of our proposed approach in comparison with other collaborative filtering based recommendation methods. The metrics RMSE is defined as:

$$RSME = \sqrt{\frac{1}{N} \sum_{i,j} I_{ij} (R_{ij} - \hat{R}_{ij})^2} \quad (14)$$

where R_{ij} denotes the ground-truth rating the user i gives to the item j , \hat{R}_{ij} denotes the corresponding predicted rating the user i gives to the item j by a method, I_{ij} is a binary matrix that indicates the ratings in the test set, and N is the total number of ratings in the test set. The metrics MAE is defined as:

$$MAE = \frac{1}{N} \sum_{i,j} I_{ij} |R_{ij} - \hat{R}_{ij}| \quad (15)$$

4.2 Benchmark Models

To demonstrate the superiority of our long-short interest model for rating prediction task, we compare our model with several benchmark models. In special, we use two popular toolkits, LibMF and LibFM, that are widely used in both academia and industry as our benchmark models.

LibMF is an open source tool for approximating an incomplete matrix using the product of two matrices in a latent space. It provides solvers for real-valued

matrix factorization, binary matrix factorization, and one-class matrix factorization, and also supports parallel computation in a multi-core machine using CPU instructions (e.g., SSE) to accelerate vector operations. Its paper [5] won the best paper award in RecSys 2013.

Factorization machines (FM) are a generic approach that allows to mimic most factorization models by feature engineering. This way, factorization machines combine the generality of feature engineering with the superiority of factorization models in estimating interactions between categorical variables of large domain. LibFM [20] is a software implementation for factorization machines that features Stochastic Gradient Descent (SGD) and Alternating Least Squares (ALS) optimization as well as Bayesian inference using Markov Chain Monte Carlo (MCMC).

4.3 Overall Results

In this section, we report the experimental results to demonstrate the effectiveness of our long-short interest model.

The model proposed in LibMF toolkit is denoted as **LibFM**, the dimension of latent factors is set to 100 and the number of iterations is set to 1000.

For the model proposed in LibFM, we denote the LibFM model optimized by SGD as **LibFM-SGD**, the learn rate is set to 0.01, the regular parameter is set to (0,0,0.01) and the stdev for initialization of 2-way factors is set to 0.1. Then we denote the LibFM model optimized by ALS as **LibFM-ALS**, the regular parameter is set to (0,0,10) and the stdev for initialization of 2-way factors is set to 0.1. Lastly, the LibFM model optimized by MCMC is denoted as **LibFM-MCMC**, the stdev for initialization of 2-way factors is set to 0.1.

Meanwhile, we denote our long-short interest model as **LSIM**, the slide window size is set to 10, the dimension of latent factors is set to 256 and the dimension of sequential features is set to 100. For feature based collaborative filtering framework, the learn rate is set to 0.005 and the regular parameter is set to 0.024 for both users and items.

We split the rating data in each dataset into random 90%-10% training-test datasets, the training dataset are used for building our proposed model and benchmark models, the remaining data are used for testing. This process is repeated five times, and we report the mean value and standard deviation of RMSE and MAE.

Table 3 shows the RSME performance of our proposed model and benchmark methods with a training ratio of 90% / 10% on three MovieLens datasets while Table 4 shows the RMSE performance of these models in the same experimental setting respectively. The best performances are marked in bold typeface. From the two tables, we can clearly observe:

1. The rating prediction performance of **LibMF** is worst among all the models, the reason should be that it focuses on accelerating the training speed of matrix factorization by parallelization, and doesn't pay enough attention on optimizing the prediction precision.

Table 3. RMSE Performance comparison of our proposed model and benchmark models with a training ratio of 90%/10% on three MovieLens datasets.

Algorithms	MovieLens-1M	MovieLens-10M	MovieLens-20M
LibMF	0.8554 ± 0.0013	0.8090 ± 0.0004	0.8023 ± 0.0004
LibFM-SGD	0.8641 ± 0.0015	0.8022 ± 0.0013	0.7945 ± 0.0023
LibFM-ALS	0.8453 ± 0.0015	0.7936 ± 0.0004	0.7860 ± 0.0004
LibFM-MCMC	0.8460 ± 0.0011	0.7866 ± 0.0004	0.7787 ± 0.0005
LSIM	0.8355 ± 0.0014	0.7740 ± 0.0013	0.7656 ± 0.0019

Table 4. MAE Performance comparison of our proposed model and benchmark models with a training ratio of 90%/10% on three MovieLens datasets.

Algorithms	MovieLens-1M	MovieLens-10M	MovieLens-20M
LibMF	0.6816 ± 0.0008	0.6311 ± 0.0003	0.6229 ± 0.0004
LibFM-SGD	0.6674 ± 0.0027	0.6127 ± 0.0034	0.6016 ± 0.0036
LibFM-ALS	0.6609 ± 0.0010	0.6068 ± 0.0002	0.5971 ± 0.0003
LibFM-MCMC	0.6661 ± 0.0008	0.6039 ± 0.0002	0.5941 ± 0.0005
LSIM	0.6539 ± 0.0010	0.5928 ± 0.0008	0.5837 ± 0.0018

2. In three **LibFM** models, MCMC optimization **LibFM-MCMC** shows a better performance in large datasets (MovieLens-10M and MovieLens-20M) while ALS optimization **LibFM-ALS** shows a better performance in small datasets (MovieLens-1M). The performance of SGD optimization **LibFM-SGD** is the worst.
3. Our **LSIM** has the best performance between among all models in both RMSE and MAE performance on all three MovieLens datasets, which shows the effectiveness of the long-short interest model. To the best of our knowledge, the best results published regarding MovieLens-1M and MovieLens-10M are reported by both [13, 26] with a final RMSE of 0.831 ± 0.003 and 0.782 ± 0.003 . These scores are obtained with a training ratio of 90%/10% and without side information. That means our proposed model can achieve the start-of-the-art performance.

4.4 Effects on Different Users

In our long-short interest model, the change of users' interest over time has been embedding into the sequential features. For these users with a lot of interactions, their long interest and short interest will be detailedly extracted via LSIM. But for the others with few interactions, LSIM only can extract marginal sequential features because their interest keeps stable. This phenomenon means that our LSIM will show a better performance on users with lots of interactions than ones with few interactions.

In order to exhibit this character of LSIM, we carry on some experiments on the MovieLens-1M dataset. First, we sort the users on MovieLens-1M according to their respective number of ratings, those users are grouped by their ranking order into 5 clusters (i.e. 0%-20%, 20%-40%, 40%-60%, 60%-80% and

80%-100%), and RMSE and MAE are computed by cluster respectively with a training ratio of 90%/10%. For instance, the first cluster contains the 20% of users with the least number of ratings and the last cluster contains the 20% of users with the highest number of ratings. The provided results are the mean reported through 5-cross validation.

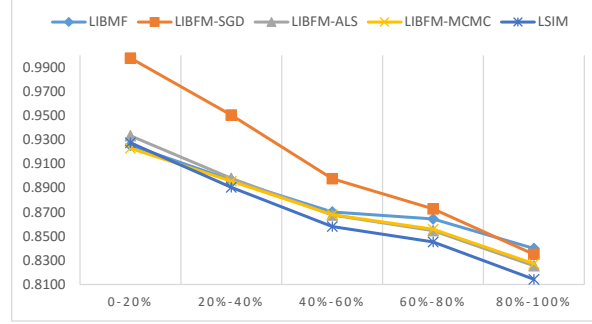


Fig. 5. RMSE computed by cluster of users sorted by their respective number of ratings on MovieLens-10M with a training ratio of 90%/10%.

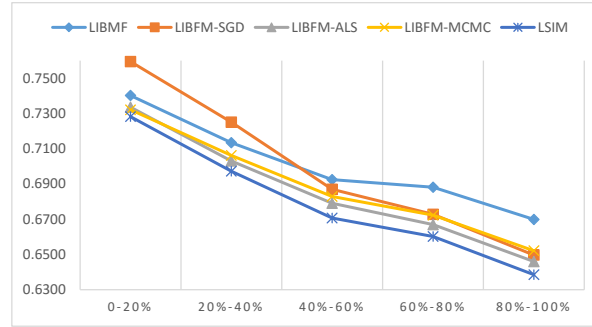


Fig. 6. MAE computed by cluster of users sorted by their respective number of ratings on MovieLens-10M with a training ratio of 90%/10%.

Figure 5 shows the RMSE performance of all recommendation models in different user clusters on MovieLens-1M dataset while Figure 6 shows the MAE performance of all recommendation models in the same experimental setting respectively. It can be clearly observed from the two tables that all these recommendation models benefit from the increase of the count of items the user interact with on both RSME and MAE metrics. Our proposed **LSIM** don't per-

form well at the first “0%-20%” user cluster, its RMSE score is 0.9275 which is less than **LibFM-MCMC**. But when the interactions become more, its performance rapidly increases, and beats the baseline methods obviously, which shows the effectiveness of our long-short interest model.

5 Conclusions and Future Works

In this study, we propose to use long-short interest model to utilize the rich sequential information on users interaction history to solve the rating prediction task in recommender system. By incorporating the sequential features into feature based collaborative filtering framework, better prediction performance can be obtained. Our thorough evaluation, using three standard MovieLens datasets, demonstrates the effectiveness of the proposed method.

Many studies remain for the future work. (1) When we carry out the long-short interest model, items given different ratings are treated equally, and ratings information is ignored. Hence, how to take advantage of the ratings information is one of our future works. (2) Side information plays an important role in recommender system, especially for solving the cold start problem. In our long-short interest model, we don’t use any side information, but it is certain that with the help of side information, users’ long-short interest can be learned more effectively. Therefore, we will focus on making use of side information in the future.

Bibliography

- [1] J. Bennett and S. Lanning. The netflix prize. In *Proceedings of KDD cup and workshop*, volume 2007, page 35, 2007.
- [2] M. Chen, Z. Xu, K. Weinberger, and F. Sha. Marginalized denoising autoencoders for domain adaptation. *arXiv preprint arXiv:1206.4683*, 2012.
- [3] T. Chen, W. Zhang, Q. Lu, K. Chen, Z. Zheng, and Y. Yu. Svdfeature: a toolkit for feature-based collaborative filtering. *The Journal of Machine Learning Research*, 13(1):3619–3622, 2012.
- [4] T. Chen, Z. Zheng, Q. Lu, W. Zhang, and Y. Yu. Feature-based matrix factorization. *arXiv preprint arXiv:1109.2271*, 2011.
- [5] W.-S. Chin, Y. Zhuang, Y.-C. Juan, and C.-J. Lin. A fast parallel stochastic gradient method for matrix factorization in shared memory systems. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 6(1):2, 2015.
- [6] C. A. Gomez-Urbe and N. Hunt. The netflix recommender system: Algorithms, business value, and innovation. *ACM Transactions on Management Information Systems (TMIS)*, 6(4):13, 2015.
- [7] G. E. Hinton, S. Osindero, and Y.-W. Teh. A fast learning algorithm for deep belief nets. *Neural computation*, 18(7):1527–1554, 2006.
- [8] G. E. Hinton and R. R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, 2006.
- [9] Y. Koren, R. Bell, and C. Volinsky. Matrix factorization techniques for recommender systems. *Computer*, (8):30–37, 2009.
- [10] N. D. Lawrence and R. Urtasun. Non-linear matrix factorization with gaussian processes. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 601–608. ACM, 2009.
- [11] Q. V. Le and T. Mikolov. Distributed representations of sentences and documents. *arXiv preprint arXiv:1405.4053*, 2014.
- [12] D. D. Lee and H. S. Seung. Algorithms for non-negative matrix factorization. In *Advances in neural information processing systems*, pages 556–562, 2001.
- [13] J. Lee, S. Kim, G. Lebanon, and Y. Singer. Local low-rank matrix approximation. In *Proceedings of The 30th International Conference on Machine Learning*, pages 82–90, 2013.
- [14] S. Li, J. Kawale, and Y. Fu. Deep collaborative filtering via marginalized denoising auto-encoder. In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*, pages 811–820. ACM, 2015.
- [15] R. Mazumder, T. Hastie, and R. Tibshirani. Spectral regularization algorithms for learning large incomplete matrices. *The Journal of Machine Learning Research*, 11:2287–2322, 2010.
- [16] T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.

- [17] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013.
- [18] M. J. Pazzani and D. Billsus. Content-based recommendation systems. In *The adaptive web*, pages 325–341. Springer, 2007.
- [19] S. Rendle. Factorization machines. In *Data Mining (ICDM), 2010 IEEE 10th International Conference on*, pages 995–1000. IEEE, 2010.
- [20] S. Rendle. Factorization machines with libfm. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 3(3):57, 2012.
- [21] 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*, pages 175–186. ACM, 1994.
- [22] F. Ricci, L. Rokach, and B. Shapira. *Introduction to recommender systems handbook*. Springer, 2011.
- [23] R. Salakhutdinov and A. Mnih. Probabilistic matrix factorization. Citeseer, 2011.
- [24] R. Salakhutdinov, A. Mnih, and G. Hinton. Restricted boltzmann machines for collaborative filtering. In *Proceedings of the 24th international conference on Machine learning*, pages 791–798. ACM, 2007.
- [25] 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*, pages 285–295. ACM, 2001.
- [26] S. Sedhain, A. K. Menon, S. Sanner, and L. Xie. Autorec: Autoencoders meet collaborative filtering. In *Proceedings of the 24th International Conference on World Wide Web Companion*, pages 111–112. International World Wide Web Conferences Steering Committee, 2015.
- [27] X. Su and T. M. Khoshgoftaar. A survey of collaborative filtering techniques. *Advances in artificial intelligence*, 2009:4, 2009.
- [28] A. Van den Oord, S. Dieleman, and B. Schrauwen. Deep content-based music recommendation. In *Advances in Neural Information Processing Systems*, pages 2643–2651, 2013.
- [29] P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning*, pages 1096–1103. ACM, 2008.
- [30] C. Wang and D. M. Blei. Collaborative topic modeling for recommending scientific articles. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 448–456. ACM, 2011.
- [31] H. Wang, N. Wang, and D.-Y. Yeung. Collaborative deep learning for recommender systems. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1235–1244. ACM, 2015.
- [32] Y. Zhou, D. Wilkinson, R. Schreiber, and R. Pan. Large-scale parallel collaborative filtering for the netflix prize. In *Algorithmic Aspects in Information and Management*, pages 337–348. Springer, 2008.