



# Practical Counterfactual Policy Learning for Top- $K$ Recommendations

Yaxu Liu<sup>\*†</sup>  
National Taiwan University  
d08944012@ntu.edu.tw

Jui-Nan Yen<sup>†</sup>  
National Taiwan University  
juinanyen@gmail.com

Bowen Yuan<sup>†‡</sup>  
Amazon  
bwyuan@amazon.com

Rundong Shi  
Meituan  
shirundong@meituan.com

Peng Yan  
Meituan  
yanpeng04@meituan.com

Chih-Jen Lin  
National Taiwan University  
cjlin@csie.ntu.edu.tw

## ABSTRACT

For building recommender systems, a critical task is to learn a policy with collected feedback (e.g., ratings, clicks) to decide which items to be recommended to users. However, it has been shown that the selection bias in the collected feedback leads to biased learning and thus a sub-optimal policy. To deal with this issue, counterfactual learning has received much attention, where existing approaches can be categorized as either value learning or policy learning approaches. This work studies policy learning approaches for top- $K$  recommendations with a large item space and points out several difficulties related to importance weight explosion, observation insufficiency, and training efficiency. A practical framework for policy learning is then proposed to overcome these difficulties. Our experiments confirm the effectiveness and efficiency of the proposed framework.

## CCS CONCEPTS

• Computing methodologies → Machine learning.

## KEYWORDS

Policy learning, Selection bias, Counterfactual learning, Recommender systems

## ACM Reference Format:

Yaxu Liu, Jui-Nan Yen, Bowen Yuan, Rundong Shi, Peng Yan, and Chih-Jen Lin. 2022. Practical Counterfactual Policy Learning for Top- $K$  Recommendations. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '22)*, August 14–18, 2022, Washington, DC, USA. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3534678.3539295>

<sup>\*</sup>Work done at Meituan as an intern.

<sup>†</sup>Contribute equally for this work.

<sup>‡</sup>Work done at National Taiwan University.

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 ACM 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](mailto:permissions@acm.org).

KDD '22, August 14–18, 2022, Washington, DC, USA

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9385-0/22/08...\$15.00

<https://doi.org/10.1145/3534678.3539295>

## 1 INTRODUCTION

In many applications, users are overwhelmed by a large number of items. To satisfy their needs more efficiently, recommender systems have been constructed to pre-select items according to their preference and some contextual information. Learning a policy to decide which items to be recommended is the core problem of building a recommender system. A widely used approach is to consider rewards (like ratings, clicks, and dwell time) for historical recommendations as labels and solve classification/regression problems to learn a model (e.g., matrix factorization [10], factorization machines [17]) as an estimator of rewards. Then the system decides future recommendations according to the estimated rewards from the learned model.

However, recent works argue that such a learning scheme may lead to biased recommendations. The reason is that the historical recommendations are decided by a previously deployed policy, which is generally referred to as the behavior policy. The mismatch between data distributions of the behavior policy and the new policy introduces the selection bias. To remove the bias, the key idea is to consider rewards for both recommended and non-recommended items. However, the realization of this idea is non-trivial. Unlike rewards for recommended items, rewards for non-recommended items, often known as counterfactual data, are generally not available. Learning an unbiased policy without counterfactual data, which is also referred to as counterfactual learning, is an emerging topic in recommender systems.

Generally speaking, existing approaches of counterfactual learning fall into two groups, value learning approaches [25, 27] and policy learning approaches [9, 19, 20]. A review of them will be given in Section 2.2. Between the two approaches, besides few works (e.g., [8]) that propose combining these approaches for recommender systems, a systematic comparison between these two paradigms is lacking. Our original purpose is to fill this gap. However, when we attempt to implement policy learning approaches for a large-scale top- $K$  recommender system, where  $K$  items are recommended for each context, we meet the following challenges.

- The propensity of a top- $K$  recommendation can be very small, resulting in the explosion of importance weights and poor performances of the learned policy.
- For recommender systems with a large item space, it is almost impossible to observe every item for every context. Thus, for some rarer contexts, the empirical reward of each item can be far away from the underlying true reward distribution. In particular, the empirical reward can be wrongly concentrated on one single

**Table 1: Main notation.**

$m, n, K$	numbers of contexts, items, and positions
$\beta, \pi$	behavior (deployed) policy and new policy
$\mathbb{A}$	set of $n$ items
$\mathbf{u}, \mathbf{c}$	random variable for the feature vector of a context and a reward vector of $n$ items
$\mathbf{V}$	feature matrix of $n$ items
$\mathbf{A}, \mathbf{A}_k$	random variable for a permutation of $K$ items and the item placed at position $k$
$\mathbf{C}$	random variable for a reward vector of items in $\mathbf{A}$

item if only a few items are observed and one of them is largely amplified by its low propensity weight. This can make the learned policy to be very sharp, a situation which has been pointed out to be harmful to robustness and generalization.

- Training efficiency is a serious issue as optimization problems of policy learning require computations over the entire item space, the size of which can reach over millions in a real-world system.

This work aims to overcome these difficulties with the following contributions.

- We propose a regularized per-item estimator to quell the first issue mentioned above.
- We propose a formulation to handle the probability mass and thus avoid the sharpness of the learned policy.
- We introduce a novel and efficient training method for policy learning. It can be effectively deployed for large-scale top- $K$  recommender systems.

The paper is organized as follows. The preliminary is in Section 2, which includes the problem setting and a review of existing approaches of counterfactual learning. We position the above-mentioned difficulties and propose the corresponding solutions in Section 3. Related works are given in Section 4. In Section 5, a series of experiments confirm our contributions. Finally, Section 6 concludes this work. A list of notations used in this work is in Table 1. Supplementary materials and programs/data used for our experiments are available at [https://www.csie.ntu.edu.tw/~cjlin/papers/counterfactual\\_topk/](https://www.csie.ntu.edu.tw/~cjlin/papers/counterfactual_topk/).

## 2 PRELIMINARY

In this section, after describing the problem settings, we briefly review existing approaches for counterfactual learning.

### 2.1 Problem Setting

In this work, we focus on top- $K$  recommender systems with  $n$  items (e.g., movies, products). Let  $\mathbb{A} = \{1, \dots, n\}$  and  $\mathbf{V} \in \mathbb{R}^{n \times D_v}$  be the associated side-feature matrix. The interactions between users and a top- $K$  recommender system can be described as the following procedures.

- When a user visits a webpage with  $K$  recommendation positions, a feature vector  $\mathbf{u} \in \mathbb{R}^{D_u}$  including the contextual information (e.g., information of a user or a webpage) is sent to the system, where  $\mathbf{u}$  follows an underlying distribution  $\Pr(\mathbf{u})$ .
- Rewards (e.g., clicks, views)  $\mathbf{c}$  for all items follow another underlying distribution  $\Pr(\mathbf{c} \mid \mathbf{u}; \mathbf{V})$ , where for each item  $j$ , the

corresponding  $c_j$  depends only on its feature vector  $\mathbf{v}_j$  in the matrix  $\mathbf{V}$  and the feature vector  $\mathbf{u}$  above. Note that  $\mathbf{v}_j$  is the  $j$ th row of  $\mathbf{V}$ .

- Given the received  $\mathbf{u}$ , a recommendation policy  $\beta$  deployed in the system recommends  $K$  distinct items according to the probability  $\beta(\mathbf{A} \mid \mathbf{u}; \mathbf{V})$ . These items form a  $K$ -element permutation of  $\mathbb{A}$  as following,

$$\mathbf{A} \in G(\mathbb{A}, K) \equiv \{(\mathbf{A}_1, \dots, \mathbf{A}_K) \mid \mathbf{A}_1, \dots, \mathbf{A}_K \in \mathbb{A}; \mathbf{A}_i \neq \mathbf{A}_j, \forall i \neq j\}. \quad (1)$$

- Once items in  $\mathbf{A}$  are recommended to users, their corresponding rewards  $\mathbf{C} = (c_{\mathbf{A}_1}, \dots, c_{\mathbf{A}_K})$  are revealed to the system, while rewards of other non-recommended items are still unknown.

With the repetition of the above procedures, the system can collect a set  $\mathbb{S}$ , which includes vast logged events  $(\mathbf{u}, \mathbf{A}, \mathbf{C})$ . One of the most crucial tasks of building a recommender system is to learn a new policy  $\pi$  so that the probability distribution  $\pi(\mathbf{A} \mid \mathbf{u}; \mathbf{V}, \theta)$  parameterized by  $\theta$  maximizes the expected cumulated reward

$$V^\pi = \mathbb{E}_{P_\pi} [r(\mathbf{c}, \mathbf{A})], \quad (2)$$

where

$$\mathbb{E}_{P_\pi} [\cdot] = \mathbb{E}_{\Pr(\mathbf{u})} \mathbb{E}_{\pi(\mathbf{A} \mid \mathbf{u}; \mathbf{V}, \theta)} \mathbb{E}_{\Pr(\mathbf{c} \mid \mathbf{u}; \mathbf{V})} [\cdot] \quad (3)$$

and

$$r(\mathbf{c}, \mathbf{A}) = \sum_{k=1}^K c_{\mathbf{A}_k} \quad (4)$$

is the cumulated reward from  $\mathbf{A}$ . To distinguish  $\pi$  from the previous policy  $\beta$ , the latter is also called the behavior policy. A natural setting to obtain a useful  $\pi$  is by using the historical events in  $\mathbb{S}$ . However, the challenge comes from the fact that directly maximizing the expected reward in (2) requires rewards  $\mathbf{c}$  for all items. Unfortunately, among all  $\frac{n!}{(n-K)!}$  permutations as the candidate of  $\mathbf{A}$ , we only observe the one recommended by  $\beta$  for each context  $\mathbf{u}$  in  $\mathbb{S}$ . Opposed to the logged events in  $\mathbb{S}$ , these events of non-recommended permutations are called counterfactual data. Conducting policy learning with partial rewards in logged events is also referred to as counterfactual learning or batch learning from logged bandit feedback [19].

### 2.2 A Review on Existing Approaches of Counterfactual learning

Existing approaches of counterfactual learning can be roughly categorized into two groups, the value learning approach and the policy learning approach. They differ on explicitly or implicitly learning the new policy  $\pi$ . Subsequently, we give a brief review of them.

**2.2.1 Value Learning Approaches.** The main idea of value learning methods is to learn a reward estimation model  $\hat{s}(\cdot)$  parametrized by  $\theta$  to model the reward probability  $\Pr(\mathbf{c} \mid \mathbf{u}; \mathbf{V})$ . To have an unbiased estimate, it has been pointed out [25] that we should minimize the following full-labeled risk

$$L(\theta) = \mathbb{E}_{\Pr(\mathbf{u})} \mathbb{E}_{\Pr(\mathbf{c} \mid \mathbf{u}; \mathbf{V})} \left[ \sum_{j=1}^n \ell(c_j, \hat{s}(\theta; \mathbf{u}, \mathbf{v}_j)) \right], \quad (5)$$

where  $\mathbf{v}_j$  is the feature vector of the  $j$ th item extracted from row  $j$  of  $\mathbf{V}$ , and  $\ell(a, b)$  is a loss function (e.g., logistic loss or squared loss). However, for each event  $(\mathbf{u}, \mathbf{A}, \mathbf{C})$ , we only have rewards contained in  $\mathbf{C}$ , which correspond to those selected items included in  $\mathbf{A}$ . It has been shown [25, 27] that because the behavior policy  $\beta(\mathbf{A} \mid \mathbf{u}; \mathbf{V})$

non-uniformly selects items with the highest reward, considering a modified form of (5) with only items included in  $\mathbf{A}$  results in a heavily biased  $\hat{s}(\cdot)$ .

For addressing the bias, a commonly used way is to apply an inverse-propensity-score (IPS) method [5]. The main idea is to re-weight the loss term of each recommended item by its inverse propensity score such that the following unbiased risk is minimized

$$L_{\text{IPS}}(\theta) = \mathbb{E}_{\text{Pr}(\mathbf{u})} \mathbb{E}_{\beta(\mathbf{A}|\mathbf{u};\mathbf{V})} \mathbb{E}_{\text{Pr}(\mathbf{c}|\mathbf{u};\mathbf{V})} \left[ \sum_{k=1}^K \frac{\ell(c_{A_k}, \hat{s}(\theta; \mathbf{u}, \mathbf{v}_{A_k}))}{z(A_k | \mathbf{u}; \mathbf{V})} \right], \quad (6)$$

where  $z(A_k | \mathbf{u}; \mathbf{V})$  is the propensity score, which is a value proportional to the probability of the item being recommended at the position. For  $z(A_k | \mathbf{u}; \mathbf{V})$ , we can infer it with the behavior policy  $\beta$  as follows

$$z(A_k | \mathbf{u}; \mathbf{V}) = \sum_{\hat{A} \in G(\mathbf{A}, K)} \beta(\hat{A} | \mathbf{u}; \mathbf{V}) \mathbb{1}[A_k = \hat{A}], \quad (7)$$

where  $\mathbb{1}[\cdot]$  is the indicator function. In some case, when the behavior policy  $\beta$  is not reserved, we need to estimate  $z(A_k | \mathbf{u}; \mathbf{V})$  from  $\mathbb{S}$ . Unfortunately, it is known that accurately estimating  $z(A_k | \mathbf{u}; \mathbf{V})$  is not an easy task. Therefore, some recent works (e.g., [25, 27]) extend (6) to minimize the following doubly robust risk [4] for mitigating the issue caused by poor estimates of propensity scores.

$$L_{\text{DR}}(\theta) = \mathbb{E}_{\text{Pr}(\mathbf{u})} \mathbb{E}_{\beta(\mathbf{A}|\mathbf{u};\mathbf{V})} \mathbb{E}_{\text{Pr}(\mathbf{c}|\mathbf{u};\mathbf{V})} \left[ \gamma \sum_{j=1}^n \bar{\ell}(\hat{c}_j, \hat{s}(\theta; \mathbf{u}, \mathbf{v}_j)) + \sum_{k=1}^K \frac{\ell(c_{A_k}, \hat{s}(\theta; \mathbf{u}, \mathbf{v}_{A_k})) - \gamma \ell(\hat{c}_{A_k}, \hat{s}(\theta; \mathbf{u}, \mathbf{v}_{A_k}))}{z(A_k | \mathbf{u}; \mathbf{V})} \right], \quad (8)$$

where  $\hat{c}_j$  is an imputed reward of item  $j$  and  $\bar{\ell}(\cdot)$  is the loss function for the imputation part. Because the imputed rewards of those non-recommended items are less reliable than recommended items, a small user-specified hyper-parameter  $\gamma$  is applied to balance the two parts.

The reward estimation model  $\hat{s}(\cdot)$  is learned from the following minimization problem

$$\min_{\theta} L_{\text{IPS}}(\theta) \text{ (or } L_{\text{DR}}(\theta)) + \lambda \Phi(\theta), \quad (9)$$

where  $\Phi(\theta)$  is a regularizer to avoid overfitting (e.g., l2 regularizer), and  $\lambda$  is a pre-specified regularization coefficient. Then, the optimal  $\pi(\mathbf{A} | \mathbf{u}; \mathbf{V}, \theta)$  can be obtained by selecting the top- $K$  items with the highest reward estimates from  $\hat{s}(\cdot)$ . However, it has been pointed out [25] that such a deterministic policy causes difficulties for optimizing future policies. Thus the  $\epsilon$ -greedy method is always applied to impose stochasticity. That is, we greedily select the item with the highest reward estimates with probability  $1 - \epsilon$ , and uniformly draw a random item with probability  $\epsilon$ .

**2.2.2 Policy Learning Approaches.** Instead of estimating rewards, policy learning approaches [4, 19] directly find an optimal  $\pi(\mathbf{A} | \mathbf{u}; \mathbf{V}, \theta)$  to maximize  $V^\pi$  in (2). Similar to value learning approaches discussed in the previous section, since  $\mathbb{S}$  is generally collected by a non-uniform  $\beta$ , we should apply the inverse propensity weighting method to correct the distribution mismatching between  $\beta$  and  $\pi$ . This leads to the following estimate of  $V^\pi$ ,

$$V_{\text{IPS}}^\pi = \mathbb{E}_{P_\beta} \left[ \mathbf{w}_A r(\mathbf{c}, \mathbf{A}) \right], \quad (10)$$

where

$$\mathbb{E}_{P_\beta}[\cdot] = \mathbb{E}_{\text{Pr}(\mathbf{u})} \mathbb{E}_{\beta(\mathbf{A}|\mathbf{u};\mathbf{V})} \mathbb{E}_{\text{Pr}(\mathbf{c}|\mathbf{u};\mathbf{V})} [\cdot] \quad (11)$$

and

$$\mathbf{w}_A = \frac{\pi(\mathbf{A} | \mathbf{u}; \mathbf{V}, \theta)}{\beta(\mathbf{A} | \mathbf{u}; \mathbf{V})}. \quad (12)$$

It is known that  $V_{\text{IPS}}^\pi$  is an unbiased estimate of  $V^\pi$  under the following assumption. The proof can be found in our supplementary materials.

**ASSUMPTION 1.** For any  $\mathbf{A}$  and  $\mathbf{u}$ ,  $\pi(\mathbf{A} | \mathbf{u}; \mathbf{V}, \theta) \neq 0$  only if  $\beta(\mathbf{A} | \mathbf{u}; \mathbf{V}) \neq 0$ .

It indicates that our new policy  $\pi$  only considers items with non-zero probability in our behavior policy  $\beta$ .

Thus we can solve the following problem to learn  $\pi$ ,

$$\max_{\theta} R(\theta) = V_{\text{IPS}}^\pi(\theta) - \lambda \Phi(\theta), \quad (13)$$

where similar to (5), a regularizer  $\Phi(\theta)$  is considered to avoid overfitting. To solve (13), the following policy gradient method is applied,

$$\theta \leftarrow \theta + \eta \nabla R(\theta), \quad (14)$$

where  $\eta$  is a specified learning rate. In (14),  $\nabla R(\theta)$  is computed as follows,

$$\begin{aligned} \nabla R(\theta) &= \mathbb{E}_{P_\beta} \left[ r(\mathbf{c}, \mathbf{A}) \frac{\nabla \pi(\mathbf{A} | \mathbf{u}; \mathbf{V}, \theta)}{\beta(\mathbf{A} | \mathbf{u}; \mathbf{V})} \right] - \lambda \nabla \Phi(\theta) \\ &= \mathbb{E}_{P_\beta} \left[ r(\mathbf{c}, \mathbf{A}) \frac{\pi(\mathbf{A} | \mathbf{u}; \mathbf{V}, \theta)}{\beta(\mathbf{A} | \mathbf{u}; \mathbf{V})} \nabla \log \pi(\mathbf{A} | \mathbf{u}; \mathbf{V}, \theta) \right] - \lambda \nabla \Phi(\theta), \end{aligned} \quad (15)$$

where the second equality is from the following chain rule,

$$\nabla \log \pi = \frac{1}{\pi} \nabla \pi. \quad (16)$$

### 3 A PRACTICAL FRAMEWORK FOR POLICY LEARNING

Except a recent work [8] that proposes a new model that combines approaches of both value learning and policy learning for recommender systems, few studied these two paradigms together in this field. Thus the original purpose of this work is to detailedly compare them for recommendation problems. However, to implement them for a large-scale top- $K$  recommender system, where both the number of logged events and the number of items are tremendous, we find some extra difficulties in policy learning approaches. In this section, after positioning a difficulty in each of the following subsections, we propose the corresponding solution. Finally, integrating these solutions leads to a practical framework of policy learning for top- $K$  recommendations.

#### 3.1 Regularized Per-item IPS Estimator

**3.1.1 Weight Explosion in IPS Estimators.** In practice, given finite historical events logged by  $\beta$ , we generally derive an unbiased estimate of  $V_{\text{IPS}}^\pi$  via Monte Carlo approximation. Expressly, assuming that the collected set includes  $m$  logged events, we maximize the following estimate of  $V_{\text{IPS}}^\pi$ ,

$$\hat{V}_{\text{IPS}}^\pi(\theta) = \frac{1}{m} \sum_{i=1}^m \mathbf{w}_A^i r(\mathbf{c}_i, \mathbf{A}_i). \quad (17)$$

The quality of  $\hat{V}_{\text{IPS}}^\pi$  estimating  $V^\pi$  deeply relies on the range of  $\mathbf{w}_A$ . To explain this, we make the following assumptions.

ASSUMPTION 2. For any  $\mathbf{A}$  and  $\mathbf{u}$ , the reward  $c_{\mathbf{A}}$  for every  $\mathbf{A} \in \mathcal{A}$  is bounded in  $[0, 1]$ .

ASSUMPTION 3. For any  $\mathbf{A}$  and  $\mathbf{u}$ , if  $\beta(\mathbf{A} | \mathbf{u}; \mathbf{V}) \neq 0$ , then

$$w_{\mathbf{A}} = \frac{\pi(\mathbf{A} | \mathbf{u}; \mathbf{V}, \theta)}{\beta(\mathbf{A} | \mathbf{u}; \mathbf{V})} \leq w_{\max}.$$

Then, we give the following theorem.<sup>1</sup>

THEOREM 1. Given a finite policy class  $\mathcal{H}$  with  $|\mathcal{H}| = N$ , we assume that Assumptions 1-3 are all satisfied. Then, with probability at least  $1 - \delta$ , we have

$$\sup_{\pi \in \mathcal{H}} |\hat{V}_{\text{IPS}}^{\pi} - V^{\pi}| \leq K w_{\max} \sqrt{\frac{2}{m} \log \frac{2N}{\delta}}.$$

From Theorem 1, we conclude that the quality of  $\hat{V}_{\text{IPS}}^{\pi}$  can suffer when the range of  $w_{\mathbf{A}}$  is large. Unfortunately, this situation commonly occurs for large-scale top- $K$  recommender systems. The leading cause of this issue is that the number of possible item permutations  $\mathbf{A}$  is an extremely large number  $n!/(n-K)!$ . As a consequence, for most item permutations, both  $\pi(\mathbf{A} | \mathbf{u}; \mathbf{V}, \theta)$  and  $\beta(\mathbf{A} | \mathbf{u}; \mathbf{V})$  can be very small. During the training procedure, any slight disagreement between  $\pi(\mathbf{A} | \mathbf{u}; \mathbf{V}, \theta)$  and  $\beta(\mathbf{A} | \mathbf{u}; \mathbf{V})$  may cause an explosion of  $w_{\mathbf{A}}$ . This conjecture has been confirmed in a real-world top- $K$  advertising system [11]. When  $K$  grows, by increasing the disagreement between  $\pi(\mathbf{A} | \mathbf{u}; \mathbf{V}, \theta)$  and  $\beta(\mathbf{A} | \mathbf{u}; \mathbf{V})$ , they demonstrate a vast gap between  $\hat{V}_{\text{IPS}}^{\pi}$  and  $V^{\pi}$ .

**3.1.2 Weight Factorization and Pruning.** The explosion of  $w_{\mathbf{A}}$  is caused by the massive number of possible permutations of  $K$  items. Thus, for avoiding this issue, an intuitive idea is to assume that  $\pi(\mathbf{A} | \mathbf{u}; \mathbf{V}, \theta)$  and  $\beta(\mathbf{A} | \mathbf{u}; \mathbf{V})$  can be respectively factorized as

$$\pi(\mathbf{A} | \mathbf{u}; \mathbf{V}, \theta) = \prod_{k=1}^K \pi(A_k | \mathbf{u}, \mathbf{A}_{1:k-1}; \mathbf{V}, \theta), \quad (18)$$

and

$$\beta(\mathbf{A} | \mathbf{u}; \mathbf{V}) = \prod_{k=1}^K \beta(A_k | \mathbf{u}, \mathbf{A}_{1:k-1}; \mathbf{V}). \quad (19)$$

This decomposes the decision of  $\mathbf{A}$  into sequential decisions of  $A_k$  in  $K$  sub-contexts. Now each sub-context consists of not only  $\mathbf{u}$  but also items  $\mathbf{A}_{1:k-1}$  that have been recommended beforehand. From (4), we can rewrite (10) as

$$V_{\text{IPS}}^{\pi} = \mathbb{E}_{P_{\beta}} \left[ w_{\mathbf{A}} \left( \sum_{k=1}^K c_{A_k} \right) \right] = \mathbb{E}_{P_{\beta}} \left[ \sum_{k=1}^K w_{\mathbf{A}} c_{A_k} \right]. \quad (20)$$

From (18) and (19), we can rewrite (12) as

$$w_{\mathbf{A}} = \prod_{k=1}^K w_k, \quad (21)$$

where

$$w_k = \frac{\pi(A_k | \hat{\mathbf{u}}_k; \mathbf{V}, \theta)}{\beta(A_k | \hat{\mathbf{u}}_k; \mathbf{V})} \text{ and } \hat{\mathbf{u}}_k = [\mathbf{u}, \mathbf{A}_{1:k-1}] \quad (22)$$

are respectively the importance weight and the feature vector of the  $k$ th sub-context.

However, an explosive  $w_{\mathbf{A}}$  may still occur as it involves the product of  $K$  unbounded  $w_k$ . To address the issue, our idea is to prune  $w_k$  included in  $w_{\mathbf{A}}$ . The following theorem shows that under certain

conditions, pruning  $\prod_{j=k+1}^K w_j$  will not change the unbiasedness of the estimate.<sup>2</sup>

THEOREM 2. Assume that

$$P_{\beta}(c_{A_k} | \mathbf{u}, \mathbf{A}_{1:k}, \mathbf{A}_{k+1:K}) = P_{\beta}(c_{A_k} | \mathbf{u}, \mathbf{A}_{1:k}).$$

The resulting estimate

$$\mathbb{E}_{P_{\beta}} \left[ \sum_{k=1}^K \left( \prod_{j=1}^k w_j \right) c_{A_k} \right]$$

after pruning  $\prod_{j=k+1}^K w_j$  is still unbiased to  $V^{\pi}$ .

According to the problem setting described in Section 2.1, we assume  $c_{A_k}$  is only dependent on  $\mathbf{u}$  and  $\mathbf{v}_{A_k}$ , so the conditions required in Theorem 2 can hold for our considered scenarios. Thus we can prune  $\prod_{j=k+1}^K w_j$  in (20) without introducing any bias. Next we discuss how to prune  $\prod_{j=1}^{k-1} w_{A_j}$ . Let

$$V_{\text{IPS}}^{\pi} = \mathbb{E}_{P_{\beta}} \left[ \sum_{k=1}^K w_k c_{A_k} \right] \quad (23)$$

be the resulting estimator after pruning  $\prod_{j=1}^{k-1} w_j$ . From (18), (19), and (21),

$$\begin{aligned} |V_{\text{IPS}}^{\pi} - V_{\text{IPS}}^{\pi}| &= \mathbb{E}_{P_{\beta}} \left[ \sum_{k=1}^K \left| \prod_{j=1}^{k-1} w_j - 1 \right| w_k c_{A_k} \right] \\ &= \mathbb{E}_{P_{\beta}} \left[ \sum_{k=1}^K \frac{|\pi(\mathbf{A}_{1:k-1} | \mathbf{u}; \mathbf{V}, \theta) - \beta(\mathbf{A}_{1:k-1} | \mathbf{u}; \mathbf{V})|}{\beta(\mathbf{A}_{1:k-1} | \mathbf{u}; \mathbf{V})} w_k c_{A_k} \right]. \end{aligned} \quad (24)$$

This then implies that the loss of pruning  $\prod_{j=1}^{k-1} w_j$  can be bounded by the difference between  $\pi$  and  $\beta$ . Intuitively, as long as  $\pi$  imitates  $\beta$ , the product can be pruned without any loss. However, any improvement of  $\pi$  over  $\beta$  can not be achieved, which is contradictory to the purpose of policy optimization. Instead, we alleviate the loss caused by pruning  $\prod_{j=1}^{k-1} w_j$  through constraining the difference between  $\pi$  and  $\beta$ . Consider Pinsker's inequality defined as follows

$$\sup |\pi - \beta| \leq \sqrt{\frac{1}{2} D_{\text{KL}}(\beta || \pi)}, \quad (25)$$

where  $D_{\text{KL}}(\beta || \pi)$  is the Kullback-Leibler divergence between  $\pi$  and  $\beta$ .

$$\begin{aligned} D_{\text{KL}}(\beta || \pi) &= -\mathbb{E}_{\text{Pr}(\mathbf{u})} \sum_{\mathbf{A}} \beta(\mathbf{A} | \mathbf{u}; \mathbf{V}) \log \left( \frac{\pi(\mathbf{A} | \mathbf{u}; \mathbf{V}, \theta)}{\beta(\mathbf{A} | \mathbf{u}; \mathbf{V})} \right) \\ &= -\mathbb{E}_{\text{Pr}(\mathbf{u})} \mathbb{E}_{\beta(\mathbf{A} | \mathbf{u}; \mathbf{V})} \log \left( \frac{\pi(\mathbf{A} | \mathbf{u}; \mathbf{V}, \theta)}{\beta(\mathbf{A} | \mathbf{u}; \mathbf{V})} \right) \\ &= -\mathbb{E}_{\text{Pr}(\mathbf{u})} \mathbb{E}_{\beta(\mathbf{A} | \mathbf{u}; \mathbf{V})} \sum_{k=1}^K \log \left( \frac{\pi(A_k | \hat{\mathbf{u}}_k; \mathbf{V}, \theta)}{\beta(A_k | \hat{\mathbf{u}}_k; \mathbf{V})} \right). \end{aligned} \quad (26)$$

The last equality in (26) follows from (18) and (19). Then, the task of alleviating the loss caused by pruning  $\prod_{j=1}^{k-1} w_j$  can be accomplished by bounding  $D_{\text{KL}}(\beta || \pi)$ . Therefore, we consider to bound  $D_{\text{KL}}(\beta || \pi)$  as an additional constraint by the method of adding a penalty function and propose the following regularized per-item IPS (RIIPS) estimator.

$$V_{\text{RIIPS}}^{\pi} = V_{\text{IPS}}^{\pi} - \alpha D_{\text{KL}}(\beta || \pi), \quad (27)$$

<sup>1</sup>See the proof in our supplementary materials.

<sup>2</sup>See the proof in our supplementary materials.

where  $\alpha$  is the penalty hyper-parameter decided by a validation procedure. Finally, for policy learning, by equipping the regularization term, we achieve the following problem.

$$\max_{\theta} \mathbb{E}_{P_{\beta}} \left[ \sum_{k=1}^K \left( \frac{\pi(A_k | \hat{\mathbf{u}}_k; \mathbf{V}, \theta)}{\beta(A_k | \hat{\mathbf{u}}_k; \mathbf{V})} c_{A_k} + \alpha \log(\pi(A_k | \hat{\mathbf{u}}_k; \mathbf{V}, \theta)) \right) \right] - \lambda \Phi(\theta), \quad (28)$$

where  $\beta$  in  $D_{\text{KL}}(\beta || \pi)$  has been omitted due to its independence to  $\theta$ . It is straightforward to solve the problem by the policy gradient method in (14). Specifically, with the trick defined in (16), the gradient of the objective function in (28) can be computed as

$$\mathbb{E}_{P_{\beta}} \left[ \sum_{k=1}^K \left( \left( \frac{\pi(A_k | \hat{\mathbf{u}}_k; \mathbf{V}, \theta)}{\beta(A_k | \hat{\mathbf{u}}_k; \mathbf{V})} c_{A_k} + \alpha \right) \nabla \log(\pi(A_k | \hat{\mathbf{u}}_k; \mathbf{V}, \theta)) \right) \right] + \lambda \nabla \Phi(\theta). \quad (29)$$

Now, it comes to how to model  $\pi$ . Let  $\hat{y}(\cdot)$  be a score function parametrized by  $\theta$ . To guarantee that the  $\mathbf{A}$  resulting from  $\pi$  is a  $K$ -element permutation of  $\mathbb{A}$ , we apply the Plackett-Luce (PL) model [15], which requires to exclude items recommended beforehand for each sub-context. Recalling (22), we can divide each sub-context  $\hat{\mathbf{u}}_k$  into two parts:  $\mathbf{u}$  for the inference of  $\hat{y}(\cdot)$  and  $\mathbf{A}_{1:k-1}$  for item exclusion. Then, with the softmax function,  $\pi$  is defined as,

$$\pi(A_k | \hat{\mathbf{u}}_k; \mathbf{V}, \theta) = \frac{e^{\hat{y}(\theta; \mathbf{u}, v_{A_k})}}{\sum_{j=1, j \notin \mathbf{A}_{1:k-1}}^n e^{\hat{y}(\theta; \mathbf{u}, v_j)}}. \quad (30)$$

According to  $\mathbf{A}_{1:k-1}$ , items recommended previously are explicitly excluded from the item space  $\mathbb{A}$ . For practical convenience, instead of (30), during training we can use the following formulation,

$$\pi(A_k | \hat{\mathbf{u}}_k; \mathbf{V}, \theta) = \frac{e^{\hat{y}(\theta; \hat{\mathbf{u}}_k, v_{A_k})}}{\sum_{j=1}^n e^{\hat{y}(\theta; \hat{\mathbf{u}}_k, v_j)}}, \quad (31)$$

which does not explicitly exclude  $\mathbf{A}_{1:k-1}$  but instead includes  $\mathbf{A}_{1:k-1}$  as additional contextual features for  $\hat{y}(\cdot)$ . The PL model is then only used for generating a  $K$ -element permutation of  $\mathbb{A}$  during the online deployment. Empirically, the models trained with (30) and (31) have similar performance in our experiments.

### 3.2 Adaptive Policy Learning

The main idea of IPS and its derivatives is to divide the observed reward  $c_{A_k}$  by the propensity  $\beta$ , so that in expectation this adjusted empirical reward will match the underlying reward distribution. As described in [18], IPS generally works when the behavior policy is stochastic and a sufficient amount of events are observed.

However, for recommender systems, the context feature  $\mathbf{u}$  is often very sparse and high-dimensional. This makes some contexts very rare while the number of items is usually very large. Consequently, it is almost impossible to go through every item for every context. In fact, it is more likely that we only observed a small fraction of items for each context.

To see why this can be bad for IPS, we give a simple example in Figure 1, which only involves two items and two independent contexts. Specifically, we assume the first context appears more frequently than the second one. Under these two contexts, the empirical rewards of the two items are observed. We show the normalized empirical reward in Figure 1(a).

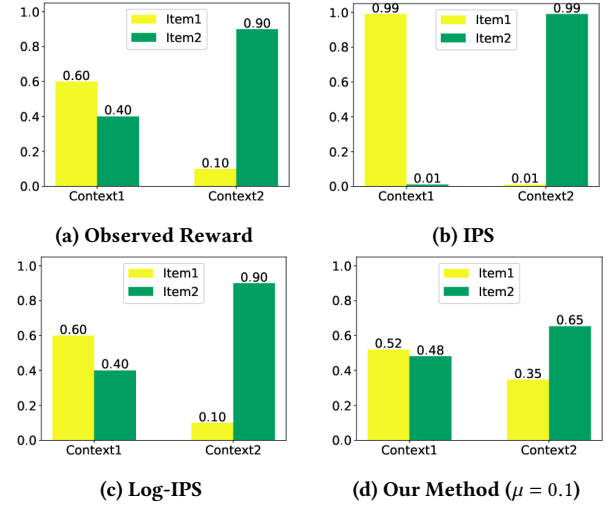


Figure 1: A toy example with two contexts and two items.

From Figure 1(b), we can see that for both contexts, the learned policy from IPS almost puts the probability mass entirely on the item with the highest empirical reward, and the other item is simply ignored. This is known as the winner-takes-all effect, and it has been pointed out [8] that the effect can be harmful to the robustness and generalization of the learned policy from IPS. The issue is particularly acute when we do not have enough observations, and as in this case, the item with the highest empirical reward might not actually be the best one.

To handle the winner-takes-all effect, [8] proposes the so-called log-IPS as below. By applying the log transformation to our RIIPS and switching to a minimization setting, (28) becomes

$$\min_{\theta} \mathbb{E}_{P_{\beta}} \left[ \sum_{k=1}^K \left( -Z_{A_k} \log(\pi(A_k | \hat{\mathbf{u}}_k; \mathbf{V}, \theta)) \right) \right] + \lambda \Phi(\theta), \quad (32)$$

where

$$Z_{A_k} = \frac{c_{A_k}}{\beta(A_k | \hat{\mathbf{u}}_k; \mathbf{V})} + \alpha. \quad (33)$$

According to [8], the log transformation in (32) moves the global optimum and thus makes the learned policy from log-IPS allocate probability mass proportional to the normalized empirical reward in Figure 1(a). This results in what is shown in Figure 1(c). Thus, if the empirical reward spreads across multiple items with enough observations, then log-IPS can alleviate the winner-takes-all effect successfully. For example, as the first context occurs frequently, both items have been observed for several times and the ratio of their empirical reward becomes even and stable. A comparison between Figure 1(c) and Figure 1(b) shows that log-IPS makes good use of this even ratio while IPS does not.

However, we argue that for contexts with few occurrences, it is likely that one item will wrongly occupy a large portion of the empirical reward. The reason behind this is that when one of the many lower propensity items is selected, its empirical reward will be largely amplified by the inverse propensity scaling. If we do not have enough observations, this will result in a high concentration of empirical rewards as the second context in Figure 1(a). In this

case, log-IPS does not help much, as the probability mass of the learned policy is still highly concentrated on one single item.

To address this problem, our idea is to introduce a hyper-parameter which can further smoothen the probability mass of the learned policy. However, we find it difficult to introduce such a hyper-parameter in (32), so an alternative formulation is needed. To this end, we connect (32) to the multi-class classification. A closer look at (32) shows that if the softmax function is used for  $\pi$  as in (31), then  $-\log \pi$  in (32) is the well-known cross-entropy loss function commonly applied for multi-class classification problems. Specifically, (31) is the probability for an instance to be in class  $A_k$  and (32) can be considered as to minimize a weighted negative log-likelihood.

A common alternative used for multi-class classification problems is the one-versus-all scheme. In this setting, a binary probability model is considered for each class. That is, the recommended item  $A_k$  is seen as positive and the non-recommended items are seen as negative. Then, for class  $A_k$ , the following probability model is considered

$$\Pr(\text{positive} \mid \hat{y}(\theta; \hat{\mathbf{u}}_k, \mathbf{v}_{A_k})) = \sigma(\hat{y}(\theta; \hat{\mathbf{u}}_k, \mathbf{v}_{A_k})), \quad (34)$$

where  $\sigma(\cdot)$  is the sigmoid function. In the one-versus-all setting, each binary problem corresponds to its own parameter  $\theta_{A_k}$ , so an independent optimization problem to minimize the negative log-likelihood of (34) is solved.

Keeping this process in mind, here we still consider one  $\theta$  for all binary models and propose the following formulation,

$$\min_{\theta} \mathbb{E}_{P_{\beta}} \sum_{k=1}^K Z_{A_k} \left[ \ell_{\log}^+(\hat{y}(\theta; \hat{\mathbf{u}}_k, \mathbf{v}_{A_k})) + \mu \sum_{j=1, j \neq A_k}^n \ell_{\log}^-(\hat{y}(\theta; \hat{\mathbf{u}}_k, \mathbf{v}_j)) \right] + \lambda \Phi(\theta), \quad (35)$$

where

$$\ell_{\log}^+(x) = \log \sigma(x) \text{ and } \ell_{\log}^-(x) = \log(1 - \sigma(x)), \quad (36)$$

and  $\mu$  is a hyper-parameter explained below. This formulation sums over all the binary classification problems in the one-versus-all scheme. To construct  $\pi$ , we must combine all  $n$  probability models in (34) into one. This issue has been well studied in the literature of the one-versus-all multi-class classification; see, for example, Section 4.1 in [7]. Here, instead of using some sophisticated settings, we use the following heuristic to construct  $\pi$ .

$$\pi(A_k \mid \hat{\mathbf{u}}; \mathbf{V}, \theta) = \frac{\sigma(\hat{y}(\theta; \hat{\mathbf{u}}, \mathbf{v}_{A_k}))}{\sum_{j=1, j \neq A_{1:k-1}}^n \sigma(\hat{y}(\theta; \hat{\mathbf{u}}, \mathbf{v}_j))}, \quad (37)$$

which also explicitly excludes  $A_{1:k-1}$  as in (30).

One major benefit of (35) is that the loss function is now decoupled for the positive and negative labels, and the hyper-parameter  $\mu$  can be naturally added in front of  $\ell_{\log}^-(\cdot)$  to control the smoothness of the final probability distribution. As we can see in Figure 1(d), if we set  $\mu = 0.1$ , this new setting can alleviate the winner-takes-all effect for both the first and the second contexts. In practice, we can use a grid search to find a suitable value for  $\mu$ , and the probability mass of the learned policy will then adapt accordingly. We thus name our approach Adaptive-RIIPS, as it can adapt the smoothness of the learned policy to address the winner-takes-all effect.

### 3.3 Efficient Training

Another critical challenge of policy learning approaches for recommender systems is the training efficiency, which is rarely discussed in past works. As we mentioned in Section 2.2.2, the policy gradient method is applied to solve a policy learning problem. For each context, this involves the calculation of all possible decisions of  $\pi$ . For (13), the number of possibilities is equal to the number of all  $K$ -element permutations of  $\mathbb{A}$ , which has the order of  $O(\frac{n!}{(n-K)!})$ . Clearly, for any top- $K$  recommender system with a tremendous number of contexts and items, solving (13) through the policy gradient method is infeasible.

As for the RIIPS estimator proposed in Section 3.1, we decompose the  $K$ -element permutation decision for each context into  $K$  sequential decisions for  $K$  sub-contexts. Thus  $\pi$  is only responsible for deciding a single item  $A_k$  from  $n$  items for each sub-context, and the decision space of  $\pi$  is remarkably reduced to  $n$ . Consequently, given a training set including  $m$  logged events, the cost of computing (29) by going through all events becomes  $O(Mn)$ , where  $M = mK$  is the number of sub-contexts after decomposition.

However, since both  $M$  and  $n$  can reach over millions in recommender systems, the training time of (29) is still prohibitive with the  $O(Mn)$  cost. This difficulty has been reported and addressed in recent works [3, 13] of policy learning with a large action (item) space. Their main idea is to reduce the cost by subsampling items from  $\mathbb{A}$ . But just as reported in these works, this kind of subsampling mechanism often degenerates the performance. To have a better trade-off between the training efficiency and the performance, we propose a novel solution in this work based on (35), which considers a non-subsampled training setting. As shown in [23], such a non-subsampled setting empirically leads to less performance loss than the subsampled one.

Define

$$\hat{\mathbf{u}}_i = [\mathbf{u}_i]_{A_{1:k-1}},$$

where  $i = (l-1) \times K + k$ , to be the  $k$ th sub-context of the  $l$ th context realization  $\mathbf{u}_l$ . Thus, for  $M = mK$  sub-contexts after the decomposition from  $m$  contexts, we derive the following empirical formulation of (35).

$$\min_{\theta} \frac{1}{m} \sum_{i=1}^M z_i \left[ \ell_{\log}^+(\hat{Y}_{ia_i}) + \mu \sum_{j=1, j \neq a_i}^n \ell_{\log}^-(\hat{Y}_{ij}) \right] + \lambda \Phi(\theta), \quad (38)$$

where  $\hat{Y}_{ij} = \hat{y}(\theta; \hat{\mathbf{u}}_i, \mathbf{v}_j)$ ,  $\mathbf{a} \in \mathbb{A}^M$  and  $\mathbf{z} \in \mathbb{R}^M$ . Here  $\mathbf{a}$  and  $\mathbf{z}$  contain realizations of  $A_k$  and  $Z_{A_k}$  for  $M$  sub-contexts, respectively. With this formulation, we transform the policy learning problem to a binary classification problem of  $Mn$  training instances. Similar to the situation of the RIIPS estimator, solving (38) by conventional optimization methods is not easy. For example, the stochastic gradient (SG) method may face difficulties in the process of sampling  $Mn$  samples [23].

To overcoming this difficulty, our idea is to connect (38) with the extreme similarity learning problem, which learns the relation between a huge number of (context, item) pairs. Efficient training has been well studied for such problems, like [23, 24, 26]. This connection relies on the following facts.



- In our experiments, we consider a two-tower structure for the model  $\hat{g}(\cdot)$ .<sup>3</sup> This setting is very prevalent in modern recommender systems.
- By applying the second-order Taylor expansion, the second term in (38) can be approximated with a weighted squared loss as  $\ell_{sq}(\hat{Y}_{ij}, \hat{Y}_{ij}) = \omega_{ij}(\hat{Y}_{ij} - \hat{Y}_{ij})^2$ , where  $\omega_{ij}$  is a cost weight associated with the loss and  $\hat{Y}_{ij}$  is an imputed value.<sup>4</sup>

As these two facts satisfy the requirements in works of extreme similarity learning, we can apply efficient optimization methods for this problem. In our implementation, we consider the Gauss-Newton method proposed in [26], which can factorize all operations involving the  $\mathcal{O}(Mn)$  cost into a series of operations with a much smaller  $\mathcal{O}(M) + \mathcal{O}(n)$  cost.

## 4 RELATED WORKS

In previous sections, we discuss the counterfactual learning in a storyline of recommender systems, where actions decided by a policy are specified as recommended items. For general scenarios, counterfactual learning has been immensely discussed in recent studies. The mainline of these studies [9, 19, 20] focus on stabilizing the learning process and consider learning a policy that decides only one single action. To the best of our knowledge, [2] is the first work extending policy learning to top-K recommender systems. Similar to our discussion in Section 3.1, they meet the difficulty caused by the vast number of item permutations. To make the problem tractable, they impose the following two constraints.

- For each context, at most one item in  $\mathbf{A}$  has a non-zero reward, which is denoted by  $A^+$ .
- For each  $\mathbf{A}$ , the  $K$  recommended items are sampled with replacement during training but without replacement during deployment.

Under the above constraints, they propose the following top-K REINFORCE estimator,

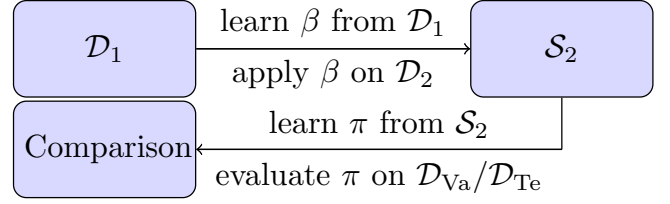
$$\hat{V}_{\text{top-K}}^{\pi}(\theta) = \mathbb{E}_{P_{\beta}} \left[ \frac{1 - (1 - \pi(A^+ | \mathbf{u}; \mathbf{V}; \theta))^K}{z(A^+ | \mathbf{u}; \mathbf{V})} c_{A^+} \right], \quad (39)$$

where  $1 - (1 - \pi(A^+ | \mathbf{u}; \mathbf{V}; \theta))^K$  is the probability of  $A^+$  being included in  $K$  items sampled by  $\pi$  with replacement, and  $z(A^+ | \mathbf{u}; \mathbf{V})$  is the propensity score the same as that in (6). For the gradient of (39), there will be an additional multiplier derived from  $\frac{1 - (1 - \pi(A^+ | \mathbf{u}; \mathbf{V}; \theta))^K}{z(A^+ | \mathbf{u}; \mathbf{V})}$ , which allocates some probability mass of the learned policy to other items of interest and thus helps to alleviate the winner-takes-all effect.

Our idea of pruning importance weights in the RIIPS estimator is inspired by reinforcement learning [16], where the pruning method is used for decreasing weights from long horizons. A similar idea appeared in the policy evaluation [12, 21] for top-K recommendations, where both  $\beta$  and  $\pi$  are pre-known and fixed. For policy learning, we are not the only work to apply this idea. A parallel work [13] also considers pruning the weight around the current action. However, different from their empirical study, we offer a theoretical analysis on the loss of the pruning in Section 3.1. This

<sup>3</sup>See the details in Appendix A.1.

<sup>4</sup>See the detailed derivations in our supplementary materials.



**Figure 2: An illustration of the supervised-to-bandit conversion, where  $\beta$  is the behavior policy, and  $\pi$  is the policy to be learnt and evaluated.**

analysis leads to the addition of a regularization term to alleviate the loss, whose effectiveness is confirmed in our empirical experiments.

On the other hand, while most works of policy learning only consider a relatively small action space (e.g., less than hundreds), [2, 13] are the pioneers to extend policy learning for tremendous actions. To tackle the training efficiency issue described in Section 3.3, which is resulted from going over all the possible actions, they pre-select actions from  $\mathbf{A}$  to reduce the cost. Specifically, in [2, 3], they perform sampled softmax [1] for training  $\pi$ . In [13], they consider a heuristic  $\rho(\cdot)$  to pre-select actions, which results in their proposed POXM estimator as follows.

$$\hat{V}_{\text{POXM}}^{\pi}(\theta) = \mathbb{E}_{P_{\beta}} \left[ \sum_{k=1}^K \frac{\pi(A_k | \hat{\mathbf{u}}_k, \rho(\mathbf{u}); \mathbf{V}, \theta)}{\beta(A_k | \hat{\mathbf{u}}_k; \mathbf{V})} (c_{A_k} - \alpha) \right], \quad (40)$$

where they only consider actions in the pre-selected set  $\rho(\mathbf{u})$  of each context such that

$$\pi(A_k | \hat{\mathbf{u}}_k, \rho(\mathbf{u}); \mathbf{V}, \theta) = \begin{cases} \frac{e^{\hat{y}(\theta; \mathbf{u}, \sigma_{A_k})}}{\sum_{j \in \rho(\mathbf{u}) \setminus \{A_1, \dots, A_{k-1}\}} e^{\hat{y}(\theta; \mathbf{u}, \sigma_j)}} & A_k \in \rho(\mathbf{u}), \\ 0 & A_k \notin \rho(\mathbf{u}). \end{cases}$$

Besides, they follow [9] to adjust  $\alpha$  in (40). In [13], they empirically construct  $\rho(\mathbf{u})$  by referring to the behavior policy  $\beta$  such that

$$\rho(\mathbf{u}) = \{A_k | A_k \text{ has the largest } \beta(A_k | \hat{\mathbf{u}}_k; \mathbf{V}), \text{ and } 1 \leq k \leq p\}, \quad (41)$$

where  $p$  is a specified number of pre-selected actions. The performance of this kind of methods highly relies on the quality of the pre-selection procedure. For example, in [3], they find the number of sampled actions significantly impacts the performance of the learned policy. For POXM, if  $\rho$  can successfully cover actions with the highest rewards with only a small action space, POXM should perform better than approaches learned from the whole action space. However, the task to construct such an effective  $\rho$  is just as challenging as learning  $\pi$ . That is to say, if for all contexts, we are already able to obtain the wanted  $\rho$ , which includes the actions with the highest rewards, then there may be no need for us to learn a new  $\pi$ . In contrast, our proposed training method in Section 3.3 does not require any pre-selection on items.

## 5 EXPERIMENTS

In this section, after presenting our experimental setup for online simulation, we conduct a series of experiments to compare our proposed framework in Section 3 with other existing state-of-the-art approaches.

**Table 2: Data statistics, where  $n$  is the number of items.  $|\cdot|$  indicates the number of elements in a set.**

Data Set	$n$	$ \mathcal{D}_1 $	$ \mathcal{D}_2 $	$ \mathcal{D}_{Va} $	$ \mathcal{D}_{Te} $
ml1m	3,513	12,066	27,153	9,051	12,077
ml10m	10,210	139,022	312,836	104,279	139,081

### 5.1 Experimental Setup for Online Simulation

To simulate the problem setting described in Section 2.1, we following the supervised-to-bandit conversion in [14, 19] to generate our data from MovieLens 1M (ml1m) and 10M (ml10m) data sets.<sup>5</sup> We first randomly split a data set into four independent subsets:  $\mathcal{D}_1$ ,  $\mathcal{D}_2$ ,  $\mathcal{D}_{Va}$  and  $\mathcal{D}_{Te}$ . Then, as shown in Figure 2, we learn a behavior policy  $\beta$  with  $\mathcal{D}_1$  and deploy it on  $\mathcal{D}_2$  by logging a  $(\mathbf{u}, \mathbf{A}, \mathbf{C})$  event for each context to form a set  $\mathcal{S}_2$ . Next we conduct different approaches to learn a new policy  $\pi$ , where  $\mathcal{D}_{Va}$  is used as the validation set for tuning hyper-parameters. Finally, we evaluate  $\pi$  on the test set  $\mathcal{D}_{Te}$  and compare the resulting performance. The whole procedure is illustrated in Figure 2, while the statistics of the generated data sets are in Table 2. More details of data preprocessing and constructing  $\beta$  are in Appendices A.2 and A.3.

For learning  $\pi$ , we compare seven approaches, which can be grouped into four categories:

- Value-IPS, Value-DR: They are two value learning approaches solving (9) respectively with (6) and (8).
- BanditNet: This approach considers the estimate in (20) without pruning  $\mathbf{w}_A$ . To alleviate the possible issue caused by any extremely large  $\mathbf{w}_A$ , it additionally applies the technique proposed in [9].
- Top- $K$  REINFORCE [2] and POXM [13]: They are two state-of-the-art competitors in the policy learning for top- $K$  decision making; see the discussion in Section 4. For POXM, we follow [13] to construct  $\rho(\cdot)$  by (41) and conduct grid searches on  $p$ .
- RIIPS, Log-RIIPS and Adaptive-RIIPS: Three approaches derived from our proposed RIIPS estimator, which respectively solve problems in (28), (32) and the approximation of (38). For Adaptive-RIIPS, we empirically set all  $\tilde{Y}_{ij} = -1$  and consider  $\omega_{ij} = \omega$  as a hyper-parameter to tune.

More details of the hyper-parameter setting and selection are in Appendix A.4.

To evaluate  $\pi$ , we simulate the online environment. Specifically, we compute the following average cumulated reward (ACR) on a data set with fully observed rewards for all items.

$$\text{ACR} = \frac{1}{\hat{m}} \sum_{i=1}^{\hat{m}} r(\mathbf{c}_i, \mathcal{A}_i), \quad (42)$$

where  $\hat{m}$  is the number of contexts included in this set,  $\mathbf{c}_i$  is the fully-observed reward vector of context  $i$ , and  $\mathcal{A}_i$  is the set of  $K$  items recommended according to  $\pi(\mathcal{A}_i | \mathbf{u}_i; \mathbf{V})$ . Here we focus on the exploitation ability of each learned policy  $\pi$ , so the decision of  $\mathcal{A}_i$  is made by selecting top  $K$  items with the highest scores predicted by  $\hat{y}(\cdot)$ . For consistency, we use the same  $K$  for constructing  $\mathcal{S}_2$  and evaluating models on  $\mathcal{D}_{Va}$  and  $\mathcal{D}_{Te}$ .

<sup>5</sup><https://grouplens.org/datasets/movielens/>

**Table 3: For two data sets, we report the ACR of all approaches on  $\mathcal{D}_{Te}$ . All scores below are scaled up by 100.**

Approach	ml1m		ml10m	
	$K = 1$	$K = 10$	$K = 1$	$K = 10$
Value-IPS	24.04	25.91	25.23	17.07
Value-DR	38.80	28.75	27.46	26.25
POXM	28.66	16.48	22.74	14.54
top- $K$ REINFORCE	34.47	29.13	30.41	25.39
BanditNet	34.51	18.21	30.38	19.14
RIIPS	35.99	28.44	30.41	24.19
Log-RIIPS	35.11	26.40	28.73	20.57
Adaptive-RIIPS	<b>39.02</b>	<b>32.27</b>	<b>34.66</b>	<b>28.27</b>
RIIPS ( $\alpha = 0$ )	34.47	23.01	30.41	23.01

### 5.2 Comparison on Various Approaches

By comparing results in Table 3, we have the following observations.

- Firstly, Value-DR is much better than Value-IPS. This result confirms what we described in Section 3.2. That is, the IPS approach may not work properly when the amount of observed events is deficient. Compared to Value-IPS, Value-DR imputes rewards for unobserved items and reweights the loss term for these items, which can be seen as a way to tune the reward distribution adaptively. This leads to a great improvement on the performance.
- RIIPS makes a significant improvement over BanditNet when  $K = 10$ . The reason is that RIIPS prunes  $\mathbf{w}_A$ . As reported in [11], the value of  $\mathbf{w}_A$  gets exploded very easily when  $K > 1$ , which leads to a huge divergence between  $\hat{V}_{IPS}^\pi$  and  $V^\pi$ . BanditNet's performance suffers without pruning  $\mathbf{w}_A$ .
- In Table 3, RIIPS is better than RIIPS ( $\alpha = 0$ ) for  $K = 10$ , which confirms the rationality of adding the regularization term  $D_{KL}(\beta || \pi)$  with pruning  $\prod_{j=1}^{K-1} \mathbf{w}_j$ .
- Compared to other approaches, top- $K$  REINFORCE is almost the second optimal approach in policy learning when  $K = 10$ . This verifies the applicability of (39) proposed in [2] and confirms the necessity for alleviating the winner-takes-all effect.
- Adaptive-RIIPS proposed in this work performs the best. As stated in Section 3.2, by tuning  $\omega$ , we can adapt the probability mass of the learned policy. This effectively addresses the winner-takes-all effect and overcomes the drawback of Log-RAIPS.
- In our data sets, POXM performs the worst because as described in Section 4, it is tough to get an effective  $\rho(\mathbf{u})$ . Even though our behavior policy  $\beta$  is trained on  $\mathcal{D}_1$  with fully-revealed  $\mathbf{c}$ , the pre-selected  $\rho$  still can not cover enough items with the highest rewards. Increasing the size of  $\rho$  may help, but it also lowers the training efficiency. This issue might limit the usability of POXM.

To further confirm the efficiency and effectiveness of the framework proposed in this work, we compare the performance and training time of Adaptive-RIIPS trained by two different optimizers: Adagrad and Gauss-Newton. In Table 4, we can observe that Adaptive-RIIPS with Gauss-Newton consistently converges faster and better than the one with Adagrad. Take ml1m with  $K = 1$  for example. Gauss-Newton not only speeds up more than 10 times for the model training but also improves the model performance by about 6% respectively when compared to Adagrad.



**Table 4: Comparison on Adaptive-RIIPS with two different optimizers: Adagrad and Gauss-Newton. Training time and performance under  $K = 1$  and 10 for two datasets are shown.**

Dataset	Optimizers	$K = 1$		$K = 10$	
		Time(s)	ACR	Time(s)	ACR
ml1m	Adagrad	739.04	36.74	6028.79	30.70
	Gauss-Newton	58.12	39.02	1422.50	32.70
ml10m	Adagrad	13826.98	34.24	60850.22	27.07
	Gauss-Newton	673.52	34.66	14500.11	28.27

## 6 CONCLUSION

In this work, we categorize the conventional approaches of counterfactual learning for recommender systems into two classes: value learning and policy learning. To do a comparison between these two kinds of approaches, we point out that some existing difficulties in the policy learning approaches for large top- $K$  recommender systems must be addressed first.

- It is likely that the extremely small propensity of a top- $K$  recommendation would lead to a poor estimation of  $\hat{V}_{\text{IPS}}^{\pi}$  and thus a sub-optimal policy.
- The robustness and generalization of policy learning approaches suffer from the limited observations of each item for each context.
- The issue of training efficiency in recommender systems with a large item space limits the usability of policy learning approaches.

To address the above difficulties, we derive a novel policy learning framework. We introduce a regularized per-item approach to balance between the variance and the bias of policy learning. Then, through decoupling the objective function and introducing an extra hyper-parameter for tuning the smoothness, we manage to improve the robustness of policy learning. For efficient training, we integrate our proposed approach with the two-tower structure and the algorithm developed in [26] to be a framework able to train large top- $K$  recommender systems. With experiments conducted on real-world data sets, we confirm the effectiveness and efficiency of our proposed framework.

## ACKNOWLEDGMENTS

This work was supported in part by MOST of Taiwan grant 110-2221-E-002-115-MY3.

## REFERENCES

- [1] Yoshua Bengio and Jean-Sébastien Senécal. 2003. Quick Training of Probabilistic Neural Nets by Importance Sampling. In *Proceedings of the Ninth International Workshop on Artificial Intelligence and Statistics*. 17–24.
- [2] 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*. 456–464.
- [3] Minmin Chen, Bo Chang, Can Xu, and Ed H. Chi. 2021. User Response Models to Improve a REINFORCE Recommender System. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining (WSDM)*. 121–129.
- [4] Miroslav Dudík, John Langford, and Lihong Li. 2011. Doubly Robust Policy Evaluation and Learning. In *Proceedings of the 28th International Conference on Machine Learning (ICML)*. 1097–1104.
- [5] Daniel G. Horvitz and Donovan J. Thompson. 1952. A Generalization of Sampling Without Replacement From a Finite Universe. *J. Amer. Statist. Assoc.* 47 (1952), 663–685.
- [6] Po-Sen Huang, Xiaodong He, Jianfeng Gao, Li Deng, Alex Acero, and Larry Heck. 2013. Learning deep structured semantic models for web search using clickthrough data. In *Proceedings of the 22nd ACM international conference on Information and Knowledge Management*. 2333–2338.
- [7] Tzu-Kuo Huang, Ruby C. Weng, and Chih-Jen Lin. 2006. Generalized Bradley-Terry Models and Multi-class Probability Estimates. *Journal of Machine Learning Research* 7 (2006), 85–115. <http://www.csie.ntu.edu.tw/~cjlin/papers/generalBT.pdf>
- [8] Olivier Jeunen, David Rohde, Flavian Vasile, and Martin Bompaire. 2020. Joint Policy-Value Learning for Recommendation. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*. 1223–1233.
- [9] Thorsten Joachims, Adith Swaminathan, and Maarten de Rijke. 2018. Deep learning with logged bandit feedback. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- [10] Yehuda Koren, Robert M. Bell, and Chris Volinsky. 2009. Matrix Factorization Techniques for Recommender Systems. *Computer* 42 (2009), 30–37.
- [11] Damien Lefortier, Adith Swaminathan, Xiaotao Gu, Thorsten Joachims, and Maarten de Rijke. 2016. Large-scale validation of counterfactual learning methods: A test-bed. In *NIPS Workshop on Inference and Learning of Hypothetical and Counterfactual Interventions in Complex Systems*.
- [12] Shuai Li, Yasin Abbasi-Yadkori, Branislav Kveton, Shan Muthukrishnan, Vishwa Vinay, and Zheng Wen. 2018. Offline evaluation of ranking policies with click models. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*. 1685–1694.
- [13] Romain Lopez, Inderjit Dhillon, and Michael I. Jordan. 2021. Learning from eXtreme Bandit Feedback. (2021).
- [14] Jiaqi Ma, Zhe Zhao, Xinyang Yi, Ji Yang, Minmin Chen, Jiayi Tang, Lichan Hong, and Ed H Chi. 2020. Off-policy Learning in Two-stage Recommender Systems. In *Proceedings of The Web Conference*. 463–473.
- [15] John I. Marden. 1995. *Analyzing and Modeling Rank Data*. Chapman & Hall, London.
- [16] Doina Precup, Richard S. Sutton, and Satinder P. Singh. 2000. Eligibility Traces for Off-Policy Policy Evaluation. In *Proceedings of the Seventeenth International Conference on Machine Learning (ICML)*. 759–766.
- [17] Steffen Rendle. 2010. Factorization machines. In *Proceedings of IEEE International Conference on Data Mining (ICDM)*. 995–1000.
- [18] Naveen Sachdeva, Yi Su, and Thorsten Joachims. 2020. Off-policy bandits with deficient support. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*. 965–975.
- [19] Adith Swaminathan and Thorsten Joachims. 2015. Batch learning from logged bandit feedback through counterfactual risk minimization. *Journal of Machine Learning Research* 16, 1 (2015), 1731–1755.
- [20] Adith Swaminathan and Thorsten Joachims. 2015. The Self-normalized Estimator for Counterfactual Learning. In *Proceedings of the 28th International Conference on Neural Information Processing Systems (NIPS)*. 3231–3239.
- [21] Adith Swaminathan, Akshay Krishnamurthy, Alekh Agarwal, Miroslav Dudík, John Langford, Damien Jose, and Imed Zitouni. 2017. Off-Policy Evaluation for Slate Recommendation. (2017), 3635–3645.
- [22] Xinyang Yi, Ji Yang, Lichan Hong, Derek Zhiyuan Cheng, Lukasz Heldt, Aditee Kumthekar, Zhe Zhao, Li Wei, and Ed Chi. 2019. Sampling-Bias-Corrected Neural Modeling for Large Corpus Item Recommendations. In *Proceedings of the 13th ACM Conference on Recommender Systems*. 269–277.
- [23] Hsiang-Fu Yu, Mikhail Bilenko, and Chih-Jen Lin. 2017. Selection of Negative Samples for One-class Matrix Factorization. In *Proceedings of SIAM International Conference on Data Mining (SDM)*. <http://www.csie.ntu.edu.tw/~cjlin/papers/one-class-mf/biased-mf-sdm-with-supp.pdf>
- [24] Hsiang-Fu Yu, Hsin-Yuan Huang, Inderjit S. Dhillon, and Chih-Jen Lin. 2017. A Unified Algorithm for One-class Structured Matrix Factorization with Side Information. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence (AAAI)*. <http://www.csie.ntu.edu.tw/~cjlin/papers/ocmf-side/biased-leml-aaai-with-supp.pdf>
- [25] Bowen Yuan, Jui-Yang Hsia, Meng-Yuan Yang, Hong Zhu, Chihyao Chang, Zhenhua Dong, and Chih-Jen Lin. 2019. Improving Ad Click Prediction by Considering Non-displayed Events. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management (CIKM)*. [http://www.csie.ntu.edu.tw/~cjlin/papers/occtr/ctr\\_oc.pdf](http://www.csie.ntu.edu.tw/~cjlin/papers/occtr/ctr_oc.pdf)
- [26] Bowen Yuan, Yu-Sheng Li, Pengrui Quan, and Chih-Jen Lin. 2021. Efficient optimization methods for extreme similarity learning with nonlinear embeddings. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*. [http://www.csie.ntu.edu.tw/~cjlin/papers/similarity\\_learning/pq.pdf](http://www.csie.ntu.edu.tw/~cjlin/papers/similarity_learning/pq.pdf)
- [27] Bowen Yuan, Yaxu Liu, Jui-Yang Hsia, Zhenhua Dong, and Chih-Jen Lin. 2020. Unbiased Ad click prediction for position-aware advertising systems. In *Proceedings of the 14th ACM Conference on Recommender Systems*. <http://www.csie.ntu.edu.tw/~cjlin/papers/debiases/debiases.pdf>
- [28] Bowen Yuan, Meng-Yuan Yang, Jui-Yang Hsia, Hong Zhu, Zhirong Liu, Zhenhua Dong, and Chih-Jen Lin. 2019. *One-class Field-aware Factorization Machines for Recommender Systems with Implicit Feedbacks*. Technical Report. National Taiwan University. [http://www.csie.ntu.edu.tw/~cjlin/papers/ocfm/imp\\_ffm.pdf](http://www.csie.ntu.edu.tw/~cjlin/papers/ocfm/imp_ffm.pdf)

## A APPENDIX

### A.1 Details of Model Setting and its Connection with Extreme Similarity Learning Problems

The connection between our proposed Adaptive-RIIPS and extreme similarity learning problems relies on the fact that the two-tower structure models are prevalent in modern recommender systems. Let  $\theta = [\theta_u]$ . The two-tower structure transforms  $\hat{\mathbf{u}}$  and  $\mathbf{v}$  into two  $d$ -dimensional embeddings respectively by learning two embedding functions  $f(\theta_u; \hat{\mathbf{u}})$  and  $g(\theta_v; \mathbf{v})$ , and then  $\hat{y}(\theta; \hat{\mathbf{u}}, \mathbf{v}) = f(\theta_u; \hat{\mathbf{u}})^\top g(\theta_v; \mathbf{v})$  is the score function of this structure. Many commonly-used recommendation models fall into the two-tower structure. For example, when  $f(\cdot)$  and  $g(\cdot)$  are linear embedding functions, it has been pointed out [28] that  $\hat{y}(\cdot)$  is a variant of the factorization machine [17]. Besides, non-linear embedding functions (e.g., neural networks) are also widely applied in two-tower recommendation models [6, 22].

Through applying the two-tower structure models for  $\hat{y}(\theta; \hat{\mathbf{u}}, \mathbf{v})$ , the problem (38) can be considered as an extreme similarity learning problem. That is, we learn the similarity between a sub-context and an item from extremely large  $Mn$  (sub-context, item) pairs. Here the similarity presents the tendency of an item being recommended by  $\pi$  under a given sub-context. To avoid any  $\mathcal{O}(Mn)$  cost in solving an extreme similarity learning problem, existing works in this field [23, 24, 26] impose the following squared loss on a certain part of the loss function. In our case, the second term in (38), which deals with the similarity between sub-contexts and non-recommended items, is replaced with  $\ell_{sq}(\tilde{Y}_{ij}, \hat{Y}_{ij}) = \omega_{ij}(\tilde{Y}_{ij} - \hat{Y}_{ij})^2$ , where  $\omega_{ij}$  is a cost associated with the loss and  $\tilde{Y}_{ij}$  is an imputed value. For each  $\omega_{ij}$  and  $\tilde{Y}_{ij}$ , it is required that they can be decomposed to multiple parts solely related to  $i$  or  $j$ , respectively.

Then, the remaining issue is how to choose  $\omega_{ij}$  and impute  $\tilde{Y}_{ij}$ . We apply the second-order Taylor expansion of  $\ell_{\log}^-(\hat{Y}_{ij})$  at our chosen point  $\tilde{y}_j$  for each item to get the following alternative loss function for non-recommended items.<sup>6</sup>

$$\frac{1}{2} \nabla^2 \ell_{\log}^-(\tilde{y}_j) \left( \hat{Y}_{ij} + \frac{\nabla \ell_{\log}^-(\tilde{y}_j) - \nabla^2 \ell_{\log}^-(\tilde{y}_j) \tilde{y}_j}{\nabla^2 \ell_{\log}^-(\tilde{y}_j)} \right)^2. \quad (43)$$

By applying (43) to (38), we attain the following problem to solve.

$$\min_{\theta} \frac{1}{m} \sum_{i=1}^M z_i \left[ \ell_{\log}^+(\hat{Y}_{ia_i}) + \omega_{ij} \sum_{j=1, j \neq a_i}^n \frac{1}{2} (\tilde{Y}_{ij} - \hat{Y}_{ij})^2 \right] + \lambda \Phi(\theta), \quad (44)$$

where for all  $i$ ,  $\omega_{ij} = \mu \nabla^2 \ell_{\log}^-(\tilde{y}_j)$ , and  $\tilde{Y}_{ij} = \frac{\nabla^2 \ell_{\log}^-(\tilde{y}_j) \tilde{y}_j - \nabla \ell_{\log}^-(\tilde{y}_j)}{\nabla^2 \ell_{\log}^-(\tilde{y}_j)}$

are pre-defined constants solely dependent to  $j$ , and  $\tilde{Y}_{ij}$  comes from the two-tower structure. As deciding  $\tilde{y}_j$  is like to select a hyper-parameter, we instead treat  $\omega_{ij}$  and  $\tilde{Y}_{ij}$  as two hyper-parameters and tune them by grid search. In our experiments, this simple transformation is found to be effective enough empirically.

<sup>6</sup>See the detailed derivations in our supplementary materials.

### A.2 Details of Data Preprocessing

For MovieLens 1M (ml1m) and 10M (ml10m) data sets, we follow [23] to binarize ratings included in original sets. We consider pairs with rating  $\geq 4$  as positive while the rest including unrated items as negative. As original sets do not include contextual information, we need to generate the contextual features  $\mathbf{u}$  and context-aware rewards  $\mathbf{c}$  respectively being realizations of  $\mathbf{u}$  and  $\mathbf{c}$ . Specifically, for each user, we group his/her positive items into a set and then randomly divide it into two subsets equally. Positive items in the first subset, like movies the user has watched, are used as the contextual information to construct  $\mathbf{u} \in \{0, 1\}^n$  where

$$u_j = \begin{cases} 1 & \text{item } j \text{ is included in the first subset,} \\ 0 & \text{otherwise.} \end{cases}$$

Besides, positive items in the second subset, like movies the user is going to watch, are used for constructing  $\mathbf{c}$  where

$$c_j = \begin{cases} 1 & \text{item } j \text{ is included in the second subset,} \\ 10^{-3} & \text{otherwise.} \end{cases}$$

We repeat the above procedure to generate ten  $(\mathbf{u}, \mathbf{c})$  pairs for each user. For  $\mathbf{V}$ , we generate a feature vector of each item by one-hot encoding with its identity number.

### A.3 Details of Behavior Policy $\beta$

For learning  $\beta$ , similar to past works [13, 19], we learn  $\hat{s}(\cdot)$  by a value learning approach with the fully-revealed  $\mathbf{c}$  included in  $\mathcal{D}_1$ . We apply the approach in [24] for this task. With the learned  $\hat{s}(\cdot)$ , we transform it into a policy  $\beta$  defined in (19) by the PL ranking model defined in (30) such that for each sub-context, we have

$$\beta(A_k | \hat{\mathbf{u}}_k; \mathbf{V}) = \frac{e^{\hat{s}(\theta; \hat{\mathbf{u}}_k, \mathbf{v}_{A_k})}}{\sum_{j=1, j \notin A_{1:k-1}}^n e^{\hat{s}(\theta; \hat{\mathbf{u}}_k, \mathbf{v}_j)}}. \quad (45)$$

### A.4 Details of Hyper-Parameter Selection

For the two-tower structure we use in all approaches, the output size of the embeddings is set to be 128. Adaptive-RIIPS uses the Gauss-Newton method proposed in [26], while all the other approaches are trained by Adagrad. For Adagrad, the learning rate is initialized to 0.05 and the batch size is set to 10% of the training data. As an exception, for ml10m with  $K = 10$ , the batch size of Adagrad is set to 1% of the training data due to the size of our RAM. For Value-DR with (8), we set  $\hat{c}_j = \log(\frac{\overline{\text{CTR}}}{1 - \overline{\text{CTR}}})$ ,  $\forall j$ , where  $\overline{\text{CTR}}$  is the average click-through-rate from an unbiased set collected by a uniform behavior policy. For convenience, we directly compute  $\overline{\text{CTR}}$  from  $\mathcal{D}_2$ , which is unbiased due to fully-observed reward vectors.

We train each approach for  $T$  epochs, retain the model at the epoch with the best validation performance, and report the test performance by applying the model to predict the test set. For Adagrad,  $T = 500$  and for Gauss-Newton,  $T = 30$ .

The following hyper-parameters are selected by a grid search. Candidates of each hyper-parameter are also listed below.

- $\lambda \in \{4^{-1}, 4^0, \dots, 4^4\}$ : the  $l_2$  regularization coefficient.
- $\alpha \in \{0, 0.01, 0.1, 1\}$ : the coefficient of  $D_{\text{KL}}(\beta || \pi)$ .
- $\gamma, \omega \in \{16^{-5}, 16^{-4}, 16^{-3}, 16^{-2}\}$ : the coefficients of the square loss term in (8) and (44).

- $\zeta \in \{0.6, 0.8, 1.0, 1.2\}$ : the Lagrange multiplier in [9].
- $\xi \in \{0.1, 0.3, 0.5\}$ : the tolerance in the relative stopping condition for Gauss-Newton.
- $p \in \{10, 20, 50\}$ : the size of  $\rho$  for POXM.