

# Online Interactive Collaborative Filtering Using Multi-armed Bandit with Dependent Arms

Qing Wang, Chunqiu Zeng, Wubai Zhou,

Tao Li, S. S. Iyengar, Larisa Shwartz, and Genady Ya. Grabarnik

**Abstract**—Online interactive recommender systems strive to promptly suggest users appropriate items (e.g., movies, news articles) according to the current context including both user and item content information. Such contextual information is often unavailable in practice, where only the users’ interaction data on items can be utilized by recommender systems. The lack of interaction records, especially for new users and items, inflames the performance of recommendation further. To address these issues, both collaborative filtering, one of the most popular recommendation techniques relying on the interaction data only, and bandit mechanisms, capable of achieving the balance between exploitation and exploration, are adopted into an online interactive recommendation setting assuming independent items (i.e., arms). This assumption rarely holds in reality, since the real-world items tend to be correlated with each other. In this paper, we study online interactive collaborative filtering problems by considering the dependencies among items. We explicitly formulate item dependencies as the clusters of arms in the bandit setting, where the arms within a single cluster share the similar latent topics. In light of topic modeling techniques, we come up with a novel generative model to generate the items from their underlying topics. Furthermore, an efficient particle-learning based online algorithm is developed for inferring both latent parameters and states of our model by taking advantage of the fully adaptive inference strategy of particle learning techniques. Additionally, our inferred model can be naturally integrated with existing multi-armed selection strategies in an interactive collaborative filtering setting. Empirical studies on two real-world applications, online recommendations of movies and news, demonstrate both the effectiveness and efficiency of our proposed approach.

**Index Terms**—Recommender systems; Interactive collaborative filtering; Topic modeling; Cold-start problem; Particle learning.

## 1 INTRODUCTION

THE overwhelming amount of data requires an efficient online interactive recommendation system where online users constantly interact with the system, and user feedback is instantly collected to improve recommendation performance. Online interactive recommender systems are challenged to immediately recommend the most proper items (e.g., movies, news articles) to users based on the

current user and item content information aiming to continuously maximize users’ satisfaction over a long run. To achieve this goal, it becomes a critical task for such recommender systems to constantly track user preferences and recommend interesting items from a large item repository.

In the process of identifying the appropriate match between the user preferences and the target items the recommender systems encounter difficulties due to several existing practical challenges. One challenge is the well-known *cold-start* problem since a number of users/items might be completely new to the system, that is, they may have no historical records at all. This problem makes recommender systems ineffective unless additional information including both items and users is collected [45], [12]. The second challenge is that most recommender systems typically assume the entire set of contextual features with respect to both users and items can be accessed to infer users’ preference. Due to a number

- 
- Q. Wang, C. Zeng, W. Zhou, T. Li and S. S. Iyengar are with the School of Computing and Information Sciences, Florida International University, Miami, FL, 33199.  
E-mail: {qwang028, czeng001, wzhou005, taoli, iyengar}@cs.fiu.edu
  - Larisa Shwartz is with the department of Cognitive Service Foundations, IBM T.J. Watson Research Center, Yorktown Heights, NY, 10598.  
E-mail: lshwart@us.ibm.com
  - Genady Ya. Grabarnik is with the department of Math & Computer Science, St. John’s University, Queens, NY, 11439.  
E-mail: grabarnig@stjohns.edu

of reasons (e.g., privacy or sampling constraints), it is challenging to obtain all relevant features ahead of time, thus rendering many factors unobservable to recommendation algorithms.

In the first challenge, an *exploration or exploitation dilemma* [6] is identified in the aforementioned setting. A tradeoff between two competing goals needs to be considered in recommender systems: maximizing user satisfaction using their consumption history, while gathering new information for improving the goodness of match between user preferences and items [24]. This dilemma is typically formulated as a multi-armed bandit problem where each arm corresponds to an item. The recommendation algorithm determines the strategies for selecting an arm to pull according to the contextual information at each trial. Pulling an arm indicates that the corresponding item is recommended. When an item matches the user preference (e.g., a recommended news article or movie is consumed), a reward is obtained; otherwise, no reward is provided. The reward information is fed back to optimize the strategies. The optimal strategy is to pull the arm with the maximum expected reward based on the historical interaction on each trial, and then to maximize the total accumulated rewards for the whole series of trials.

Collaborative filtering (CF) is widely applied in recommender systems [34], [7], [22] to address the second challenge. CF has gained its popularity due to its advantage over other recommendation techniques, where CF requires no extra information about items or users for recommendation but only users' historical ratings on items [20], [21]. Further, considering both aforementioned challenges simultaneously aggravates the difficulties when recommending items. Recently, an online interactive collaborative filtering system has been suggested [21], [47] adopting both techniques, multi-armed bandit and collaborative filtering. Typically, one collaborative filtering task is formulated as a matrix factorization problem. Matrix factorization derives latent features for both users and items from the historical interaction records. It assumes that a user's preference (i.e., rating) on a given item can be predicted considering items' and users' latent feature vectors. Based on this assumption, multi-armed bandit policies make use of the predicted reward (i.e., user preference) for arm (i.e., item) selection. The feedback occurring between the current user and arm is used to update the user's and arm's latent vectors, without impacting the inference of other arms' latent vectors assuming arms are inde-

pendent from each other. However, the assumption about the independency among arms rarely holds in real-world applications. For example, in the movie recommendation scenario, each movie corresponds to an arm. The dependent arms (i.e., movies) typically share similar latent topics (e.g., science fiction movies, action movies, etc.), and are likely to receive similar rewards (i.e., ratings or feedback) from users. Intuitively, the dependencies among arms can be utilized for reward prediction improvement and further facilitated the maximization of users' satisfaction in a long run.

In this paper, we introduce an interactive collaborative topic regression model that utilizes bandit algorithms with dependent arms to recommend appropriate items for target users. A sequential online inference method is proposed to learn the latent parameters and infer the latent states. We adopt a generative process based on topic model to explicitly formulate the arm dependencies as the clusters on arms, where dependent arms are assumed to be generated from the same cluster. Every time an arm is pulled, the feedback is not only used for inferring the involved user and item latent vectors, but it is also employed to update the latent parameters with respect to the arm's cluster. The latent cluster parameters further help with reward prediction for other arms in the same cluster. The fully adaptive online inference strategy of particle learning [11] allows our model to effectively capture arm dependencies. In addition, the learnt parameters can be naturally integrated into existing multi-arm selection strategies, such as *UCB* and *Thompson sampling*. We conduct empirical studies on two real-world applications, movie and news recommendations. The experimental results demonstrate the effectiveness of our proposed approach.

The rest of this paper is organized as follows. In Section 2, we provide a brief summary of prior work relevant to collaborative filtering, collaborative topic model, multi-armed bandit and the online inference with particle learning. We formulate the underlying problem in Section 3. The solution to the problem is presented in Section 4. Extensive evaluation results are reported in Section 5. Finally, Section 6 concludes the paper.

## 2 RELATED WORK

This section highlights the existing works in the literature that are closely relevant to our work.

## 2.1 Interactive Collaborative Filtering

In recommender systems, collaborative filtering (CF) has gained an extensive popularity in recent decades due to its capability of identifying the user preference from the historical interactions between users and items [21], [28], [31], [32], [33], [34], [43]. However, it is still an immense challenge to effectively predict preferences for new users. This challenge typically referred to as the *cold-start* problem [3], [10], [35]. A straightforward solution to address this issue involves two separated stages, where it first explicitly figures out the user profile, then makes further recommendation based on the established user profile [29], [30]. By contrast, some preliminary works, referred to as interactive collaborative filtering (ICF), have recently emerged as an alternative way to deal with the *cold-start* issue [21], [47]. These works do not explicitly fulfill the two stages separately, but formulate the recommendation problem as a multi-armed bandit problem, and then naturally integrate the two stages together by striking a balance between exploration and exploitation. Our work is primarily relevant to this research area addressing the ICF problem.

The ICF problem is first introduced in [47], where several multi-armed bandit algorithms (e.g., Thompson sampling [13], UCB [9]) are used for item recommendation in light of the user-item rating prediction with the probabilistic matrix factorization (PMF) framework [31]. However, the proposed method in [47] does not work in a completely online interactive mode since the multi-armed bandit algorithms partially rely on the latent item feature vector distributions, which are learnt with the offline Gibbs sampling in advance. In [21], an efficient Thompson sampling algorithm named particle Thompson sampling (PTS) addresses the ICF problem with Bayesian probabilistic matrix factorization (BPMF) [32] in a completely online mode. To reduce the reward prediction uncertainty, Wang et al. [39] incorporated the contextual features into the learned latent feature vectors for ICF problem. But these methods assume the latent item feature vectors in the ICF setting are independent. Although the work in [26] formulates the arm dependencies as an arm clustering problem, it fails to present an efficient online method to learn arm dependencies. By comparison, we explicitly learn the dependent arms with a generative topic model in the ICF setting and develop an efficient online solution capable of tracking the dependencies between arms as well as addressing the online recommendation.

Some recent studies explore the bandit dependencies for a group recommendation delivery by assuming that users in the same group react with similar feedback to the same recommended item [18], [37], [40], [42], [44]. Most existing works utilize the context information for users or predefined social network to build the user dependencies. Wu et al. [44] exploit social information to find dependency among users for improving the accuracy of reward prediction. Wang et al. [40] propose a context-aware collaborative bandit model, which could incorporate mutual influence among users directly for matrix completion. In [42], an interactive social recommendation model is proposed to predict the target user’s preference using a weighted combination of a user’s preferences and his/her friends’ preferences. A context-dependent clustering of bandits algorithm [17] is investigated, where the clusters over users are based on the current item content. Our work is orthogonal to those studies since we investigate the arm (item) dependencies in a bandit model rather than the dependencies among users. Wang et al. [41] come up with hierarchical multi-armed bandit algorithms leveraging the explicit taxonomy information of items for online recommendation. Our proposed method is capable of instantly learning the item dependencies during the online interactive recommendation process without explicit context information provided.

Multi-armed bandits are widely adopted in diverse applications such as online advertising [27], [45], web content optimization [1], [24], and robotic routing [5]. The core task of bandit problem is to balance the tradeoff between exploration and exploitation. A series of algorithms have been proposed to deal with this problem including  $\epsilon$ -greedy [38], UCB [9], [25], EXP3 [4], Thompson sampling [2]. In this paper, we model the ICF problem as a multi-armed bandit problem with dependent arms and proposes online recommendation methodologies based on UCB and Thompson sampling.

## 2.2 Sequential Online Inference

Our model leverages topic modeling [8] to formulate arm dependencies and sequential online inference to infer the latent states and learn the unknown parameters. Popular sequential learning methods include sequential monte carlo sampling [15], [19] and particle learning [11], [46].

Sequential Monte Carlo (SMC) sampling consists of a set of Monte Carlo methodologies to solve the filtering problem [16]. These methodologies allow inference of full posterior distributions

in general state space models, which may be both nonlinear and non-Gaussian.

Particle learning provides state filtering, sequential parameter learning and smoothing in a general class of state space models [11]. Particle learning is used to approximate the sequence of filtering and smoothing distributions in light of parameter uncertainty for a wide class of state space models. The central idea behind particle learning is to create a particle directly from the approximation to the joint posterior distribution of states and conditional sufficient statistics of fixed parameters in a fully-adapted resample-propagate framework. In this paper, we leverage the idea of particle learning for both latent state inference and parameter learning.

### 3 PROBLEM FORMULATION

In this section, we provide a mathematical formulation of the interactive collaborative filtering (ICF) problem into a bandit setting. Then we introduce a new generative model describing the dependency among arms (i.e. items). A glossary of notations mentioned in this paper is summarized in Table 1.

TABLE 1: Important Notations

Notation	Description
$M, N$	number of rows (users) and columns (items).
$\mathbf{R} \in \mathcal{R}^{M \times N}$	the rating matrix.
$\mathbb{S}(t)$	the sequence of $(n(t-1), r_{m,n(t-1)})$ observed until time $t$ .
$n(t)$	the recommended item index in the $t$ -th iteration.
$r_{m,t}$	the rating (reward) of the $m$ -th user by pulling the given item in the $t$ -th iteration.
$y_{m,t}$	the predicted rating for the $m$ -th user over given item in the $t$ -th iteration.
$\pi$	the policy to recommend items sequentially.
$R_\pi$	the cumulative rating (reward) of the policy $\pi$ .
$K$	the number of topics and the number of dimensions for latent vectors.
$\mathbf{p}_m \in \mathcal{R}^K$	the latent feature vector of the $m$ -th user.
$\mathbf{q}_n \in \mathcal{R}^K$	the latent feature vector of the $n$ -th item.
$\Phi_k \in \mathcal{R}^N$	the item distribution of the $k$ -th topic.
$\mathcal{P}_{m,n(t-1)}$	the set of particles for the item $n(t-1)$ given user $m$ at time $t-1$ .
$z_{m,t}$	the latent topic of the $m$ -th user in the $t$ -th iteration.
$x_{m,t}$	the selected item of the $m$ -th user in the $t$ -th iteration.
$\lambda$	Dirichlet priors over topics for topic model.
$\eta$	Dirichlet priors over items for topic model.
$\sigma_n^2$	the variance of rating prediction.
$\alpha, \beta$	the hyper parameters determine the distribution of $\sigma_n^2$ .
$\mu_q, \Sigma_q$	the hyper parameters determine the gaussian distribution of $\mathbf{q}_n$ .
$\xi$	the observation noise of the rating

#### 3.1 Basic Concepts and Terminologies

Assume that there are  $M$  users and  $N$  items in the system. The preferences of the users for the items are recorded by a partially observable matrix  $\mathbf{R} = \{r_{m,n}\} \in \mathcal{R}^{M \times N}$ , where the rating score  $r_{m,n}$  indicates how user  $m$  would like item  $n$ . The basic collaborative filtering task is to predict the unknown rating score in light of the observed rating scores in  $\mathbf{R}$ . However, it is very challenging to fulfill the

task in practice due to the high dimensionality and sparsity of the rating matrix. Matrix factorization addresses this challenge by mapping each user  $m$  and item  $n$  to the latent feature vectors  $\mathbf{p}_m \in \mathcal{R}^K$  and  $\mathbf{q}_n \in \mathcal{R}^K$  in a shared low  $K$ -dimension space (typically,  $K \ll M, N$ ). It assumes that the rating  $r_{m,n}$  can be predicted by

$$\hat{r}_{m,n} = \mathbf{p}_m^\top \mathbf{q}_n. \quad (1)$$

Therefore, the latent features  $\{\mathbf{p}_m\}$  and  $\{\mathbf{q}_n\}$  can be learned by minimizing the prediction error for all observed ratings in  $\mathbf{R}$ , while each unobserved rating value can be estimated using Equation (1) with its corresponding latent features learned by matrix factorization. In practice, since the feedback (i.e., rating scores) from users is received over time, the system is required to address the collaborative filtering problem in an interactive mode, which is referred to as an interactive recommender system.

In an interactive recommender system, user  $m$  constantly arrives to interact with the system over time. At each time  $t = [1, \dots, T]$ , the system, according to the observed rating history, recommends an item  $n(t)$  to the corresponding user  $m$ . After consuming item  $n(t)$ , the feedback (i.e., rating)  $r_{m,n(t)}$  from user  $m$  is collected by the system and further utilized to update the model for the next item delivery. The interactive recommendation process involves a series of decisions over a finite but possibly unknown time horizon  $T$ . Accordingly, such an interactive recommendation process is modeled as a multi-armed bandit problem, where each item corresponds to an arm. Pulling an arm indicates that its corresponding item is being recommended and the rating score is considered as the reward received after pulling the corresponding arm.

Let  $\mathbb{S}(t)$  be the available information at time  $t$  collected by the system for the target user  $m$ ,

$$\mathbb{S}(t) = \{(n(1), r_{m,n(1)}), \dots, (n(t-1), r_{m,n(t-1)})\}. \quad (2)$$

A policy  $\pi$  is defined as a function and used to select an arm based on the current cumulative information  $\mathbb{S}(t)$ ,

$$n(t) = \pi(\mathbb{S}(t)). \quad (3)$$

The total rewards received by the policy  $\pi$  after  $T$  iterations is

$$R_\pi = \sum_{t=1}^T r_{m,\pi(\mathbb{S}(t))}. \quad (4)$$

The optimal policy  $\pi^*$  is defined as the one with maximum accumulated expected reward after  $T$

iterations,

$$\pi^* = \arg \max_{\pi} \mathbb{E}(R_{\pi}) = \arg \max_{\pi} \sum_{t=1}^T \mathbb{E}(r_{m,\pi(\mathbb{S}(t))}|t). \quad (5)$$

Therefore, our goal is to identify an optimal policy for maximizing the total rewards. Herein we use reward instead of regret to express the objective function, since maximization of the cumulative rewards is equivalent to minimization of regret during the  $T$  iterations [47]. Before selecting one arm at time  $t$ , a policy  $\pi$  typically learns a model to predict the reward for every arm according to the historical accumulated information  $\mathbb{S}(t)$ . The reward prediction helps the policy  $\pi$  make decisions to increase the total rewards.

In the latent factor model [31], [32], the rating is estimated by a product of user and item feature vectors  $\mathbf{p}_m$  and  $\mathbf{q}_n$  in Equation (1). From the probabilistic perspective, PMF introduces an observation noise  $\xi$ , a zero-mean Gaussian noise with variance  $\sigma^2$  (i.e.,  $\xi \sim \mathcal{N}(0, \sigma^2)$ ), to the rating prediction function given in Equation (1). The derived rating prediction is as follows:

$$r_{m,n} = \mathbf{p}_m^T \mathbf{q}_n + \xi. \quad (6)$$

In this setting, Equation (5) can be re-formulated as:

$$\pi^* = \arg \max_{\pi} \sum_{t=1}^T \mathbb{E}_{\mathbf{p}_m, \mathbf{q}_n(\mathbb{S}(t))} (\mathbf{p}_m^T \mathbf{q}_{\pi(\mathbb{S}(t))}|t). \quad (7)$$

Consequently, the goal of an interactive recommender system is reduced to the optimization of the objective function in Equation (7).

**Thompson Sampling**, one of earliest heuristics for the bandit problem [13], belongs to the probability matching family. Its main idea is to randomly allocate the pulling chance according to the probability that an arm gives the largest expected reward at a particular time  $t$ . Based on the objective function in Equation (7), the probability of pulling arm  $n$  can be expressed as follows:

$$p(n(t) = n) = \int \mathbb{I}[\mathbb{E}(r_{m,n}|\mathbf{p}_m, \mathbf{q}_n) = \max_i \mathbb{E}(r_{m,i}|\mathbf{p}_m, \mathbf{q}_i)] p(\mathbf{p}_m, \mathbf{q}_n|t) d\mathbf{p}_m d\mathbf{q}_n. \quad (8)$$

At each time  $t$ , Thompson sampling samples both the user and item feature vectors together from their corresponding distributions, and then selects the item that leads to the largest reward expectation. Therefore, using Thompson sampling, the item selection function can be defined as:

$$n(t) = \arg \max_n (\tilde{\mathbf{p}}_m^T \tilde{\mathbf{q}}_n|t), \quad (9)$$

where  $\tilde{\mathbf{p}}_m$  and  $\tilde{\mathbf{q}}_n$  denote the sampled feature vectors for user  $m$  and item  $n$ , respectively.

To accomplish Thompson sampling, it is critical to model the random variable  $\mathbf{p}_m$  and  $\mathbf{q}_n$  using distributions, where the latent feature vectors can be easily sampled and the feedback at every time can be reasonably integrated. Most of the previous studies suppose a Gaussian prior for both user and item feature vectors, with an assumption that items are independent from each other [21], [47]. However, this assumption rarely holds in real applications. In the following section, we explicitly formulate the dependent arms with a generative model.

### 3.2 Modeling the Arm Dependency

Based on the fact that similar items (i.e., arms) are likely to receive similar feedback (i.e., rewards), we assume that a dependency exists among similar items. The dependencies among items can be further leveraged to improve the users' preferences inference on a particular item even if the item has little historical interaction data in the system. The challenge here lies in how to sequentially infer the arms' dependencies as well as the users' preferences simultaneously, providing the feedback over time.

In our work, the arms' dependencies are expressed in the form of the clusters of arms, where the dependent arms fall into the same one. In order to explore the dependencies in the bandit setting, Latent Dirichlet Allocation (LDA) [8], a generative statistic model for topic modeling, is adopted to construct the arms' clusters. We propose the ICTR (Interactive Collaborative Topic Regression) model to infer the clusters of arms as well as the arm selection.

The main idea of our model is to treat an item  $n$  as a word, while consider a user  $m$  as a document. All the items rated by a user indicate the hidden preferences of the user, analogous to the scenario in topic modeling where the words contained in a document imply its latent topics. Specifically, let  $K$  be the number of latent aspects (i.e., topics or clusters) the users concern when consuming items. We assume that  $\mathbf{p}_m \in \mathcal{R}^K$  corresponds to the latent vector for user  $m$ , where the  $k$ -th component of  $\mathbf{p}_m$  indicates the user's preference over the  $k$ -th aspect of items. Further,  $\mathbf{q}_n \in \mathcal{R}^K$  is supposed to be the latent vector for item  $n$ , and the  $k$ -th component value of  $\mathbf{q}_n$  represents that it belongs to the  $k$ -th cluster. The rating score  $r_{m,n}$ , given by user  $m$  after consuming item  $n$ , is assumed to be the inner product of  $\mathbf{p}_m$  and  $\mathbf{q}_n$ . By linking to the topic model,

a generative process for user ratings is accordingly introduced and presented in Figure 1.

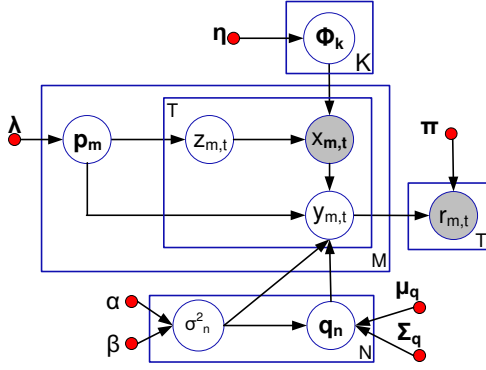


Fig. 1: The graphic model for the ICTR model. Random variable is denoted as a circle. The circle with filled color denotes the observed random variable. Red dot represents a hyper parameter.

Based on the above description, the user latent vector  $\mathbf{p}_m$  is assumed to follow a Dirichlet prior distribution with a predefined hyper parameter  $\lambda$ , shown in Equation (10).

$$\mathbf{p}_m | \lambda \sim \text{Dir}(\lambda). \quad (10)$$

As presented in Equation (6), we denote  $\sigma^2$  as the variance of the noise for reward prediction and assume  $\sigma_n^2$  is drawn from the Inverse Gamma ( $\mathcal{IG}$ ) distribution shown in the following.

$$p(\sigma_n^2 | \alpha, \beta) = \mathcal{IG}(\alpha, \beta), \quad (11)$$

where  $\alpha$  and  $\beta$  are predefined hyper parameters for  $\mathcal{IG}$  distribution.

Given  $\sigma_n^2$ , the item latent vector  $\mathbf{q}_n$  is generated by a Gaussian prior distribution as follows:

$$\mathbf{q}_n | \mu_q, \Sigma_q, \sigma_n^2 \sim \mathcal{N}(\mu_q, \sigma_n^2 \Sigma_q), \quad (12)$$

where  $\mu_q$  and  $\Sigma_q$  are predefined hyper parameters.

Further, let  $\Phi_k \in \mathcal{R}^N$  be the item distribution for topic  $k$ . Similar to  $\mathbf{p}_m$ , Dirichlet distribution is specified as the prior of  $\Phi_k$  presented in Equation (13).

$$\Phi_k | \eta \sim \text{Dir}(\eta), \quad (13)$$

where  $\eta \in \mathcal{R}^N$  is the hyper parameter.

When user  $m$  arrives to interact with the system at time  $t$ , one of  $K$  topics, denoted as  $z_{m,t}$ , is first selected according to the user's latent preference  $\mathbf{p}_m$ , indicating that the user  $m$  shows interest in the topic  $z_{m,t}$  at this moment. Accordingly,  $z_{m,t}$  is supposed to follow a multinomial distribution governed by  $\mathbf{p}_m$  as follows,

$$z_{m,t} | \mathbf{p}_m \sim \text{Mult}(\mathbf{p}_m). \quad (14)$$

W.L.O.G, we assume  $z_{m,t} = k$ , and then the item distribution for topic  $k$  (i.e.,  $\Phi_k$ ) is for generating

the item  $x_{m,t}$  recommended to the user  $m$  at time  $t$ . We assume the random variable  $x_{m,t}$  follows the multinomial distribution ruled by  $\Phi_k$ , i.e.,

$$x_{m,t} | \Phi_k \sim \text{Mult}(\Phi_k). \quad (15)$$

W.L.O.G, item  $n$  is assumed to be selected by user  $m$  at time  $t$  (i.e.,  $x_{m,t} = n$ ) where the latent vector corresponding to item  $n$  is  $\mathbf{q}_n$ . Let  $y_{m,t}$  be the predicted reward (i.e., rating), given by user  $m$  at time  $t$ . The predicted reward  $y_{m,t}$  can be inferred by

$$y_{m,t} \sim \mathcal{N}(\mathbf{p}_m^T \mathbf{q}_n, \sigma_n^2). \quad (16)$$

By Equation (16), the rewards of different items are predicted. Based on the predicted rewards, the policy  $\pi$  selects an item and recommends it to user  $m$ , considering the tradeoff between exploitation and exploration. After consuming the recommended item, the system receives the actual reward  $r_{m,t}$  from user  $m$ . The objective of the model is to maximize the expected accumulative rewards in a long run as described in Equation (5).

In this section, taking the clusters of arms into account, we formally introduced our ICTR model, which integrates matrix factorization with topic modeling in the bandit setting. We develop our solution to infer ICTR model from a Bayesian perspective in the following section.

## 4 METHODOLOGY AND SOLUTION

In this section, we present the methodology for online inferences of ICTR model.

The posterior distribution inference involves five random variables, i.e.,  $\mathbf{p}_m$ ,  $z_{m,t}$ ,  $\Phi_k$ ,  $\mathbf{q}_n$ , and  $\sigma_n^2$ . According to the graphical model in Figure 1, the five random variables belong to two categories: parameter random variable and latent state random variable.  $\Phi_k$ ,  $\mathbf{p}_m$ ,  $\mathbf{q}_n$ , and  $\sigma_n^2$  are parameter random variables since they are assumed to be fixed but unknown, and their values do not change with time. Instead,  $z_{m,t}$  is referred to as a latent state random variable since it is not observable and its value is time dependent. After pulling arm  $n(t)$ , where  $n(t) = x_{m,t}$  according to Equation (15) at time  $t$ , a reward is observed as  $r_{m,t}$ . Thus,  $x_{m,t}$  and  $r_{m,t}$  are referred to as observed random variables.

Our goal is to infer both latent parameter variables and latent state random variables to sequentially fit the observed data at time  $t - 1$ , and predict the rewards for arm selection with respect to the incoming user at time  $t$ . However, since the inference of our model cannot be conducted by a simple closed-form solution, we adopt the sequential sampling-based inference strategy that is widely

used in sequential Monte Carlo sampling [36], particle filtering [14], and particle learning [11] to learn the distribution of both parameter and state random variables. Specifically, particle learning that allows both state filtering and sequential parameter learning simultaneously is a perfect solution to our proposed model inference. In order to develop the solution based on particle learning, we first define a particle as follows.

**Definition 1 (Particle).** A particle for predicting the reward  $y_{m,t}$  is a container that maintains the current status information for both user  $m$  and item  $x_{m,t}$ . The status information comprises of random variables such as  $\mathbf{p}_m$ ,  $\sigma_n^2$ ,  $\Phi_k$ ,  $\mathbf{q}_n$ , and  $z_{m,t}$ , as well as the hyper parameters of their corresponding distributions, such as  $\lambda$ ,  $\alpha$ ,  $\beta$ ,  $\eta$ ,  $\mu_q$  and  $\Sigma_q$ .

In particle learning, each particle corresponds to a sample for modeling inference status information. At each time stamp, particles are re-sampled according to their fitness to the current observable data. Then, the re-sampled particles are propagated to new particles and obtain the status information for the next time stamp. In the following subsections, we develop our solution based on particle learning.

#### 4.1 Re-sample Particles with Weights

At time  $t - 1$ , a fixed-size set of particles is maintained for the reward prediction for each arm  $n(t-1)$  given user  $m$ . We denote the particle set at time  $t - 1$  as  $\mathcal{P}_{m,n(t-1)}$  and assume the number of particles in  $\mathcal{P}_{m,n(t-1)}$  is  $B$ . Let  $\mathcal{P}_{m,n(t-1)}^{(i)}$  be the  $i^{\text{th}}$  particles given both the user  $m$  and the item  $n(t-1)$  at time  $t - 1$ , where  $1 \leq i \leq B$ . Each particle  $\mathcal{P}_{m,n(t-1)}^{(i)}$  has a weight, denoted as  $\rho^{(i)}$ , indicating its fitness for the new observed data at time  $t$ . Note that  $\sum_{i=1}^B \rho^{(i)} = 1$ . The fitness of each particle  $\mathcal{P}_{m,n(t-1)}^{(i)}$  is defined as the likelihood of the observed data  $x_{m,t}$  and  $r_{m,t}$ . Therefore,

$$\rho^{(i)} \propto p(x_{m,t}, r_{m,t} | \mathcal{P}_{m,n(t-1)}^{(i)}). \quad (17)$$

Further,  $y_{m,t}$  is the predicted value of  $r_{m,t}$ . The distribution of  $y_{m,t}$ , determined by  $\mathbf{p}_m$ ,  $\mathbf{q}_n$ ,  $z_{m,t}$ ,  $\Phi_k$ , and  $\sigma_n^2$ , has been described in Section 3.2.

Therefore, we can compute  $\rho^{(i)}$  as proportional to the density value given  $y_{m,t} = r_{m,t}$  and  $x_{m,t} = n$ . Thus, we obtain

$$\rho^{(i)} \propto \sum_{z_{m,t}=1}^K \{ \mathcal{N}(r_{m,t} | (\mathbf{p}_m^T \mathbf{q}_n, \sigma_n^2)) \bullet p(z_{m,t} = k, x_{m,t} = n | \mathcal{P}_{m,n(t-1)}^{(i)}) \},$$

where

$$\begin{aligned} p(z_{m,t} = k, x_{m,t} = n | \mathcal{P}_{m,n(t-1)}^{(i)}) &= \iint_{\mathbf{p}_m, \Phi_k} p(z_{m,t} = k, x_{m,t} = n, \mathbf{p}_m, \Phi_k | \lambda, \eta) d\mathbf{p}_m d\Phi_k \\ &= \int_{\mathbf{p}_m} \text{Mult}(z_{m,t} = k | \mathbf{p}_m) \text{Dir}(\mathbf{p}_m | \lambda) d\mathbf{p}_m \\ &\bullet \int_{\Phi_k} \text{Mult}(x_{m,t} = n | \Phi_k) \text{Dir}(\Phi_k | \eta) d\Phi_k \\ &= E(\mathbf{p}_{m,k} | \lambda) \bullet E(\Phi_{k,n} | \eta). \end{aligned} \quad (18)$$

Thus, we have:

$$\rho^{(i)} \propto \sum_{z_{m,t}=1}^K \{ \mathcal{N}(r_{m,t} | (\mathbf{p}_m^T \mathbf{q}_n, \sigma_n^2)) \bullet E(\mathbf{p}_{m,k} | \lambda) \bullet E(\Phi_{k,n} | \eta) \}, \quad (19)$$

where  $E(\mathbf{p}_{m,k} | \lambda)$  and  $E(\Phi_{k,n} | \eta)$  represent the conditional expectations of  $\mathbf{p}_{m,k}$  and  $\Phi_{k,n}$  given the observed reward  $\lambda$  and  $\eta$  of  $\mathcal{P}_{m,n(t-1)}^{(i)}$ . The expectations can be inferred by  $E(\mathbf{p}_{m,k} | \lambda) = \frac{\lambda_k}{\sum_{k=1}^K \lambda_k}$  and  $E(\Phi_{k,n} | \eta) = \frac{\eta_{k,n}}{\sum_{n=1}^N \eta_{k,n}}$ .

Before updating any parameters, a re-sampling process is conducted. We replace the particle set  $\mathcal{P}_{m,n(t-1)}$  with a new set  $\mathcal{P}_{m,n(t)}$ , where  $\mathcal{P}_{m,n(t)}$  is generated from  $\mathcal{P}_{m,n(t-1)}$  using sampling with replacement based on the weights of particles. Then sequential parameter updating is based on  $\mathcal{P}_{m,n(t)}$ .

#### 4.2 Latent State Inference

Provided with the new observation  $x_{m,t}$  and  $r_{m,t}$  at time  $t$ , the random state  $z_{m,t}$  can be one of  $K$  topics and the posterior distribution of  $z_{m,t}$  is shown as follows:

$$z_{m,t} | x_{m,t}, r_{m,t}, \mathcal{P}_{m,n(t-1)}^{(i)} \sim \text{Mult}(\theta), \quad (20)$$

where  $\theta \in \mathcal{R}^K$  is the parameter specifying the multinomial distribution. According to Equation (18), since

$$p(z_{m,t} | x_{m,t}, r_{m,t}, \lambda, \eta) \propto p(z_{m,t}, x_{m,t} | r_{m,t}, \lambda, \eta),$$

$\theta$  can be computed by  $\theta_k \propto E(\mathbf{p}_{m,k} | r_{m,t}, \lambda) \bullet E(\Phi_{k,n} | r_{m,t}, \eta)$ . Further,  $E(\mathbf{p}_{m,k} | r_{m,t}, \lambda)$  and  $E(\Phi_{k,n} | r_{m,t}, \eta)$  can be obtained as follows,

$$\begin{aligned} E(\mathbf{p}_{m,k} | r_{m,t}, \lambda) &= \frac{\mathcal{I}(z_{m,t} = k) r_{m,t} + \lambda_k}{\sum_{k=1}^K [\mathcal{I}(z_{m,t} = k) r_{m,t} + \lambda_k]}, \\ E(\Phi_{k,n} | r_{m,t}, \eta) &= \frac{\mathcal{I}(x_{m,t} = n) r_{m,t} + \eta_{k,n}}{\sum_{n=1}^N [\mathcal{I}(x_{m,t} = n) r_{m,t} + \eta_{k,n}]} \end{aligned} \quad (21)$$

where  $\mathcal{I}(\bullet)$ , an indicator function, returns 1 when the input boolean expression is true, otherwise returns 0. Specifically, if  $r_{m,t} \in \{0, 1\}$ , the value of  $r_{m,t}$  indicates whether  $x_{m,t}$  should be included in the

preferred item list of the user  $m$ . If  $r_{m,t} \in [0, +\infty)$ , the value of  $r_{m,t}$  implies how the user  $m$  likes the item  $x_{m,t}$ . Therefore, our solution can effectively handle the non-negative rating score at different scales.

### 4.3 Parameter Statistics Inference

At time  $t - 1$ , the sufficient statistics for the parameter random variables ( $\mathbf{q}_n, \sigma_n^2, \mathbf{p}_m, \Phi_k$ ) are  $(\mu_{\mathbf{q}}, \Sigma_{\mathbf{q}}, \alpha, \beta, \lambda, \eta)$ . Assume  $\mu'_{\mathbf{q}}, \Sigma'_{\mathbf{q}}, \alpha', \beta', \lambda',$  and  $\eta'$  are the sufficient statistics at time  $t$ , which are updated on the basis of both the sufficient statistics at time  $t - 1$  and the new observation data (i.e,  $x_{m,t}$  and  $r_{m,t}$ ). The sufficient statistics for parameters are updated as follows:

$$\begin{aligned} \Sigma'_{\mathbf{q}_n} &= (\Sigma_{\mathbf{q}_n}^{-1} + \mathbf{p}_m \mathbf{p}_m^\top)^{-1}, \quad \mu'_{\mathbf{q}_n} = \Sigma'_{\mathbf{q}_n} (\Sigma_{\mathbf{q}_n}^{-1} \mu_{\mathbf{q}_n} + \mathbf{p}_m r_{m,t}) \\ \alpha' &= \alpha + \frac{1}{2} \\ \beta' &= \beta + \frac{1}{2} (\mu_{\mathbf{q}_n}^\top \Sigma_{\mathbf{q}_n}^{-1} \mu_{\mathbf{q}_n} + r_{m,t}^\top r_{m,t} - \mu_{\mathbf{q}_n}^\top \Sigma_{\mathbf{q}_n}^{-1} \mu'_{\mathbf{q}_n}) \\ \lambda'_k &= \mathcal{I}(z_{m,t} = k) r_{m,t} + \lambda_k, \quad \eta'_{k,n} = \mathcal{I}(x_{m,t} = n) r_{m,t} + \eta_{k,n} \end{aligned} \quad (22)$$

At time  $t$ , the sampling process for the parameter random variables  $\sigma_n^2, \mathbf{q}_n, \mathbf{p}_m$  and  $\Phi_k$  is summarized as below:

$$\begin{aligned} \sigma_n^2 &\sim \mathcal{IG}(\alpha', \beta'), \quad \mathbf{q}_n | \sigma_n^2 \sim \mathcal{N}(\mu'_{\mathbf{q}_n}, \sigma_n^2 \Sigma'_{\mathbf{q}_n}), \\ \mathbf{p}_m &\sim \text{Dir}(\lambda'), \quad \Phi_k \sim \text{Dir}(\eta'). \end{aligned} \quad (23)$$

### 4.4 Integration with Policies

In our ICTR model, when user  $m$  arrives at time  $t$ , reward  $r_{m,t}$  is unknown since it is not observed until one of arms  $x_{m,t}$  is pulled. Without observed  $x_{m,t}$  and  $r_{m,t}$ , the particle re-sampling, latent state inference, and parameter statistics inference for time  $t$  cannot be conducted. Therefore, we utilize the latent vectors  $\mathbf{p}_m$  and  $\mathbf{q}_n$ , sampled from their corresponding posterior distributions by Equation (23) at time  $t - 1$ , to predict the reward for each arm. In this section, two policies based on Thompson sampling and UCB for ICF are integrated with our model.

In the model, every item has  $B$  independent particles given the user  $m$ . Each particle  $i$  contains its latent variables and parameters, and produces an independent reward prediction  $r_{m,t}^{(i)}$ . Specifically, according to Thompson sampling discussed in Section 3.1, we predict the reward of pulling arm  $n$  with the average value of rewards from  $B$  particles. The policy based on Thompson sampling selects an arm  $n(t)$  based on the following equation,

$$n(t) = \arg \max_n (\bar{r}_{m,n}), \quad (24)$$

where  $\bar{r}_{m,n}$  denotes the average reward, i.e.,

$$\bar{r}_{m,n} = \frac{1}{B} \sum_{i=1}^B \mathbf{p}_m^{(i)\top} \mathbf{q}_n^{(i)}.$$

Moreover, UCB policy selects an arm based on the upper bound of the predicted reward. Assuming that  $r_{m,t}^{(i)} \sim \mathcal{N}(\mathbf{p}_m^{(i)\top} \mathbf{q}_n^{(i)}, \sigma^{(i)2})$ , the UCB-based policy is developed by the mean and variance of predicted reward, i.e.,

$$n(t) = \arg \max_n (\bar{r}_{m,n} + \gamma \sqrt{\nu}), \quad (25)$$

where  $\gamma \geq 0$  is a predefined threshold and the variance is expressed as  $\nu = \frac{1}{B} \sum_i \sigma^{(i)2}$ .

### 4.5 Algorithm

Putting all the aforementioned inference together, an algorithm for ICTR model is provided below.

Online inference for ICF problem starts with MAIN procedure presented in Algorithm 1. As user  $m$  arrives at time  $t$ , EVAL procedure computes a score for each arm, where we define the score as the average reward. The arm with the highest score is selected to pull. After receiving a reward by pulling an arm, the new feedback is used to update ICTR model by UPDATE procedure. Especially in UPDATE procedure, we use the *resample-propagate* strategy in particle learning [11] rather than the *propagate-resample* strategy in particle filtering [14]. With the *resample-propagate* strategy, particles are re-sampled by taking  $\rho^{(i)}$  as the  $i^{\text{th}}$  particle's weight, where the  $\rho^{(i)}$  indicates the fitness of the observation at time  $t$  given the particle at time  $t - 1$ . The *resample-propagate* strategy is considered as an optimal and fully adapted strategy avoiding an importance sampling step.

In addition, existing algorithms [21], [47] consider all the arms independently, while our model takes the clusters of arms into account by learning the topic-related random variables (e.g.,  $\Phi_k$ ), which are shared among all the arms.

## 5 EMPIRICAL STUDY

To demonstrate the efficiency of our proposed algorithm, we conduct our experimental study over two popular real-world dataset: Yahoo! Today News and MovieLens (10M). First, we outline the general implementation of the baselines. Second, we start with a brief description of the datasets and evaluation method. Finally, we show and discuss the comparative experimental results of both the proposed algorithms and the baselines, and a case study



**Algorithm 1** The algorithm for ICTR model

---

```

1: procedure MAIN( $B$ ) ▷ main entry
2:   Initialize  $B$  particles, i.e.,  $\mathcal{P}_{m,n(0)}^{(1)} \dots \mathcal{P}_{m,n(0)}^{(B)}$ .
3:   for  $t \leftarrow 1, T$  do
4:     User  $m$  arrives for item recommendation.
5:      $n(t) = \arg \max_{n=1, N} \text{EVAL}(m, n)$  ▷ by Equation (24) or
      Equation (25).
6:     Receive  $r_{m,t}$  by rating item  $n(t)$ .
7:     UPDATE( $m, n(t), r_{m,t}$ ).
8:   end for
9: end procedure

10: procedure EVAL( $m, n$ ) ▷ get a rating score for item  $n$ , given user  $m$ .
11:   for  $i \leftarrow 1, B$  do ▷ Iterate on each particle.
12:     Get the user latent vector  $\mathbf{p}_m^{(i)}$ .
13:     Get the item latent vector  $\mathbf{q}_n^{(i)}$ .
14:     Predict  $i^{\text{th}}$  reward  $r_{m,t}^{(i)}$ .
15:   end for
16:   Compute the average reward as the final reward  $r_{m,t}$ .
17:   return the score.
18: end procedure

19: procedure UPDATE( $m, n(t), r_{m,t}$ ) ▷ update the inference.
20:   for  $i \leftarrow 1, B$  do ▷ Compute weights for each particle.
21:     Compute weight  $\rho^{(i)}$  of particle  $\mathcal{P}_{m,n(t)}^{(i)}$  by Equation (17).
22:   end for
23:   Re-sample  $\mathcal{P}'_{m,n(t)}$  from  $\mathcal{P}_{m,n(t)}$  according to the weights
       $\rho^{(i)}$ s.
24:   for  $i \leftarrow 1, B$  do ▷ Update statistics for each particle.
25:     Update the sufficient statistics for  $z_{m,t}$  by Equation (21).
26:     Sample  $z_{m,t}$  according to Equation (20).
27:     Update the statistics for  $\mathbf{q}_n, \sigma_n^2, \mathbf{p}_m, \Phi_k$  by Equation (22).
28:     Sample  $\mathbf{q}_n, \sigma_n^2, \mathbf{p}_m, \Phi_k$  by Equation (23).
29:   end for
30: end procedure

```

---

on movie topic distribution analysis of MovieLens (10M). Methods in [18], [40], [44], [48] are excluded from the baselines since our work is orthogonal to those methods.

## 5.1 Baseline Algorithms

In the experiment, we demonstrate the performance of our methods by comparing them with the following baseline algorithms:

- 1) Random: it randomly selects an item recommending to the target user.
- 2)  $\epsilon$ -greedy( $\epsilon$ ): it randomly selects an item with probability  $\epsilon$  and selects the item of the largest predicted reward with probability  $1 - \epsilon$ , where  $\epsilon$  is a predefined parameter.
- 3) UCB( $\lambda$ ): it picks the item  $j(t)$  with the highest rewards at time  $t$  as follows:

$$j(t) = \arg \max_{j=1, \dots, N} (\hat{\mu}_j + \lambda \sqrt{\frac{2 \ln(t)}{n_j(t)}})$$

where  $n_j(t)$  is the number of times that item  $n_j$  has been recommended until time  $t$ .

- 4) TS( $S_i(t), F_i(t)$ ): Thompson sampling described in Section 3.1, randomly draws the expected reward from the Beta posterior distribution,

and selects the item with the largest predicted reward.  $S_i(t)/F_i(t)$  is the number of positive/negative feedback on item  $i$  until time  $t$ .

- 5) PTS( $d, p$ ): particle Thompson sampling for matrix factorization approximates the posterior of latent feature vectors by updating a set of particles. Here  $d$  is the dimension of latent feature vector and  $p$  is the number of particles.

Our methods proposed in this paper include:

- 1) ICTRTS( $d, p$ ): it denotes our proposed interactive collaborative topic regression model with TS. Here  $d$  is the dimension of latent feature vector and  $p$  is the number of particles.
- 2) ICTRUCB( $d, p, \gamma$ ): it indicates our proposed model with UCB. Similar to UCB,  $\gamma$  is given. Here  $d$  is the dimension of latent feature vector and  $p$  represents the number of particles.

All algorithms are implemented using Java 1.8. All empirical experiments are running on Linux 2.6.32. The server is equipped with Intel(R) Xeon(R) CPU with 24 cores running at speed of 2.50GHZ. The total volume of memory is 158GB.

## 5.2 Datasets Description

We use two real-world datasets shown in Table 2 to evaluate our proposed algorithms.

TABLE 2: Description of Datasets.

Dataset	Yahoo News	MovieLens (10M)
#users	226,710	71,567
#items	652	10,681
#ratings	280,410,150	10,000,054

**Yahoo! Today News:** The core task of personalized news recommendation is to display appropriate news articles on the web page for users. The system often takes the user’s instant feedback into account to improve the prediction of his/her preferences, where the user feedback is about whether he/she clicks the recommended article or not. Here, we formulate the personalized news recommendation problem as an instance of bandit problem, where each arm corresponds to a news article. The experimental dataset is a collection based on a sample of anonymized user interaction on the news feeds published by Yahoo! Research Lab<sup>1</sup>. The dataset contains 15 days’ visit events of user-news item interaction data by randomly selecting news articles for recommendation. Besides, user’s information (e.g., demographic information) is provided for each visit event and represented as the user

1. <http://webscope.sandbox.yahoo.com/catalog.php>

identification, where users with the same information are identified as one user. In our experiments, the visit events of the first day are utilized for selecting proper parameters of ICTR model, while two million of the remaining are for the evaluation. Each interactive record in the historical logs consists of user ID, news article ID, rating feedback and a timestamp.

**MovieLens (10M):** Online movie recommender service aims to maximize the customer’s satisfaction by recommending the proper movies to the target users according to their preferences. Specifically, several movies are selected out of a movie set and displayed to users, and then users’ feedback on displayed movies are collected for improving the user satisfaction. Thereby, the problem of movie recommendation can be formulated as a bandit problem where an arm is a movie, a pull is regarded as a movie selection, and the reward is indicated by the user’s rating on the recommended movie. In our experiments, each rating associates user ID, movie ID, and a timestamp. In order to use the *replayer* evaluation method, we assume that the rating data is produced by the users when the movies are randomly recommended. The rating score in the dataset ranges from 1 to 5. Additionally, we choose the top-N (N=100) popular movies to form a movie set, from which one movie is recommended to a user by algorithms in every trial.

### 5.3 Evaluation Method and Metrics

The evaluation methods for traditional non-interactive recommender systems assume the independence among the items at different time stamps once the offline model is built. In an online interactive recommender system, the recommended items at previous time stamps are used to update the recommendation model, and then further affect the recommendation items at current time stamp.

We apply the *replayer* method to evaluate our proposed algorithms on the aforementioned two datasets. The *replayer* method, first introduced in [23], provides an unbiased offline evaluation for multi-armed bandit algorithms via historical logs, where the logs are assumed to be generated by random recommendation. The main idea of *replayer* is to replay each user visit in the historical logs to the algorithm under evaluation. If the recommended item by the testing algorithm is identical to the one in the historical log, this visit is considered as a match between the historical recommendation and the testing recommendation algorithm. The *replayer* method only counts those matched visits in

for the accumulated reward computation. Since the recommendation algorithms may result in different numbers of matched visits, the average reward (i.e., the accumulated rewards divided by the number of matched visits) is adopted for evaluation.

Particularly, in the scenario of news article recommendation, a matched visit corresponds to an impression, and a reward of one is obtained by a click, so the average reward also represents the average CTR (Click Through Rate). In the scenario of movie recommendation, we set the reward of one if the rating score of the recommended movie is no less than four, indicating that the user likes the recommended movie. If the rating is less than four, a reward value of zero is obtained. Thus, the average reward in this scenario indicates the success rate of movie recommendation. To sum it up, in our setting, the reward is one if the recommended article (movie) is clicked (liked), otherwise it is zero.

TABLE 3: Evaluation Metric Computation for *Replayer*.

		Item in Random Recommendation Logs	
		Clicked/Liked	Not Clicked/Not Liked
Recommended Item	Matched	TP	FP
	Not Matched	N/A	N/A

Consider a matched visit shown in Table 3. If the item in the logs is clicked or liked by a user, the recommended item is referred to as a true positive (TP), otherwise it is referred to as a false positive (FP). The average reward is computed as  $\frac{TP*1+FP*0}{TP+FP} = \frac{TP}{TP+FP}$ , corresponding to the formula for the precision for matched visits. However, for the unmatched visits, we can not determine whether an item is false negative or true negative since no ground truth is provided. Therefore, the computation of the average reward in this case as the recall, relying on the false negative, is not feasible.

### 5.4 Recommendation Evaluation

In this section we first conduct the *replayer* evaluation method for each algorithm with different parameter settings. The aforementioned average reward is used as the performance metric in the experiments.

All baseline algorithms are configured with different parameter settings provided in Table 4. The settings of all algorithms with the highest average reward are highlighted in bold. Our algorithm **ICTRUCB(2, 10, 0.1)** achieves the best performance among all algorithms on Yahoo! Today News, and the performance comparisons among

different algorithms along different time buckets are illustrated in Figure 2. For MovieLens (10M), ICTRTS (3, 10) outperforms all others and the corresponding performance comparisons are shown in Figure 3.

Our proposed algorithms outperform the baseline algorithms using independent arms because ICTR model can leverage the dependencies among items by clustering items (arms) using items' latent aspects. The feedback received after recommending an item is not only used to update the model parameters related to this item, but also utilized to refine the parameters for the item's cluster. As a result, the updated cluster parameters further influence the model's parameter inference for other items within the same cluster. The effect of the clustering is illustrated in more details in the next section.

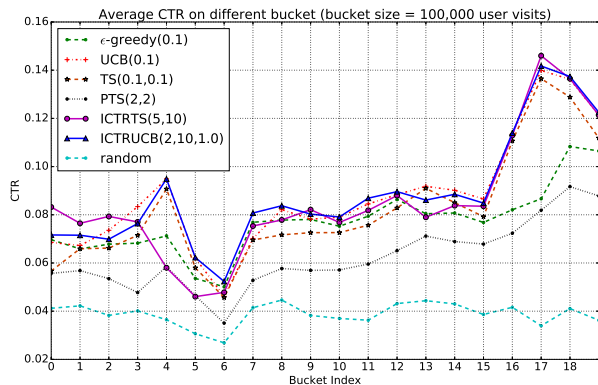


Fig. 2: The average CTR of Yahoo! Today News data is given along each time bucket. All algorithms shown here are configured with their best parameter settings.

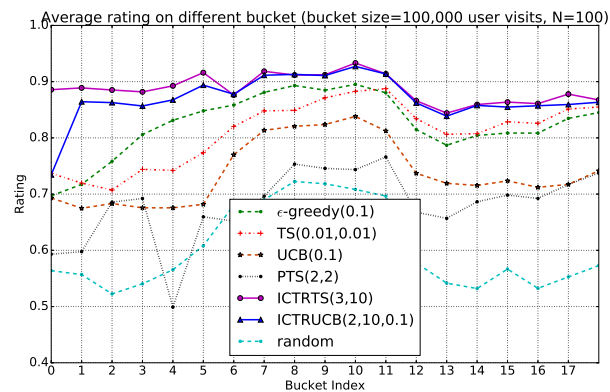


Fig. 3: The average rating of MovieLens (10M) data is given along each time bucket. All algorithms shown here are configured with their best parameter settings.

## 5.5 A Case Study: Topic Distribution Analysis on MovieLens (10M)

We conduct an experiment to demonstrate that our model can effectively capture the dependency between items, i.e., finding the latent topics among movies and clustering similar movies together. In this experiment, top-N ( $N=8$ ) popular movies are selected and topic number ( $K=2$ ) is set for our model. After millions of training iterations, the learned latent movie feature vectors will represent each movie's topic distribution over the two latent topics, in which the  $i$ -th dimension of the feature vector encodes the probability that the movie belongs to the  $i$ -th movie topic cluster. We separately choose four movies with the highest value of the first element and the second element of these latent feature vectors, and list their IDs, names, and movie types in Table 5, which clearly proves our assumption that the model is able to capture the dependency between items and cluster similar movies together.

## 5.6 Time Cost

The cumulative time cost of each algorithm on both datasets is presented in Figure 4a and Figure 4b, where all algorithms are configured with their best parameter settings. Our proposed algorithms have higher running time since they need to learn the latent features for arms. However, the computational complexity of both  $\text{ICTRUCB}(1, 1, 1.0)$  and  $\text{ICTRTS}(1, 1)$  is comparable to the baselines'. We also evaluate the time costs of  $\text{ICTRTS}$  and  $\text{ICTRUCB}$  with different number of particles and latent feature vector dimensions on the two datasets (see Figure 4c and Figure 4d). It shows that the time cost grows linearly with the number of particles and dimensions of latent feature vector.

The observations can be summarized as follows: (1) MovieLens (10M) requires much more time than Yahoo! Today New due to a larger amount of items and users. (2) In general, UCB-based algorithms (e.g.,  $\text{ICTRUCB}$ ,  $\text{UCB}$ ) are faster than TS-based ones (e.g.,  $\text{ICTRTS}$ ,  $\text{PTS}$ ) since the TS-based algorithms highly depend on the sampling process.

## 6 CONCLUSION AND FUTURE WORK

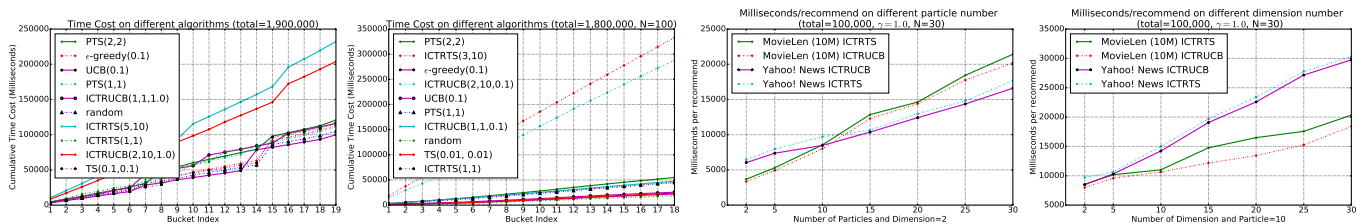
In this paper, we propose an interactive collaborative topic regression model that adopts a generative process based on topic model to explicitly formulate the arm dependencies as the clusters on arms, where dependent arms are assumed to be

TABLE 4: Average CTR/Rating on two real world datasets

Algorithm	Yahoo! Today News				MovieLens (10M)			
	mean	std	min	max	mean	std	min	max
$\epsilon$ -greedy (0.01)	0.06916	0.00312	0.06476	0.07166	0.70205	0.06340	0.60752	0.78934
$\epsilon$ -greedy (0.1)	<b>0.07566</b>	0.00079	0.07509	0.07678	<b>0.82038</b>	0.01437	0.79435	0.83551
$\epsilon$ -greedy (0.3)	0.07006	0.00261	0.06776	0.07372	0.80447	0.01516	0.77982	0.82458
$\epsilon$ -greedy (1.0)	0.03913	0.00051	0.03842	0.03961	0.60337	0.00380	0.59854	0.60823
UCB (0.01)	0.05240	0.00942	0.04146	0.06975	0.62133	0.10001	0.45296	0.73369
UCB (0.1)	<b>0.08515</b>	0.00021	0.08478	0.08544	<b>0.73537</b>	0.07110	0.66198	0.85632
UCB (0.5)	0.05815	0.00059	0.05710	0.05893	0.71478	0.00294	0.63623	0.64298
UCB (1.0)	0.04895	0.00036	0.04831	0.04932	0.63909	0.00278	0.60324	0.61296
TS (0.01, 0.01)	0.07853	0.00058	0.07759	0.07921	<b>0.83585</b>	0.00397	0.82927	0.84177
TS (0.1, 0.1)	<b>0.07941</b>	0.00040	0.07869	0.07988	0.83267	0.00625	0.82242	0.84001
TS (0.5, 0.5)	0.07914	0.00106	0.07747	0.08041	0.82988	0.00833	0.81887	0.84114
TS (1.0, 1.0)	0.07937	0.00079	0.07788	0.08044	0.83493	0.00798	0.82383	0.84477
PTS (2, 2)	<b>0.06069</b>	0.00575	0.05075	0.06470	<b>0.70484</b>	0.03062	0.64792	0.74610
PTS (2, 10)	0.05699	0.00410	0.05130	0.06208	0.65046	0.01124	0.63586	0.66977
PTS (5, 10)	0.05778	0.00275	0.05589	0.06251	0.63777	0.00811	0.62971	0.65181
PTS (5, 20)	0.05726	0.00438	0.05096	0.06321	0.62289	0.00714	0.61250	0.63567
PTS (10, 20)	0.05490	0.00271	0.05179	0.05839	0.61819	0.01044	0.60662	0.63818
ICTRTS (2, 5)	0.06888	0.00483	0.06369	0.07671	0.70386	0.15772	0.48652	0.85596
ICTRTS (2, 10)	0.06712	0.01873	0.03731	0.08487	0.56643	0.10242	0.42974	0.67630
ICTRTS (3, 10)	0.06953	0.00783	0.05857	0.07804	<b>0.88512</b>	0.00052	0.88438	0.88553
ICTRTS (5, 10)	<b>0.08321</b>	0.08236	0.08492	0.06292	0.55748	0.14168	0.38715	0.73404
ICTRTS (7, 10)	0.05066	0.00885	0.04229	0.06423	0.517826	0.07120	0.42297	0.59454
ICTRTS (7, 20)	0.04925	0.00223	0.04672	0.05285	0.61414	0.12186	0.44685	0.73365
ICTRUCB (2, 10, 0.01)	0.06673	0.01233	0.04588	0.08112	0.44650	0.06689	0.38678	0.53991
ICTRUCB (2, 10, 1.0)	<b>0.08597</b>	0.00056	0.08521	0.08675	<b>0.86411</b>	0.01528	0.85059	0.88547
ICTRUCB (3, 10, 0.05)	0.07250	0.00426	0.06799	0.07694	0.54757	0.13265	0.43665	0.73407
ICTRUCB (3, 10, 1.0)	0.08196	0.00296	0.07766	0.08530	0.57805	0.08716	0.46453	0.67641
ICTRUCB (5, 10, 0.01)	0.07009	0.00722	0.06411	0.08244	0.62282	0.02572	0.59322	0.65594
ICTRUCB (5, 10, 1.0)	0.08329	0.00140	0.08098	0.08481	0.80038	0.24095	0.9625	0.88554

TABLE 5: Movie Topic Distribution of MovieLen (10M)

Topic Cluster I			Topic Cluster II		
MovieId	MovieName	MovieType	MovieId	MovieName	MovieType
32	12 Monkeys	Sci-Fi,Thriller	344	Pet Detective	Comedy
50	Usual Suspects	Crime,Mystery,Thriller	588	Aladdin	Children,Animation,Comedy
590	Dances with wolves	Adventure,Drama,Western	595	Beauty and the Beast	Animation,Children,Musical
592	Batman	Action,Crime,Sci-Fi,Thriller	2857	Yellow Submarine	Adventure,Animation,Comedy,Musical



(a) Cumulative time cost of Yahoo! Today News is given along each time bucket. (b) Cumulative time cost of MovieLens (10M) is given along each time bucket. (c) Time cost are given with different number of particles. (d) Time cost are given with different number of latent feature vector dimensions.

Fig. 4: Time cost comparison on both two datasets.

generated from the same cluster. Every time an arm is pulled, the feedback is not only used for inferring the involved user and item latent vectors, but also employed to update the latent parameters with respect to the arm's cluster. The latent cluster parameters further help with the reward prediction for other arms in the same cluster. We conduct

empirical studies on two real-world applications, including movie and news recommendation, and the experimental results demonstrate the effectiveness of our proposed approach.

Individual preferences on news and movies usually evolve over time. One possible research direction is to extend our model considering the

time-varying property in user preferences for better online personal recommendation [45]. In addition, we would like to provide a comprehensive regret analysis [17] of our model in the future work.

## REFERENCES

- [1] D. Agarwal, B.-C. Chen, and P. Elango. Explore/exploit schemes for web content optimization. In *Data Mining, 2009. ICDM'09. Ninth IEEE International Conference on*, pages 1–10. IEEE, 2009.
- [2] S. Agrawal and N. Goyal. Thompson sampling for contextual bandits with linear payoffs. In *International Conference on Machine Learning*, pages 127–135, 2013.
- [3] H. J. Ahn. A new similarity measure for collaborative filtering to alleviate the new user cold-starting problem. *Information Sciences*, 178(1):37–51, 2008.
- [4] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM journal on computing*, 32(1):48–77, 2002.
- [5] B. Awerbuch and R. Kleinberg. Online linear optimization and adaptive routing. *Journal of Computer and System Sciences*, 74(1):97–114, 2008.
- [6] A. Barraza-Urbina. The exploration-exploitation trade-off in interactive recommender systems. In *Proceedings of the Eleventh ACM Conference on Recommender Systems*, pages 431–435. ACM, 2017.
- [7] J. Bennett, S. Lanning, et al. The netflix prize. In *Proceedings of KDD cup and workshop*, volume 2007, page 35. New York, NY, USA, 2007.
- [8] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.
- [9] D. Bounieffouf, A. Bouzeghoub, and A. L. Gançarski. A contextual-bandit algorithm for mobile context-aware recommender system. In *International Conference on Neural Information Processing*, pages 324–331. Springer, 2012.
- [10] K. Buza and L. Peska. Aladin: A new approach for drug-target interaction prediction. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 322–337. Springer, 2017.
- [11] C. Carvalho, M. S. Johannes, H. F. Lopes, and N. Polson. Particle learning and smoothing. *Statistical Science*, 25(1):88–106, 2010.
- [12] S. Chang, J. Zhou, P. Chubak, J. Hu, and T. S. Huang. A space alignment method for cold-start tv show recommendations. In *Proceedings of the 24th International Conference on Artificial Intelligence*, pages 3373–3379. AAAI Press, 2015.
- [13] O. Chapelle and L. Li. An empirical evaluation of thompson sampling. In *Advances in neural information processing systems*, pages 2249–2257, 2011.
- [14] P. M. Djurić, J. H. Kotecha, J. Zhang, Y. Huang, T. Ghirmai, M. F. Bugallo, and J. Miguez. Particle filtering. *Signal Processing Magazine, IEEE*, 20(5):19–38, 2003.
- [15] A. Doucet, N. De Freitas, and N. Gordon. An introduction to sequential monte carlo methods. In *Sequential Monte Carlo methods in practice*, pages 3–14. Springer, 2001.
- [16] A. Doucet, S. Godsill, and C. Andrieu. On sequential monte carlo sampling methods for bayesian filtering. *Statistics and computing*, 10(3):197–208, 2000.
- [17] C. Gentile, S. Li, P. Kar, A. Karatzoglou, G. Zappella, and E. Etrúe. On context-dependent clustering of bandits. In *International Conference on Machine Learning*, pages 1253–1262, 2017.
- [18] C. Gentile, S. Li, and G. Zappella. Online clustering of bandits. In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pages 757–765, 2014.
- [19] J. H. Halton. Sequential monte carlo. In *Mathematical Proceedings of the Cambridge Philosophical Society*, volume 58, pages 57–78. Cambridge Univ Press, 1962.
- [20] J. L. Herlocker, J. A. Konstan, A. Borchers, and J. Riedl. An algorithmic framework for performing collaborative filtering. In *Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval*, pages 230–237. ACM, 1999.
- [21] J. Kawale, H. H. Bui, B. Kveton, L. Tran-Thanh, and S. Chawla. Efficient thompson sampling for online matrix-factorization recommendation. In *Advances in Neural Information Processing Systems*, pages 1297–1305, 2015.
- [22] Y. Koren, R. Bell, and C. Volinsky. Matrix factorization techniques for recommender systems. *Computer*, 42(8), 2009.
- [23] L. Li, W. Chu, J. Langford, T. Moon, and X. Wang. An unbiased offline evaluation of contextual bandit algorithms with generalized linear models. *JMLR*, 26:19–36, 2012.
- [24] L. Li, W. Chu, J. Langford, and R. E. Schapire. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, pages 661–670. ACM, 2010.
- [25] D. K. Mahajan, R. Rastogi, C. Tiwari, and A. Mitra. Logucb: an explore-exploit algorithm for comments recommendation. In *Proceedings of the 21st ACM international conference on Information and knowledge management*, pages 6–15. ACM, 2012.
- [26] S. Pandey, D. Chakrabarti, and D. Agarwal. Multi-armed bandit problems with dependent arms. In *Proceedings of the 24th international conference on Machine learning*, pages 721–728. ACM, 2007.
- [27] S. Pandey and C. Olston. Handling advertisements of unknown quality in search advertising. In *Advances in neural information processing systems*, pages 1065–1072, 2007.
- [28] L. Peska, K. Buza, and J. Koller. Drug-target interaction prediction: A bayesian ranking approach. *Computer methods and programs in biomedicine*, 152:15–21, 2017.
- [29] A. M. Rashid, I. Albert, D. Cosley, S. K. Lam, S. M. McNee, J. A. Konstan, and J. Riedl. Getting to know you: learning new user preferences in recommender systems. In *Proceedings of the 7th international conference on Intelligent user interfaces*, pages 127–134. ACM, 2002.
- [30] A. M. Rashid, G. Karypis, and J. Riedl. Learning preferences of new users in recommender systems: an information theoretic approach. *ACM SIGKDD Explorations Newsletter*, 10(2):90–100, 2008.
- [31] R. Salakhutdinov and A. Mnih. Probabilistic matrix factorization. In *Nips*, volume 1, pages 2–1, 2007.
- [32] R. Salakhutdinov and A. Mnih. Bayesian probabilistic matrix factorization using markov chain monte carlo. In *Proceedings of the 25th international conference on Machine learning*, pages 880–887. ACM, 2008.
- [33] 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.
- [34] J. B. Schafer, D. Frankowski, J. Herlocker, and S. Sen. Collaborative filtering recommender systems. In *The adaptive web*, pages 291–324. Springer, 2007.
- [35] A. I. Schein, A. Popescul, L. H. Ungar, and D. M. Pennock. Methods and metrics for cold-start recommendations. In *Proceedings of the 25th annual international ACM SIGIR con-*

ference on Research and development in information retrieval, pages 253–260. ACM, 2002.

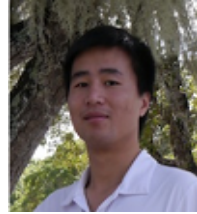
- [36] A. Smith, A. Doucet, N. de Freitas, and N. Gordon. *Sequential Monte Carlo methods in practice*. Springer Science & Business Media, 2013.
- [37] L. Song, C. Tekin, and M. van der Schaar. Online learning in large-scale contextual recommender systems. *IEEE Transactions on Services Computing*, 9(3):433–445, 2016.
- [38] M. Tokic. Adaptive  $\epsilon$ -greedy exploration in reinforcement learning based on value differences. In *KI 2010: Advances in Artificial Intelligence*, pages 203–210. Springer, 2010.
- [39] H. Wang, Q. Wu, and H. Wang. Learning hidden features for contextual bandits. In *Proceedings of the 25th ACM International Conference on Information and Knowledge Management*, pages 1633–1642. ACM, 2016.
- [40] H. Wang, Q. Wu, and H. Wang. Factorization bandits for interactive recommendation. In *AAAI*, pages 2695–2702, 2017.
- [41] Q. Wang, T. Li, S. Iyengar, L. Shwartz, and G. Y. Grabarnik. Online it ticket automation recommendation using hierarchical multi-armed bandit algorithms. In *Proceedings of the 2018 SIAM International Conference on Data Mining*, pages 657–665. SIAM, 2018.
- [42] X. Wang, S. C. Hoi, C. Liu, and M. Ester. Interactive social recommendation. In *Proceedings of the 2017 ACM Conference on Information and Knowledge Management*, pages 357–366. ACM, 2017.
- [43] L. Wu, C.-J. Hsieh, and J. Sharpnack. Large-scale collaborative ranking in near-linear time. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 515–524. ACM, 2017.
- [44] Q. Wu, H. Wang, Q. Gu, and H. Wang. Contextual bandits in a collaborative environment. In *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*, pages 529–538. ACM, 2016.
- [45] C. Zeng, Q. Wang, S. Mokhtari, and T. Li. Online context-aware recommendation with time varying multi-armed bandit. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 2025–2034. ACM, 2016.
- [46] C. Zeng, Q. Wang, W. Wang, T. Li, and L. Shwartz. Online inference for time-varying temporal dependency discovery from time series. In *Big Data (Big Data), 2016 IEEE International Conference on*, pages 1281–1290. IEEE, 2016.
- [47] X. Zhao, W. Zhang, and J. Wang. Interactive collaborative filtering. In *Proceedings of the 22nd ACM international conference on Conference on information and knowledge management*, pages 1411–1420. ACM, 2013.
- [48] L. Zhou and E. Brunskill. Latent contextual bandits and their application to personalized recommendations for new users. *arXiv preprint arXiv:1604.06743*, 2016.



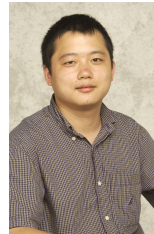
**Qing Wang** is currently a Ph.D. student in the School of Computer Science at Florida International University. She received a B.S. and M.S. degree in Computer Science from Zhengzhou University and Xidian University in 2009 and 2013, respectively. Her research interests include interactive recommender system, multi-armed bandit and large scale data mining.



**Chunqiu Zeng** is currently a Ph.D. student in the School of Computer Science at Florida International University. He received a B.S. and M.S. degree in Computer Science from Sichuan University in 2006 and 2009, respectively. His research interests include event mining, system management and large scale data mining.



**Wubai Zhou** is currently a Ph.D. student in the School of Computer Science at Florida International University. He received a B.S. degree in Computer Science and Technology from Wuhan University in 2012. His research interests include system oriented data mining, system management and machine learning.



**Tao Li** was an Eminent Scholar Professor in the School of Computer Science at Florida International University deceased in Dec. 2017. He received his Ph.D. in Computer Science in 2004 from the University of Rochester. His research interests lie in data mining and machine learning studying both on the algorithmic and application issues, computing system management, and information retrieval. He is a recipient of the USA NSF CAREER Award and multiple IBM Faculty Research Awards.



**S. S. Iyengar** currently is the Director and Ryder Professor in the School of Computer Science at Florida International University. He is a Member of the European Academy of Sciences, a Fellow of the Institute of Electrical and Electronics Engineers (IEEE), a Fellow of the Association of Computing Machinery (ACM), a Fellow of the American Association for the Advancement of Science (AAAS), and Fellow of the Society for Design and Process Science (SDPS). Dr. S. S. Iyengar is a leading researcher in the fields of distributed sensor networks, computational robotics, and oceanographic applications, and is perhaps best known for introducing novel data structures and algorithmic techniques for large-scale computations in sensor technologies and image processing applications.



**Larisa Shwartz** (M'05) is currently a Researcher at IBM T.J. Watson Research Center, Yorktown Heights, NY, USA. She received Ph.D. in mathematics from UNISA University. Dr. Shwartz has research experience in mathematics and computer science, now she is focusing on IT service management technologies for service delivery. She has more than 55 publications

and 52 patents.

**Genady Ya. Grabarnik** currently teaches at Math and CS Department, St John's University. He is a trained mathematician and authored over 80 papers. He spent 10 years at IBM T.J. Watson Research Center where his work was celebrated with a number of awards including Outstanding Technical Achievement Award and Research Achievement Awards. He is a prolific inventor with over 65 US patents. His interests include research in functional analysis, inventions, and research in computer science and artificial intelligence.