

Alleviating Matthew Effect of Offline Reinforcement Learning in Interactive Recommendation

Chongming Gao
chongming.gao@gmail.com
University of Science and Technology
of China

Kexin Huang
huangkx@mail.ustc.edu.cn
University of Science and Technology
of China

Jiawei Chen*
sleepyhunt@zju.edu.cn
Zhejiang University
Hangzhou, China

Yuan Zhang
yuanz.pku@gmail.com
Kuaishou Technology Co., Ltd.

Biao Li
libiao@kuaishou.com
Kuaishou Technology Co., Ltd.

Peng Jiang
jiangpeng@kuaishou.com
Kuaishou Technology Co., Ltd.

Shiqi Wang
shiqi@cqu.edu.cn
Chongqing University
Chongqing, China

Zhong Zhang
zhongzhang@std.ustc.edu.cn
University of Electronic Science and
Technology of China

Xiangnan He*
xiangnanhe@gmail.com
University of Science and Technology
of China

ABSTRACT

Offline reinforcement learning (RL), a technology that offline learns a policy from logged data without the need to interact with online environments, has become a favorable choice in decision-making processes like interactive recommendation. Offline RL faces the value overestimation problem. To address it, existing methods employ conservatism, e.g., by constraining the learned policy to be close to behavior policies or punishing the rarely visited state-action pairs. However, when applying such offline RL to recommendation, it will cause a severe Matthew effect, i.e., the rich get richer and the poor get poorer, by promoting popular items or categories while suppressing the less popular ones. It is a notorious issue that needs to be addressed in practical recommender systems.

In this paper, we aim to alleviate the Matthew effect in offline RL-based recommendation. Through theoretical analyses, we find that the conservatism of existing methods fails in pursuing users' long-term satisfaction. It inspires us to add a penalty term to relax the pessimism on states with high entropy of the logging policy and indirectly penalizes actions leading to less diverse states. This leads to the main technical contribution of the work: *Debiased model-based Offline RL* (DORL) method. Experiments show that DORL not only captures user interests well but also alleviates the Matthew effect. The implementation is available via <https://github.com/chongminggao/DORL-codes>.

CCS CONCEPTS

• **Information systems** → **Recommender systems**.

*Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGIR '23, July 23–27, 2023, Taipei, Taiwan

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-9408-6/23/07...\$15.00
<https://doi.org/10.1145/3539618.3591636>

KEYWORDS

Offline Reinforcement Learning; Recommendation; Matthew effect

ACM Reference Format:

Chongming Gao, Kexin Huang, Jiawei Chen, Yuan Zhang, Biao Li, Peng Jiang, Shiqi Wang, Zhong Zhang, and Xiangnan He. 2023. Alleviating Matthew Effect of Offline Reinforcement Learning in Interactive Recommendation. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '23)*, July 23–27, 2023, Taipei, Taiwan. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3539618.3591636>

1 INTRODUCTION

Recommender systems, a powerful tool for helping users select preferred items from massive items, are continuously investigated by e-commerce companies. Previously, researchers tried to dig up static user interests from historical data by developing supervised learning-based recommender models. With the recent development of deep learning and the rapid growth of available data, fitting user interests is not a bottleneck for now.

A desired recommendation policy should be able to satisfy users for a long time [45]. Therefore, it is natural to involve Reinforcement Learning (RL) which is a type of Machine Learning concerned with how an intelligent agent can take actions to pursue a long-term goal [1, 4, 5, 35]. In this setting, the recommendation process is formulated as a sequential decision process where the recommender interacts with users and receives users' online feedback (i.e., rewards) to optimize users' long-term engagement, rather than fitting a model on a set of samples based on supervised learning [7, 53, 62]. However, it is expensive and impractical to learn a policy from scratch with real users, which becomes the main obstacle that impedes the deployment of RL to recommender systems.

One remedy is to leverage historical interaction sequences, i.e., recommendation logs, to conduct offline RL (also called batch RL) [33]. The objective is to learn an online policy that makes counterfactual decisions to perform better than the behavior policies induced by the offline data. However, without real-time feedback, directly employing conventional online RL algorithms in offline

scenarios will result in poor performance due to the value overestimation problem in offline RL. The problem is induced when the function approximator of the agent tries to extrapolate values (e.g., Q-values in Q-learning [47]) for the state-action pairs that are not well-covered by logged data. More specifically, since the RL model usually maximizes the expected value or trajectory reward, it will intrinsically prefer overestimated values induced by the extrapolation error, and the error will be compounded in the bootstrapping process when estimating Q-values, which results in unstable learning and divergence [30]. In recommendation, this may lead to an overestimation of user preferences for items that infrequently appear in the offline logs.

This is the core challenge for offline RL algorithms because of the inevitable mismatch between the offline dataset and the learned policy [43]. To solve this problem, offline RL algorithms incorporate conservatism into the policy design. Model-free offline RL algorithms directly incorporate conservatism by constraining the learned policy to be close to the behavior policy [11, 30], or by penalizing the learned value functions from being over-optimistic upon out-of-distribution (OOD) decisions [31]. Model-based offline RL algorithms learn a pessimistic model as a proxy of the environment, which results in a conservative policy [27, 55]. This philosophy guarantees that offline RL models can stick to offline data without making OOD actions, which has been proven to be effective in lots of domains, such as robotic control [9] or games [2].

However, applying conservatism to recommender systems gives rise to a severe Matthew effect [36], which can be summarized as “the rich get richer and the poor get poorer”. In recommendation, it means that the popular items or categories in previous data will get larger opportunities to be recommended later, whereas the unpopular ones get neglected. This is catastrophic since users desire diverse recommendations and the repetition of certain contents will incur the filter bubble issue, which in turn hurts users’ satisfaction even though users favored them before [12, 16, 41, 58]. We will show the Matthew effect in the existing offline RL-based recommender (Fig. 3), and analyze how users’ satisfaction will be hurt (Fig. 2).

In this paper, we embrace the model-based RL paradigm. The basic idea is to learn a user model (i.e., world model) that captures users’ preferences, then use it as a pseudo-environment (i.e., simulated users) to produce rewards to train a recommendation policy. Compared to model-free RL, model-based RL has several advantages in recommendation. First and foremost, model-based RL is much more sample efficient [55]. That it needs significantly fewer samples makes it more suitable for the highly sparse recommendation data. Second, explicitly learning the user model simplifies the problem and makes it easier to incorporate expert knowledge. For example, the user model can be implemented as any state-of-the-art recommendation model (e.g., DeepFM [19] in this work) or sophisticated generative adversarial frameworks [53, 59]. Although some works have adopted this paradigm in their recommender systems [12, 22, 23, 63], they did not explicitly consider the value overestimation problem in offline RL, not to mention the Matthew effect in the solutions.

To address the value overestimation problem while reducing the Matthew effect, we propose a *Debiased model-based Offline RL* (DORL) method for recommendation. By theoretically analyzing the mismatch between real users’ long-term satisfaction and the

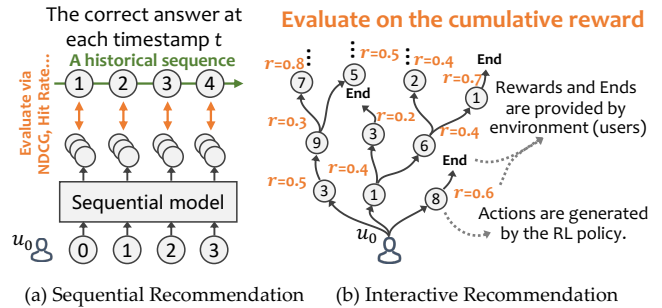


Figure 1: Illustration of evaluations in the traditional sequential recommendation and interactive recommendation settings.

preferences estimated from the offline data, DORL adds a penalty term that relaxes the pessimism on states with high entropy of the logging policy and indirectly penalizes actions leading to less diverse states. By introducing such a counterfactual exploration mechanism, DORL can alleviate the Matthew effect in final recommendations. Our contributions are summarized as:

- We point out that conservatism in offline RL can incur the Matthew effect in recommendation. We show this phenomenon in existing methods and how it hurts user satisfaction.
- After theoretically analyzing how existing methods fail in recommendation, we propose the DORL model that introduces a counterfactual exploration in offline data.
- We demonstrate the effectiveness of DORL in an interactive recommendation setting, where alleviating the Matthew effect increases users’ long-term experience.

2 RELATED WORK

Here, we briefly review the Matthew effect in recommendation. We introduce the interactive recommendation and offline RL.

2.1 Matthew Effect in Recommendation

Liu and Huang [36] confirmed the existence of the Matthew effect in YouTube’s recommendation system, and Wang et al. [42] gave a quantitative analysis of the Matthew effect in collaborative filtering-based recommenders. A common way of mitigating the Matthew effect in recommendation is to take into account diversity [3, 20, 34, 60]. Another perspective on this problem is to remove popularity bias [61]. We consider the Matthew effect in offline RL-based recommendation systems. We will analyze why this problem occurs and provide a novel way to address it.

2.2 Interactive Recommendation

The interactive recommendation is a setting where a model interacts with a user online [13, 32, 48]. The model recommends items to the user and receives the user’s real-time feedback. This process is repeated until the user quits. The model will update its policy with the goal to maximize the cumulative satisfaction over the whole interaction process (instead of learning on I.I.D. samples). This setting well reflects the real-world recommendation scenarios, for example, a user will continuously watch short videos and leave feedback (e.g. click, add to favorite) until he chooses to quit.

Here, we emphasize the most notable difference between the interactive recommendation setting and traditional sequential recommendation settings [50, 51]. Fig. 1 illustrates the learning and evaluation processes in sequential and interactive recommendation settings. Sequential recommendation uses the philosophy of supervised learning, i.e., evaluating the top- k results by comparing them with a set of “correct” answers in the test set and computing metrics such as Precision, Recall, NDCG, and Hit Rate. By contrast, interactive recommendation evaluates the results by accumulating the rewards along the interaction trajectories. There is no standard answer in interactive recommendation, which is challenging [8].

Interactive recommendation requires offline data of high quality, which hampers the development of this field for a long time. We overcome this problem by using the recently-proposed datasets that support interactive learning and off-policy evaluation [18].

2.3 Offline Reinforcement Learning

Recently, many offline RL models have been proposed to overcome the value overestimation problem. For model-free methods, BCQ [11] uses a generative model to constrain probabilities of state-action pairs the policy utilizes, thus avoiding using rarely visited data to update the value network; CQL [31] contains a conservative strategy to penalize the overestimated Q-values for the state-action pairs that have not appeared in the offline data; GAIL [21] utilizes a discriminator network to distinguish between expert policies with others for imitation learning. IQL [29] enables the learned policy to improve substantially over the best behavior in the data through generalization, without ever directly querying a Q-function with unseen actions. For model-based methods, MOPO [55] learns a pessimistic dynamics model and use it to learn a conservative estimate of the value function; COMBO [54] learns the value function based on both the offline dataset and data generated via model rollouts, and it suppresses the value function on OOD data generated by the model. Almost all offline RL methods have a similar philosophy: to introduce conservatism or pessimism in the learned policy [33].

There are efforts to conduct offline RL in recommendation [6, 12, 23, 24, 49]. However, few works explicitly discuss the Matthew effect in recommendation. Chen et al. [6] mentioned this effect in their experiment section, but their method is not tailored to overcome this issue.

3 EMPIRICAL STUDY ON MATTHEW EFFECT

We conduct empirical studies in recommendation to show how the Matthew effect affects user satisfaction. When the Matthew effect is amplified, the recommender will repeatedly recommend the items with dominant categories. To illustrate the long-term effect on user experience, we explore the logs of the KuaiRand-27K video dataset¹ [15] and the LFM-1b music dataset² [38]. KuaiRand-27K contains a 23 GB log recording 27,285 users’ 322,278,385 interactions on 32,038,725 videos with 62 categories, which are collected from April 8th, 2022 to May 8th, 2022. LFM-1b contains a 40GB log recording 120,322 users’ 1,088,161,692 listening events on 32,291,134 tracks with 3,190,371 artists, which are fetched from Last.FM in the range from January 2013 to August 2014.

¹<https://kuairand.com/>

²<http://www.cp.jku.at/datasets/LFM-1b/>

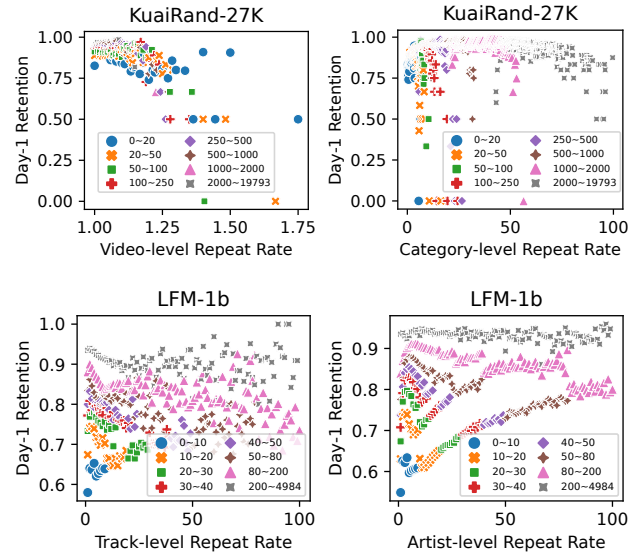


Figure 2: An increase in the repeat rate (i.e., the Matthew effect) negatively impacts user experience. In nearly all activity level groups marked by different colors and maker types, Day-1 Retention of users declines as the repeat rate rises.

Both the two datasets provide the timestamp of each event, hence we can assess the long-term effect of overexposure by investigating the change of Day-1 Retention. Day-1 Retention is defined as the probability of a user who returns to the app tomorrow after finishing today’s viewing/listening. This metric is more convincing than real-time signals (e.g., click, adding to favorite) in regard to reflecting the long-term effect on user satisfaction. We consider the item-level and category/artist-level repeat rates as the metrics to measure the Matthew effect. The item-level (or category-level) repeat rate of a user viewing videos on a certain day is defined as:

$$\frac{\text{the number of viewing events}}{\text{the number of unique videos (or unique categories)}}$$

For example, if a user views 5 unique videos (which belong to 3 unique categories) 20 times in a day, then the item-level repeat rate is $20/5=4.0$ and the category-level repeat rate is $20/3=6.67$. The item-level and artist-level repeat rates for music listening are defined in a similar way. Note that in KuaiRand, video-level overexposure rarely appears because of the rule of video recommendation, i.e., the same video will not be recommended twice.

In this definition, user activity can become a confounder. For instance, a user who views 100 videos a day can be more active than a user who views only 10 videos a day, and thus is more likely to revisit the App the next day. Therefore, we control for this confounder by splitting the users w.r.t. the number of their daily viewing events. Groups with different user activity levels are marked by different colors and marker types. The results are shown in Fig. 2.

In short, Day-1 Retention reduces when the repeat rate increases in each group with a user activity level. This phenomenon can

be observed for both the item-level and category/artist-level repeat rates within the video dataset and music dataset. The results show that users' satisfaction will be hurt when the Matthew effect becomes severer.

4 PRELIMINARY ON MODEL-BASED RL

We introduce the basics of RL and model-based offline RL.

4.1 Basics of Reinforcement Learning

Reinforcement learning (RL) is the science of decision making. We usually formulate the problem as a Markov decision process (MDP): $M = (\mathcal{S}, \mathcal{A}, T, r, \gamma)$, where \mathcal{S} and \mathcal{A} represent the state space and action space, $T(s, a, s') = P(s_{t+1} = s' | s_t = s, a_t = a)$ is the transition probability from (s, a) to s' , $r(s, a)$ is the reward of taking action a at state s , and γ is the discount factor. Accordingly, the offline MDP can be denoted as $\widehat{M} = (\mathcal{S}, \mathcal{A}, \widehat{T}, \widehat{r}, \gamma)$, where \widehat{T} and \widehat{r} are the transition probability and reward function predicted by an offline model. In offline RL, the policy is trained on an offline dataset \mathcal{D} which was collected by a behavior policy π_β running in online environment M . By modifying the offline MDP \widehat{M} to be conservative for overcoming the value overestimation issue, we will derive a modified MDP $\widetilde{M} = (\mathcal{S}, \mathcal{A}, \widetilde{T}, \widetilde{r}, \gamma)$, where the modified reward \widetilde{r} is modified from the predicted reward \widehat{r} .

Since RL considers long-term utility, we can define the value function as $V_M^\pi(s) = \mathbb{E}_{\pi, T}[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | s_0 = s]$, denoting the cumulative reward gain by policy π after state s in MDP M . Let $P_{T,t}^\pi$ be the probability of the agent's being in state s at time t , if the agent uses policy π and transits with T . Defining $\rho_T^\pi(s, a) = (1 - \gamma)\pi(a|s) \sum_{t=0}^{\infty} \gamma^t P_{T,t}^\pi(s)$ as the discounted distribution of state-action pair (s, a) for policy π over T , we can derive another form of the policy's accumulated reward as $\eta_M(\pi) = \mathbb{E}_{(s,a) \sim \rho_T^\pi} [r(s, a)]$.

4.2 Model-Based Offline RL Framework

In this paper, we follow a state-of-the-art general Model-based Offline Policy Optimization framework, MOPO³ [55]. The basic idea is to learn a dynamics model \widehat{T} which captures the state transition $(s, a) \rightarrow s'$ of the environment and estimates reward $\widehat{r}(s, a)$ given state s and action a . For addressing the distributional shift problem where values $V_M^\pi(s)$ are usually over-optimistically estimated, MOPO introduces a penalty function $p(s, a)$ on the estimated reward $\widehat{r}(s, a)$ as:

$$\widetilde{r}(s, a) = \widehat{r}(s, a) - \lambda p(s, a). \quad (1)$$

On the modified reward $\widetilde{r}(s, a)$, the offline MDP \widehat{M} will be modified to be a conservative MDP: $\widetilde{M} = (\mathcal{S}, \mathcal{A}, \widetilde{T}, \widetilde{r}, \gamma)$. MOPO learns its policy in this MDP: \widetilde{M} . By defining $\epsilon_p(\pi) = \mathbb{E}_{(s,a) \sim \rho_T^\pi} [p(s, a)]$, MOPO has the following theoretical guarantee:

THEOREM 4.1. *If the penalizer $p(s, a)$ meets:*

$$\lambda \mathbb{E}_{(s,a) \sim \rho_T^\pi} [p(s, a)] \geq |\eta_{\widetilde{M}}(\pi) - \eta_M(\pi)|, \quad (2)$$

then the best offline policy $\widehat{\pi}$ trained in \widetilde{M} satisfies:

$$\eta_M(\widehat{\pi}) \geq \sup_{\pi} \{\eta_M(\pi) - 2\lambda \epsilon_p(\pi)\}. \quad (3)$$

³We consider an RL framework to be general if it doesn't require domain-related prior knowledge or any specific algorithms. Satisfying our demands, MOPO is shown to be one of the best-performing model-based offline RL frameworks.

The proof can be found in [55]. Eq. (2) requires the penalty to be a measurement of offline and online mismatch, thus $\epsilon_p(\pi)$ can be interpreted as how much policy π will be affected by the offline extrapolation error. Eq. (3) is considered to be a theoretical guarantee for reward penalty in model-based offline RL. For example, with π^* denoting optimal policy in online MDP M , we have $\eta_M(\widehat{\pi}) \geq \eta_M(\pi^*) - 2\lambda \epsilon_p(\pi^*)$.

Remark: Through learning $\widehat{\pi}$ offline in the conservative MDP \widetilde{M} with Eq. (1), we can obtain the result that will not deviate too much from the result of learning an optimal policy π^* online in the ground-truth MDP M . The deviation will not exceed $2\lambda \epsilon_p(\pi^*)$.

However, there was no sufficient analysis on how to properly choose the penalty term $p(s, a)$. Next, We introduce how to adapt this framework to recommendation and reformulate the $p(s, a)$ according to the characteristic of the recommendation scenario.

5 METHOD

We implement the model-based offline RL framework in recommendation. we redesign the penalty to alleviate the accompanied Matthew effect. Then, we introduce the proposed DORL model.

5.1 Model-based RL in Recommendation

In recommendation, we cannot directly obtain a state from the environment, we have to model the state by capturing the interaction context and the user's mood. Usually, a state $s \in \mathcal{S}$ is defined as the vector extracted from the user's previously interacted items and corresponding feedback. After the system recommends an item as action $a \in \mathcal{A}$, the user will give feedback as a scale reward signal $r := R(s, a)$. For instance, $r \in \{0, 1\}$ indicates whether the user clicks the item, or $r \in \mathbb{R}^+$ reflects a user's viewing time for a video. The state transition function (i.e., state encoder) T can be written as $s' := f_\omega(s, a, r)$, where $f_\omega(s, a, r)$ autoregressively outputs the next state s' and can be implemented as any sequential models.

When learning offline, we cannot obtain users' reactions to the items that are not covered by the offline dataset. we address this problem by using a user model (or reward model) $\widehat{R}(s, a)$ to learn users' static interests. This model can be implemented as any state-of-the-art recommender such as DeepFM [19]. The user model will generate an estimated reward $\widehat{r} = \widehat{R}(s, a)$ representing a user's intrinsic interest in an item. The transition function \widehat{T} will be written as $s' := f_\omega(s, a, \widehat{r})$. The offline MDP is defined as $\widehat{M} = (\mathcal{S}, \mathcal{A}, \widehat{T}, \widehat{r}, \gamma)$. Since the estimated reward \widehat{r} can deviate from the ground-truth value r , we follow MOPO to use Eq. (1) to get the modified reward $\widetilde{r}(s, a) = \widehat{r}(s, a) - \lambda p(s, a)$. Afterward, we can train the recommendation policy on the modified reward \widetilde{r} by treating the user model as simulated users.

Now, the problem turns into designing the penalty term $p(s, a)$. To begin with, we extend the mismatch function in [37, 55]. We use R and \widehat{R} as the shorthand for $R(s, a)$ and $\widehat{R}(s, a)$, respectively.

LEMMA 5.1. *Define the mismatch function $G_M^\pi(s, a)$ of a policy π on the ground truth MDP M and the estimated MDP \widehat{M} as:*

$$\begin{aligned} G_M^\pi(s, a) &:= \mathbb{E}_{s' \sim \widehat{T}, \widehat{r} \sim \widehat{R}} [\gamma V_M^\pi(s') + \widehat{r}] - \mathbb{E}_{s' \sim T, r \sim R} [\gamma V_M^\pi(s') + r] \\ &= \mathbb{E}_{\widehat{r} \sim \widehat{R}} [\gamma V_M^\pi(f_\omega(s, a, \widehat{r})) + \widehat{r}] - \mathbb{E}_{r \sim R} [\gamma V_M^\pi(f_\omega(s, a, r)) + r]. \end{aligned} \quad (4)$$

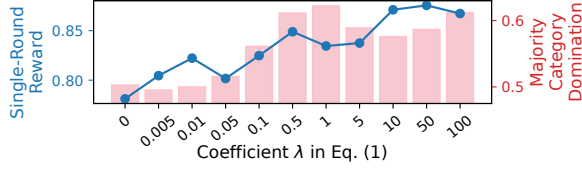


Figure 3: Effect of conservatism in recommendation.

It satisfies:

$$\mathbb{E}_{(s,a) \sim p_T^\pi} [G_M^\pi(s, a)] = \eta_{\hat{M}}(\pi) - \eta_M(\pi). \quad (5)$$

The mismatch function $G_M^\pi(s, a)$ extends the definition presented in [55], with the key distinction being the separation of state transition function into state s and reward r . This is due to the fact that, in the context of recommendation systems, the stochastic nature of state probabilities arises solely from the randomness associated with their reward signals r . Consequently, when integrating along the state transition, it is essential to explicitly express the impact of reward r . The proof of Lemma 5.1 can be adapted from the proof procedure of the telescoping lemma in [37].

Following the philosophy of conservatism, we add a penalty term $p(s, a)$ according to the mismatch function $G_M^\pi(s, a)$ by assuming:

$$\lambda p(s, a) \geq |G_M^\pi(s, a)|. \quad (6)$$

By combining Eq. (5) and Eq. (6), the condition in Eq. (2) is met, which provides the theoretical guarantee for the recommendation policy π learned in the conservative MDP: \tilde{M} .

Remark: Lemma 5.1 provides a perspective for designing the penalty term $p(s, a)$ that satisfies the theoretical guarantee in Theorem 4.1. According to Eq. (6), the problem of defining $p(s, a)$ turns into analyzing $G_M^\pi(s, a)$, which will be described in Section 5.3.

The original MOPO model uses the uncertainty of the dynamics model P_U as the penalty, i.e., $p(s, a) = P_U$. However, penalizing uncertainty will encourage the model to pay more attention to items that are frequently recommended while neglecting the rarely recommended ones. This will accelerate the Matthew effect.

5.2 Matthew Effect

To quantify the Matthew effect in the results of recommendation, we use a metric: majority category domination (MCD), which is defined as the percentage of the recommended items that are labeled as the dominated categories in training data⁴.

We show the effect of conservatism of MOPO by varying the coefficient λ of Eq. (1). The results on the KuaiRec dataset are shown in Fig. 3. With increasing λ , the model receives a higher single-round reward (the blue line), which means the policy captures users' interest more accurately. On the other hand, MCD also increases (the red bars), which means the recommended items tend to be the most popular categories (those cover 80% items) in training data. I.e., **the more conservative the policy is, the stronger the Matthew effect becomes.**

⁴The dominated categories are the most popular categories that cover 80% items in the training set. There are 13 (out of 46) dominated categories in KuaiRand, and 12 (out of 31) dominated categories in KuaiRec.

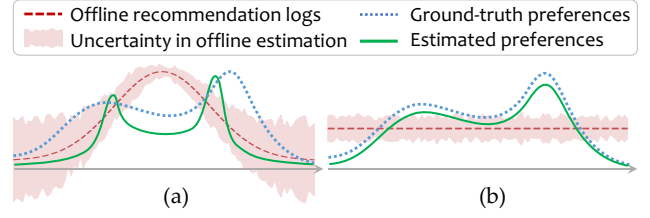


Figure 4: Estimating user preferences using behavior policies with a Gaussian distribution and a uniform distribution.

When the results narrow down to these categories, users' satisfaction will be hurt and the interaction process terminates early, which results in low cumulative rewards over the interaction sequence. More details will be described in Section 6.

5.3 Solution: Re-design the Penalty

To address this issue, we consider a more sophisticated manner to design the penalty term $p(s, a)$ in Eq. (1). We dissect the mismatch function in Eq. (4) as:

$$|G_M^\pi(s, a)| \leq \gamma |\mathbb{E}_{\hat{r} \sim \hat{R}} [V_M^\pi(f_\omega(s, a, \hat{r}))] - \mathbb{E}_{r \sim R} [V_M^\pi(f_\omega(s, a, r))]| + |\mathbb{E}_{\hat{r} \sim \hat{R}} r - \mathbb{E}_{r \sim R} r| := \gamma d_V(\hat{R}, R) + d_1(\hat{R}, R), \quad (7)$$

where $d_1(\hat{R}, R)$ represents the deviation of estimated reward \hat{R} from true reward R , $d_V(\hat{R}, R)$ measures the difference between the value functions V_M^π of next state calculated offline (via \hat{R}) and online (via R). Both of them can be seen as specific metrics measuring the distance between \hat{R} and R . While $d_1(\hat{R}, R)$ is straightforward, $d_V(\hat{R}, R)$ considers the long-term effect on offline learning and is hard to estimate. Based on the aforementioned analysis, a pessimistic reward model \hat{R} in MOPO will amplify the Matthew effect that reduces long-term satisfaction, thus resulting in a large $d_V(\hat{R}, R)$.

An intuitive way to solve this dilemma is to introduce exploration on states with high entropy of the logging policy. Without access to online user feedback, we can only conduct the counterfactual exploration in the offline data.

• **An illustrative example.** To illustrate the idea, we give an example in Fig. 4. The goal is to estimate a user's preferences given the logged data induced by a behavior policy. In reality, the distribution of the logged data is dependent on the policies of previous recommenders. For convenience, we use a Gaussian distribution as the behavior policy in Fig. 4(a). Since previous recommenders cannot precisely reflect users' ground-truth preferences, there is always a deviation between the behavior policy (the red line) and users' ground-truth preferences (the blue line).

Besides, as the items are not equally exposed in the behavior policy, there will be high uncertainty in estimating the rarely appeared items (as shown in the filled area). Offline RL methods emphasize conservatism in estimation and penalize the uncertain samples, which results in the distribution that narrows down the preferences to these dominated items (the green line). This is how the Matthew effect appears.

By contrast, using a uniform distribution to collect data can prevent biases and reduce uncertainty (Fig. 4(b)). Ideally, a policy

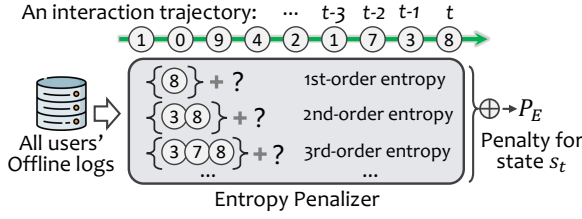


Figure 5: Illustration of the mechanism of entropy penalizer.

learned on sufficient data collected uniformly can capture unbiased user preferences and produce recommendations without the Matthew effect, i.e., $\gamma d_V(\widehat{R}, R) + d_1(\widehat{R}, R)$ can reduce to 0.

Therefore, an intuitive way to design penalty term $p(s, a)$ is to add a term: the discrepancy of behavior between the uniform distribution $\pi_u(\cdot|s)$ and the behavior policy $\pi_\beta(\cdot|s)$ given state s . We use the Kullback–Leibler divergence $D_{KL}(\pi_\beta(\cdot|s)||\pi_u(\cdot|s))$ to measure the distance, which can be written as:

$$\begin{aligned} P_E &:= -D_{KL}(\pi_\beta(\cdot|s)||\pi_u(\cdot|s)) \\ &= -\mathbb{E}_{a \sim \pi_\beta(\cdot|s)} [\log(\pi_\beta(a|s)) - \log(\pi_u(a|s))] \quad (8) \\ &= \mathcal{H}(\pi_\beta(\cdot|s)) - \log(|\mathcal{A}|), \end{aligned}$$

where $|\mathcal{A}|$ is a constant representing the number of items. Hence, the term P_E depends on the entropy of the behavior policy $\pi_\beta(\cdot|s)$ given state s . The modified penalty term can be written as $p(s, a) = P_U + P_E$, and the modified reward model will be formulated as:

$$\tilde{r}(s, a) = \hat{r}(s, a) - \lambda_1 P_U + \lambda_2 P_E. \quad (9)$$

Except for penalizing high uncertainty areas, the new model also penalizes policies with a low entropy at state s . Intuitively, if the behavior policy $\pi_\beta(\cdot|s)$ recommended only a few items at state s , then the true user preferences at state s may be unrevealed. Under such circumstances, the entropy term $\mathcal{H}(\pi_\beta(\cdot|s))$ is low, hence we penalize the estimated reward $\hat{r}(s, a)$ by a large P_E .

The entropy penalizer does not depend on the chosen action but only on the state in which the agent is. Which means the effect of this penalty will be indirect and penalize actions that lead to less diverse states, because of the long-term optimization. Hence, The learned policy achieves the counterfactual exploration in the offline data, which in turn counteracts the Matthew effect in offline RL.

5.4 The DORL Method

Now, we provide a practical implementation motivated by the analysis above. The proposed model is named Debiased model-based Offline RL (DORL) model, whose framework is illustrated in Fig. 6.

• **Penalty on entropy.** We introduce how to compute P_E in the entropy penalizer module. Fig. 5 shows a trajectory of the interaction process, where the current action at time t is to recommend item 8. We define P_E in Eq. (9) to be the summation of k -order entropy ($k = 1, 2, \dots$). For example, when $k = 3$, we search all users' recommendation logs to collect all continuous sub-sequences with pattern $[\{3, 7, 8\}, ?]$, where “?” can match any item, and $\{3, 7, 8\}$ is a sorted set that can cover all of its enumeration, e.g., $[8, 3, 7]$ or $[7, 3, 8]$. On these sub-sequences, we can count the frequencies of action “?” to estimate the entropy of behavior policy π_β given the

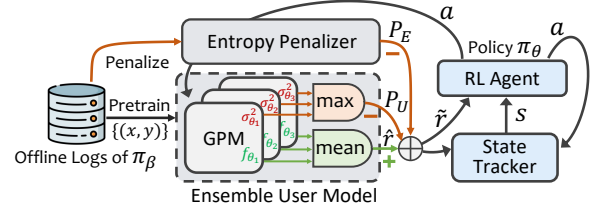


Figure 6: Illustration of Debiased Model-based Offline RL-based Recommendation (DORL) framework.

previous three recommended items. Without losing generality, we normalize the entropy to range (0, 1].

• **Penalty on uncertainty.** We penalize both the epistemic uncertainty of the reward model and the aleatoric uncertainty of offline data. We use the variance of K ensemble reward models $\{\widehat{R}_{\theta_k}, k = 1, 2, \dots, K\}$ to capture the epistemic uncertainty, which is commonly used to capture the uncertainty of the model in offline RL [33]. Aleatoric uncertainty is data-dependent [26]. By formulating the user model as a Gaussian probabilistic model (GPM), we can directly predict the variance of the reward and take this predicted variance as aleatoric uncertainty. For the k -th model \widehat{R}_{θ_k} , the loss function is:

$$\mathcal{L}(\theta_k) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2\sigma_{\theta_k}^2(x_i)} \|y_i - f_{\theta_k}(x_i)\|^2 + \frac{1}{2} \log \sigma_{\theta_k}^2(x_i), \quad (10)$$

where N is the number of samples, $f_{\theta_k}(x_i)$ and $\sigma_{\theta_k}^2(x_i)$ are the predicted mean and variance of sample x_i , respectively. By combining epistemic uncertainty and aleatoric uncertainty, we formulate the uncertainty penalizer P_U in Eq. (9) as: $P_U := \max_{k \in \{1, 2, \dots, K\}} \sigma_{\theta_k}^2$. We define the fitted reward \hat{r} as the mean of K ensemble models: $\hat{r}(s, a) = \frac{1}{K} \sum_k f_{\theta_k}(s, a)$. The final modified reward $\tilde{r}(s, a)$ will be computed by Eq. (9).

The framework of the proposed DORL model is illustrated in Fig. 6. Without losing generality, we use DeepFM [19] as the backbone for the user model and implement the actor-critic method [28] as the RL policy.

The state tracker $f_\omega(s, a, r)$ is a network modeling the transition function $T(s, a, s') = P(s_{t+1} = s' | s_t = s, a_t = a)$. It can be implemented as any sequential model such as recurrent neural network (RNN)-based models [44], Convolutional models [40, 56], Transformer-based methods [12, 25, 51]. Huang et al. [22] investigated the performances of different state encoders in RL-based recommenders. We use a naive average layer as the state tracker since it requires the least training time but nonetheless outperforms many complex encoders [22]. It can be written as:

$$\vec{s}_{t+1} := \frac{1}{N} \sum_{n=t-N+1}^t [\vec{e}_{a_n} \oplus \tilde{r}_n], \quad (11)$$

where \oplus is the concatenation symbol, \vec{s}_{t+1} is the vector representing the state at time $t + 1$, \vec{e}_{a_n} is the embedding vector of action a_n . \tilde{r}_n is the reward value calculated by Eq. (9) and we normalize it to range (0,1] here. N is the window size reflecting how many previous item-reward pairs are calculated.

6 EXPERIMENTS

We introduce how we evaluate the proposed DORL model in the interactive recommendation setting. We want to investigate the following questions:

- **(RQ1)** How does DORL perform compared to state-of-the-art offline RL methods in the interactive recommendation setting?
- **(RQ2)** To what extent can DORL alleviate the Matthew effect and pursue long-term user experience?
- **(RQ3)** How does DORL perform in different environments with different user tolerance to repeated content?

6.1 Experimental Setup

We introduce the experimental settings with regard to environments and state-of-the-art offline RL methods.

6.1.1 Recommendation Environments. As mentioned in Section 2.2, in the interactive recommendation setting, we are interested in users’ long-term satisfaction rather than users’ fitting capabilities [8]. **Traditional recommendation datasets are too sparse or lack necessary information** (e.g., timestamps, explicit feedback, item categories) to evaluate the interactive recommender systems. We create two recommendation environments on two recently-proposed datasets, KuaiRec and KuaiRand-Pure, which contain high-quality logs.

KuaiRec [14] is a video dataset that contains a fully-observed user-item interaction matrix where 1,411 users have viewed all 3,327 videos and left feedback. By taking the fully-observed matrix as users’ true interest, we can give a reward for the model’s every recommendation (without missing entries like other datasets). We use the normalized viewing time (i.e., the ratio of viewing time to the video length) as the online reward.

KuaiRand-Pure [15] is a video dataset that inserted 1,186,059 random recommendations involving 7,583 items into 27,285 users’ standard recommendation streams. These randomly exposed data can reflect users’ unbiased preferences, from which we can complete the matrix to emulate the fully-observed matrix in KuaiRec. This is an effective way to evaluate RL-based recommendation [22, 23]. We use the “is_click” signal to indicate users’ ground-truth interest, i.e., as the online reward.

In Section 3, we have shown that users’ experience can be hurt by the Matthew effect. To let the environments reflect this phenomenon, we follow [12, 52] to introduce a quit mechanism: when the model recommends more than M items with the same category in previous N rounds, the interaction terminates. Note that the same item will not be recommended twice in an interaction sequence. Since we evaluate the model via the cumulative rewards $\sum_t r_t$ over the interaction trajectory, quitting early (due to the Matthew effect) will lead to inferior performances. For now, the two environments can play the same role as the online users. Therefore, we can evaluate the model as the process shown in Fig. 1 (b).

The evaluation environments are used for assessing models and they are not available in the training stage. For the training purpose, both KuaiRec and KuaiRand provide additional recommendation logs. The statistics of the training data are illustrated in Table 1.

6.1.2 Baselines. We select two naive bandit-based algorithms, four model-free offline RL methods, and four model-based offline RL

Table 1: Statistics of two datasets.

Datasets	Usage	#Users	#Items	#Interactions	#Categories
KuaiRec	Train	7,176	10,728	12,530,806	31
	Test	1,411	3,327	4,676,570	31
KuaiRand	Train	27,285	7,551	1,436,609	46
	Test	27,285	7,583	1,186,059	46

methods (including ours) in evaluation. We use the DeepFM model [19] as the backbone in the two bandit methods and four model-based methods. These baselines are:

- **ϵ -greedy**, a naive bandit-based policy that outputs a random result with probability ϵ or outputs the deterministic results of DeepFM with probability $1 - \epsilon$.
- **UCB**, a naive bandit-based policy that maintains an upper confidence bound for each item and follows the principle of optimism in the face of uncertainty.
- **SQN**, or Self-Supervised Q-learning [50], contains two output layers (heads): one for the cross-entropy loss and the other for RL. We use the RL head to generate final recommendations.
- **BCQ**, or Batch-Constrained deep Q-learning [11], adapts the conventional deep Q-learning to batch RL. We use the discrete-action version [10], whose core idea is to reject these uncertain data and update the policy using only the data of high confidence.
- **CQL**, or Conservative Q-Learning [31], is a model-free RL method that adds a Q-value regularizer on top of an actor-critic policy.
- **CRR**, or Critic Regularized Regression [46], is a model-free RL method that learns the policy by avoiding OOD actions.
- **MBPO**, a vanilla model-based policy optimization method that uses DeepFM as the user model to train an actor-critic policy.
- **IPS** [39] is a well-known statistical technique adjusting the target distribution by re-weighting each sample in the collected data. We implement IPS in a DeepFM-based user model, then learn the policy using an actor-critic method.
- **MOPO**, a model-based offline policy optimization method [55] that penalizes the uncertainty of the DeepFM-based user model and then learns an actor-critic policy.

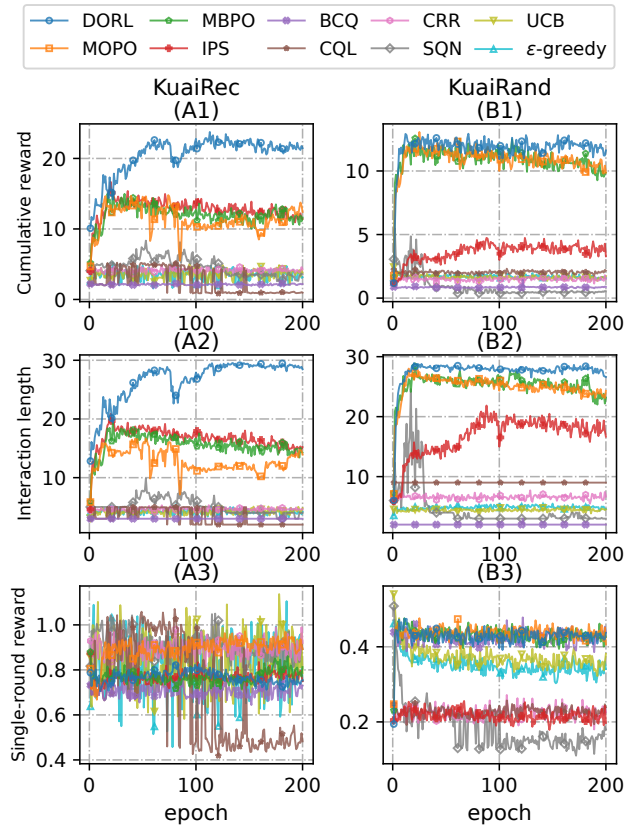
6.2 Overall Performance Comparison (RQ1)

We evaluate all methods in two environments. For the four model-based RL methods (MBPO, IPS, MOPO, and our DORL), we use the same DeepFM model as the user model and fixed its parameters to make sure the difference comes only from the policies. We use the grid search technique on the key parameters to tune all methods in the two environments. For DORL, we search the combination of two key parameters λ_1 and λ_2 in Eq. (9). Both of them are searched in $\{0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1.0, 5, 10, 50, 100\}$. We report the results with $\lambda_1 = 0.01, \lambda_2 = 0.05$ for KuaiRand and $\lambda_1 = 0.05, \lambda_2 = 5$ for KuaiRec. All methods in two environments are evaluated with the quit parameters: $M = 0, N = 4$, and the maximum round is set to 30. The results are the average metrics of 100 interaction trajectories.

The results are shown in Fig. 7, where all policies are learned with 200 epochs. After learning in each epoch, we will evaluate all methods with 100 episodes (i.e., interaction trajectories) in the two interactive environments. The first row shows the cumulative

Table 2: Average Results of all methods in the two environments. (Bold: Best; underline: runner-up).

Methods	KuaiRec				KuaiRand			
	R _{tra}	R _{each}	Length	MCD	R _{tra}	Reach	Length	MCD
UCB	3.606 ± 0.609	0.853 ± 0.114	4.219 ± 0.389	0.811 ± 0.058	1.651 ± 0.152	0.372 ± 0.028	4.431 ± 0.212	0.789 ± 0.024
ϵ -greedy	3.515 ± 0.731	0.828 ± 0.129	4.219 ± 0.405	0.823 ± 0.048	1.711 ± 0.126	0.351 ± 0.025	4.880 ± 0.270	0.773 ± 0.024
SQN	4.673 ± 1.215	0.913 ± 0.055	5.111 ± 1.288	0.686 ± 0.093	0.912 ± 0.929	0.182 ± 0.058	4.601 ± 3.712	0.621 ± 0.187
CRR	4.163 ± 0.253	<u>0.895 ± 0.037</u>	4.654 ± 0.215	0.865 ± 0.017	1.481 ± 0.124	0.226 ± 0.015	6.561 ± 0.352	0.733 ± 0.019
CQL	2.506 ± 1.767	0.684 ± 0.228	3.224 ± 1.365	0.386 ± 0.385	2.032 ± 0.107	0.226 ± 0.012	9.000 ± 0.000	0.778 ± 0.000
BCQ	2.123 ± 0.081	0.708 ± 0.027	3.000 ± 0.000	0.667 ± 0.000	0.852 ± 0.052	0.425 ± 0.016	2.005 ± 0.071	0.998 ± 0.024
MBPO	12.043 ± 1.312	0.770 ± 0.029	15.646 ± 1.637	<u>0.362 ± 0.047</u>	10.933 ± 0.946	<u>0.431 ± 0.021</u>	25.345 ± 1.819	0.306 ± 0.040
IPS	<u>12.833 ± 1.353</u>	0.767 ± 0.023	<u>16.727 ± 1.683</u>	0.215 ± 0.064	3.629 ± 0.676	0.216 ± 0.014	16.821 ± 3.182	0.201 ± 0.116
MOPO	11.427 ± 1.750	0.892 ± 0.051	12.809 ± 1.850	0.479 ± 0.062	<u>10.934 ± 0.963</u>	0.437 ± 0.019	25.002 ± 1.891	0.343 ± 0.029
Ours	20.494 ± 2.671	0.767 ± 0.026	26.712 ± 3.419	0.379 ± 0.015	11.850 ± 1.036	0.428 ± 0.022	27.609 ± 2.121	<u>0.296 ± 0.036</u>

**Figure 7: Results of all methods in two environments.**

reward, which directly reflects the long-term satisfaction in our interactive recommendation setting. The second row and third rows dissect the cumulative reward into two parts: the length of the interaction trajectory and the single-round reward, respectively. For a better comparison, we average the results in 200 epochs and

show them in Table 2. Besides the three metrics, we also report the majority category domination (MCD) in Table 2.

From the results, we observe that the four model-based RL methods (MBPO, IPS, MOPO, and DORL) significantly outperform the four model-free RL methods (SQN, CRR, CQL, and BCQ) with respect to trajectory length and cumulative reward. This is because model-based RL is much more sample efficient than model-free RL. In recommendation, the training data is highly sparse. Model-free RL learns directly from the recommendation logs that we have split into different sequences according to the exit rule described above. However, it is extremely difficult to capture the exit mechanism from the sparse logs. By contrast, the model-based RL can leverage the user model to construct as many interaction sequences as possible during training, which guarantees that the policy can distill useful knowledge from the limited offline samples. That is why we embrace model-based RL in recommendation.

For model-based RL methods, MOPO shows an obvious improvement compared to the vanilla method MBPO in terms of single-round reward. It is because MBPO does not consider the OOD actions in the offline data that will incur extrapolation errors in the policy. MOPO introduces the uncertainty penalizer to make the policy pay more attention to the samples of high confidence, which in turn makes the policy capture users' interest more precisely. However, MOPO sacrifices many unpopular items because they appear less frequently and are considered uncertain samples. Therefore, the average length decreases, which in turn reduces the cumulative reward. Our method DORL overcomes this problem. From Fig. 7, we observe that DORL attains the maximal average cumulative reward after several epochs in both KuaiRec and KuaiRand due to that it reaches the largest interaction length. Compared to MOPO and MBPO, DORL sacrifices a little bit of the single-round reward due to its counterfactual exploration philosophy, meanwhile, this greatly improves the diversity and enlarges the length of interactions due. Therefore, it achieves the goal of maximizing users' long-term experiences. After enhancing the vanilla MBPO with the IPS technique, the learned user model gives adjustments to the distribution of training data by re-weighting all items. IPS obtains

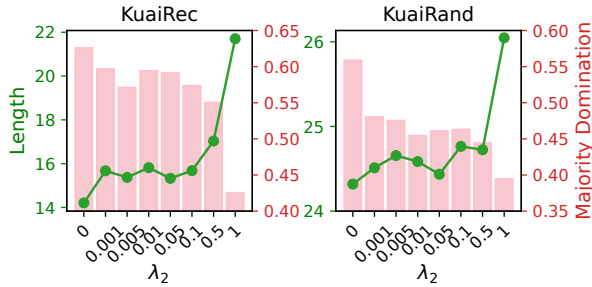


Figure 8: Effect of penalizing entropy.

a satisfactory performance in KuaiRec but receives abysmal performances in KuaiRand. This is due to its well-known high variance issue can incur estimation errors. Compared to IPS’s hard debiased mechanism, DORL’s soft debiased method is more suitable for model-based RL in recommendation.

For model-free RL methods, As discussed above, four model-free methods fail in two datasets due to limited offline samples. Though they can capture users’ interest by returning a high single-round reward (e.g., SQN and CRR in KuaiRec, and BCQ in KuaiRand), they cannot maintain a long interaction trajectory. For example, BCQ updates its policy only on those samples with high confidence, which results in a severe Matthew effect in the recommendation results (reflected by high MCD and short length). SQN’s performance oscillates with the largest magnitude since its network is updated by two heads. The RL head serves as a regularizer to the self-supervised head. When the objectives of the two heads conflict with each other, the performance becomes unstable. Therefore, these methods are not suitable for recommendation where offline data are sparse.

As for the naive bandit methods, UCB and ϵ -greedy, they are designed to explore and exploit the optimal actions for the independent and identically distributed (IID) data. They do not even possess the capability to optimize long-term rewards at all. Therefore, they are inclined to recommend the same items when the model finishes exploring offline data, which leads to high MCD and short interactions. These naive policies are not suitable for pursuing long-term user experiences in recommendation.

6.3 Results on alleviating Matthew effect (RQ2)

We have shown in Fig. 3 that penalizing uncertainty will result in the Matthew effect in recommendation. More specifically, increasing λ_1 can make the recommended items to be the most dominant ones in the training set, which results in a high MCD value. Here, we show how the introduced “counterfactual” exploration mechanism helps alleviate this effect. We conduct the experiments for different combinations of (λ_1, λ_2) as described above, then we average the results along the λ_1 to show the influence of λ_2 alone.

The results are shown in Fig. 8. Obviously, increasing λ_2 can lengthen the interaction process and reduce majority category domination. I.e., When we penalize the entropy of behavior policy hard, (1) the recommender does not repeat the items with the same categories; (2) the recommended results will be diverse instead of

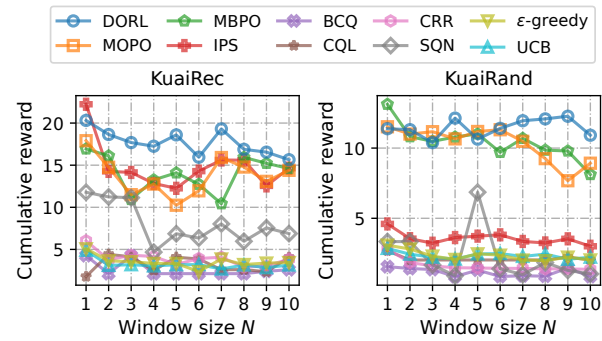


Figure 9: Results under different ending conditions

focusing on dominated items. The results show the effectiveness of penalizing entropy in DORL in alleviating the Matthew effect.

6.4 Results with different environments (RQ3)

To validate that DORL can work robustly in different environment settings, we vary the window size N in the exit mechanism and fix $M = 10$ during the evaluation. The results are shown in Fig. 9. We only visualize the most important metric: the cumulative reward. When N is small ($N = 1$), other model-based methods can surpass our DORL. When N gets larger ($N > 3$), users’ tolerance for similar content (i.e., items with the same category) becomes lower, and the interaction process comes to be easier to terminate. Under such a circumstance, DORL outperforms all other policies, which demonstrates the robustness of DORL in different environments.

7 CONCLUSION

We point out that conservatism in offline RL can incur the Matthew effect in recommendation. We conduct studies to show that the Matthew effect hurts users’ long-term experiences in both the music and video datasets. Through theoretical analysis of the model-based RL framework, we show that the reason for amplifying the Matthew effect is the philosophy of suppressing uncertain samples. It inspires us to add a penalty term to make the policy emphasize the data induced by the behavior policies with high entropy. This will reintroduce the exploration mechanism that conservatism has suppressed, which alleviates the Matthew effect.

In the future, when fitting user interests is not a bottleneck anymore, researchers could consider higher-level goals, such as pursuing users’ long-term satisfaction [57] or optimizing social utility [17]. With the increase in high-quality offline data, we believe that offline RL can be better adapted to recommender systems to achieve these goals. During this process, many interesting yet challenging issues (such as the Matthew effect in this work) will be raised. After addressing these issues, we can create more intelligent recommender systems that benefit society.

ACKNOWLEDGEMENTS

This work is supported by the National Key Research and Development Program of China (2021YFF0901603), the National Natural Science Foundation of China (61972372, U19A2079, 62121002), and the CCCD Key Lab of Ministry of Culture and Tourism.

REFERENCES

- [1] M Mehdi Afsar, Trafford Crump, and Behrouz Far. 2022. Reinforcement Learning based Recommender Systems: A Survey. *Comput. Surveys* 55, 7 (2022), 1–38.
- [2] Rishabh Agarwal, Dale Schuurmans, and Mohammad Norouzi. 2020. An optimistic perspective on offline reinforcement learning. In *International Conference on Machine Learning (ICML '20)*. PMLR, 104–114.
- [3] Ashton Anderson, Lucas Maystre, Ian Anderson, Rishabh Mehrotra, and Mounia Lalmas. 2020. Algorithmic Effects on the Diversity of Consumption on Spotify. In *Proceedings of The Web Conference 2020 (Taipei, Taiwan) (WWW '20)*. 2155–2165.
- [4] Qingpeng Cai, Shuchang Liu, Xueliang Wang, Tianyou Zuo, Wentao Xie, Bin Yang, Dong Zheng, Peng Jiang, and Kun Gai. 2023. Reinforcing User Retention in a Billion Scale Short Video Recommender System. *arXiv preprint arXiv:2302.01724* (2023).
- [5] Qingpeng Cai, Zhenghai Xue, Chi Zhang, Wanqi Xue, Shuchang Liu, Ruohan Zhan, Xueliang Wang, Tianyou Zuo, Wentao Xie, Dong Zheng, et al. 2023. Two-Stage Constrained Actor-Critic for Short Video Recommendation. *arXiv preprint arXiv:2302.01680* (2023).
- [6] Minmin Chen, Alex Beutel, Paul Covington, Sagar Jain, Francois Belletti, and Ed H. Chi. 2019. Top-K Off-Policy Correction for a REINFORCE Recommender System. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining (Melbourne VIC, Australia) (WSDM '19)*. 456–464.
- [7] Minmin Chen, Can Xu, Vince Gatto, Devanshu Jain, Aviral Kumar, and Ed Chi. 2022. Off-Policy Actor-Critic for Recommender Systems. In *Proceedings of the 16th ACM Conference on Recommender Systems (Seattle, WA, USA) (RecSys '22)*. 338–349.
- [8] Romain Deffayet, Thibaut Thonet, Jean-Michel Renders, and Maarten de Rijke. 2023. Offline Evaluation for Reinforcement Learning-based Recommendation: A Critical Issue and Some Alternatives. *arXiv preprint arXiv:2301.00993* (2023).
- [9] Frederik Ebert, Chelsea Finn, Sudeep Dasari, Annie Xie, Alex Lee, and Sergey Levine. 2018. Visual foresight: Model-based deep reinforcement learning for vision-based robotic control. *arXiv preprint arXiv:1812.00568* (2018).
- [10] Scott Fujimoto, Edoardo Conti, Mohammad Ghavamzadeh, and Joelle Pineau. 2019. Benchmarking batch deep reinforcement learning algorithms. *arXiv preprint arXiv:1910.01708* (2019).
- [11] Scott Fujimoto, David Meger, and Doina Precup. 2019. Off-Policy Deep Reinforcement Learning without Exploration. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). 2052–2062.
- [12] Chongming Gao, Wenqiang Lei, Jiawei Chen, Shiqi Wang, Xiangnan He, Shijun Li, Biao Li, Yuan Zhang, and Peng Jiang. 2022. CIRS: Bursting Filter Bubbles by Counterfactual Interactive Recommender System. *arXiv preprint arXiv:2204.01266* (2022).
- [13] Chongming Gao, Wenqiang Lei, Xiangnan He, Maarten de Rijke, and Tat-Seng Chua. 2021. Advances and Challenges in Conversational Recommender Systems: A Survey. *AI Open* 2 (2021), 100–126.
- [14] Chongming Gao, Shijun Li, Wenqiang Lei, Jiawei Chen, Biao Li, Peng Jiang, Xiangnan He, Jiaxin Mao, and Tat-Seng Chua. 2022. KuaiRec: A Fully-Observed Dataset and Insights for Evaluating Recommender Systems. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management (Atlanta, GA, USA) (CIKM '22)*. 540–550. <https://doi.org/10.1145/3511808.3557220>
- [15] Chongming Gao, Shijun Li, Yuan Zhang, Jiawei Chen, Biao Li, Wenqiang Lei, Peng Jiang, and Xiangnan He. 2022. KuaiRand: An Unbiased Sequential Recommendation Dataset with Randomly Exposed Videos. In *Proceedings of the 31st ACM International Conference on Information and Knowledge Management (Atlanta, GA, USA) (CIKM '22)*. 3953–3957. <https://doi.org/10.1145/3511808.3557624>
- [16] Zhaolin Gao, Tianshu Shen, Zheda Mai, Mohamed Reda Bouadjenek, Isaac Waller, Ashton Anderson, Ron Bodkin, and Scott Sanner. 2022. Mitigating the Filter Bubble While Maintaining Relevance: Targeted Diversification with VAE-Based Recommender Systems. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval (Madrid, Spain) (SIGIR '22)*. 2524–2531.
- [17] Yingqiang Ge, Xiaoting Zhao, Lucia Yu, Saurabh Paul, Diane Hu, Chu-Cheng Hsieh, and Yongfeng Zhang. 2022. Toward Pareto Efficient Fairness-Utility Trade-off in Recommendation through Reinforcement Learning. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining (Virtual Event, AZ, USA) (WSDM '22)*. 316–324.
- [18] Alexandre Gilotte, Clément Calauzènes, Thomas Nedelec, Alexandre Abraham, and Simon Dollé. 2018. Offline A/B Testing for Recommender Systems. In *WSDM '18*. 198–206.
- [19] Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. 2017. DeepFM: A Factorization-Machine Based Neural Network for CTR Prediction. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence (Melbourne, Australia) (IJCAI'17)*. 1725–1731.
- [20] Christian Hansen, Rishabh Mehrotra, Casper Hansen, Brian Brost, Lucas Maystre, and Mounia Lalmas. 2021. Shifting Consumption towards Diverse Content on Music Streaming Platforms. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining (Virtual Event, Israel) (WSDM '21)*. 238–246.
- [21] Jonathan Ho and Stefano Ermon. 2016. Generative Adversarial Imitation Learning. In *Advances in Neural Information Processing Systems (NeurIPS '16, Vol. 29)*, D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett (Eds.).
- [22] Jin Huang, Harrie Oosterhuis, Bunyamin Cetinkaya, Thijs Rood, and Maarten de Rijke. 2022. State Encoders in Reinforcement Learning for Recommendation: A Reproducibility Study. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval (Madrid, Spain) (SIGIR '22)*. 2738–2748.
- [23] Jin Huang, Harrie Oosterhuis, Maarten de Rijke, and Herke van Hoof. 2020. Keeping Dataset Biases out of the Simulation: A Debiased Simulator for Reinforcement Learning Based Recommender Systems. In *RecSys '20*. 190–199.
- [24] Olivier Jeunen and Bart Goethals. 2021. Pessimistic Reward Models for Off-Policy Learning in Recommendation. In *Proceedings of the 15th ACM Conference on Recommender Systems (Amsterdam, Netherlands) (RecSys '21)*. 63–74.
- [25] Wang-Cheng Kang and Julian McAuley. 2018. Self-attentive sequential recommendation. In *International Conference on Data Mining (ICDM '18)*. IEEE, 197–206.
- [26] Alex Kendall and Yarin Gal. 2017. What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision?. In *Proceedings of the 31st International Conference on Neural Information Processing Systems (Long Beach, California, USA) (NeurIPS '17)*. 5580–5590.
- [27] Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. 2020. MOREL: Model-Based Offline Reinforcement Learning. In *Advances in Neural Information Processing Systems (NeurIPS '20, Vol. 33)*, H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (Eds.). 21810–21823.
- [28] Vijay Konda and John Tsitsiklis. 1999. Actor-critic algorithms. *Advances in neural information processing systems* 12 (1999).
- [29] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. 2022. Offline Reinforcement Learning with Implicit Q-Learning. In *International Conference on Learning Representations (ICLR '22)*.
- [30] Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. 2019. Stabilizing off-policy q-learning via bootstrapping error reduction. *Advances in Neural Information Processing Systems* 32 (2019).
- [31] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. 2020. Conservative Q-Learning for Offline Reinforcement Learning. In *Advances in Neural Information Processing Systems (NeurIPS '20, Vol. 33)*, H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (Eds.). 1179–1191.
- [32] Wenqiang Lei, Chongming Gao, and Maarten de Rijke. 2021. RecSys 2021 Tutorial on Conversational Recommendation: Formulation, Methods, and Evaluation (*RecSys '21*). Association for Computing Machinery, New York, NY, USA, 842–844. <https://doi.org/10.1145/3460231.3473325>
- [33] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. 2020. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643* (2020).
- [34] Yile Liang, Tiejun Qian, Qing Li, and Hongzhi Yin. 2021. Enhancing Domain-Level and User-Level Adaptivity in Diversified Recommendation. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval (Virtual Event, Canada) (SIGIR '21)*. 747–756.
- [35] Shuchang Liu, Qingpeng Cai, Bowen Sun, Yuhao Wang, Ji Jiang, Dong Zheng, Kun Gai, Peng Jiang, Xiangyu Zhao, and Yongfeng Zhang. 2023. Exploration and Regularization of the Latent Action Space in Recommendation. *arXiv preprint arXiv:2302.03431* (2023).
- [36] Ying Chieh Liu and Min Qi Huang. 2021. Examining the Matthew Effect on YouTube Recommendation System. In *International Conference on Technologies and Applications of Artificial Intelligence (TAAI '21)*. 146–148. <https://doi.org/10.1109/TAAI54685.2021.00035>
- [37] Yiping Luo, Huazhe Xu, Yuanzhi Li, Yuandong Tian, Trevor Darrell, and Tengyu Ma. 2018. Algorithmic framework for model-based deep reinforcement learning with theoretical guarantees. *arXiv preprint arXiv:1807.03858* (2018).
- [38] Markus Schedl. 2016. The LFM-1b Dataset for Music Retrieval and Recommendation. In *Proceedings of the 2016 ACM on International Conference on Multimedia Retrieval (New York, New York, USA) (ICMR '16)*. 103–110.
- [39] Adith Swaminathan and Thorsten Joachims. 2015. Counterfactual Risk Minimization: Learning from Logged Bandit Feedback. In *International Conference on Machine Learning (ICML '15)*. PMLR, 814–823.
- [40] Jiayi Tang and Ke Wang. 2018. Personalized Top-N Sequential Recommendation via Convolutional Sequence Embedding. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining (Marina Del Rey, CA, USA) (WSDM '18)*. 565–573.
- [41] Matus Tomlein, Branislav Pecher, Jakub Simko, Ivan Srba, Robert Moro, Elena Stefanova, Michal Kompan, Andrea Hrcokova, Juraj Podrouzek, and Maria Bielikova. 2021. An Audit of Misinformation Filter Bubbles on YouTube: Bubble Bursting and Recent Behavior Changes. In *RecSys '21*. 1–11.
- [42] Hao Wang, Zonghu Wang, and Weishi Zhang. 2018. Quantitative analysis of Matthew effect and sparsity problem of recommender systems. In *2018 IEEE 3rd International Conference on Cloud Computing and Big Data Analysis (ICCCBDA)*. 78–82. <https://doi.org/10.1109/ICCCBDA.2018.8386490>

- [43] Ruosong Wang, Dean P Foster, and Sham M Kakade. 2020. What are the statistical limits of offline RL with linear function approximation? *arXiv preprint arXiv:2010.11895* (2020).
- [44] Shiqi Wang, Chongming Gao, Min Gao, Junliang Yu, Zongwei Wang, and Hongzhi Yin. 2022. Who Are the Best Adopters? User Selection Model for Free Trial Item Promotion. *IEEE Transactions on Big Data* (2022).
- [45] Yuyan Wang, Mohit Sharma, Can Xu, Sriraj Badam, Qian Sun, Lee Richardson, Lisa Chung, Ed H. Chi, and Minmin Chen. 2022. Surrogate for Long-Term User Experience in Recommender Systems. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining* (Washington DC, USA) (KDD '22). 4100–4109.
- [46] Ziyu Wang, Alexander Novikov, Konrad Zolna, Josh S Merel, Jost Tobias Springenberg, Scott E Reed, Bobak Shahriari, Noah Siegel, Caglar Gulcehre, Nicolas Heess, et al. 2020. Critic Regularized Regression. *Advances in Neural Information Processing Systems* 33 (2020), 7768–7778.
- [47] Christopher JCH Watkins and Peter Dayan. 1992. Q-learning. *Machine learning* 8, 3 (1992), 279–292.
- [48] Junda Wu, Zhihui Xie, Tong Yu, Handong Zhao, Ruiyi Zhang, and Shuai Li. 2022. Dynamics-Aware Adaptation for Reinforcement Learning Based Cross-Domain Interactive Recommendation (SIGIR '22). 290–300.
- [49] Teng Xiao and Donglin Wang. 2021. A general offline reinforcement learning framework for interactive recommendation. In *Proceedings of the AAAI Conference on Artificial Intelligence* (AAAI '21, Vol. 35). 4512–4520.
- [50] Xin Xin, Alexandros Karatzoglou, Ioannis Arapakis, and Joemon M. Jose. 2020. Self-Supervised Reinforcement Learning for Recommender Systems. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval* (SIGIR '20). 931–940.
- [51] Xin Xin, Tiago Pimentel, Alexandros Karatzoglou, Pengjie Ren, Konstantina Christakopoulou, and Zhaochun Ren. 2022. Rethinking Reinforcement Learning for Recommendation: A Prompt Perspective. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Madrid, Spain) (SIGIR '22). 1347–1357.
- [52] Shuyuan Xu, Juntao Tan, Zuohui Fu, Jianchao Ji, Shelby Heinecke, and Yongfeng Zhang. 2022. Dynamic Causal Collaborative Filtering. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management* (Atlanta, GA, USA) (CIKM '22). 2301–2310.
- [53] Wanqi Xue, Qingpeng Cai, Ruohan Zhan, Dong Zheng, Peng Jiang, and Bo An. 2023. ResAct: Reinforcing Long-term Engagement in Sequential Recommendation with Residual Actor (ICLR '23).
- [54] Tianhe Yu, Aviral Kumar, Rafael Rafailov, Aravind Rajeswaran, Sergey Levine, and Chelsea Finn. 2021. Combo: Conservative offline model-based policy optimization. *Advances in Neural Information Processing Systems* 34 (2021), 28954–28967.
- [55] Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea Finn, and Tengyu Ma. 2020. MOPO: Model-based Offline Policy Optimization. In *Advances in Neural Information Processing Systems (NeurIPS '20, Vol. 33)*, H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (Eds.). 14129–14142.
- [56] Fajie Yuan, Alexandros Karatzoglou, Ioannis Arapakis, Joemon M. Jose, and Xiangnan He. 2019. A Simple Convolutional Generative Network for Next Item Recommendation. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining* (Melbourne VIC, Australia) (WSDM '19). 582–590.
- [57] Qihua Zhang, Junning Liu, Yuzhuo Dai, Yiyang Qi, Yifan Yuan, Kunlun Zheng, Fan Huang, and Xianfeng Tan. 2022. Multi-Task Fusion via Reinforcement Learning for Long-Term User Satisfaction in Recommender Systems. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining* (Washington DC, USA) (KDD '22). 4510–4520.
- [58] Yuan Zhang, Xue Dong, Weijie Ding, Biao Li, Peng Jiang, and Kun Gai. 2023. Dive and Conquer: Towards Better Embedding-based Retrieval for Recommender Systems From a Multi-task Perspective. *arXiv preprint arXiv:2302.02657* (2023).
- [59] Xiangyu Zhao, Long Xia, Lixin Zou, Hui Liu, Dawei Yin, and Jiliang Tang. 2021. UserSim: User Simulation via Supervised Generative Adversarial Network. In *WWW '21*. 3582–3589.
- [60] Yu Zheng, Chen Gao, Liang Chen, Depeng Jin, and Yong Li. 2021. DGCN: Diversified Recommendation with Graph Convolutional Networks. In *Proceedings of the Web Conference 2021* (Ljubljana, Slovenia) (WWW '21). 401–412.
- [61] Yu Zheng, Chen Gao, Xiang Li, Xiangnan He, Yong Li, and Depeng Jin. 2021. Disentangling User Interest and Conformity for Recommendation with Causal Embedding. In *Proceedings of the Web Conference 2021* (Ljubljana, Slovenia) (WWW '21). 2980–2991.
- [62] Lixin Zou, Long Xia, Zhuoye Ding, Jiaying Song, Weidong Liu, and Dawei Yin. 2019. Reinforcement Learning to Optimize Long-Term User Engagement in Recommender Systems. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (Anchorage, AK, USA) (KDD '19). 2810–2818.
- [63] Lixin Zou, Long Xia, Pan Du, Zhuo Zhang, Ting Bai, Weidong Liu, Jian-Yun Nie, and Dawei Yin. 2020. Pseudo Dyna-Q: A Reinforcement Learning Framework for Interactive Recommendation. In *WSDM '20*. 816–824.