

Off-Policy Evaluation of Ranking Policies under Diverse User Behavior

Haruka Kiyohara
Hanjuku-Kaso Co., Ltd.
kiyohara@hanjuku-kaso.com

Masatoshi Uehara
Cornell University
mu223@cornell.edu

Yusuke Narita
Yale University
yusuke.narita@yale.edu

Nobuyuki Shimizu
Yahoo Japan Corporation
nobushim@yahoo-corp.jp

Yasuo Yamamoto
Yahoo Japan Corporation
yasyamam@yahoo-corp.jp

Yuta Saito
Cornell University
ys552@cornell.edu

ABSTRACT

Ranking interfaces are everywhere in online platforms. There is thus an ever growing interest in their *Off-Policy Evaluation* (OPE), aiming towards an accurate performance evaluation of ranking policies using logged data. A de-facto approach for OPE is *Inverse Propensity Scoring* (IPS), which provides an unbiased and consistent value estimate. However, it becomes extremely inaccurate in the ranking setup due to its high variance under large action spaces. To deal with this problem, previous studies assume either independent or cascade user behavior, resulting in some ranking versions of IPS. While these estimators are somewhat effective in reducing the variance, all existing estimators apply a single universal assumption to every user, causing excessive bias and variance. Therefore, this work explores a far more general formulation where user behavior is diverse and can vary depending on the user context. We show that the resulting estimator, which we call *Adaptive IPS* (AIPS), can be unbiased under any complex user behavior. Moreover, AIPS achieves the minimum variance among all unbiased estimators based on IPS. We further develop a procedure to identify the appropriate user behavior model to minimize the mean squared error (MSE) of AIPS in a data-driven fashion. Extensive experiments demonstrate that the empirical accuracy improvement can be significant, enabling effective OPE of ranking systems even under diverse user behavior.

CCS CONCEPTS

• **Information systems** → **Retrieval models and ranking; Evaluation of retrieval results.**

KEYWORDS

off-policy evaluation, ranking policy, inverse propensity score.

ACM Reference Format:

Haruka Kiyohara, Masatoshi Uehara, Yusuke Narita, Nobuyuki Shimizu, Yasuo Yamamoto, and Yuta Saito. 2023. Off-Policy Evaluation of Ranking Policies under Diverse User Behavior. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '23), August*

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.

KDD '23, August 6–10, 2023, Long Beach, CA, USA

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0103-0/23/08...\$15.00

<https://doi.org/10.1145/3580305.3599447>

6–10, 2023, Long Beach, CA, USA. ACM, New York, NY, USA, 20 pages.
<https://doi.org/10.1145/3580305.3599447>

1 INTRODUCTION

Ranking interfaces serve as a crucial component in many real-world search and recommender systems, where rankings (as actions) are often optimized through contextual bandit processes [8, 19, 24]. As these ranking policies interact with the environment, they collect logged data valuable for *Off-Policy Evaluation* (OPE) [27, 28]. The goal of OPE is to accurately evaluate the performance of new policies using only the logged data without deploying them in the field, providing a safe alternative to online A/B testing [11, 18].

A popular approach for OPE is *Inverse Propensity Scoring* (IPS) [25, 34], which provides an unbiased estimate of policy performance. Although several estimators are developed on top of IPS in standard OPE [35, 40], they can severely degrade in large action spaces [29]. In particular, IPS becomes vulnerable when applied to ranking policies where the number of actions grows exponentially [22, 24]. An existing popular approach to deal with this variance issue is to make some assumptions about user behavior and define the corresponding IPS estimators [19, 22, 24]. For instance, *Independent IPS* (IIPS) [22] assumes that a user interacts with the items in a ranking completely independently. In contrast, *Reward interaction IPS* (RIPS) [24] assumes that a user interacts with the items sequentially from top to bottom, namely the cascade assumption [12]. These estimators provide some variance reduction, while remaining unbiased when the corresponding assumption is satisfied.

Although this approach has been shown to improve IPS in some ranking applications [19, 24], a critical limitation is that all existing estimators model every user's behavior based on a single, universal assumption (such as independence and cascade). However, it is widely acknowledged that real user behavior can be much more diverse [23, 41, 43] and heterogeneous [42]. With such diverse user behavior, the existing approach can suffer from significant bias and variance. For instance, consider a scenario with two groups of users, one following independent user behavior and the other following a cascade model. In this situation, IIPS is no longer unbiased, and RIPS is sub-optimal in terms of variance. An ideal strategy would arguably be to apply IIPS and RIPS to each group adaptively.

We thus explore a much more general formulation assuming that user behavior is sampled from a distribution conditional on the user context to capture possibly diverse behavior. On top of our general formulation, we propose a new estimator called *Adaptive IPS* (AIPS), which applies different importance weighting schemes

to different users based on their behavior, namely *adaptive importance weighting*. We show that AIPS is unbiased for virtually any distribution of user behavior and that AIPS achieves the minimum variance among all unbiased IPS estimators. We also analyze the bias-variance tradeoff of AIPS in the case of unknown user behavior, which interestingly implies that the true behavior model may not result in an optimal OPE. Thus, we develop a strategy to *optimize* the behavior model from the logged data in a way that minimizes the MSE of AIPS rather than merely trying to estimate the true behavior model. Experiments on synthetic and real-world data demonstrate that AIPS provides a significant gain in MSE over existing methods particularly when the user behavior is diverse.

Our contributions can be summarized as follows.

- We propose a novel formulation and estimator for OPE of ranking policies capturing diverse user behavior.
- We show that AIPS is unbiased for any distributions of user behavior and that it achieves the minimum variance.
- We develop a non-parametric procedure to minimize the MSE of AIPS through optimizing (rather than estimating) the behavior model from the logged data.
- We empirically demonstrate that AIPS enables much more accurate OPE particularly under diverse user behavior.

2 PRELIMINARIES

This section formulates OPE of ranking policies.

2.1 Off-Policy Evaluation of Ranking Policies

We use $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^d$ to denote a context vector (e.g., user demographics) and \mathcal{A} to denote a finite set of discrete actions. Let then $\mathbf{a} = (a_1, a_2, \dots, a_k, \dots, a_K)$ denote a ranking action vector of length K (e.g., a ranked list of songs). We call a function $\pi : \mathcal{X} \rightarrow \Delta(\mathcal{A}^K)$ a *factored* policy. Given context \mathbf{x} , it chooses an action at each position (a_k) independently, where $\pi(\mathbf{a} | \mathbf{x}) = \prod_{k=1}^K \pi(a_k | \mathbf{x})$ is the probability of choosing a specific ranking action \mathbf{a} . In contrast, we call $\pi : \mathcal{X} \rightarrow \Delta(\Pi_K(\mathcal{A}))$ a *non-factored* policy, where $\Pi_K(\mathcal{A})$ is a set of K -permutation of \mathcal{A} . Note that a *factored* policy may choose the same action more than once in a ranking, whereas a *non-factored* policy selects a ranking action without replacement (i.e., $\forall 1 \leq k < l \leq K, a_k \neq a_l$). In addition, let $\mathbf{r} = (r_1, r_2, \dots, r_k, \dots, r_K)$ denote a reward vector with r_k being a random reward observed at the k -th position (e.g., clicks, conversions, dwell time).

In OPE of ranking policies, we are interested in estimating the following *policy value* of evaluation policy π as a measure of its effectiveness [19, 24]:

$$\begin{aligned} V(\pi) &:= \mathbb{E}_{p(\mathbf{x})\pi(\mathbf{a}|\mathbf{x})} \left[\sum_{k=1}^K \alpha_k q_k(\mathbf{x}, \mathbf{a}) \right] \\ &= \sum_{k=1}^K \alpha_k \underbrace{\mathbb{E}_{p(\mathbf{x})\pi(\mathbf{a}|\mathbf{x})} [q_k(\mathbf{x}, \mathbf{a})]}_{V_k(\pi)}, \end{aligned} \quad (1)$$

where $q_k(\mathbf{x}, \mathbf{a}) := \mathbb{E}[r_k | \mathbf{x}, \mathbf{a}]$ is the position-wise expected reward function given context \mathbf{x} and ranking action \mathbf{a} . $V_k(\pi)$ is the *position-wise* policy value and α_k is a non-negative weight assigned to position k . Our definition of the policy value in Eq. (1) captures a wide variety of information retrieval metrics. For example, when

$\alpha_k := 1/\log_2(k+1)$, $V(\pi)$ becomes identical to the discounted cumulative gain (DCG) [14] under policy π . Throughout this paper, we focus on estimating the position-wise policy value $V_k(\cdot)$, as estimating $V(\cdot)$ is straightforward given an estimate of $V_k(\cdot)$.

For performing an OPE, we can leverage logged bandit data collected under the *logging policy* π_0 , i.e., $\mathcal{D} := \{(\mathbf{x}_i, \mathbf{a}_i, \mathbf{r}_i)\}_{i=1}^n$ where \mathbf{a}_i is a vector of discrete variables that indicate which ranking action is chosen by π_0 for individual i . \mathbf{x}_i and \mathbf{r}_i denote the context and reward vectors observed for i . To sum, a logged bandit dataset is generated in the following process:

$$\{(\mathbf{x}_i, \mathbf{a}_i, \mathbf{r}_i)\}_{i=1}^n \sim \prod_{i=1}^n p(\mathbf{x}_i)\pi_0(\mathbf{a}_i | \mathbf{x}_i)p(\mathbf{r}_i | \mathbf{x}_i, \mathbf{a}_i).$$

Note that we assume that the logging policy provides full support over the ranking action space. The accuracy of an estimator \hat{V} is measured by its MSE, i.e., $\text{MSE}(\hat{V}) := \mathbb{E}_{\mathcal{D}}[(V(\pi) - \hat{V}(\pi; \mathcal{D}))^2]$, which can be decomposed into squared bias and variance of \hat{V} .

2.2 Existing Estimators

Here, we summarize some notable existing estimators for OPE in the ranking setup and their statistical properties.

Inverse Propensity Scoring. IPS uses the *ranking-wise* importance weight to provide an unbiased and consistent estimate as follows.

$$\hat{V}_k^{\text{IPS}}(\pi; \mathcal{D}) := \frac{1}{n} \sum_{i=1}^n \frac{\pi(\mathbf{a}_i | \mathbf{x}_i)}{\pi_0(\mathbf{a}_i | \mathbf{x}_i)} r_{i,k}.$$

IPS does not impose any particular user behavior model, and thus it is generally unbiased and consistent under standard identification assumptions. However, it suffers from extremely high variance when the action space ($|\mathcal{A}^K|$ or $|\Pi_K(\mathcal{A})|$) is large [29], which is particularly problematic in the ranking setup [19, 22, 24, 36].

Independent IPS. IIPS assumes that a user interacts with the actions in a ranking independently, which is known as the independence assumption or item-position model [22]. This assumption posits that the reward observed at each position depends solely on the action chosen at that particular position, not on the other actions presented in the same ranking. Under this independence assumption, it is sufficient to condition only on a_k to characterize the corresponding position-wise expected reward, i.e., $q_k(\mathbf{x}, \mathbf{a}) = \mathbb{E}[r_k | \mathbf{x}, a_k]$. Based on this assumption, IIPS defines the *position-wise* importance weight as follows.

$$\hat{V}_k^{\text{IIPS}}(\pi; \mathcal{D}) := \frac{1}{n} \sum_{i=1}^n \frac{\pi(a_{i,k} | \mathbf{x}_i)}{\pi_0(a_{i,k} | \mathbf{x}_i)} r_{i,k}.$$

where $\pi(a_k | \mathbf{x}) := \sum_{\mathbf{a}'} \pi(\mathbf{a}' | \mathbf{x}) \mathbb{I}\{a'_k = a_k\}$ is the marginal action choice probability at position k under policy π . IIPS substantially reduces the variance of IPS while remaining unbiased under the independence assumption. However, since the independence assumption is overly restrictive to describe real user behavior, IIPS often suffers from severe bias [19, 24].

Reward interaction IPS. RIPS leverages a weaker assumption called the cascade assumption, which assumes that a user interacts with the actions in a ranking sequentially from top to bottom [12]. Hence, the reward observed at each position (r_k) is influenced only by the actions observed at higher positions ($\mathbf{a}_{1:k}$). Since the cascade

Table 1: Correspondence among user behavior assumptions, estimators, and relevant set of actions.

assumption	estimator	relevant actions $\Phi_k(\mathbf{a}, \mathbf{c})$
no assumption	IPS	\mathbf{a}
cascade	RIPS	$\mathbf{a}_{1:k}$
independence	IIPS	a_k
adaptive	AIPS (ours)	$\Phi_k(\mathbf{a}, \mathbf{c}), \mathbf{c} \sim p(\cdot \mathbf{x})$

model assumes that r_k is independent of lower positions, it is sufficient to condition on $\mathbf{a}_{1:k}$ to identify the position-wise expected reward, i.e., $q_k(\mathbf{x}, \mathbf{a}) = \mathbb{E}[r_k | \mathbf{x}, \mathbf{a}_{1:k}]$. Based on this assumption, RIPS applies the *top-k* importance weight as

$$\hat{V}_k^{\text{RIPS}}(\pi; \mathcal{D}) := \frac{1}{n} \sum_{i=1}^n \frac{\pi(\mathbf{a}_{i,1:k} | \mathbf{x}_i)}{\pi_0(\mathbf{a}_{i,1:k} | \mathbf{x}_i)} r_{i,k}.$$

where $\mathbf{a}_{i,k_1:k_2} := (a_{i,k_1}, a_{i,k_1+1}, \dots, a_{i,k_2})$. RIPS is unbiased under the cascade assumption, while reducing the variance of IPS [24]. However, when the cascade assumption does not hold true, it may produce a large bias. Furthermore, RIPS can suffer from high variance when the ranking size is large [19].

Limitation of the existing estimators. We have so far seen that existing estimators have tried to control the bias-variance trade-off by leveraging some assumptions on user behavior – a stronger assumption reduces the variance more but introduces a larger bias. Although this approach has shown some success, a critical limitation is that existing estimators apply a single universal assumption to the entire population, while real user behavior can often be much more diverse and heterogeneous [23, 41–43]. In such realistic scenarios, imposing a single assumption can result in highly sub-optimal estimations. For example, a strong assumption (e.g., independence) produces a large bias in a subpopulation following more complex behavior models, while a weak assumption (e.g., cascade) produces unnecessary variance in another subpopulation following simpler behavior models. This limitation motivates the development of a new estimator that can better exploit the potentially diverse and heterogeneous user behavior to substantially improve OPE of ranking policies.

3 THE ADAPTIVE IPS ESTIMATOR

Our key idea in deriving a new estimator is to take into account various user behaviors by refining the typical formulation of OPE of ranking systems. More specifically, here we introduce an *action-reward interaction* matrix denoted by $\mathbf{c} \in \{0, 1\}^{K \times K}$ whose (k, l) element $(c_{k,l})$ indicates whether r_k is affected by a_l . Given \mathbf{c} , the position-wise expected reward can be expressed as follows.

$$q_k(\mathbf{x}, \mathbf{a}, \mathbf{c}) = \mathbb{E}[r_k | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c})],$$

where $\Phi_k(\mathbf{a}, \mathbf{c}) := \{a_l \in \mathcal{A} | c_{k,l} = 1\}$ is a set of *relevant* or *sufficient* actions needed to identify the expected reward function at the k -th position. The elements of the matrix are considered to be sampled from some unknown probability distribution $p(\mathbf{c} | \mathbf{x})$, which is

conditioned on the context \mathbf{x} to capture potentially diverse and heterogeneous user behavior. Table 1 describes how our formulation generalizes the assumptions used by existing estimators.

Leveraging the action-reward interaction matrix, our *Adaptive IPS* (AIPS) estimator is defined as

$$\hat{V}_k^{\text{AIPS}}(\pi; \mathcal{D}) := \frac{1}{n} \sum_{i=1}^n \frac{\pi(\Phi_k(\mathbf{a}_i, \mathbf{c}_i) | \mathbf{x}_i)}{\pi_0(\Phi_k(\mathbf{a}_i, \mathbf{c}_i) | \mathbf{x}_i)} r_{i,k}.$$

At a high level, AIPS applies *adaptive* importance weighting based on the context-aware behavior model \mathbf{c} . Specifically, when estimating the position-wise policy value at the k -th position, AIPS considers only the actions that affect the reward observed at that particular position ($\Phi_k(\mathbf{a}_i, \mathbf{c}_i)$) to define the importance weight. In this way, AIPS is able not only to deal with potential bias due to diverse user behavior but also to avoid producing unnecessary variance. The following sections show how AIPS enables a much more effective OPE of ranking policies compared to existing estimators.

3.1 Theoretical Analysis

This section provides some key statistical properties of AIPS assuming that the user behavior model \mathbf{c} is observable. Then, we analyze the bias of AIPS when using an estimated user behavior model $\hat{\mathbf{c}}$. Finally, we present an algorithm to *optimize* the behavior model in a way that minimizes the MSE of the resulting estimator. Note that all proofs omitted from the main text are provided in Appendix B.

First, we show that AIPS can be unbiased under any (context-dependent) distribution of user behavior $p(\mathbf{c} | \mathbf{x})$.

PROPOSITION 3.1. *If the user behavior model \mathbf{c} is observed, AIPS is unbiased, i.e., $\mathbb{E}_{\mathcal{D}}[\hat{V}_k^{\text{AIPS}}(\pi; \mathcal{D})] = V_k(\pi)$ for any π and $p(\mathbf{c} | \mathbf{x})$.*

Proposition 3.1 suggests that AIPS is unbiased in a far more general situation about user behavior compared to that of existing work, which relies on a particular behavior model. Next, we show that the variance reduction of AIPS from IPS can be substantial.

THEOREM 3.2. *(Variance Reduction of AIPS over IPS) Compared to IPS, AIPS reduces the variance by the following amount.*

$$\begin{aligned} & n \left(\mathbb{V}_{\mathcal{D}}(\hat{V}_k^{\text{IPS}}(\pi; \mathcal{D})) - \mathbb{V}_{\mathcal{D}}(\hat{V}_k^{\text{AIPS}}(\pi; \mathcal{D})) \right) \\ &= \mathbb{E} \left[\left(\frac{\pi(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})} \right)^2 \mathbb{V}_{\Phi_k^c(\mathbf{a}, \mathbf{c})} \left(\frac{\pi(\Phi_k^c(\mathbf{a}, \mathbf{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))}{\pi_0(\Phi_k^c(\mathbf{a}, \mathbf{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))} \right) \right. \\ & \quad \left. \cdot \mathbb{E}[r_k^2 | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c})] \right], \end{aligned}$$

where the outer expectation is taken over $p(\mathbf{x})p(\mathbf{c}|\mathbf{x})\pi_0(\Phi_k(\mathbf{a}, \mathbf{c})|\mathbf{x})$, and $\Phi_k^c(\mathbf{a}, \mathbf{c})$ is the complement of $\Phi_k(\mathbf{a}, \mathbf{c})$.

Theorem 3.2 ensures that AIPS always provides a non-negative variance reduction over IPS. Moreover, Theorem 3.2 suggests that variance reduction becomes substantial when \mathbf{c} is sparse and the importance weight about irrelevant actions ($\Phi_k^c(\mathbf{a}, \mathbf{c})$) is large.

The following also shows that AIPS achieves the minimum variance among all IPS-based unbiased estimators.

THEOREM 3.3. *(Variance Optimality of AIPS) Let*

$$\hat{V}_k(\pi; \mathcal{D}, \tilde{\mathbf{c}}) := \frac{1}{n} \sum_{i=1}^n \frac{\pi(\Phi_k(\mathbf{a}_i, \tilde{\mathbf{c}}) | \mathbf{x}_i)}{\pi_0(\Phi_k(\mathbf{a}_i, \tilde{\mathbf{c}}) | \mathbf{x}_i)} r_{i,k},$$

Table 2: A toy example illustrating the possible benefit of strategic variance reduction with an incorrect behavior model. AIPS with an incorrect (but optimized) behavior model produces much smaller variance while introducing some small bias, resulting in a smaller MSE than AIPS with the true behavior model.

	bias	variance	MSE (= bias ² + variance)
AIPS with the true behavior model \mathbf{c}	0.0	0.5	0.50 (= (0.0) ² + 0.5)
AIPS with an incorrect (but optimized) behavior model $\hat{\mathbf{c}}$	0.1	0.3	0.31 (= (0.1) ² + 0.3)

so that $\mathbb{E}[\hat{V}_k(\pi; \mathcal{D}, \hat{\mathbf{c}})] = V_k(\pi)$. Then, for any $\hat{\mathbf{c}}$ (s.t. $\mathbf{c} \subseteq \hat{\mathbf{c}}$) and π , we have

$$\mathbb{V}_{\mathcal{D}}(\hat{V}_k^{\text{AIPS}}(\pi; \mathcal{D})) \leq \mathbb{V}_{\mathcal{D}}(\hat{V}_k(\pi; \mathcal{D}, \hat{\mathbf{c}})).$$

Theorem 3.3 ensures that AIPS guarantees the minimum variance among all unbiased IPS estimators under any distribution of user behavior, suggesting that AIPS is the optimal unbiased estimator.¹

Although we have shown above that AIPS can exhibit favorable statistical properties under general user behavior, it should be noted that we currently assume that the true user behavior \mathbf{c} is observable. Since this is generally not the case, the following investigates the bias of AIPS when given is an estimated user behavior $\hat{\mathbf{c}}$.

THEOREM 3.4. (Bias of AIPS with an estimated user behavior) When an estimated user behavior $\hat{\mathbf{c}}$ is used, AIPS has the following bias.

$\text{Bias}(\hat{V}_k^{\text{AIPS}}; \hat{\mathbf{c}}) = \mathbb{E}_{p(\mathbf{x})p(\mathbf{c}|\mathbf{x})\pi(\mathbf{a}|\mathbf{x})} [(\Delta w_k(\mathbf{a}, \mathbf{c}, \hat{\mathbf{c}}) - 1) q_k(\mathbf{x}, \mathbf{a}, \mathbf{c})]$, where

$$\Delta w_k(\mathbf{a}, \mathbf{c}, \hat{\mathbf{c}}) := \frac{\pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) \setminus \Phi_k(\mathbf{a}, \hat{\mathbf{c}}) | \mathbf{x}, \Phi_k(\mathbf{a}, \hat{\mathbf{c}}))}{\pi(\Phi_k(\mathbf{a}, \mathbf{c}) \setminus \Phi_k(\mathbf{a}, \hat{\mathbf{c}}) | \mathbf{x}, \Phi_k(\mathbf{a}, \hat{\mathbf{c}}))}.$$

Theorem 3.4 suggests that AIPS remains unbiased when the true model is a subset of the estimated model ($\mathbf{c} \subseteq \hat{\mathbf{c}}$). Furthermore, we can see that the bias of AIPS is characterized by the overlap between the true \mathbf{c} and estimated user behavior $\hat{\mathbf{c}}$, i.e., when there is a large overlap between \mathbf{c} and $\hat{\mathbf{c}}$, the cardinality of $\Phi_k(\mathbf{a}, \mathbf{c}) \setminus \Phi_k(\mathbf{a}, \hat{\mathbf{c}})$ becomes small, resulting in a smaller bias for AIPS.

Controlling the bias-variance tradeoff. Theorems 3.2 and 3.4 suggest that the bias-variance tradeoff of AIPS is mainly characterized by $\hat{\mathbf{c}}$.² When $\hat{\mathbf{c}}$ is dense, the bias of the resulting AIPS estimator will be very small, but the variance can be high. On the other hand, a sparse $\hat{\mathbf{c}}$ can substantially reduce the variance of AIPS while introducing some bias. This suggests that true user behavior \mathbf{c} may not necessarily minimize the MSE of AIPS, and that there exists an interesting strategy to intentionally utilize an incorrect model $\hat{\mathbf{c}}$ to further improve the accuracy of the downstream estimation. Table 2 provides a toy example illustrating a situation where AIPS with an incorrect behavior model can achieve a lower MSE than that with the true behavior model. In this example, the minimum variance among all unbiased estimators is 0.5, which is achieved by the true behavior model as per Theorem 3.3. However, a lower MSE can be realized by intentionally using an incorrect (overly sparse) model. This is because we can gain a large variance reduction (-0.2)

¹Note that AIPS may not be optimal if we also take some *biased* estimators into consideration. This motivates our idea of intentionally leveraging an incorrect behavior model to further improve the MSE, which is the sum of the squared bias and variance.

²Note that the variance of AIPS with an estimated behavior model can immediately be obtained by replacing \mathbf{c} with $\hat{\mathbf{c}}$ in Theorem 3.2.

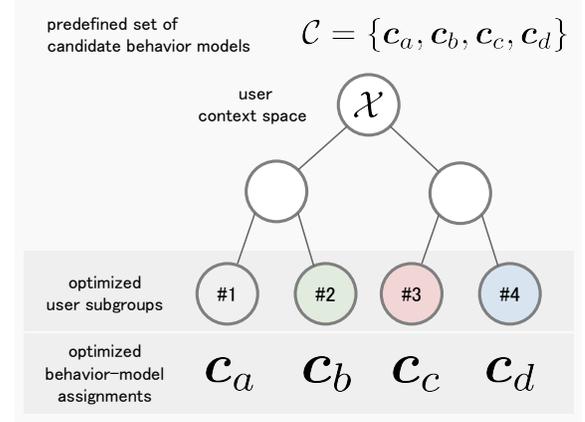


Figure 1: Tree-based optimization of user behavior model, which optimizes the partition in the context space and assignments of the user behavior model of each subgroup (from #1 to #4) so the resulting MSE of AIPS is minimized.

by allowing only a small squared bias (+0.01), and hence using the true user behavior model does not result in the optimal MSE of AIPS. Therefore, instead of discussing how to *estimate* the true user behavior, the following section describes a data-driven approach to *optimize* $\hat{\mathbf{c}}$ in a way that minimizes the MSE of AIPS.

3.2 Optimizing the User Behavior Model

Our goal here is to *optimize* the context-aware user behavior model $\hat{\mathbf{c}}(\mathbf{x})$ to minimize the MSE of AIPS, rather than merely identifying the true model. We achieve this by optimizing the behavior model at a *subgroup level*, inspired by the subgroup identification techniques proposed in treatment effect estimation [1, 17].

Due to its flexibility in handling even non-differentiable objectives, we employ a non-parametric tree-based model to optimize the partition in the context space (user subgroups) and the assignments of user behavior models to each subgroup simultaneously. Specifically, we recursively partition the feature space (\mathcal{X}) and assign an appropriate behavior model to each node in a way that minimizes the MSE of AIPS. This means that we follow the classification and regression tree (CART) algorithm [2] and set the MSE of AIPS as its objective function. More specifically, we first define a candidate set of user behavior models $C := \{\hat{\mathbf{c}}_0, \dots, \hat{\mathbf{c}}_m\}$.³ Then, at each parent

³In general, set of candidate models C should be as large as possible, but a larger C might be infeasible due to intensive computation. In experiments, we show that AIPS remains effective even when C is not large if we include a diverse set of models so that the tree model can find an appropriate behavior model for each user subgroup.

node $l \subseteq \mathcal{X}$, the tree partitions it into child nodes $(l_{(l^*)}, l_{(r^*)}) \subseteq l$ where $l_{(l^*)} \cup l_{(r^*)} = l$ and $l_{(l^*)} \cap l_{(r^*)} = \emptyset$ and assigns behavior models $(\hat{c}_{(l^*)}, \hat{c}_{(r^*)}) \in \mathcal{C}$ to these nodes by the following criterion.

$$(l_{(l^*)}, l_{(r^*)}, \hat{c}_{(l^*)}, \hat{c}_{(r^*)}) := \arg \min_{(l_{(l^*)}, l_{(r^*)}, \hat{c}_{(l^*)}, \hat{c}_{(r^*)})} \widehat{MSE}(\hat{c}_{(l^*)}, \hat{c}_{(r^*)}; l_{(l^*)}, l_{(r^*)}) \quad (2)$$

where $\hat{c}_{(l)} \in \mathcal{C}$ is a candidate behavior model assigned to node l , and $\widehat{MSE}(\hat{c}_{(l^*)}, \hat{c}_{(r^*)}; l_{(l^*)}, l_{(r^*)})$ is an estimated MSE when $\hat{c}_{(l^*)}, \hat{c}_{(r^*)} \in \mathcal{C}$ are assigned to the left and right nodes, respectively. Algorithm 1 in the appendix provides the complete optimization procedure.

Compared to existing subgroup identification procedures [1, 17], our algorithm is unique in that it directly optimizes the MSE in OPE rather than minimizing some prediction loss for treatment effect estimation. It is important to note that our algorithm is agnostic to the method used for estimating the MSE. For example, we can estimate the MSE by following existing methods from [35] or [37]. We thus consider the MSE estimation task as an independent research topic and do not propose specific approaches to estimate the MSE from the logged data. Instead, our experiments will demonstrate that AIPS with our data-driven procedure for behavior model optimization performs reasonably well across a variety of experiment settings, even with a noisy MSE estimate and with an MSE estimated via an existing method from Udagawa et al. [37] that uses only the observed logged data.

4 SYNTHETIC EXPERIMENTS

This section empirically compares the proposed estimator with existing estimators (IPS, IIPS, and RIPS) on synthetic ranking data. Our experiment is implemented on top of *OpenBanditPipeline* [27]⁴, a modular Python package for OPE. Our experiment code is available at <https://github.com/aiueola/kdd2023-aips>. Other experiment details and additional results are provided in Appendix A.

4.1 Setup

Basic setting. To generate synthetic data, we randomly sample five-dimensional context ($d = 5$) from the standard normal distribution. Then, for each position k , we sample continuous rewards from a normal distribution as $r_k \sim \mathcal{N}(q_k(\mathbf{x}, \mathbf{a}, \mathbf{c}), \sigma^2)$, where we use $\sigma = 0.5$. The following describes how to define the expected reward function $q_k(\mathbf{x}, \mathbf{a}, \mathbf{c})$ and user behavior distribution $p(\mathbf{c} | \mathbf{x})$.

Position-wise expected reward function. Following Kiyohara et al. [19], we first define the following position-wise **base** reward function $\tilde{q}_k(\mathbf{x}, a_k)$, which depends only on the action presented at the corresponding position (a_k) rather than the entire ranking.

$$\tilde{q}_k(\mathbf{x}, a_k) = \theta_{a_k}^\top \mathbf{x} + b_{a_k},$$

where θ_{a_k} is a parameter vector sampled from the standard normal distribution, and b_{a_k} is a bias term that corresponds to action a_k .

Then, we define the position-wise expected reward function given a particular user behavior model \mathbf{c} as follows.

$$q_k(\mathbf{x}, \mathbf{a}, \mathbf{c}) = c_{k,k} \tilde{q}_k(\mathbf{x}, a_k) + \sum_{l \neq k} c_{k,l} \mathbb{W}(a_k, a_l)$$

where $c_{k,l} \in \{0, 1\}$ is the (k, l) element of \mathbf{c} , which indicates whether a_l affects r_k . \mathbb{W} is a $|\mathcal{A}| \times |\mathcal{A}|$ matrix whose elements are sampled from a uniform distribution with range $[0, 1]$. This matrix defines how the co-occurrence of a pair of actions affects $q_k(\mathbf{x}, \mathbf{a}, \mathbf{c})$.

Distribution of user behavior. Next, the following defines the three basic user behavior models used in existing work [19, 24].

- **standard (S):** $c_S(k, l) = 1, \forall l \in [K]$.
- **cascade (C):** $c_C(k, l) = 1, \forall l \leq k$, and $c_C(k, l) = 0$, otherwise.
- **independence (I):** $c_I(k, k) = 1$ and $c_I(k, l) = 0$, otherwise.

for each $k \in [K]$. To introduce more diverse behaviors beyond the above basic models, we define the following *h-neighbor perturbation*:

$$c_{neighbor,h}(k, l) = 1, \forall |l - k| \leq h, \\ \text{and } c_{neighbor,h}(k, l) = 0, \text{ otherwise.}$$

where h is the number of neighboring items that perturb the basic model. By applying this perturbation to the basic models, we define the following more complex behavior models.

- **C1:** $c_{C1}(k, l) = c_C(k, l) + c_{neighbor,1}(k, l)$
- **C2:** $c_{C2}(k, l) = c_C(k, l) + c_{neighbor,2}(k, l)$
- **I1:** $c_{I1}(k, l) = c_I(k, l) + c_{neighbor,1}(k, l)$

We also define two additional user behaviors by applying *random perturbation* to the independence model as follows.

- **R3:** $c_{R3}(k, l) = c_I(k, l) + c_{random,3}(k, l)\mathbb{I}\{k \neq l\}$
- **R6:** $c_{R6}(k, l) = c_I(k, l) + c_{random,6}(k, l)\mathbb{I}\{k \neq l\}$

where $c_{random,h}(k, \cdot) = 1$ only for randomly chosen h positions for each $k \in [K]$.

To study how the estimators work under diverse and heterogeneous user behaviors, we use **{S, R6, R3, C2, C1, I1}** and sample them from the following distribution given a user context:

$$p(\mathbf{c}_z | \mathbf{x}) := \text{softmax}(\lambda_z \cdot |\theta_z^\top \mathbf{x}|) = \frac{\exp(\lambda_z \cdot |\theta_z^\top \mathbf{x}|)}{\sum_{z'} \exp(\lambda_{z'} \cdot |\theta_{z'}^\top \mathbf{x}|)},$$

where $z \in \{S, R6, R3, C2, C1, I1\}$ is the index of each user behavior. θ_z is a parameter vector sampled from the standard uniform distribution, and λ_z is some weight parameter. By assigning different values of λ_z to different user behaviors, we can control the distribution of user behavior. In particular, we define λ_z as follows.

$$\lambda_z := \exp((2\delta - 1) \cdot \gamma_z),$$

where γ_z is some coefficient value defined for each user behavior as $\{\gamma_S, \gamma_{R6}, \gamma_{R3}, \gamma_{C2}, \gamma_{C1}, \gamma_{I1}\} = \{1.5, 0.9, 0.3, -0.3, -0.9, -1.5\}$, and $\delta \in [0, 1]$ is an experiment parameter called the ‘‘user behavior distribution parameter’’, which controls the entropy of the behavior distribution. For example, all user behavior will be uniformly distributed when $\delta = 0.5$, as $\gamma_z = 0, \forall z$. In contrast, $\delta < 0.5$ samples user behaviors having negative values of γ_z more frequently, while $\delta > 0.5$ prioritizes those having positive values of γ_z . In particular, under our definition of $\{\gamma_z\}$, a smaller value of δ leads to simpler user behavior, while a larger value leads to more complex behavior in general.⁵

⁴<https://github.com/st-tech/zr-obp>

⁵Figures 6 and 7 in Appendix A show how changes in the value of distribution parameters (δ and λ) control the distribution of user behavior $p(\mathbf{c}_z | \mathbf{x})$.

Logging and evaluation policies. We define a factored logging policy to generate synthetic logged data as follows.

$$\pi_0(\mathbf{a} | \mathbf{x}) = \prod_{k=1}^K \pi_0(a_k | \mathbf{x}) = \prod_{k=1}^K \frac{\exp(f_0(\mathbf{x}, a_k))}{\sum_{a' \in \mathcal{A}} \exp(f_0(\mathbf{x}, a'))},$$

where $f_0(\mathbf{x}, a) = \theta_a^\top \mathbf{x} + b_a$. We sample θ_a and b_a from the standard uniform distribution. Then, we define the evaluation policy by applying the following transformation to the logging policy.

$$\pi(\mathbf{a} | \mathbf{x}) = \prod_{k=1}^K \left((1 - \epsilon) \mathbb{I}\{a_k = \arg \min_{a' \in \mathcal{A}} f_0(\mathbf{x}, a')\} + \epsilon / |\mathcal{A}| \right), \quad (3)$$

where $\epsilon \in [0, 1]$ is an experiment parameter that determines the stochasticity of π . Specifically, a small value of ϵ leads to a near-deterministic policy, while a large value leads to a near-uniform policy. We use $\epsilon = 0.3$ throughout our synthetic experiment.

Compared estimators. We compare AIPS against IPS, RIPS, and IIPS. We also report the results of AIPS (true), which uses the true user behavior \mathbf{c} and thus is infeasible in practice. However, this provides a useful reference to investigate the effectiveness of our strategic variance reduction method from Section 3.2.

Note that AIPS uses the following surrogate MSE as the objective function when performing user behavior optimization.

$$\widehat{MSE}(\hat{V}_k^{\text{AIPS}}(\pi; \mathcal{D}, \tilde{\mathbf{c}})) = \widehat{Bias}(\hat{V}_k^{\text{AIPS}}(\pi; \mathcal{D}, \tilde{\mathbf{c}}))^2 + \hat{V}(\hat{V}_k^{\text{AIPS}}(\pi; \mathcal{D}, \tilde{\mathbf{c}})),$$

where $\hat{V}(\cdot)$ is the sample variance. To control the accuracy of the bias estimation, we use its noisy estimate \widehat{Bias} . Specifically, we first estimate the bias based on an on-policy estimate of the policy value, which is denoted as $\widehat{Bias}_{\text{on}}(\cdot)$, and then add some Gaussian noise as $\widehat{Bias} \sim \mathcal{N}(\widehat{Bias}_{\text{on}}, \sigma_\Delta^2)$ where $\sigma_\Delta = 0.3 \times |\widehat{Bias}_{\text{on}}|$. By doing so, we can simulate a practical situation where AIPS relies on some noisy estimate of MSE. This procedure also enables us to evaluate the robustness of AIPS to the varying accuracies of MSE estimation, as demonstrated in Appendix A.

4.2 Results and Discussion

We run the OPE simulations 1000 times with different random seeds. We report the MSE, bias, and variance of the estimators normalized by the true policy value $V(\pi)$. Note that we use $n = 8k$, $K = 8$, and $\delta = 0.6$ as default experiment parameters.⁶ In all figures, the solid lines indicate the performance metrics averaged over the simulation runs and the shaded regions show their 95% confidence intervals.

RQ (1): How do the estimators perform with varying data sizes? Figure 2 compares the estimators' MSEs (normalized by the true value $V(\pi)$) with varying data sizes. The result clearly suggests that AIPS (ours) achieves the best (lowest) MSE in a range of logged data sizes, while the existing estimators fail drastically in some specific cases. First, we observe that IIPS and RIPS fail to improve their MSE even with increasing logged data sizes. We attribute this to their high bias due to the mismatch between their behavior assumption (independence or cascade) and the true user

behavior (which is diverse and context-dependent in our experiment). Second, we can see that IPS enables an unbiased estimation,⁷ however, it suffers from extreme variance, particularly when the data size is small. This is because IPS often applies unnecessarily large importance weights regardless of the true user behavior.

In contrast, AIPS (true) deals with the bias-variance issues of the existing estimators by leveraging adaptive importance weighting based on prior knowledge about the true user behavior (which is unavailable in practice). Specifically, Figure 2 demonstrates that AIPS (true) is unbiased and thus performs better than IIPS and RIPS when the data size is large where AIPS (true) becomes increasingly accurate with a reduced variance while IIPS and RIPS remain highly biased. Moreover, AIPS (true) has a much lower variance than IPS by applying importance weighting to only the relevant set of actions for each given context. However, it should be noted that the variance of AIPS (true) can still be high, particularly when the data size is extremely small. In particular, AIPS (true) exhibits a worse MSE than RIPS when $n \leq 2k$, which interestingly implies that naive use of true user behavior when performing importance weighting is not optimal in terms of MSE.

Our AIPS estimator performs much better than all existing estimators and even overcomes the limitations of AIPS (true) by *optimizing* the user behavior model rather than merely exploiting the true model. More specifically, AIPS further improves the MSE of AIPS (true) by greatly reducing the variance at the cost of introducing only a small amount of bias. Note here that this is achieved even though we impose some estimation error in the MSE estimation, suggesting that the subgroup optimization procedure from Section 3.2 is robust to the estimation error of the MSE.⁸ These empirical results demonstrate that AIPS is able to adaptively optimize the user behavior model in a way that improves the MSE and thus enables a more reliable OPE in a range of logged data sizes particularly under diverse user behavior and even without the true knowledge of the behavior model.

RQ (2): How do the estimators perform with varying lengths of ranking? Next, we compare the performance of the estimators with varying lengths of ranking (K) in Figure 3. The overall trend and qualitative comparison are similar to the previous arguments made in RQ (1) – AIPS (ours) works stably well across a range of settings, while the existing estimators fail for some specific values of K . Specifically, when $K \geq 10$, IPS and AIPS (true) produce extremely high variance due to excessive importance weights, while IIPS and RIPS produce substantial bias due to their strong assumption about user behavior. In contrast, AIPS (ours) leads to a much better bias-variance tradeoff by *optimizing* the user behavior. In particular, we observe that AIPS (ours) prioritizes reducing bias when $K \leq 8$, while it puts more priority on variance reduction when $K \geq 10$, resulting in its superior performance against existing methods as well as AIPS (true) in a range of ranking sizes.

RQ (3): How do the estimators perform with various user behavior distributions? Figure 4 shows how the accuracy of the

⁷Note that the squared bias of IPS is not exactly zero even though this estimator is always theoretically unbiased. This is due to the fact that we estimate the squared bias based on the simulation results where there is some small variance.

⁸We also observe the similar results and superior behavior of AIPS with varying amounts of noise on the MSE estimate, which is reported in Appendix A.

⁶We use 1k, 2k, ... to denote 1000, 2000, ...

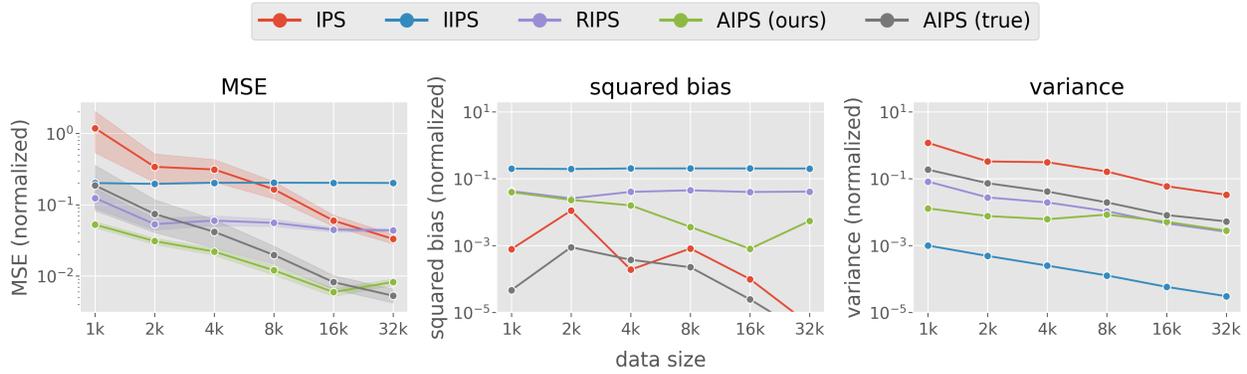


Figure 2: Comparison of the estimators' MSE (normalized by the true value $V(\pi)$) with varying data sizes (n)

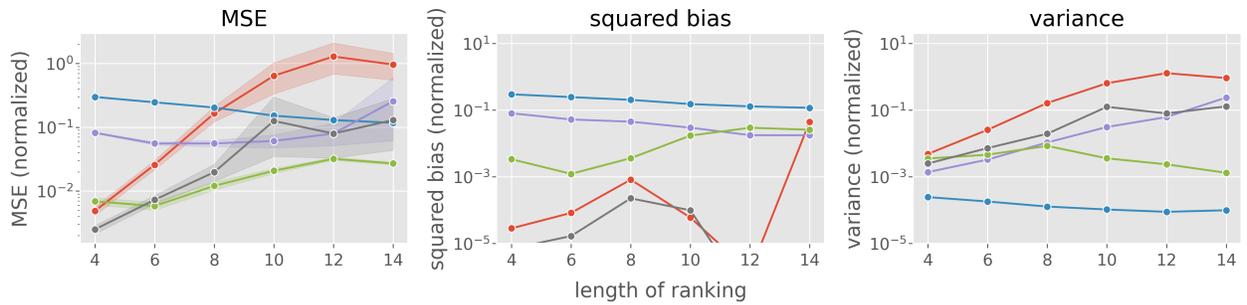


Figure 3: Comparison of the estimators' MSE (normalized by the true value $V(\pi)$) with varying lengths of ranking (K)

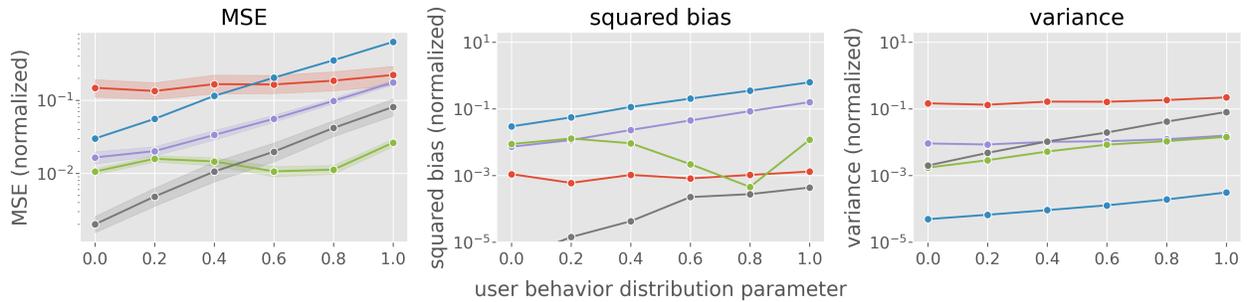


Figure 4: Comparison of the estimators' MSE (normalized by the true value $V(\pi)$) with varying behavior distributions (δ)

estimators changes as the user behavior distribution shifts from a simple user behavior ($\delta = 0.0$) to a more realistic, complex one ($\delta = 1.0$). First, the result demonstrates that IPS produces inaccurate OPE across various behavior distributions, as its variance is consistently high. In contrast, IIPS and RIPS are accurate when user behavior is simple ($\delta = 0.0$). However, as user behavior gradually becomes more complex, IIPS and RIPS produce larger bias because their assumptions become increasingly incorrect. Similarly, AIPS (true) enables an accurate estimation, particularly when the user behavior is simple, but its MSE gradually becomes worse as the user behavior becomes more complex. Specifically, AIPS (true) is accurate when $\delta \leq 0.4$ due to its optimal variance, however, it suffers from extremely high variance due to large importance weights

and shows substantial accuracy deterioration in the presence of complex user behaviors. Finally, we observe that AIPS (ours) consistently achieves a much more accurate estimation compared to IPS, IIPS, and RIPS across various behavior distributions, particularly under the challenging cases of complex user behaviors ($\delta \geq 0.6$). Moreover, AIPS (ours) is even better than AIPS (true) when $\delta \geq 0.6$, because AIPS (ours) optimizes the user behavior model and thus avoids the excessive variance of AIPS (true) as long as this strategic variance reduction does not introduce considerable bias. In the case of simple behaviors ($\delta \leq 0.4$), it becomes more important to reduce the bias by leveraging the true behavior model, and thus AIPS (true) performs the best in these cases. However, the results clearly demonstrate the benefit of AIPS against existing methods

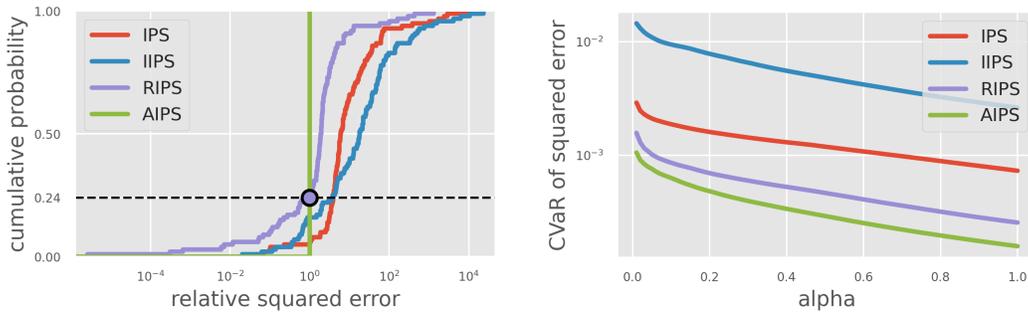


Figure 5: Estimators’ performance comparison in the real-world experiment: (Left) Cumulative distribution function (CDF) of the estimators’ squared error (relative to that of AIPS). (Right) Conditional Value at Risk (CVaR) of the estimators’ squared error with varying values of α .

(IPS, IIPS, RIPS) and that of the idea of behavior model optimization in practical situations where the user behavior is highly complex.

5 REAL-WORLD EXPERIMENT

This section demonstrates the effectiveness of AIPS using the logged data collected on a real-world ranking system.

Setup. To evaluate and compare the estimators in a more practical situation, we collect some logged bandit data by running an A/B test of two (factored) ranking policies π_A and π_B on a real e-commerce platform whose aim is to optimize a ranking of modules (which showcase a set of products inside) to maximize the number of clicks. Our A/B test produces two sets of logged data \mathcal{D}_A and \mathcal{D}_B where $|\mathcal{D}_A| = 1,979$ and $|\mathcal{D}_B| = 1,954$. Note that, in this application, \mathbf{x} is a five-dimensional user context, \mathbf{a} is a ranking of modules where $|\mathcal{A}| = 2$, $K = 6$, and r_k is a binary click indicator.

To perform an OPE experiment, we regard π_A as a logging policy and π_B as an evaluation policy. We use \mathcal{D}_A to estimate the value of π_B by estimators, while we use \mathcal{D}_B to approximate the ground truth value of π_B by on-policy estimation. Then, we calculate the squared error (SE) of an estimator as $\text{SE}(\hat{V}) := (\hat{V}_{\text{on}}(\pi_B; \mathcal{D}_B) - \hat{V}(\pi_B; \mathcal{D}_A))^2$. We run the experiment 100 times using different bootstrapped samples of \mathcal{D}_A and report the cumulative distribution function (CDF) of SE of IPS, IIPS, RIPS, and AIPS, relative to that of AIPS. To evaluate the worst case performance of the estimators, we also report the conditional value at risk (CVaR) of SE, which measures the average performance of the worst $\alpha \times 100\%$ trials for each estimator. AIPS uses $C = \{c_S, c_C, c_I, c', c''\}$ as the candidate set of behavior models where c_S , c_C , and c_I are defined in Section 4, and c' and c'' are defined in Appendix A. When performing user behavior optimization of AIPS, its MSE is estimated by PAS-IF [37] using only the observable logged data. Note that we cannot implement AIPS (true) in this section, since we do not know the true user behavior in the real-world dataset.

Result. Figure 5 (Left) shows the estimators’ CDF of relative SE and demonstrates that AIPS performs the best in 76% of the trials. Moreover, in Figure 5 (Right), we observe that AIPS improves the CVaR of SE more than 30% compared to RIPS for a range of α . These results suggest that AIPS enables a more accurate and stable OPE than previous estimators in the real-world situation.

Summary of empirical findings. In summary, AIPS achieves far more accurate OPE than all existing estimators (IPS, IIPS, and RIPS) in both synthetic and real-world experiments via leveraging adaptive importance weighting with an optimized user behavior model. Specifically, AIPS has a much lower bias than IIPS and RIPS by identifying more appropriate behavior models that have a sufficient overlap with the true user behavior. Moreover, AIPS substantially reduces the variance of IPS by avoiding unnecessarily large importance weights. As a result, AIPS shows a superior performance particularly in realistic situations where the ranking size is large and user behavior is diverse and complex.⁹ Moreover, we observe that AIPS performs even better than AIPS (true) in many cases, implying that strategically leveraging an incorrect behavior model can lead to a better MSE. We thus conclude that AIPS enables a more effective OPE of ranking systems and that we should consider optimizing the behavior model to improve the MSE rather than being overly sensitive to its correct estimation.

6 RELATED WORK

Off-Policy Evaluation. OPE is of great practical relevance in search and recommender systems, as it enables the performance evaluation of counterfactual policies without interacting with the actual users [18, 20, 28, 32]. In particular, OPE in the single action setting has been studied extensively, producing many estimators with good theoretical guarantees [11, 16, 21, 29]. Among them, IPS is often considered a benchmark estimator [25], which uses the importance sampling technique to correct the distribution shift between different policies. IPS is unbiased under some identification assumptions such as full support and unconfoundedness, but it often suffers from high variance [9]. Doubly Robust (DR) [9] reduces the variance of IPS by using an estimated reward function as a control variate. However, DR can still struggle with high variance when the action space is extremely large [29] such as in the ranking setup.

Beyond the standard OPE, there has also been a growing interest in OPE of ranking systems due to its much practical relevance. In the ranking setting where a policy chooses a ranked list of items to present to the users, OPE faces the critical variance issue due to

⁹Appendix A provides additional experiment results demonstrating that AIPS is more robust to reward noise and changes in user behavior distribution compared to baseline estimators (IPS, IIPS, and RIPS).

combinatorial action spaces. To address this variance issue, existing work has introduced some assumptions about user behavior. In particular, IIPS [22] assumes that a user interacts with the actions independently across positions. Under this assumption, the reward observed at each position depends only on the action presented at the same position, leading to a significant variance reduction compared to IPS. Although IIPS is unbiased when the independence assumption holds true, it can have a large bias when users follow a more complicated behavior [19, 24]. RIPS [24] assumes a more reasonable assumption, called the cascade assumption, which requires that a user interacts with the actions sequentially from top to bottom [12]. Therefore, the reward observed at each position depends only on the actions presented at higher positions. Leveraging the cascade assumption, RIPS can somewhat reduce the variance of IPS while being unbiased in more realistic cases compared to IIPS. To further improve the variance of RIPS, Kiyohara et al. [19] propose the Cascade-DR estimator, leveraging the recursive structure of the cascade assumption and a control variate. Although the above approach has shown some empirical success, the critical issue is that all the above estimators rely on a single assumption (independence or cascade) applied to every user, which can cause large bias and unnecessary variance. Therefore, we were based on a more general formulation by assuming that user behavior is sampled from some unknown context-dependent distribution. As a result, AIPS provides an unbiased estimation even under arbitrarily diverse user behavior and achieves the minimum variance among the class of IPS estimators that are unbiased. Moreover, we developed a method to *optimize* the user behavior model rather than accurately estimating it given the theoretical observations that the true user behavior is not optimal in terms of MSE.

Note that there is another estimator called the Pseudo Inverse (PI) estimator [35, 36, 38] in the slate recommendation setting. This estimator considers a situation where only the slate-wise reward is observed (i.e., the position-wise rewards are unobservable). Since PI is not able to leverage position-wise rewards, it is often highly sub-optimal in our setup where position-wise rewards are observable, as empirically verified in McInerney et al. [24].

Click Models. The click models aim to formulate how users interact with a list of documents [3, 5, 10, 12, 15, 26, 30, 33, 39], and it has typically been studied based on the following *examination hypothesis*: $p(c_k = 1 | \mathbf{x}, \mathbf{a}) = p(o_k = 1 | \mathbf{x}, \mathbf{a}) \cdot p(r_k = 1 | \mathbf{x}, a_k)$, where c_k is a click indicator while r_k is a relevance indicator of the document presented at the k -th position. $p(o_k | \mathbf{x}, \mathbf{a})$ is the probability that a user examines the k -th document in a ranking. When the user examines the document (i.e., $o_k = 1$), the click probability is assumed identical to the probability of relevance. Much research has been done to better parameterize the examination probability to explain finer details of the real-world examination behavior. For example, the Position-based model assumes that the examination probability depends only on the position in a ranking, while the Cascade model [6, 12] assumes that the examination probability at the k -th position depends on the relevance of the documents shown at higher positions.

In contrast, the user behavior models utilized in OPE focus more on modeling the dependencies among actions and rewards rather than modeling the examination probability [19, 24]. As already

discussed, the critical drawback of the previous methods is that only a single assumption is assumed to model every user’s behavior. In the information retrieval literature, Chen et al. [4] considered context-dependent click models, which assume that the examination behavior may change depending on the search query. Moreover, several studies have indicated the need to incorporate some context information in building and estimating click models such as devices [23], user browsing history [7], and user intention [13]. In this work, we deal with potentially diverse user behaviors by formulating them via a context-dependent distribution for the first time in OPE of ranking policies. Note, however, that our motivation is substantially different from that of the click modeling literature. That is, we aim to develop an accurate OPE estimator in terms of MSE while click modeling aims to estimate the true user behavior as accurately as possible. This difference motivates our unique strategy to intentionally rely on an incorrect behavior model to further improve the MSE of our estimator as discussed in Section 3.2.

7 CONCLUSION AND FUTURE WORK

This paper studied OPE of ranking systems under diverse user behavior. When the user behavior is diverse and depends on the user context, all existing estimators can be highly sub-optimal because they apply a single assumption to the entire population. To achieve an effective OPE even under much more diverse user behavior, we propose the *Adaptive IPS* estimator based on a new formulation where the user behavior is assumed to be sampled from a *context-dependent* distribution. We began by theoretically characterizing the bias and variance of AIPS assuming known user behaviors, showing that it can be unbiased under any distribution of user behavior and that it achieves the optimal variance among unbiased IPS estimators. Interestingly, though, our analysis also indicates that myopically using the true user behavior in OPE might not be optimal in terms of MSE. Therefore, we provided a data-driven procedure to *optimize* the user behavior model to minimize the MSE of the resulting AIPS estimator rather than trying to *estimate* the true behavior, which tends to be sub-optimal in OPE. Experiments demonstrate that AIPS provides a substantial gain in MSE against existing methods in a range of OPE situations.

Our work also raises several intriguing research questions for future studies. First, it would be valuable to develop an accurate way to estimate the MSE of an OPE estimator beyond existing methods [35, 37] to better optimize the user behavior model to further improve AIPS. Second, OPE of ranking policies can still become extremely difficult when the number of unique actions ($|\mathcal{A}|$) is large. Therefore, it would be interesting to leverage the recent action embedding approach [29, 31] to overcome this critical limitation in the ranking setup. Besides, as a practical, yet simple extension, adding a control variate to AIPS is expected to further improve its variance and outperform Cascade-DR of Kiyohara et al. [19], which assumes the cascade assumption. Finally, this work only studied the statistical problem of estimating the value of a fixed new policy, so it would be interesting to use our estimator to enable more efficient off-policy learning in ranking systems.

REFERENCES

- [1] Susan Athey and Guido Imbens. 2016. Recursive Partitioning for Heterogeneous Causal Effects. *Proceedings of the National Academy of Sciences* 113, 27 (2016),

- 7353–7360.
- [2] Leo Breiman, Jerome H Friedman, Richard A Olshen, and Charles J Stone. 2017. *Classification and regression trees*. Routledge.
 - [3] Olivier Chapelle and Ya Zhang. 2009. A Dynamic Bayesian Network Click Model for Web Search Ranking. In *Proceedings of the 18th International Conference on World Wide Web*. 1–10.
 - [4] Jia Chen, Jiaxin Mao, Yiqun Liu, Min Zhang, and Shaoping Ma. 2020. A Context-Aware Click Model for Web Search. In *Proceedings of the 13th International Conference on Web Search and Data Mining*. 88–96.
 - [5] Aleksandr Chuklin, Ilya Markov, and Maarten de Rijke. 2015. Click Models for Web Search. *Synthesis lectures on information concepts, retrieval, and services* 7, 3 (2015), 1–115.
 - [6] Nick Craswell, Onno Zoeter, Michael Taylor, and Bill Ramsey. 2008. An Experimental Comparison of Click Position-Bias Models. In *Proceedings of the 2008 international conference on web search and data mining*. 87–94.
 - [7] Lifang Deng, Jin Niu, Angulia Yang, Qidi Xu, Xiang Fu, Jiandong Zhang, and Anxiang Zeng. 2020. Hybrid Interest Modeling for Long-tailed Users. *arXiv preprint arXiv:2012.14770* (2020).
 - [8] Maria Dimakopoulou, Nikos Vlassis, and Tony Jebara. 2019. Marginal Posterior Sampling for Slate Bandits. In *IJCAI*. 2223–2229.
 - [9] Miroslav Dudík, John Langford, and Lihong Li. 2011. Doubly Robust Policy Evaluation and Learning. In *Proceedings of the 28th International Conference on International Conference on Machine Learning* (Bellevue, Washington, USA) (ICML'11). Omnipress, Madison, WI, USA, 1097–1104.
 - [10] Georges E Dupret and Benjamin Piwowarski. 2008. A User Browsing Model to Predict Search Engine Click Data from Past Observations. In *Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. 331–338.
 - [11] Alexandre Gilotte, Clément Calauzènes, Thomas Nedelec, Alexandre Abraham, and Simon Dollé. 2018. Offline A/B Testing for Recommender Systems. In *Proceedings of the 11th ACM International Conference on Web Search and Data Mining*. 198–206.
 - [12] Fan Guo, Chao Liu, and Yi Min Wang. 2009. Efficient Multiple-Click Models in Web Search. In *Proceedings of the 2nd ACM International Conference on Web Search and Data Mining*. 124–131.
 - [13] Botao Hu, Yuchen Zhang, Weizhu Chen, Gang Wang, and Qiang Yang. 2011. Characterizing Search Intent Diversity into Click Models. In *Proceedings of the 20th International Conference on World Wide Web*. 17–26.
 - [14] Kalervo Järvelin and Jaana Kekäläinen. 2002. Cumulated Gain-Based Evaluation of IR Techniques. *ACM Transactions on Information Systems (TOIS)* 20, 4 (2002), 422–446.
 - [15] Thorsten Joachims, Adith Swaminathan, and Tobias Schnabel. 2017. Unbiased Learning-to-Rank with Biased Feedback. In *Proceedings of the 10th ACM International Conference on Web Search and Data Mining*. 781–789.
 - [16] Nathan Kallus, Yuta Saito, and Masatoshi Uehara. 2021. Optimal Off-Policy Evaluation from Multiple Logging Policies. In *Proceedings of the 38th International Conference on Machine Learning*, Vol. 139. PMLR, 5247–5256.
 - [17] Ramtin Keramati, Omer Gottesman, Leo Anthony Celi, Finale Doshi-Velez, and Emma Brunskill. 2022. Identification of Subgroups With Similar Benefits in Off-Policy Policy Evaluation. In *Proceedings of the Conference on Health, Inference, and Learning*, Vol. 174. 397–410.
 - [18] Haruka Kiyohara, Kosuke Kawakami, and Yuta Saito. 2021. Accelerating Offline Reinforcement Learning Application in Real-Time Bidding and Recommendation: Potential Use of Simulation. *arXiv preprint arXiv:2109.08331* (2021).
 - [19] Haruka Kiyohara, Yuta Saito, Tatsuya Matsuhiro, Yusuke Narita, Nobuyuki Shimizu, and Yasuo Yamamoto. 2022. Doubly Robust Off-Policy Evaluation for Ranking Policies under the Cascade Behavior Model. In *Proceedings of the 15th ACM International Conference on Web Search and Data Mining*. 487–497.
 - [20] 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).
 - [21] Lihong Li, Wei Chu, John Langford, and Xuanhui Wang. 2011. Unbiased Offline Evaluation of Contextual-bandit-based News Article Recommendation Algorithms. In *Proceedings of the Fourth ACM International Conference on Web Search and Data Mining*. *WSDM*, 297–306.
 - [22] Shuai Li, Yasin Abbasi-Yadkori, Branislav Kveton, S 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*. 1685–1694.
 - [23] Jiaxin Mao, Cheng Luo, Min Zhang, and Shaoping Ma. 2018. Constructing Click Models for Mobile Search. In *The 41st International ACM SIGIR Conference on Research and Development in Information Retrieval*. 775–784.
 - [24] James McInerney, Brian Brost, Praveen Chandar, Rishabh Mehrotra, and Benjamin Carterette. 2020. Counterfactual Evaluation of Slate Recommendations with Sequential Reward Interactions. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 1779–1788.
 - [25] Doña Precup, Richard S. Sutton, and Satinder P. Singh. 2000. Eligibility Traces for Off-Policy Policy Evaluation. In *Proceedings of the 17th International Conference on Machine Learning*. 759–766.
 - [26] Yuta Saito. 2020. Unbiased Pairwise Learning from Biased Implicit Feedback. In *Proceedings of the 2020 ACM SIGIR on International Conference on Theory of Information Retrieval*. 5–12.
 - [27] Yuta Saito, Shunsuke Aihara, Megumi Matsutani, and Yusuke Narita. 2020. Open Bandit Dataset and Pipeline: Towards Realistic and Reproducible Off-Policy Evaluation. *arXiv preprint arXiv:2008.07146* (2020).
 - [28] Yuta Saito and Thorsten Joachims. 2021. Counterfactual Learning and Evaluation for Recommender Systems: Foundations, Implementations, and Recent Advances. In *Proceedings of the 15th ACM Conference on Recommender Systems*. 828–830.
 - [29] Yuta Saito and Thorsten Joachims. 2022. Off-Policy Evaluation for Large Action Spaces via Embeddings. In *Proceedings of the 39th International Conference on Machine Learning*. 19089–19122.
 - [30] Yuta Saito, Gota Morishta, and Shota Yasui. 2020. Dual learning algorithm for delayed conversions. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1849–1852.
 - [31] Yuta Saito, Qingyang Ren, and Thorsten Joachims. 2023. Off-Policy Evaluation for Large Action Spaces via Conjunct Effect Modeling. *arXiv preprint arXiv:2305.08062* (2023).
 - [32] Yuta Saito, Takuma Udagawa, Haruka Kiyohara, Kazuki Mogi, Yusuke Narita, and Kei Tateno. 2021. Evaluating the Robustness of Off-Policy Evaluation. In *Proceedings of the 15th ACM Conference on Recommender Systems*. 114–123.
 - [33] Yuta Saito, Suguru Yaginuma, Yuta Nishino, Hayato Sakata, and Kazuhide Nakata. 2020. Unbiased Recommender Learning from Missing-Not-At-Random Implicit Feedback. In *Proceedings of the 13th International Conference on Web Search and Data Mining*. 501–509.
 - [34] Alex Strehl, John Langford, Lihong Li, and Sham M Kakade. 2010. Learning from Logged Implicit Exploration Data. In *Advances in Neural Information Processing Systems*, Vol. 23. 2217–2225.
 - [35] Yi Su, Maria Dimakopoulou, Akshay Krishnamurthy, and Miroslav Dudík. 2020. Doubly Robust Off-Policy Evaluation with Shrinkage. In *Proceedings of the 37th International Conference on Machine Learning*, Vol. 119. PMLR, 9167–9176.
 - [36] Adith Swaminathan, Akshay Krishnamurthy, Alekh Agarwal, Miro Dudík, John Langford, Damien Jose, and Imed Zitouni. 2017. Off-Policy Evaluation for Slate Recommendation. In *Advances in Neural Information Processing Systems*, Vol. 30. 3632–3642.
 - [37] Takuma Udagawa, Haruka Kiyohara, Yusuke Narita, Yuta Saito, and Kei Tateno. 2022. Policy-Adaptive Estimator Selection for Off-Policy Evaluation. *arXiv preprint arXiv:2211.13904* (2022).
 - [38] Nikos Vlassis, Fernando Amat Gil, and Ashok Chandrashekar. 2021. Off-Policy Evaluation of Slate Policies under Bayes Risk. *arXiv preprint arXiv:2101.02553* (2021).
 - [39] Xuanhui Wang, Michael Bendersky, Donald Metzler, and Marc Najork. 2016. Learning to Rank with Selection Bias in Personal Search. In *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*. 115–124.
 - [40] Yu-Xiang Wang, Alekh Agarwal, and Miroslav Dudík. 2017. Optimal and Adaptive Off-policy Evaluation in Contextual Bandits. In *Proceedings of the 34th International Conference on Machine Learning*. *ICML*, 3589–3597.
 - [41] Danqing Xu, Yiqun Liu, Min Zhang, Shaoping Ma, and Liyun Ru. 2012. Incorporating Revisiting Behaviors into Click Models. In *Proceedings of the 15th ACM International Conference on Web Search and Data Mining*. 303–312.
 - [42] Junqi Zhang, Yiqun Liu, Jiaxin Mao, Xiaohui Xie, Min Zhang, Shaoping Ma, and Qi Tian. 2022. Global or Local: Constructing Personalized Click Models for Web Search. In *Proceedings of the ACM Web Conference 2022*. 213–223.
 - [43] Ruizhe Zhang, Xiaohui Xie, Jiaxin Mao, Yiqun Liu, Min Zhang, and Shaoping Ma. 2021. Constructing a Comparison-based Click Model for Web Search. In *Proceedings of the Web Conference 2021*. 270–283.

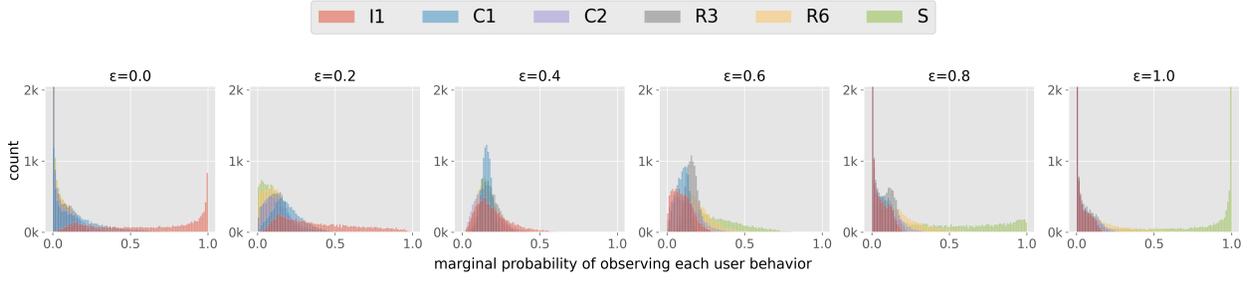
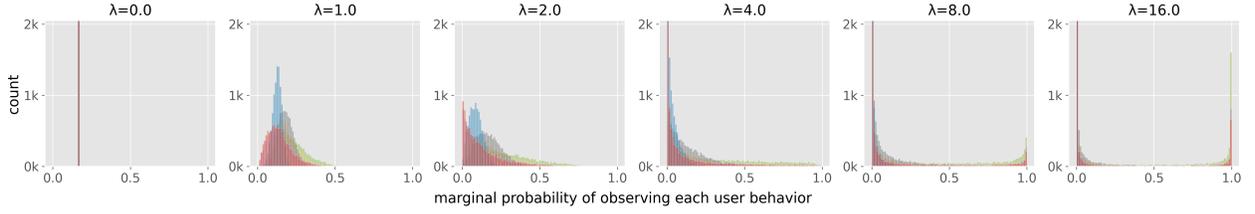
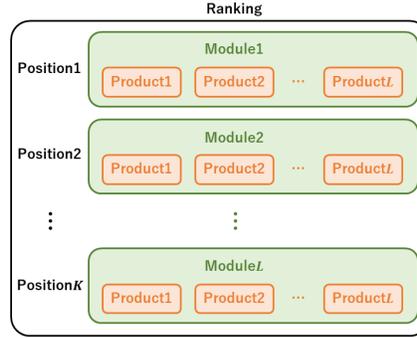
Figure 6: User behavior distribution with varying values of δ Figure 7: User behavior distribution with varying values of λ 

Figure 8: This figure illustrates a ranking of modules in the e-commerce platform used in our real-world experiment where each "Module" corresponds to an action indicating a category of products, such as "Recommended items" or "Campaign information".

A ADDITIONAL EXPERIMENT DETAILS AND RESULTS

A.1 Experimental Details

Distributions of user behavior in the synthetic experiment. In the synthetic experiment, we control the distribution of user behavior by varying the values of δ (user behavior distribution parameter) – a small value of δ increases the probability of observing simple user behaviors, while a large value of δ increases the probability of observing complex user behaviors. Figure 6 demonstrates how different values of δ control the distribution of user behavior, which we estimate with randomly sampled 10,000 user contexts.

Platform's ranking interface in the real-world experiment. Figure 8 illustrates the ranking interface of the e-commerce platform used in the real-world experiment. The two factored policies, π_A and π_B , choose which module as an action to present at each position in a ranking to maximize the sum of observed clicks during the data collection experiment.

The candidate set of behavior models for AIPS in the real-world experiment. In the real-world experiment, AIPS uses $\mathcal{C} = \{c_S, c_C, c_I, c', c''\}$ as the candidate set of behavior models when performing user behavior optimization. c_S , c_C , and c_I are defined in Section 4. c' and c'' are defined specifically as

- c' : $c'(k, 1) = c'(k, 2) = c'(k, k) = 1$, otherwise, $c'(k, l) = 0$, $\forall l \in [K]$

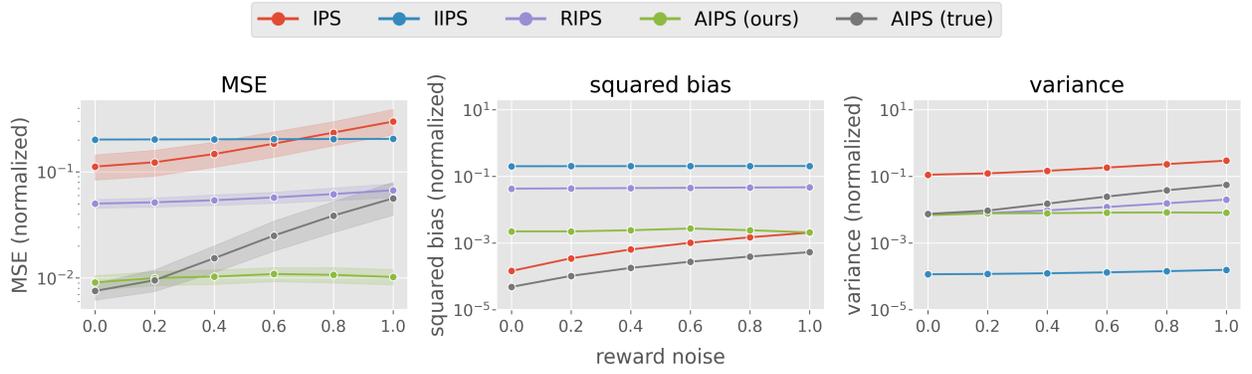


Figure 9: Comparison of the estimators' MSE (normalized by the true value $V(\pi)$) with varying reward noise levels (σ)

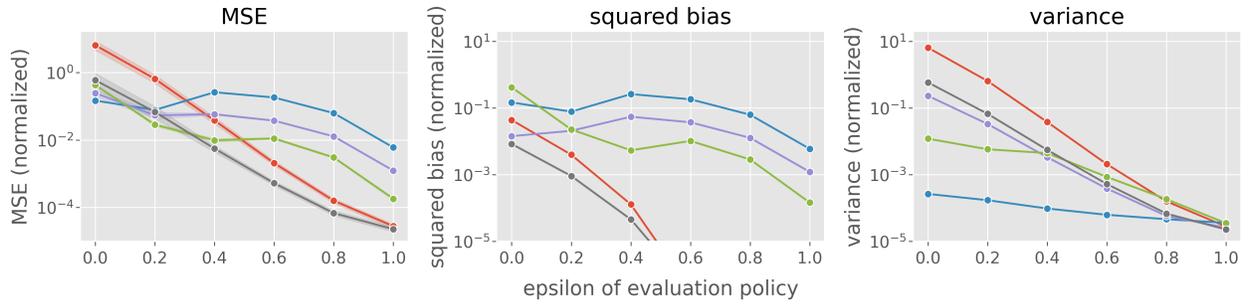


Figure 10: Comparison of the estimators' MSE (normalized by the true value $V(\pi)$) with varying evaluation policies (ϵ)

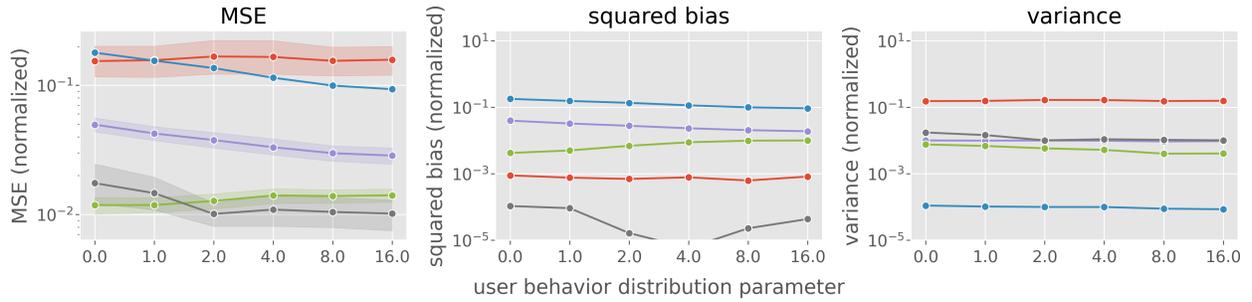


Figure 11: Comparison of the estimators' MSE (normalized by the true value $V(\pi)$) with varying behavior distributions (λ)

- c'' : $c''(k, l) = 1$ if $l \leq k$, otherwise, $c''(k, l) = 0$, $\forall l \in [K]$

for each position $k \in [K]$.

A.2 Additional Results on Synthetic Data

This section explores four additional research questions regarding: (A1) reward noise level (σ), (A2) evaluation policy (ϵ), (A3) identifiability of user behavior (λ), and (A4) estimation error in the bias estimate (σ_Δ) used in AIPS. Note that we set $n = 8k$ (data size), $K = 8$ (length of ranking), $\sigma = 0.5$ (reward noise level), $\epsilon = 0.3$ (evaluation policy parameter), $\delta = 0.6$ (user behavior setting), and $\sigma_\Delta = 0.3$ (error level in bias estimate) as default, and run OPE simulations over 300 different logged data replicated with different random seeds. The following reports and discusses the MSE, squared bias, and variance of the estimators normalized by the true policy value of the evaluation policy $V(\pi)$.

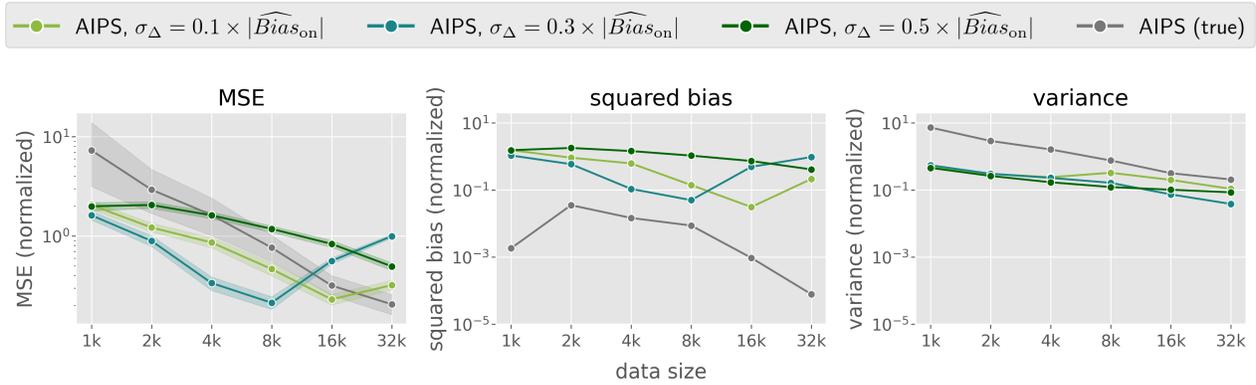


Figure 12: Comparison of AIPS's performance with varying data sizes (n) and estimation error (σ_{Δ})

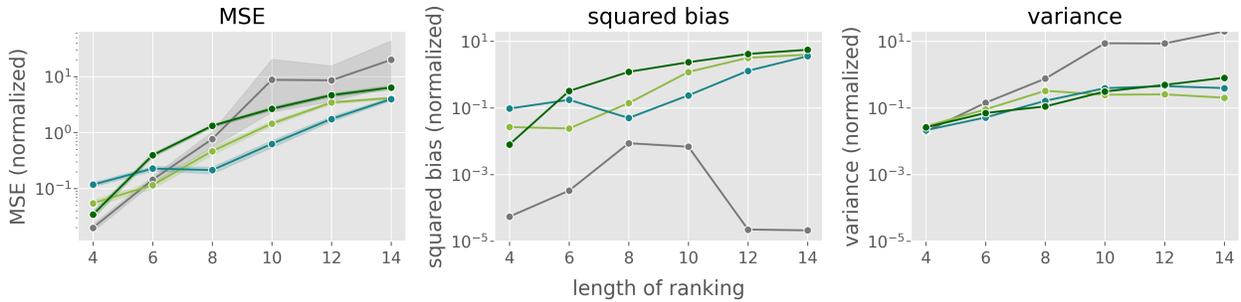


Figure 13: Comparison of AIPS's performance with varying lengths of ranking (K) and estimation error (σ_{Δ})

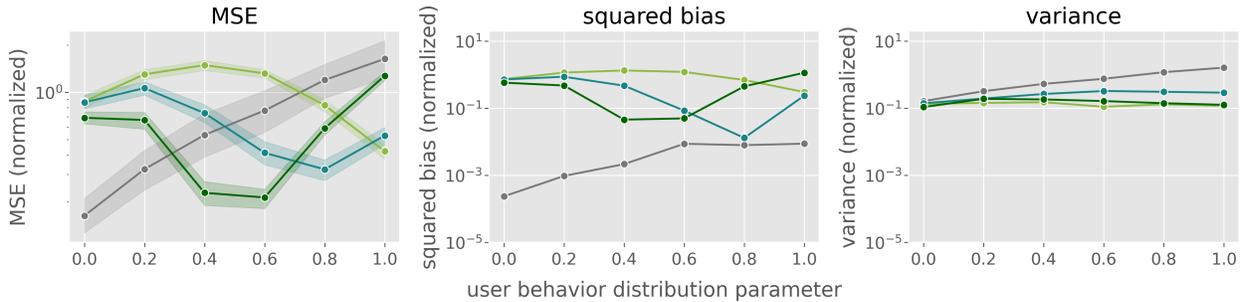


Figure 14: Comparison of AIPS's performance with varying user behavior distributions (δ) and estimation error (σ_{Δ})

RQ (A1): How do the estimators perform with varying reward noise levels? Figure 9 compares estimators' performance with varying reward noise levels $\sigma \in \{0.0, 0.2, \dots, 1.0\}$. We observe that IPS and AIPS (true) suffer from increasingly high variance as the noise level becomes higher. In contrast, our AIPS demonstrates its clear robustness to the increase in reward noise.

RQ (A2): How do the estimators perform with varying evaluation policies? Next, Figure 10 compares the estimators with various evaluation policies $\epsilon \in \{0.0, 0.2, \dots, 1.0\}$. The figure indicates that AIPS (true) and IPS are quite accurate when the evaluation policy is near-uniform and does not deviate from the logging policy greatly ($\epsilon = 0.8, 1.0$). This is because AIPS (true) and IPS are unbiased and do not suffer from high variance when the evaluation policy is highly stochastic. However, under more practical situations where the evaluation policy is more deterministic, AIPS becomes superior due to its favorable variance property while IPS performs the worst due to its extreme variance. Note that AIPS (true) also works well for a range of evaluation policies in the default setting (i.e., $n = 8k, K = 8$). However, it may

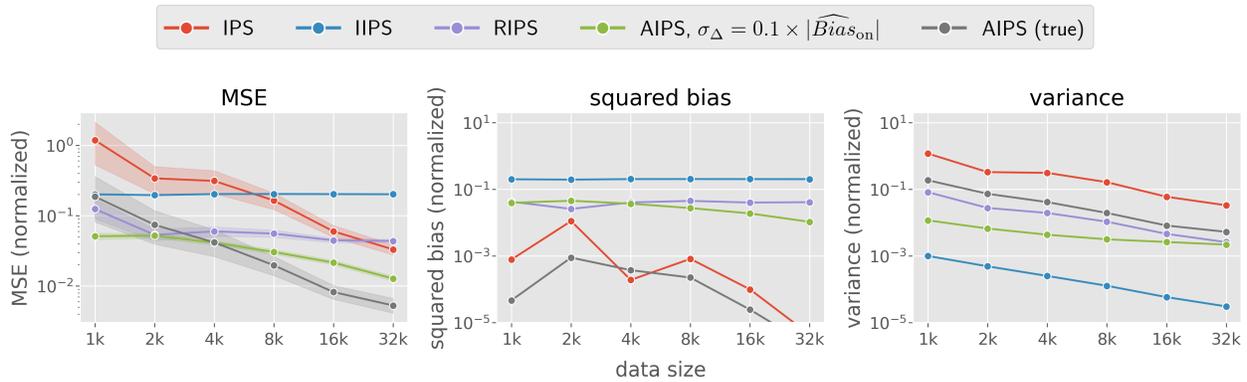


Figure 15: Comparison of the estimators' MSE (normalized by the true value $V(\pi)$) with varying data sizes (n) when $\sigma_\Delta = 0.1 \times |\widehat{Bias}_{on}|$

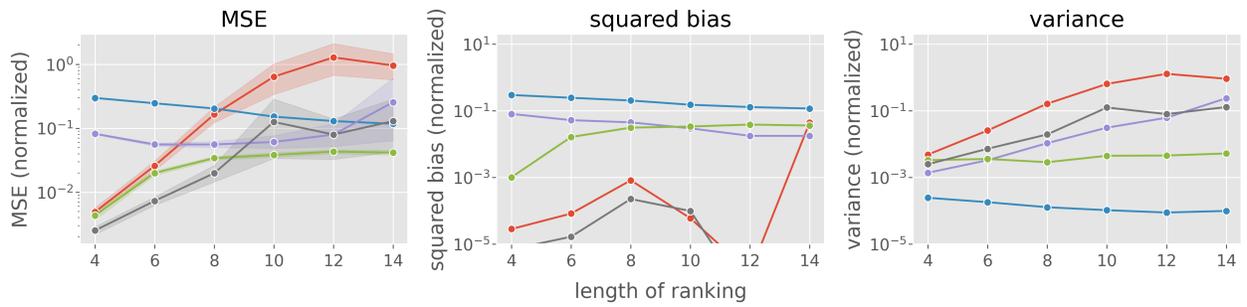


Figure 16: Comparison of the estimators' MSE (normalized by the true value $V(\pi)$) with varying lengths of ranking (K) when $\sigma_\Delta = 0.1 \times |\widehat{Bias}_{on}|$

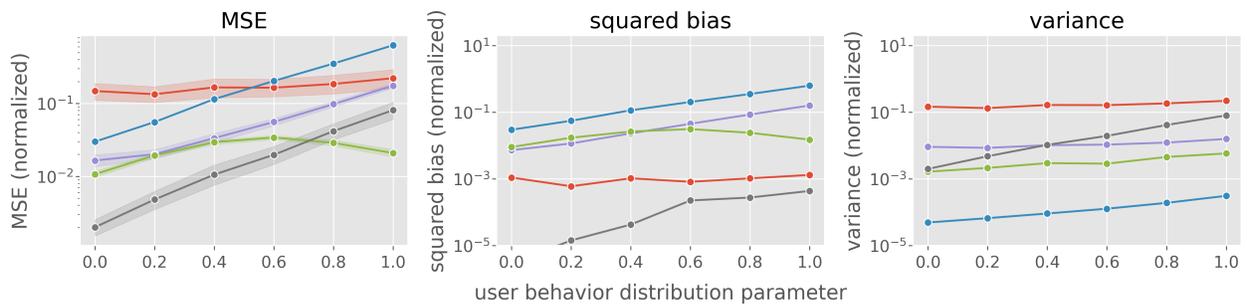


Figure 17: Comparison of the estimators' MSE (normalized by the true value $V(\pi)$) with varying user behavior distributions (δ) when $\sigma_\Delta = 0.1 \times |\widehat{Bias}_{on}|$

suffer from a higher variance when the ranking size is larger ($K > 8$) and the evaluation policy is more deterministic ($\epsilon = 0.0, 0.2$), as already shown in Figure 3.

RQ (A3): How do the estimators perform under different levels of identifiability of user behavior? Here, we investigate how the estimators perform with varying levels of identifiability of user behavior. Specifically, we control the identifiability of user behavior by varying the values of λ_z . Recall here that we sample user behavior from the following conditional distribution.

$$p(c_z | \mathbf{x}) := \text{softmax}(\lambda_z \cdot |\theta_z^\top \mathbf{x}|) = \frac{\exp(\lambda_z \cdot |\theta_z^\top \mathbf{x}|)}{\sum_{z'} \exp(\lambda_{z'} \cdot |\theta_{z'}^\top \mathbf{x}|)}.$$



Figure 18: Comparison of the estimators' MSE (normalized by the true value $V(\pi)$) with varying data sizes (n) and $\sigma_{\Delta} = 0.5 \times |\widehat{Bias}_{on}|$

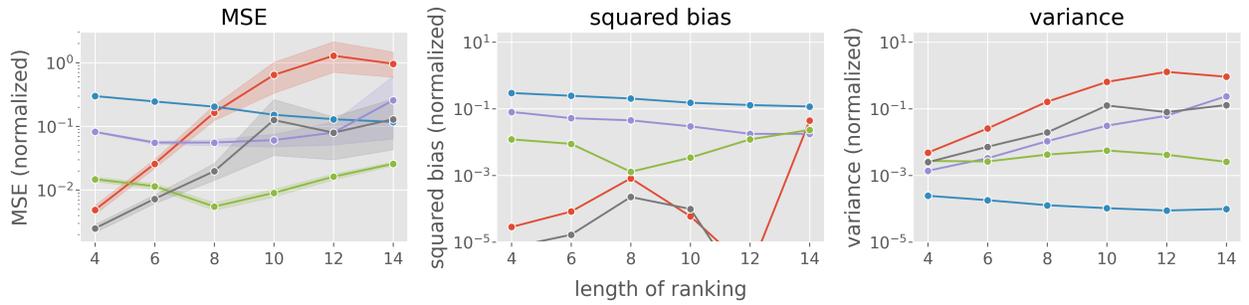


Figure 19: Comparison of the estimators' MSE (normalized by the true value $V(\pi)$) with varying lengths of ranking (K) and $\sigma_{\Delta} = 0.5 \times |\widehat{Bias}_{on}|$

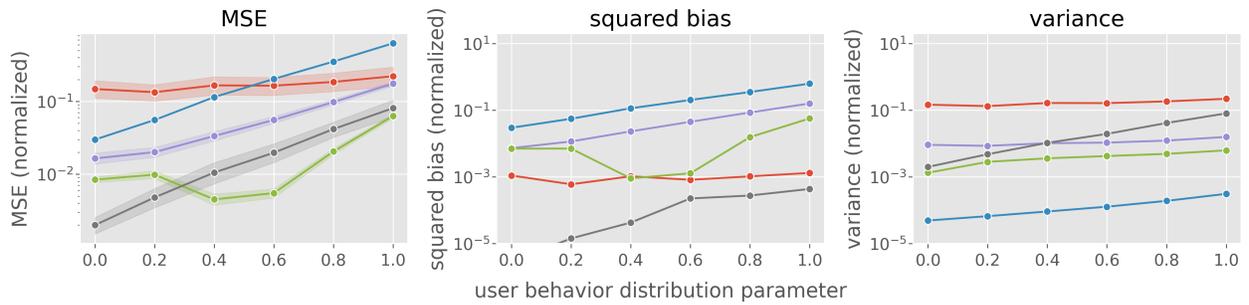


Figure 20: Comparison of the estimators' MSE (normalized by the true value $V(\pi)$) with varying user behavior distributions (δ) and $\sigma_{\Delta} = 0.5 \times |\widehat{Bias}_{on}|$

Thus, by definition, when we set $\lambda_z = \lambda, \forall z$ (constant value), user behavior will be uniformly distributed. We can increase identifiability of user behavior by using a large value of λ where user behavior becomes near-deterministic and easily identifiable from the context. By contrast, when we use a small value of λ , user behavior will be almost context-independent and may not be easily identifiable from the observed user context. We thus vary $\lambda \in \{0.0, 1.0, 2.0, \dots, 16.0\}$ to see how the estimators' performance changes with different levels of identifiability of user behavior (Figure 7 illustrates the distributions of user behavior with varying values of λ).

Somewhat surprisingly, we observe in Figure 11 that identifiability of user behavior has almost no impact on the estimators' performances. We conjecture that AIPS is robust to the change in the level of identifiability because it does not aim to precisely estimate user behavior,

Table 3: Runtime comparison between AIPS and IPS with varying data sizes (n)

data size (n)	1,000	2,000	4,000	8,000	16,000	32,000
IPS	0.6022 (\pm 0.008)	1.198 (\pm 0.021)	2.379 (\pm 0.041)	4.752 (\pm 0.072)	9.505 (\pm 0.149)	18.98 (\pm 0.251)
AIPS	23.43 (\pm 0.534)	44.22 (\pm 0.9206)	85.24 (\pm 1.502)	167.2 (\pm 2.556)	329.6 (\pm 4.162)	653.5 (\pm 7.611)
(relative)	38.90	36.98	35.83	35.18	34.68	34.42

Note: We report mean (\pm std) of the runtime (sec) of IPS and AIPS in the synthetic experiment over 100 random seeds. (relative) reports the runtime of AIPS divided by that of IPS.

Table 4: Runtime comparison between AIPS and IPS with varying lengths of ranking (K)

length of ranking (K)	4	6	8	10	12	14
IPS	2.313 (\pm 0.021)	3.390 (\pm 0.025)	4.690 (\pm 0.037)	6.192 (\pm 0.039)	7.924 (\pm 0.066)	9.878 (\pm 0.068)
AIPS	68.89 (\pm 0.977)	114.3 (\pm 1.212)	166.2 (\pm 1.822)	221.5 (\pm 2.384)	284.2 (\pm 3.136)	351.5 (\pm 3.059)
(relative)	29.78	33.70	35.44	35.76	35.87	35.59

Note: We report mean (\pm std) of the runtime (sec) of IPS and AIPS in the synthetic experiment over 100 random seeds. (relative) reports the runtime of AIPS divided by that of IPS.

but it rather aims to choose the most suitable user behavior model to minimize its MSE given the logged data. This observation provides a further argument for the applicability and robustness of the user behavior optimization procedure of AIPS.

RQ (A4): How does AIPS perform with varying estimation errors in the bias estimate? Finally, we evaluate how robust AIPS is to the estimation error in the bias estimate by varying the values of $\sigma_\Delta \in \{0.1, 0.3, 0.5\} \times |\widehat{Bias}_{on}|$. More specifically, we compare AIPS' performance with the three different values of σ_Δ across varying data sizes $n \in \{1k, 2k, \dots, 32k\}$ (Figure 12), lengths of ranking $K \in \{4, 6, \dots, 14\}$ (Figure 13), and user behavior distribution parameters $\delta \in \{0.0, 0.2, \dots, 1.0\}$ (Figure 14) as done in RQs (1)-(3) in the main text. We also compare AIPS against IPS, IIPS, and RIPS with $\sigma_\Delta = 0.1 \times |\widehat{Bias}_{on}|$ in Figures 15-17 and with $\sigma_\Delta = 0.5 \times |\widehat{Bias}_{on}|$ in Figures 18-20.

Overall, the results indicate that AIPS is robust to the estimation error in the bias estimate. In particular, the trends observed in Figures 15-17 and Figures 18-20 are quite similar to those observed in Figures 2-4 in the main text, suggesting that AIPS effectively balances the bias-variance tradeoff even in the presence of severe estimation error in the bias estimate. Interestingly, we also observe that a larger estimation error in the bias estimate does not necessarily lead to a larger MSE. Specifically, in Figure 12 and Figure 13, the accuracy of AIPS only becomes worse with larger estimation error (σ_Δ) when the data size is large ($n = 16k, 32k$) and the ranking size is small ($K = 4, 6$). This implies that the estimation error of the bias estimate may be slightly problematic only when the bias is dominant in the MSE.

Runtime analysis. One potential concern of AIPS is the additional computation overhead introduced by its user behavior optimization procedure. Tables 3 and 4 compare the computation time of AIPS against that of IPS with varying data sizes (n) and lengths of ranking (K). The result shows that the whole estimation process ends only in 653 seconds (< 11 mins) even in the largest sample size. Moreover, we can see that the relative computation time of AIPS compared to IPS does not grow with the sample size and lengths of ranking.

B OMITTED PROOFS

B.1 Proof of Proposition 3.1

PROOF. For any given $\mathbf{x} \sim p(\mathbf{x})$ and $\mathbf{c} \sim p(\mathbf{c} | \mathbf{x})$, we have

$$\begin{aligned}
\mathbb{E}_{\pi_0(\mathbf{a}|\mathbf{x})p(\mathbf{r}|\mathbf{x},\mathbf{a},\mathbf{c})} \left[\frac{\pi(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})} r_k \right] &= \sum_{\mathbf{a}} \pi_0(\mathbf{a} | \mathbf{x}) \frac{\pi(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})} q_k(\mathbf{x}, \mathbf{a}, \mathbf{c}) \\
&= \sum_{\Phi_k(\mathbf{a}, \mathbf{c})} \sum_{\Phi_k^c(\mathbf{a}, \mathbf{c})} \pi_0(\mathbf{a} | \mathbf{x}) \frac{\pi(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})} q_k(\mathbf{x}, \mathbf{a}, \mathbf{c}) \\
&= \sum_{\Phi_k(\mathbf{a}, \mathbf{c})} \frac{\pi(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})} q_k(\mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c})) \sum_{\Phi_k^c(\mathbf{a}, \mathbf{c})} \pi_0(\mathbf{a} | \mathbf{x}) \\
&= \sum_{\Phi_k(\mathbf{a}, \mathbf{c})} \frac{\pi(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})} q_k(\mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c})) \underbrace{\sum_{\Phi_k^c(\mathbf{a}, \mathbf{c})} \pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) \cup \Phi_k^c(\mathbf{a}, \mathbf{c}) | \mathbf{x})}_{=\pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})} \\
&= \sum_{\Phi_k(\mathbf{a}, \mathbf{c})} \pi(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x}) q_k(\mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c})) \\
&= \sum_{\Phi_k(\mathbf{a}, \mathbf{c})} q_k(\mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c})) \sum_{\Phi_k^c(\mathbf{a}, \mathbf{c})} \pi(\Phi_k(\mathbf{a}, \mathbf{c}) \cup \Phi_k^c(\mathbf{a}, \mathbf{c}) | \mathbf{x}) \\
&= \sum_{\mathbf{a}} \pi(\mathbf{a} | \mathbf{x}) q_k(\mathbf{x}, \mathbf{a}, \mathbf{c}) \\
&= \mathbb{E}_{\pi(\mathbf{a}|\mathbf{x})p(\mathbf{r}|\mathbf{x},\mathbf{a},\mathbf{c})} [r_k]
\end{aligned} \tag{4}$$

where $q_k(\mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c})) := \mathbb{E}_{p(\mathbf{r}|\mathbf{x},\Phi_k(\mathbf{a},\mathbf{c}))} [r_k]$, $q_k(\mathbf{x}, \mathbf{a}, \mathbf{c}) := \mathbb{E}_{p(\mathbf{r}|\mathbf{x},\mathbf{a},\mathbf{c})} [r_k]$, and $q_k(\mathbf{x}, \mathbf{a}, \mathbf{c}) = q_k(\mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))$.

Then, we have

$$\begin{aligned}
\mathbb{E}_{\mathcal{D}} \left[\hat{V}_k^{\text{AIPS}}(\pi; \mathcal{D}) \right] &= \mathbb{E}_{\mathcal{D}} \left[\frac{1}{n} \sum_{i=1}^n \frac{\pi(\Phi_k(\mathbf{a}_i, \mathbf{c}_i) | \mathbf{x}_i)}{\pi_0(\Phi_k(\mathbf{a}_i, \mathbf{c}_i) | \mathbf{x}_i)} r_{i,k} \right] \\
&= \mathbb{E}_{p(\mathbf{x})p(\mathbf{c}|\mathbf{x})\pi_0(\mathbf{a}|\mathbf{x})p(\mathbf{r}|\mathbf{x},\mathbf{a},\mathbf{c})} \left[\frac{\pi(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})} r_k \right] \\
&= \mathbb{E}_{p(\mathbf{x})p(\mathbf{c}|\mathbf{x})} \left[\mathbb{E}_{\pi_0(\mathbf{a}|\mathbf{x})p(\mathbf{r}|\mathbf{x},\mathbf{a},\mathbf{c})} \left[\frac{\pi(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})} r_k \right] \right] \\
&= \mathbb{E}_{p(\mathbf{x})p(\mathbf{c}|\mathbf{x})} \left[\mathbb{E}_{\pi(\mathbf{a}|\mathbf{x})p(\mathbf{r}|\mathbf{x},\mathbf{a},\mathbf{c})} [r_k] \right] \quad \because \text{Eq. (4)} \\
&= V_k(\pi)
\end{aligned}$$

□

B.2 Proof of Theorems 3.2 and 3.3

PROOF. To prove Theorems 3.2 and 3.3, we quantify the difference between the variance of AIPS ($\hat{V}_k^{\text{AIPS}}(\pi; \mathcal{D})$) and that of an arbitrary unbiased estimator ($\hat{V}_k(\pi; \mathcal{D}, \tilde{\mathbf{c}})$) defined in Theorem 3.3. For brevity of exposition, the following uses $\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{\mathbf{c}}) := \Phi_k(\mathbf{a}, \tilde{\mathbf{c}}) \setminus \Phi_k(\mathbf{a}, \mathbf{c})$. We will also use the fact that $\mathbf{c} \subseteq \tilde{\mathbf{c}}$ always holds true for any $\tilde{\mathbf{c}}$ when $\hat{V}_k(\pi; \mathcal{D}, \tilde{\mathbf{c}})$ is unbiased. This directly follows from Theorem 3.4, which is

later proved in Appendix B.3

$$\begin{aligned}
& n \left(\mathbb{V}_{\mathcal{D}}(\hat{V}_k(\pi; \mathcal{D}, \tilde{c})) - \mathbb{V}_{\mathcal{D}}(\hat{V}_k^{\text{AIPS}}(\pi; \mathcal{D})) \right) \\
&= n \left(\mathbb{V}_{\mathcal{D}} \left(\frac{\pi(\Phi_k(\mathbf{a}, \tilde{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \tilde{c}) | \mathbf{x})} r_k \right) - \mathbb{V}_{\mathcal{D}} \left(\frac{\pi(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})} r_k \right) \right) \\
&= \mathbb{E}_{p(\mathbf{x})p(\mathbf{c}|\mathbf{x})\pi_0(\mathbf{a}|\mathbf{x})p(\mathbf{r}|\mathbf{x}, \mathbf{a}, \mathbf{c})} \left[\left(\frac{\pi(\Phi_k(\mathbf{a}, \tilde{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \tilde{c}) | \mathbf{x})} r_k \right)^2 \right] - \underbrace{\left(\mathbb{E}_{p(\mathbf{x})p(\mathbf{c}|\mathbf{x})\pi_0(\mathbf{a}|\mathbf{x})p(\mathbf{r}|\mathbf{x}, \mathbf{a}, \mathbf{c})} \left[\frac{\pi(\Phi_k(\mathbf{a}, \tilde{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \tilde{c}) | \mathbf{x})} r_k \right] \right)^2}_{=V_k(\pi)} \\
&\quad - \underbrace{\left(\mathbb{E}_{p(\mathbf{x})p(\mathbf{c}|\mathbf{x})\pi_0(\mathbf{a}|\mathbf{x})p(\mathbf{r}|\mathbf{x}, \mathbf{a}, \mathbf{c})} \left[\frac{\pi(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})} r_k \right] \right)^2}_{=V_k(\pi)} \\
&= \mathbb{E}_{p(\mathbf{x})p(\mathbf{c}|\mathbf{x})\pi_0(\mathbf{a}|\mathbf{x})p(\mathbf{r}|\mathbf{x}, \mathbf{a}, \mathbf{c})} \left[\left(\left(\frac{\pi(\Phi_k(\mathbf{a}, \tilde{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \tilde{c}) | \mathbf{x})} \right)^2 - \left(\frac{\pi(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})} \right)^2 \right) r_k^2 \right] \\
&= \mathbb{E}_{p(\mathbf{x})p(\mathbf{c}|\mathbf{x})\pi_0(\mathbf{a}|\mathbf{x})p(\mathbf{r}|\mathbf{x}, \mathbf{a}, \mathbf{c})} \left[\left(\frac{\pi(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})} \right)^2 \left(\left(\frac{\pi(\Phi_k(\mathbf{a}, \tilde{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \tilde{c}) | \mathbf{x})} \frac{\pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})}{\pi(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})} \right)^2 - 1 \right) r_k^2 \right] \\
&= \mathbb{E}_{p(\mathbf{x})p(\mathbf{c}|\mathbf{x})\pi_0(\mathbf{a}|\mathbf{x})p(\mathbf{r}|\mathbf{x}, \mathbf{a}, \mathbf{c})} \left[\left(\frac{\pi(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})} \right)^2 \left(\left(\frac{\pi(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))}{\pi_0(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))} \right)^2 - 1 \right) r_k^2 \right] \\
&= \mathbb{E}_{p(\mathbf{x})p(\mathbf{c}|\mathbf{x})\pi_0(\mathbf{a}|\mathbf{x})} \left[\left(\frac{\pi(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})} \right)^2 \left(\left(\frac{\pi(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))}{\pi_0(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))} \right)^2 - 1 \right) \mathbb{E}_{p(\mathbf{r}|\mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))} [r_k^2] \right] \\
&= \mathbb{E}_{p(\mathbf{x})p(\mathbf{c}|\mathbf{x})\pi_0(\Phi_k(\mathbf{a}, \mathbf{c})|\mathbf{x})} \left[\left(\frac{\pi(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})} \right)^2 \mathbb{E}_{\pi_0(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))} \left[\left(\frac{\pi(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))}{\pi_0(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))} \right)^2 - 1 \right] \mathbb{E}_{p(\mathbf{r}|\mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))} [r_k^2] \right] \\
&= \mathbb{E}_{p(\mathbf{x})p(\mathbf{c}|\mathbf{x})\pi_0(\Phi_k(\mathbf{a}, \mathbf{c})|\mathbf{x})} \left[\left(\frac{\pi(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})}{\pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})} \right)^2 \mathbb{V}_{\pi_0(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))} \left[\frac{\pi(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))}{\pi_0(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))} \right] \mathbb{E}_{p(\mathbf{r}|\mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))} [r_k^2] \right], \quad (5)
\end{aligned}$$

where we use $\frac{\pi(\Phi_k(\mathbf{a}, \tilde{c}) | \mathbf{x})}{\pi(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})} = \pi(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))$ and $\pi_0(\mathbf{a} | \mathbf{x}) = \pi_0(\Phi_k(\mathbf{a}, \mathbf{c}) | \mathbf{x})\pi_0(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))$. Moreover, in Eq. (5), we use the following trick:

$$\begin{aligned}
& \mathbb{E}_{\pi_0(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))} \left[\left(\frac{\pi(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))}{\pi_0(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))} \right)^2 - 1 \right] \\
&= \mathbb{E}_{\pi_0(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))} \left[\left(\frac{\pi(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))}{\pi_0(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))} \right)^2 \right] - \underbrace{\left(\mathbb{E}_{\pi_0(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))} \left[\frac{\pi(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))}{\pi_0(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))} \right] \right)^2}_{=1} \\
&= \mathbb{V}_{\pi_0(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))} \left[\frac{\pi(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))}{\pi_0(\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) | \mathbf{x}, \Phi_k(\mathbf{a}, \mathbf{c}))} \right]
\end{aligned}$$

We can see that Eq. (5) is always non-negative, which means that the variance of AIPS is never larger than that of any unbiased IPS estimator defined by $\hat{V}_k(\pi; \mathcal{D}, \tilde{c})$ with $\mathbf{c} \subseteq \tilde{c}$. Hence, Theorem 3.3 is proved. Furthermore, we can derive Theorem 3.2 by replacing $\Phi_k(\mathbf{a}, \tilde{c})$ with \mathbf{a} (in this case, $\Phi_k^d(\mathbf{a}, \mathbf{c}, \tilde{c}) = \Phi_k^c(\mathbf{a}, \mathbf{c})$). \square

B.3 Proof of Theorem 3.4

PROOF. First, we calculate the bias of AIPS with an estimated user behavior \hat{c} below.

Algorithm 1 The procedure to optimize user behavior assignments in AIPS (detailed in Section 3.2)

Input: logged data \mathcal{D} , a loss function to minimize MSE $\mathcal{L}(\cdot)$, a set of candidate user behavior models $C = \{\hat{c}^0, \dots, \hat{c}^m\}$, the base user behavior model \hat{c}_{base} , a set of random states \mathcal{S}

Output: dictionary containing \hat{c} for each partition \mathbb{C}

- 1: Initialize node sets to partition $\mathbb{L} \leftarrow \{\mathcal{D}\}$ and dictionary containing \hat{c} for each partition $\mathbb{C} \leftarrow \emptyset$
- 2: Initialize the number of user partition $g \leftarrow 0$
- 3: **while** $\mathbb{L} \neq \emptyset$ **do**
- 4: Remove node l from \mathbb{L} as $\mathbb{L} \leftarrow \mathbb{L} \setminus \{l\}$ and set l as the parent node
- 5: Initialize the minimum loss $\widehat{MSE}^{(-)} \leftarrow \widehat{MSE}(\hat{c}_{(l)}; l)$ where $\hat{c}_{(l)} := \arg \min_{\hat{c}} \widehat{MSE}(\hat{c}; l)$
- 6: Initialize the best subset $(\hat{c}_{(l(l))}, \hat{c}_{(l(r))}, \mathcal{D}_{(l(l))}, \mathcal{D}_{(l(r))}) \leftarrow \emptyset$
- 7: **for** $s \in \mathcal{S}$ **do**
- 8: Randomly generate partitions in the feature space (\mathcal{X}) and create two subsets of the data ($\mathcal{D}_{(l^*)}, \mathcal{D}_{(r^*)}$)
 (e.g., data that satisfy $\|\mathbf{x}\|_2 \leq 1$ are deemed as the subset indicating the left node ($\mathcal{D}_{(l^*)}$), while others
 as the subset of the right node ($\mathcal{D}_{(r^*)}$))
- 9: Identify the best behavior model for each subset as
 $(\hat{c}_{(l^*)}, \hat{c}_{(r^*)}) := \arg \min_{(\hat{c}_{(l)}, \hat{c}_{(r)})} \widehat{MSE}(\hat{c}_{(l)}, \hat{c}_{(r)}; l^*, r^*)$
- 10: **if** $\widehat{MSE}^{(-)} > \widehat{MSE}(\hat{c}_{(l^*)}, \hat{c}_{(r^*)}; l^*, r^*)$ **then**
- 11: Update the best partition as
 $\widehat{MSE}^{(-)} \leftarrow \widehat{MSE}(\hat{c}_{(l^*)}, \hat{c}_{(r^*)}; l^*, r^*)$
 $(\hat{c}_{(l(l))}, \hat{c}_{(l(r))}, \mathcal{D}_{(l(l))}, \mathcal{D}_{(l(r))}) \leftarrow (\hat{c}_{(l^*)}, \hat{c}_{(r^*)}, \mathcal{D}_{(l^*)}, \mathcal{D}_{(r^*)})$
- 12: **end if**
- 13: **end for**
- 14: **if** $(\mathcal{D}_{(l(l))}, \mathcal{D}_{(l(r))}) = \emptyset$ **then**
- 15: // end of the optimization procedure
- 16: Add the parent partition and the corresponding user behavior model to \mathbb{C} as
 $\mathbb{C}[g] \leftarrow (\mathcal{D}_{(l)}, \hat{c}_{(l)}, l), \quad g \leftarrow g + 1$
- 17: **else**
- 18: // continue the optimization procedure
- 19: Add children nodes to the tree as $\mathbb{L} \leftarrow \mathbb{L} \cup \{\mathcal{D}_{(l(l))}, \mathcal{D}_{(l(r))}\}$
- 20: **end if**
- 21: **end while**
