
Uncertainty-Aware Instance Reweighting for Off-Policy Learning

Xiaoying Zhang¹ Junpu Chen² Hongning Wang³ Hong Xie⁴ Yang Liu¹
John C.S. Lui⁵ Hang Li¹

¹ByteDance Research ²ChongQing University ³Tsinghua University
⁴Chongqing Institute of Green and Intelligent Technology, Chinese Academy of Science
⁵The Chinese University of Hong Kong
{zhangxiaoying.xy, yang.liu01, lihang.lh}@bytedance.com
{jumpchan98, hongx87, wang.hongn}@gmail.com
cslui@cse.cuhk.edu.hk

Abstract

Off-policy learning, referring to the procedure of policy optimization with access only to logged feedback data, has shown importance in various real-world applications, such as search engines and recommender systems. While the ground-truth logging policy is usually unknown, previous work simply employs its estimated value for the off-policy learning, ignoring the negative impact from both high bias and high variance resulted from such an estimator. And such impact is often magnified on samples with small and inaccurately estimated logging probabilities. The contribution of this work is to explicitly model the uncertainty in the estimated logging policy, and propose an Uncertainty-aware Inverse Propensity Score estimator (UIPS) for improved off-policy learning, with a theoretical convergence guarantee. Experiment results on the synthetic and real-world recommendation datasets demonstrate that UIPS significantly improves the quality of the discovered policy, when compared against an extensive list of state-of-the-art baselines.

1 Introduction

In many real-world applications, including search engines [2], online advertisements [35], recommender systems [8, 22], only logged data is available for subsequent policy learning. For example, in recommender systems, various complex recommendation policies are optimized over logged user interactions (e.g., clicks or stay time) with items recommended by previous recommendation policies (referred to as the *logging policy*) [51, 14]. However, such logged data is often known to be biased, since the feedback on items where the logging policy did not take is unknown. This inevitably distorts the evaluation and optimization of a new policy when it differs from the logging policy.

Off-policy learning [41, 27] thus emerges as a preferred way to learn an improved policy only from the logged data, by addressing the mismatch between the learning and logging policies. One of the most commonly used off-policy learning methods is the Inverse Propensity Scoring (IPS) [8, 25], which assigns per-sample importance weight (i.e., propensity score) to the training objective on the logged data, so as to get an unbiased optimization objective in expectation. The importance weight in IPS is the probability ratio of taking an action between the learning and logging policies.

Unfortunately, the ground-truth logging policy is oftentimes unavailable to the learner in practice, due to reasons like legacy issues, i.e., it was not recorded in the data. Additionally, in specific situations like the healthcare domain [28] or two-stage recommender systems [8], access to the ground-truth logging policy is not feasible. One common treatment by many previous studies [35, 22, 8, 24] is to first estimate the logging policy using a supervised learning method (e.g., logistic regression,

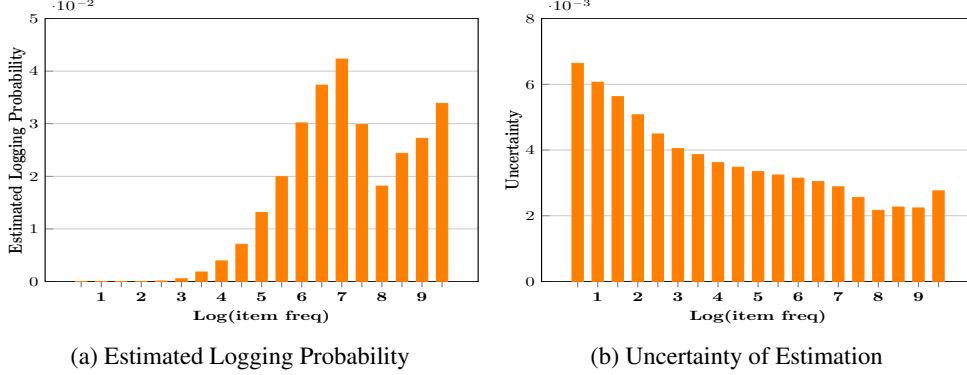


Figure 1: Estimated logging policy and its uncertainty under different item frequency on KuaiRec.

neural networks, etc.), and then employ the estimated logging policy for off-policy learning. In this work, we first show that such an approximation results in a biased estimator which is sensitive to data with small estimated logging probabilities. Worse still, small estimated logging probabilities usually suggest there are limited related samples in the logged data, whose estimations can have high uncertainties, i.e., being wrong with a high probability. Figure 1 shows a piece of empirical evidence from a large-scale recommendation benchmark KuaiRec dataset [12], where items with lower frequencies in the logged dataset have lower estimated logging probabilities (via a neural network estimator) and higher uncertainties at the same time. The high bias and variance caused by these samples can greatly hinder the performance of subsequent off-policy learning. We defer detailed discussions of this result in Section 2.

In this work, we explicitly take the uncertainty of the estimated logging policy into consideration and design an Uncertainty-aware Inverse Propensity Score estimator (UIPS) for off-policy learning. UIPS reweights the propensity score of each logged sample to control its impact on policy optimization, and learns an improved policy by alternating between: (1) Find the optimal weight that makes the estimator as accurate as possible, based on the uncertainty of the estimated logging policy; (2) Improve the policy by optimizing the resulting objective function. The optimal weight for each sample is obtained by minimizing the upper bound of the mean squared error (MSE) to the ground-truth policy evaluation, with a closed-form solution. Furthermore, UIPS ensures that off-policy learning converges to a stationary point where the true policy gradient is zero; while convergence may not be guaranteed when directly using the estimated logging policy. Extensive experiments on a synthetic and three real-world recommendation datasets against a rich set of state-of-the-art baselines demonstrate the power of UIPS. All data and code can be found in <https://github.com/Xiaoyinggit/UIPS.git>.

2 Preliminary: off-policy learning

We focus on the standard contextual bandit setup to explain the key concepts in UIPS. Following the convention [16, 29, 36], let $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^d$ be a d -dimensional context vector drawn from an unknown distribution $p(\mathbf{x})$. Each context is associated with a finite set of actions denoted by \mathcal{A} , where $|\mathcal{A}| < \infty$. Let $\pi : \mathcal{X} \rightarrow \Delta(\mathcal{A})$ denote a stochastic policy, such that $\pi(a|\mathbf{x})$ is the probability of selecting action a under context \mathbf{x} and $\sum_{a \in \mathcal{A}} \pi(a|\mathbf{x}) = 1$. Under a given context \mathbf{x} , the reward $r_{\mathbf{x},a}$ is only observed when action a is chosen, i.e., bandit feedback. Without loss of generality, we assume $r_{\mathbf{x},a} \in [0, 1]$. Let $V(\pi)$ denote the expected reward of the policy π :

$$V(\pi) = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}), a \sim \pi(a|\mathbf{x})} [r_{\mathbf{x},a}]. \quad (1)$$

We look for a policy $\pi(a|\mathbf{x})$ to maximize $V(\pi)$. In the rest, we denote $\mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}), a \sim \pi(a|\mathbf{x})}[\cdot]$ as $\mathbb{E}_{\pi}[\cdot]$.

In off-policy learning, one can only access a set of logged feedback data $D := \{(\mathbf{x}_n, a_n, r_{\mathbf{x}_n, a_n})\}_{n=1}^N \subseteq \mathcal{X} \times \mathcal{A} \times [0, 1]$. Given \mathbf{x}_n , the action a_n was generated by a stochastic logging policy β , i.e., $a_n \sim \beta(\cdot|\mathbf{x}_n)$, which is usually different from the learning policy $\pi(a|\mathbf{x})$ [24, 40, 8]. The actions a_1, \dots, a_N and their corresponding rewards $r_{\mathbf{x}_1, a_1}, \dots, r_{\mathbf{x}_N, a_N}$ are generated independently given β . The main challenge is then to address the distributional discrepancy between $\beta(a|\mathbf{x})$ and $\pi(a|\mathbf{x})$, when optimizing $\pi(a|\mathbf{x})$ to maximize $V(\pi)$ with access only to the logged dataset D .

One of the most widely used methods to address the distribution shift between $\pi(a|\mathbf{x})$ and $\beta(a|\mathbf{x})$ is the Inverse Propensity Scoring (IPS) [8, 25]. One can easily get that:

$$V(\pi) = \mathbb{E}_\beta \left[\frac{\pi(a|\mathbf{x})}{\beta(a|\mathbf{x})} r_{\mathbf{x},a} \right],$$

yielding the following empirical estimator of $V(\pi)$:

$$\hat{V}_{\text{IPS}}(\pi) = \frac{1}{N} \sum_{n=1}^N \frac{\pi(a_n|\mathbf{x}_n)}{\beta(a_n|\mathbf{x}_n)} r_{\mathbf{x}_n, a_n}, \quad (2)$$

where $\pi(a_n|\mathbf{x}_n)/\beta(a_n|\mathbf{x}_n)$ is referred to as the propensity score. Various algorithms can be readily used for policy optimization under $\hat{V}_{\text{IPS}}(\pi)$, including value-based methods [33] and policy-based methods [19, 31, 42]. In this work, we adopt a well-known policy gradient algorithm, REINFORCE [42]. Assume the policy $\pi(a|\mathbf{x})$ is parameterized by $\boldsymbol{\vartheta}$, via the ‘‘log-trick’’, the gradient of $\hat{V}_{\text{IPS}}(\pi_{\boldsymbol{\vartheta}})$ with respect to $\boldsymbol{\vartheta}$ can be readily derived as,

$$\nabla_{\boldsymbol{\vartheta}} \hat{V}_{\text{IPS}}(\pi_{\boldsymbol{\vartheta}}) = \frac{1}{N} \sum_{n=1}^N \frac{\pi_{\boldsymbol{\vartheta}}(a_n|\mathbf{x}_n)}{\beta(a_n|\mathbf{x}_n)} r_{\mathbf{x}_n, a_n} \nabla_{\boldsymbol{\vartheta}} \log(\pi_{\boldsymbol{\vartheta}}(a_n|\mathbf{x}_n)).$$

Approximation with an unknown logging policy. In many real-world applications, the ground-truth logging probabilities, i.e., $\beta(a|\mathbf{x})$ of each observation (\mathbf{x}, a) in D , are unknown. As a typical walk-around, previous work employs supervised learning methods such as logistic regression [30] and neural networks [8] to estimate the logging policy, and replaces $\beta(a|\mathbf{x})$ with its estimated value $\hat{\beta}(a|\mathbf{x})$ to get the following BIPS estimator for policy learning:

$$\hat{V}_{\text{BIPS}}(\pi_{\boldsymbol{\vartheta}}) = \frac{1}{N} \sum_{n=1}^N \frac{\pi_{\boldsymbol{\vartheta}}(a_n|\mathbf{x}_n)}{\hat{\beta}(a_n|\mathbf{x}_n)} r_{\mathbf{x}_n, a_n}. \quad (3)$$

However, as shown in the following proposition, inaccurate $\hat{\beta}(a|\mathbf{x})$ leads to high bias and variance in BIPS. Worse still, smaller and inaccurate $\hat{\beta}(a|\mathbf{x})$ further enlarges this bias and variance.

Proposition 2.1. *The bias and variance of $\hat{V}_{\text{BIPS}}(\pi_{\boldsymbol{\vartheta}})$ can be derived as follows:*

$$\begin{aligned} \text{Bias} \left(\hat{V}_{\text{BIPS}}(\pi_{\boldsymbol{\vartheta}}) \right) &= \mathbb{E}_D \left[\hat{V}_{\text{BIPS}}(\pi_{\boldsymbol{\vartheta}}) - V(\pi_{\boldsymbol{\vartheta}}) \right] = \mathbb{E}_{\pi_{\boldsymbol{\vartheta}}} \left[r_{\mathbf{x},a} \left(\frac{\beta(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} - 1 \right) \right] \\ N \text{Var}_D \left(\hat{V}_{\text{BIPS}}(\pi_{\boldsymbol{\vartheta}}) \right) &= \text{Var}_{\pi_{\boldsymbol{\vartheta}}} \left(\frac{\beta(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} r_{\mathbf{x},a} \right) + \mathbb{E}_{\pi_{\boldsymbol{\vartheta}}} \left[\left(\frac{\pi_{\boldsymbol{\vartheta}}(a|\mathbf{x})}{\beta(a|\mathbf{x})} - 1 \right) \frac{\beta(a|\mathbf{x})^2}{\hat{\beta}(a|\mathbf{x})^2} r_{\mathbf{x},a}^2 \right] \end{aligned}$$

However, smaller $\hat{\beta}(a|\mathbf{x})$ usually implies less number of related training samples in the logged data, and thus $\hat{\beta}(a|\mathbf{x})$ can be inaccurate with a higher probability. To make it more explicit, let us revisit the empirical results shown in Figure 1. We followed the method introduced in [8] to estimate the logging policy on KuaiRec dataset [12] and plotted the estimated $\hat{\beta}(a|\mathbf{x})$ and its corresponding uncertainties on items of different observation frequencies in the logged dataset. We adopted the method in [45] to measure the confidence interval of $\hat{\beta}(a|\mathbf{x})$ on each instance. A wider confidence interval, i.e., higher uncertainty in estimation, implies that with a high probability the true value may be further away from the empirical mean estimate. We can observe in Figure 1 that as item frequency decreases, the estimated logging probability also decreases, but the estimation uncertainty increases. This implies that a smaller $\hat{\beta}(a|\mathbf{x})$ is usually 1) more inaccurate and 2) associated with a higher uncertainty.

As a result, with high bias and variance caused by inaccurate $\hat{\beta}(a|\mathbf{x})$, it is erroneous to learn $\pi_{\boldsymbol{\vartheta}}(a|\mathbf{x})$ by simply optimizing $\hat{V}_{\text{BIPS}}(\pi_{\boldsymbol{\vartheta}})$. Furthermore, this approach may also hinder the convergence of off-policy learning, as discussed later in Section 3.2.

3 Uncertainty-aware off-policy learning

Our idea is to consider the uncertainty of the estimated logging policy by incorporating per-sample weight $\phi_{\mathbf{x},a}$, and perform policy learning by optimizing the following empirical estimator:

$$\hat{V}_{\text{UIPS}}(\pi_{\boldsymbol{\vartheta}}) = \frac{1}{N} \sum_{n=1}^N \frac{\pi_{\boldsymbol{\vartheta}}(a_n|\mathbf{x}_n)}{\hat{\beta}(a_n|\mathbf{x}_n)} \phi_{\mathbf{x}_n, a_n} r_{\mathbf{x}_n, a_n}. \quad (4)$$

Intuitively, one should assign lower weights to samples whose $\hat{\beta}(a/j\mathbf{x})$ is small and far away from the ground-truth $\beta(a/j\mathbf{x})$. We then divide off-policy optimization into two iterative steps:

- **Deriving the optimal instance weight:** Find the optimal $\phi_{\mathbf{x},a}$ to make $\hat{V}_{\text{UIPS}}(\pi_{\vartheta})$ approach its ground-truth $V(\pi)$ as closely as possible, so as to facilitate policy learning. The derived optimal weight is denoted as $\phi_{\mathbf{x},a}$ (see Theorem 3.2).
- **Policy improvement:** Update the policy $\pi_{\vartheta}(a/j\mathbf{x})$ using the following gradient:

$$\nabla_{\vartheta} \hat{V}_{\text{UIPS}}(\pi_{\vartheta}) = \frac{1}{N} \sum_{n=1}^N \frac{\pi_{\vartheta}(a_n/j\mathbf{x}_n)}{\hat{\beta}(a_n/j\mathbf{x}_n)} \phi_{\mathbf{x}_n, a_n} (r_{\mathbf{x}_n, a_n} - \pi_{\vartheta}(a_n/j\mathbf{x}_n)) \quad (5)$$

The whole algorithm framework and its computational cost, as well as important notations are summarized in Appendix 7.1.

3.1 Derive the optimal uncertainty-aware instance weight

We expect to find the optimal weight $\phi_{\mathbf{x},a}$ to make the empirical estimator $\hat{V}_{\text{UIPS}}(\pi_{\vartheta})$ as accurate as possible, taking into account the uncertainty in estimated logging probabilities. Intuitively, a high accuracy of the estimator is crucial for determining the correct direction of policy learning. We follow previous work [36, 29] and measure the mean squared error (MSE) of $\hat{V}_{\text{UIPS}}(\pi_{\vartheta})$ to the ground-truth policy value $V(\pi_{\vartheta})$, which captures both the bias and variance of an estimator. A lower MSE indicates a more accurate estimator.

In UIPS, instead of directly minimizing the MSE, which is intractable, we find $\phi_{\mathbf{x},a}$ to minimize the upper bound of MSE. As we show later, the optimal $\phi_{\mathbf{x},a}$ has a closed-form solution which relates to both the value of $\pi_{\vartheta}(a/j\mathbf{x})/\hat{\beta}(a/j\mathbf{x})$ and the estimation uncertainty of $\hat{\beta}(a/j\mathbf{x})$.

Theorem 3.1. *The mean squared error (MSE) between $\hat{V}_{\text{UIPS}}(\pi_{\vartheta})$ and ground-truth estimator $V(\pi_{\vartheta})$ is upper bounded as follows:*

$$\begin{aligned} \text{MSE}(\hat{V}_{\text{UIPS}}(\pi_{\vartheta})) &= \mathbb{E}_D \left[\left(\hat{V}_{\text{UIPS}}(\pi_{\vartheta}) - V(\pi_{\vartheta}) \right)^2 \right] = \text{Bias}(\hat{V}_{\text{UIPS}}(\pi_{\vartheta}))^2 + \text{Var}(\hat{V}_{\text{UIPS}}(\pi_{\vartheta})) \\ &= \mathbb{E}_{\pi_{\vartheta}} \left[r_{\mathbf{x},a}^2 \frac{\pi_{\vartheta}(a/j\mathbf{x})}{\beta(a/j\mathbf{x})} \right] - \mathbb{E}_{\beta} \left[\left(\frac{\beta(a/j\mathbf{x})}{\hat{\beta}(a/j\mathbf{x})} \phi_{\mathbf{x},a} - 1 \right)^2 \right] + \mathbb{E}_{\beta} \left[\frac{\pi_{\vartheta}(a/j\mathbf{x})^2}{\hat{\beta}(a/j\mathbf{x})^2} \phi_{\mathbf{x},a}^2 \right]. \end{aligned}$$

As the first expectation term $\mathbb{E}_{\pi_{\vartheta}} \left[r_{\mathbf{x},a}^2 \frac{\pi_{\vartheta}(a/j\mathbf{x})}{\beta(a/j\mathbf{x})} \right]$ is a non-negative constant, we denote it as $\lambda \geq [0, 1)$ when searching for $\phi_{\mathbf{x},a}$. To minimize this upper bound of MSE, the optimal $\phi_{\mathbf{x},a}$ for each sample (\mathbf{x}, a) should minimize the following,

$$\lambda \left(\frac{\beta(a/j\mathbf{x})}{\hat{\beta}(a/j\mathbf{x})} \phi_{\mathbf{x},a} - 1 \right)^2 + \frac{\pi_{\vartheta}(a/j\mathbf{x})^2}{\hat{\beta}(a/j\mathbf{x})^2} \phi_{\mathbf{x},a}^2. \quad (6)$$

An interesting observation is that setting $\phi_{\mathbf{x},a} = \frac{\hat{\beta}(a/j\mathbf{x})}{\beta(a/j\mathbf{x})}$, i.e., turning $\frac{\pi(a/j\mathbf{x})}{\hat{\beta}(a/j\mathbf{x})} \phi_{\mathbf{x},a}$ into $\frac{\pi(a/j\mathbf{x})}{\beta(a/j\mathbf{x})}$ does not result in the optimal solution of Eq.(6). This is because such a setting only reduces bias (i.e., the first term of Eq.(6)), but fails to control the second term, which is related to the variance. Moreover, we cannot directly minimize Eq.(6) due to the unknown $\beta(a/j\mathbf{x})$. But it is possible to obtain a confidence interval which contains $\beta(a/j\mathbf{x})$ with a high probability, when $\hat{\beta}(a/j\mathbf{x})$ is obtained via a specific estimator, e.g., (generalized) linear model or kernel methods.

Following previous work [23, 16, 22], we adopt the realizable assumption that $\beta(a/j\mathbf{x})$ can be represented by a softmax function applied over a parametric function $f_{\theta^*}(\mathbf{x}, a)$. Moreover, the universal approximation theorem [18] states that a parametric function with sufficient capacity, when combined with a softmax function, can approximate any distribution. Then we have:

$$\beta(a/j\mathbf{x}) \propto \exp(f_{\theta^*}(\mathbf{x}, a)), \hat{\beta}(a/j\mathbf{x}) \propto \exp(f_{\theta}(\mathbf{x}, a)), \quad (7)$$

where $f_{\theta}(\mathbf{x}, a)$ is an estimate of $f_{\theta^*}(\mathbf{x}, a)$. Following the conventional definition of confidence interval [20], we define γ and $U_{\mathbf{x},a}$ such that $|f_{\theta^*}(\mathbf{x}, a) - f_{\theta}(\mathbf{x}, a)| \leq \gamma U_{\mathbf{x},a}$ holds with probability

at least $1-\delta$, where γ is a function of δ (typically the smaller δ is, the larger γ is). Then $\gamma U_{\mathbf{x},a}$ measures the width of confidence interval of $f_{\theta}(\mathbf{x}, a)$ against its ground-truth $f_{\theta^*}(\mathbf{x}, a)$. As derived in Appendix 7.2, with probability at least $1-\delta$, we have $\beta(a/\mathbf{x}) \geq \mathbf{B}_{\mathbf{x},a}$ and

$$\mathbf{B}_{\mathbf{x},a} = \left[\frac{\hat{Z} \exp(-\gamma U_{\mathbf{x},a})}{Z} \hat{\beta}(a/\mathbf{x}), \frac{\hat{Z} \exp(\gamma U_{\mathbf{x},a})}{Z} \hat{\beta}(a/\mathbf{x}) \right],$$

where $Z = \sum_{a^o} \exp(f_{\theta}(a^o/\mathbf{x}))$ and $\hat{Z} = \sum_{a^o} \exp(f_{\theta}(a^o/\mathbf{x}))$.

As $\beta(a/\mathbf{x})$ can be any value in $\mathbf{B}_{\mathbf{x},a}$ with high probability, we aim to find the optimal $\phi_{\mathbf{x},a}$ that minimizes the worst case of Eq.(6), thereby ensuring that $\hat{V}_{\text{UIPS}}(\pi_{\vartheta})$ approaches its ground-truth $V(\pi_{\vartheta})$ under the sense of MSE, even in the worst possible scenarios. This ensures the subsequent policy improvement direction will not be much worse with high probability. Thus, we formulate the following optimization problem:

$$\min_{\phi_{\mathbf{x},a}} \max_{\beta_{\mathbf{x},a} \geq \mathbf{B}_{\mathbf{x},a}} \lambda \left(\frac{\beta_{\mathbf{x},a}}{\hat{\beta}(a/\mathbf{x})} \phi_{\mathbf{x},a} - 1 \right)^2 + \frac{\pi_{\vartheta}(a/\mathbf{x})^2}{\hat{\beta}(a/\mathbf{x})^2} \phi_{\mathbf{x},a}^2. \quad (8)$$

The following theorem derives a closed-form formula for the optimal solution of Eq.(8).

Theorem 3.2. *Let $\eta \geq [\exp(-\gamma U_{\mathbf{x}}^{\max}), \exp(\gamma U_{\mathbf{x}}^{\max})]$, where $U_{\mathbf{x}}^{\max} = \max_a U_{\mathbf{x},a}$. The optimization problem in Eq.(8) has a closed-form solution:*

$$\phi_{\mathbf{x},a} = \min \left(\lambda / \left[\frac{\lambda}{\eta} \exp(-\gamma U_{\mathbf{x},a}) + \frac{\eta \pi_{\vartheta}(a/\mathbf{x})^2}{\hat{\beta}(a/\mathbf{x})^2 \exp(\gamma U_{\mathbf{x},a})} \right], 2\eta / \left[\exp(\gamma U_{\mathbf{x},a}) + \exp(-\gamma U_{\mathbf{x},a}) \right] \right).$$

The following corollary demonstrates the advantage of UIPS. The detailed proof of Theorem 3.2 and Corollary 3.3 can be found in Appendix 7.8.

Corollary 3.3. *With $\phi_{\mathbf{x},a}$ derived in Theorem 3.2, $\hat{V}_{\text{UIPS}}(\pi_{\vartheta})$ in Eq.(4) achieves a smaller upper bound of MSE than $\hat{V}_{\text{BIPS}}(\pi_{\vartheta})$ in Eq. (3).*

Insights about $\phi_{\mathbf{x},a}$. The detailed analysis of the effect of $\phi_{\mathbf{x},a}$ can be found in Lemma 7.1 in Appendix 7.8. In summary, we have the following key findings,

- For samples whose largest possible propensity score is under control: i.e., $\frac{\pi_{\vartheta}(a/\mathbf{x})}{\min_{\mathbf{B}_{\mathbf{x},a}} \beta} < \frac{P^-}{\lambda}$, higher uncertainty implies smaller values of $\pi/\hat{\beta}$. This suggests samples of this type with positive rewards are underestimated, and the extend of underestimation increases with the estimation uncertainty. UIPS thus chooses to increase $\phi_{\mathbf{x},a}$ with uncertainty, to emphasize these long-tail positive samples.
- Conversely, for samples with large propensity scores, UIPS decreases $\phi_{\mathbf{x},a}$ as the uncertainty increases, so as to prevent their distortion in policy learning.

Uncertainty estimation. Now we describe how to calculate $U_{\mathbf{x},a}$, i.e., the uncertainty of the estimated $\hat{\beta}(a/\mathbf{x})$. In this work, we choose to estimate $\beta(a/\mathbf{x})$ using a neural network, because 1) its representation learning capacity has been proved in numerous studies, and 2) various ways [11, 45] can be leveraged to perform the uncertainty estimation in a neural network. We adopt [45] due to its computational efficiency and theoretical soundness. Following the proof of Theorem 4.4 in [45], given the logged dataset D , we can get with a high probability that there exists γ such that:

$$|f_{\theta}(\mathbf{x}_n, a_n) - f_{\theta^*}(\mathbf{x}_n, a_n)| \leq \gamma \sqrt{\mathbf{g}(\mathbf{x}_n, a_n)^\top \mathbf{M}_D^{-1} \mathbf{g}(\mathbf{x}_n, a_n)}$$

where $\mathbf{g}(\mathbf{x}_n, a_n)$ is the gradient of $f_{\theta}(\mathbf{x}_n, a_n)$ with respect to the neural network's last layer's parameter $\theta_w = \theta$, i.e., $\mathbf{g}(\mathbf{x}_n, a_n) = \nabla_{\theta_w} f_{\theta}(\mathbf{x}_n, a_n)$. And $\mathbf{M}_D = \sum_{n=1}^N \mathbf{g}(\mathbf{x}_n, a_n) \mathbf{g}(\mathbf{x}_n, a_n)^\top$, implying $U_{\mathbf{x}_n, a_n} = \sqrt{\mathbf{g}(\mathbf{x}_n, a_n)^\top \mathbf{M}_D^{-1} \mathbf{g}(\mathbf{x}_n, a_n)}$.

3.2 Convergence of policy learning under UIPS

The following theorem provides the convergence result for UIPS, which converges to a stationary point of the expected reward function. The proof is provided in Appendix 7.9.

Theorem 3.4. Denote G_{\max} and Φ as the maximum value of $k \frac{\partial \pi_{\vartheta}(a|\mathbf{x})}{\partial \vartheta} k$ and $E_{\beta} \left[\frac{\pi_{\vartheta}^2(a|\mathbf{x})}{\beta^2(a|\mathbf{x})} (\phi_{\mathbf{x},a})^2 \right]$ respectively, i.e., $k \frac{\partial \pi_{\vartheta}(a|\mathbf{x})}{\partial \vartheta} k \leq G_{\max}$ and $E_{\beta} \left[\frac{\pi_{\vartheta}^2(a|\mathbf{x})}{\beta^2(a|\mathbf{x})} (\phi_{\mathbf{x},a})^2 \right] \leq \Phi$. And denote V_{\max} as the finite maximum expected reward that can be achieved, and $\varphi_{\max} = \max_{\mathbf{x},a} \left\{ \left| \frac{\beta(a|\mathbf{x})}{\beta(a|\mathbf{x})} \phi_{\mathbf{x},a} - 1 \right| \right\}$. Assume that the expected reward of π_{ϑ} , i.e., $V(\pi_{\vartheta})$, is a differentiable and L -smooth function w.r.t ϑ . Denote the policy parameters obtained by Eq.(5) at iteration $k \in [K]$ as ϑ_k , then $\varphi_{\max} \in (0, 1)$ and

$$\frac{1}{K} \sum_{k=1}^K E[k \nabla V(\pi_{\vartheta_k}) k]^2 \leq \frac{2LV_{\max}}{K(1 - \varphi_{\max})} + \left(L + \frac{2V_{\max}}{(1 - \varphi_{\max})} \right) \frac{G_{\max} \Phi}{K},$$

where $\nabla V(\pi_{\vartheta})$ is the true policy gradient under ground-truth logging probability, i.e., $\nabla V(\pi_{\vartheta}) = E_{\beta} \left[\frac{\pi_{\vartheta}(a|\mathbf{x})}{\beta(a|\mathbf{x})} r_{\mathbf{x},a} \nabla_{\vartheta} \log(\pi_{\vartheta}(a|\mathbf{x})) \right]$.

Theorem 3.4 shows that, as $K \rightarrow \infty$ and with $1/(1 - \varphi_{\max})$ and Φ being controlled, UIPS leads policy update to converge to a stationary point where the true policy gradient $\nabla V(\pi_{\vartheta_k})$ is zero. And fortunately, UIPS is effective in controlling both $1/(1 - \varphi_{\max})$ and Φ . Specifically, we denote $\varphi_{\mathbf{x},a} = \left| \frac{\beta(a|\mathbf{x})}{\beta(a|\mathbf{x})} \phi_{\mathbf{x},a} - 1 \right|$ and $\Phi_{\mathbf{x},a} = \frac{\pi_{\vartheta}^2(a|\mathbf{x})}{\beta^2(a|\mathbf{x})} (\phi_{\mathbf{x},a})^2$. It is clear to note that $\lambda \varphi_{\mathbf{x},a}^2 + \Phi_{\mathbf{x},a}$ corresponds to the objective in Eq.(6) for deriving $\phi_{\mathbf{x},a}$ for each sample (\mathbf{x}, a) . In other words, UIPS selects $f_{\phi_{\mathbf{x},a},g}$ to minimize $\varphi_{\max} = \max f_{\phi_{\mathbf{x},a},g}$ and $\Phi = E_{\beta} [\Phi_{\mathbf{x},a}]$, which directly accelerate the policy converge to a stationary point with the true policy gradient being zero.

In the case of BIPS in Eq.(3), we have $\phi_{\mathbf{x},a} = 1$. Although Φ may be large due to small logging probabilities, the more concerning issue is that the requirement $\varphi_{\max} \in (0, 1)$ is no longer satisfied when $\beta(a|\mathbf{x}) = 2\hat{\beta}(a|\mathbf{x})$, which may happen with a non-negligible probability. Hence, the convergence of policy learning under \hat{V}_{BIPS} is no better than that under UIPS.

4 Empirical evaluations

We evaluate UIPS on both synthetic data and three real-world datasets with unbiased collection. We compare UIPS with the following baselines, which can be grouped into five categories:

- **Cross-Entropy (CE)**: A supervised learning method with the cross-entropy loss over its softmax output. No off-policy correction is performed in this method.
- **BIPS-Cap** [8]: The off-policy learning solution under the BIPS estimator in Eq.(3). The estimated propensity scores are further suppressed to control variance, i.e., taking $\min \left(c, \frac{\pi_{\vartheta}(a|\mathbf{x})}{\beta(a|\mathbf{x})} \right)$ as the propensity score. Setting c to a small value can reduce variance, but introduces bias.
- **MinVar & stableVar** [46], **Shrinkage** [36]: This line of work improves off-policy evaluation by reweighing each sample. For example, MinVar and stableVar reweigh each sample by $\frac{h_{\mathbf{x},a}}{\sum_{a^0} h_{\mathbf{x},a^0}}$ with $h_{\mathbf{x},a} = \frac{\hat{\beta}(a|\mathbf{x})}{\pi_{\vartheta}(a|\mathbf{x})^2}$ and $h_{\mathbf{x},a} = \frac{\hat{\beta}(a|\mathbf{x})}{\pi_{\vartheta}(a|\mathbf{x})}$ respectively, since they find that $\pi_{\vartheta}(a|\mathbf{x})^2 / \hat{\beta}(a|\mathbf{x})$ is directly related to policy evaluation variance. Su et al. [36] propose to shrink the propensity score by $\lambda / (\lambda + \frac{\pi_{\vartheta}(a|\mathbf{x})^2}{\hat{\beta}(a|\mathbf{x})^2})$, which is a special case of our UIPS with $U_{\mathbf{x},a} = 0$ and $\eta = 1$. All these methods simply treat $\hat{\beta}(a|\mathbf{x})$ as $\beta(a|\mathbf{x})$, and none of them consider the uncertainty of $\hat{\beta}(a|\mathbf{x})$.
- **SNIPS** [39], **BanditNet** [16], **POEM** [38], **POXM** [23], **Adaptive** [22]: This line of work aims for more stable and accurate policy learning. For example, SNIPS normalizes the estimator by the sum of propensity scores in each batch. BanditNet extends SNIPS and leverages an additional Lagrangian term to normalize the estimator by an approximated sum of propensity scores of all samples. POEM jointly optimizes the estimator and its variance. POXM controls estimation variance by pruning samples with small logging probabilities. Adaptive proposes a new formulation to utilize negative samples.
- **ApproxKNN** [5] and **IPS-C-TS**: The line of work improves off-policy learning by applying calibration to estimated logging probabilities. ApproxKNN utilizes the K-Nearest Neighbor algorithm for calibration, which exhibits the lowest calibration error in [5]. IPS-C-TS, on the other hand, employs temperature scaling, a widely recognized and effective calibration method for probability distribution [13].

Table 1: Experiment results on synthetic datasets. The best and second best results are highlighted with **bold** and underline respectively. The p -value under the t-test between UIPS and the best baseline on each dataset is also provided.

Algorithm	$\tau = 0.5$			$\tau = 1$			$\tau = 2$		
	P@5	R@5	NDCG@5	P@5	R@5	NDCG@5	P@5	R@5	NDCG@5
IPS-GT	0.5589 $1e^{-3}$	0.1582 $6e^{-4}$	0.6093 $1e^{-3}$	0.5526 $2e^{-3}$	0.1565 $6e^{-4}$	0.6007 $1e^{-3}$	0.5531 $2e^{-3}$	0.1557 $7e^{-4}$	0.6037 $1e^{-3}$
CE	0.5553 $6e^{-4}$	0.1573 $2e^{-4}$	0.6037 $5e^{-4}$	0.5510 $6e^{-4}$	0.1561 $2e^{-4}$	0.5995 $4e^{-4}$	0.5386 $2e^{-3}$	0.1524 $7e^{-4}$	0.5874 $2e^{-3}$
BIPS-Cap	0.5515 $2e^{-3}$	0.1553 $8e^{-4}$	0.6031 $2e^{-3}$	0.5526 $2e^{-3}$	0.1561 $6e^{-4}$	0.6016 $1e^{-3}$	0.5409 $3e^{-3}$	0.1529 $9e^{-4}$	0.5901 $2e^{-3}$
MinVar	0.5340 $2e^{-3}$	0.1509 $6e^{-4}$	0.5857 $2e^{-3}$	0.5282 $2e^{-3}$	0.1491 $7e^{-4}$	0.5791 $2e^{-3}$	0.5036 $4e^{-3}$	0.1415 $1e^{-3}$	0.5543 $3e^{-3}$
stableVar	0.4577 $5e^{-3}$	0.1310 $1e^{-3}$	0.5111 $2e^{-3}$	0.5373 $3e^{-3}$	0.1523 $9e^{-4}$	0.5866 $3e^{-3}$	0.5279 $3e^{-3}$	0.1492 $8e^{-4}$	0.5781 $3e^{-3}$
Shrinkage	0.5526 $2e^{-3}$	0.1562 $7e^{-4}$	0.6024 $1e^{-3}$	0.5499 $4e^{-3}$	0.1545 $1e^{-3}$	0.6040 $3e^{-3}$	0.5347 $2e^{-3}$	0.1513 $6e^{-4}$	0.5824 $2e^{-3}$
SNIPS	0.2616 $6e^{-2}$	0.0749 $2e^{-2}$	0.3150 $7e^{-2}$	0.3538 $5e^{-2}$	0.0987 $1e^{-2}$	0.4144 $6e^{-2}$	0.4379 $3e^{-2}$	0.1226 $9e^{-3}$	0.5177 $3e^{-2}$
BanditNet	0.4011 $3e^{-2}$	0.1131 $8e^{-3}$	0.4830 $2e^{-2}$	0.3894 $4e^{-2}$	0.1095 $1e^{-2}$	0.4741 $3e^{-2}$	0.4122 $3e^{-2}$	0.1153 $8e^{-3}$	0.4934 $3e^{-2}$
POEM	0.5480 $2e^{-3}$	0.1539 $8e^{-4}$	0.6008 $2e^{-3}$	0.5502 $2e^{-3}$	0.1551 $6e^{-4}$	0.6000 $2e^{-3}$	0.5399 $2e^{-3}$	0.1526 $8e^{-4}$	0.5893 $2e^{-3}$
POXM	0.4006 $3e^{-2}$	0.1130 $8e^{-3}$	0.4828 $2e^{-2}$	0.3616 $4e^{-2}$	0.1019 $1e^{-2}$	0.4522 $4e^{-2}$	0.3816 $4e^{-2}$	0.1069 $1e^{-2}$	0.4680 $4e^{-2}$
Adaptive	0.3831 $2e^{-2}$	0.1050 $4e^{-3}$	0.4382 $2e^{-2}$	0.4734 $4e^{-3}$	0.1325 $1e^{-3}$	0.5326 $3e^{-3}$	0.3936 $1e^{-2}$	0.1097 $4e^{-3}$	0.4368 $2e^{-2}$
ApproxKNN	0.5576 $1e^{-3}$	0.1580 $4e^{-4}$	0.6059 $2e^{-3}$	0.5527 $9e^{-4}$	0.1567 $1e^{-4}$	0.6010 $1e^{-3}$	0.5409 $2e^{-3}$	0.1532 $6e^{-4}$	0.5890 $1e^{-3}$
IPS-C-TS	0.5565 $1e^{-3}$	0.1577 $3e^{-4}$	0.6048 $7e^{-4}$	0.5517 $6e^{-4}$	0.1563 $2e^{-4}$	0.6002 $6e^{-4}$	0.5393 $1e^{-3}$	0.1526 $5e^{-4}$	0.5879 $1e^{-3}$
UIPS-P	0.4019 $3e^{-2}$	0.1131 $1e^{-2}$	0.4831 $3e^{-2}$	0.3904 $4e^{-2}$	0.1096 $1e^{-2}$	0.4749 $3e^{-2}$	0.4109 $3e^{-2}$	0.1149 $1e^{-2}$	0.4922 $3e^{-2}$
UIPS-O	0.4135 $4e^{-2}$	0.1167 $1e^{-2}$	0.4954 $4e^{-2}$	0.3896 $4e^{-2}$	0.1096 $1e^{-2}$	0.4739 $3e^{-2}$	0.4519 $3e^{-2}$	0.1268 $8e^{-3}$	0.5296 $2e^{-2}$
UIPS	0.5608 $2e^{-3}$	0.1589 $8e^{-4}$	0.6113 $3e^{-3}$	0.5572 $2e^{-3}$	0.1571 $8e^{-4}$	0.6074 $2e^{-3}$	0.5432 $3e^{-3}$	0.1534 $8e^{-4}$	0.5946 $2e^{-3}$
p-value	$3e^{-3}$	$1e^{-2}$	$4e^{-5}$	$2e^{-5}$	$2e^{-1}$	$4e^{-10}$	$1e^{-1}$	$5e^{-1}$	$4e^{-2}$

Table 2: Performance under different uncertainties.

Algorithm	Actions on Samples with High Uncertainty			Actions on Samples with Low Uncertainty		
	P@5(RI)	R@5(RI)	NDCG@5(RI)	P@5(RI)	R@5(RI)	NDCG@5(RI)
CE	0.5190	0.1231	0.5526	0.5913	0.1915	0.6549
BIPS-Cap	0.5117 (-1.41%)	0.1202 (-2.33%)	0.5488 (-0.68%)	0.5913 (+0.00%)	0.1903 (-0.64%)	0.6574 (+0.39%)
Shrinkage	0.5158 (-0.62%)	0.1217 (-1.11%)	0.5505 (-0.37%)	0.5892 (-0.35%)	0.1905 (-0.55%)	0.6546 (-0.05%)
UIPS	0.5222 (+0.61%)	0.1237 (+0.50%)	0.5568 (+0.77%)	0.5994 (+1.38%)	0.1940 (+1.28%)	0.6658 (+1.66%)

- **UIPS-P** and **UIPS-O**: These are two variants of our UIPS with different ways of leveraging uncertainties. UIPS-P directly penalizes samples whose estimated logging probabilities are of high uncertainty, i.e., taking $\phi_{x,a} = 1.0 / \exp(\gamma U_{x,a})$, which follows previous work on offline reinforcement learning [43, 4]. UIPS-O adversarially uses the worst propensity score for policy learning, i.e., $\phi_{x,a} = 1.0 / \exp(-\gamma U_{x,a})$.

4.1 Synthetic data

Data generation. Following previous work [24, 23], we generate a synthetic dataset by a supervise-to-bandit conversion on Wiki10-31K dataset [7], which is an extreme multi-label classification dataset. The Wiki10-31K dataset contains approximately 20K samples. Each sample is associated with a feature vector \tilde{x} of 101,938 dimensions and a label vector $\mathbf{y}_{\tilde{x}}$ of 31K classes with more than one positive class. Let $\mathbf{y}_{\tilde{x},a}$ denote the label of class a under \tilde{x} , and we take each class as an action. The huge action space creates great challenges in off-policy learning, e.g., sparse observations, and therefore better evaluates different methods.

We split the dataset into train, validation and test sets with size 11K:3K:6K. The test set is from the official split. Since the original feature vector \tilde{x} is too sparse, for ease of learning, we first embedded it to dimension d by $\mathbf{x} = \mathbf{W}\tilde{x}$, and synthesized the ground-truth logging policy $\beta(a|\mathbf{x})$ by:

$$\beta(a|\mathbf{x}) \propto \exp(\mathbf{x}^\top \boldsymbol{\theta}_a / \tau), \quad (9)$$

where \mathbf{W} and $f\boldsymbol{\theta}_a g$ are pre-trained parameters by applying a logistic regression model on the train set, τ is a positive hyper-parameter that controls the skewness of logging distribution. A small value of τ leads to a near-deterministic logging policy, while a larger τ makes it flatter. More implementation details can be found in Appendix 7.3.

Evaluation metrics. To evaluate the learned policy $\pi_\vartheta(a|\mathbf{x})$, we calculate Precision@K (P@K), Recall@K (R@K) and NDCG@K following previous work [23, 24]. Higher P@K, R@K and NDCG@K imply a better policy.

Effectiveness of policy learning. Table 1 shows the average performance and standard deviations of all algorithms under 10 random seeds on three synthetic datasets generated under different τ . As the ground-truth logging policy is accessible on the synthetic datasets, we included a new baseline IPS-GT, which uses the IPS estimator with the ground-truth logging probabilities. We calculated p -value under t-test between UIPS and the best baseline to investigate the significance of improvement.

First, we can observe that UIPS achieved similar and even better performance than IPS-GT when $\tau = 0.5$ and $\tau = 1$, but performed worse than IPS-GT when $\tau = 2$. Despite using ground-truth

logging probabilities, IPS-GT still suffered from high variance caused by samples with small logging probabilities, which is the main cause of its worse performance when $\tau = 0.5$ and $\tau = 1$. In contrast, UIPS effectively controlled the negative impact of these high-variance samples, resulting in a better bias-variance trade-off.

With an increasing τ , suggesting a decrease in the probability of selecting positive actions, most algorithms experienced a drop in performance. However, UIPS consistently outperformed all other algorithms across all three datasets and metrics. Interestingly, as τ decreases, the performance improvement of UIPS became even more pronounced, despite SNIPS, BanditNet, and POXM being designed to handle small logging probabilities of positive actions.

ApproxKNN and IPS-C-TS generally achieved better performance than BIPS-Cap, implying the effectiveness of calibration of estimated logging probabilities. However, UIPS still consistently outperformed both ApproxKNN and IPS-C-TS. The main reason is that calibration primarily focuses on adjusting the estimated probabilities to ensure on *average* the model’s predictions are reliable and accurate. In contrast, UIPS specifically handles the impact from each *individual* sample in policy learning.

UIPS also consistently outperformed Shrinkage (a special case of UIPS with uncertainties always being zero) on all three datasets, demonstrating the benefits of considering the estimation uncertainty. Finally, blindly reweighing through uncertainties, regardless of their impact on the accuracy of the resulting estimator and the learned policy, ultimately resulted in poor performance, as demonstrated by UIPS-P and UIPS-O.

Performance under different uncertainty levels. As shown in Figure 1, low-frequency samples in the logged dataset suffer higher uncertainties in their propensity estimation. Thus, we divided the test set into two subsets according to the average frequency of associated actions, where the uncertainty in the subset associated with low-frequency actions is on average 8% higher than that in high-frequency actions. Table 2 shows the results on these two subsets when $\tau = 0.5$. In addition, we include the results of the top three baselines that directly utilize the estimated logging policy. Table 2 clearly demonstrates that only UIPS performed better than CE on the test set with low-frequency actions, implying the distortion of inaccurately estimated logging probabilities and the effectiveness of UIPS in efficiently handling them.

Off-policy Evaluation. We further inspected whether \hat{V}_{UIPS} in Eq.(4) leads to more accurate off-policy evaluation. Following previous work [29, 46, 36], we evaluated the following ϵ -greedy policy: $\pi(a|\mathbf{x}) = \frac{1}{jM_x} \mathbb{1}_{\{a=j\}} + \epsilon/jA$, where M_x contains all positive actions associated with instance \mathbf{x} . For each \mathbf{x} in the test set, we randomly sample 100 actions following the logging policy in Eq.(9) to generate the logged dataset. Table 3 shows the MSE of the estimators to the ground-truth policy value under 20 different random seeds. From Table 3, one can observe that: 1) IPS-GT with a skewer logging policy (i.e., smaller τ) leads to higher MSE, consistent with previous findings [29, 46, 36]; 2) inaccurate logging probabilities result in high bias and variance, leading to much larger MSE of BIPS compared to IPS-GT. Furthermore, this distortion is particularly pronounced when the ground-truth logging policy is skewed ($\tau = 0.5$); and 3) although all using the estimated logging policy, \hat{V}_{UIPS} yields the smallest MSE, comparing to other baselines that are designed to improve over BIPS.

Hyper-parameter Tuning. Discussions about hyper-parameter tuning and performance of UIPS under different hyper-parameters can also be found in Appendix 7.3.1.

4.2 Real-world data

To demonstrate the effectiveness of UIPS in real-world scenarios, we evaluate it on three recommendation datasets: (1) Yahoo! R3¹; (2) Coat²; (3) KuaiRec [12], for music, fashion and short-video recommendations respectively. All these datasets contain an unbiased test set collected from a randomized controlled trial where items are randomly selected. The statistics of the three datasets and implementation details, e.g., model architectures and dataset splits, can be found in Appendix 7.3.2.

Following [10], we take $K = 5$ on Yahoo! R3 and Coat datasets, and $K = 50$ on KuaiRec dataset. The p -value under the t-test between UIPS and the best baseline on each dataset is also reported to investigate the significance of improvement.

¹<https://webscope.sandbox.yahoo.com/>

²<https://www.cs.cornell.edu/~schnabts/mnar/>

Table 3: MSE of different off-policy estimators. A lower MSE indicates a more accurate estimator.

Algorithm	IPS-GT	BIPS	minVar	stableVar	Shrinkage	UIPS
$\tau = 0.5$	0.0875 $4e^{-4}$	15.786 1.51	0.9021 $7e^{-13}$	0.8612 $5e^{-8}$	0.0718 $5e^{-6}$	0.0210 $2e^{-6}$
$\tau = 1.0$	0.0209 $8e^{-5}$	0.5510 0.388	0.9019 $8e^{-12}$	0.8578 $2e^{-7}$	0.1978 $2e^{-5}$	0.0093 $1e^{-6}$
$\tau = 2.0$	0.0020 $6e^{-6}$	0.5669 0.013	0.9015 $5e^{-15}$	0.8342 $5e^{-7}$	0.2952 $3e^{-5}$	0.0043 $4e^{-7}$

Table 4: Experimental results on real-world datasets. The best and second best results are highlighted with **bold** and underline respectively. The p-value under the t-test between UIPS and the best baseline on each dataset is also provided.

Algorithm	Yahoo			Coat			KuaiRec		
	P@5	R@5	NDCG@5	P@5	R@5	NDCG@5	P@50	R@50	NDCG@50
CE	0.2819 $2e^{-3}$	0.7594 $6e^{-3}$	0.6073 $7e^{-3}$	0.2799 $5e^{-3}$	0.4618 $1e^{-2}$	0.4529 $7e^{-3}$	0.8802 $2e^{-3}$	0.0240 $8e^{-5}$	0.8810 $6e^{-3}$
BIPS-Cap	0.2808 $2e^{-3}$	0.7576 $5e^{-3}$	0.6099 $8e^{-3}$	0.2758 $6e^{-3}$	0.4582 $7e^{-3}$	0.4399 $9e^{-3}$	0.8750 $3e^{-3}$	0.0238 $7e^{-5}$	0.8788 $5e^{-3}$
MinVar	0.2843 $4e^{-3}$	<u>0.7685</u> $1e^{-2}$	0.6168 $1e^{-2}$	0.2813 $3e^{-3}$	<u>0.4668</u> $9e^{-3}$	0.4414 $8e^{-3}$	0.8827 $1e^{-3}$	0.0240 $5e^{-5}$	0.8886 $2e^{-3}$
stableVar	0.2787 $2e^{-3}$	0.7499 $7e^{-3}$	0.5919 $7e^{-3}$	0.2840 $3e^{-3}$	0.4662 $5e^{-3}$	0.4393 $7e^{-3}$	0.8524 $7e^{-3}$	0.0231 $2e^{-4}$	0.8570 $4e^{-3}$
Shrinkage	0.2843 $3e^{-3}$	0.7654 $8e^{-3}$	0.6204 $7e^{-3}$	0.2790 $5e^{-3}$	0.4636 $4e^{-3}$	0.4464 $1e^{-2}$	0.8744 $3e^{-3}$	0.0238 $9e^{-5}$	0.8771 $6e^{-3}$
SNIPS	0.2222 $4e^{-3}$	0.5828 $1e^{-2}$	0.4357 $1e^{-2}$	0.2643 $7e^{-3}$	0.4287 $1e^{-2}$	0.4009 $9e^{-3}$	0.8411 $6e^{-3}$	0.0228 $2e^{-4}$	0.8431 $6e^{-3}$
BanditNet	0.2413 $8e^{-3}$	0.6442 $2e^{-2}$	0.4988 $2e^{-2}$	0.2781 $8e^{-3}$	0.4527 $1e^{-2}$	0.4251 $1e^{-2}$	0.8758 $5e^{-3}$	0.0239 $2e^{-4}$	0.8810 $4e^{-3}$
POEM	0.2732 $3e^{-3}$	0.7357 $1e^{-2}$	0.5880 $1e^{-2}$	0.2791 $4e^{-3}$	0.4566 $6e^{-3}$	0.4375 $6e^{-3}$	0.7785 $1e^{-2}$	0.0210 $2e^{-4}$	0.7779 $6e^{-3}$
POXM	0.2250 $5e^{-3}$	0.5940 $1e^{-2}$	0.4542 $2e^{-2}$	0.2663 $6e^{-3}$	0.4308 $9e^{-3}$	0.4006 $1e^{-2}$	0.8962 $1e^{-2}$	0.0245 $4e^{-4}$	0.9041 $1e^{-2}$
Adaptive	0.2762 $3e^{-3}$	0.7451 $9e^{-3}$	0.5919 $8e^{-3}$	0.2830 $3e^{-3}$	0.4634 $5e^{-3}$	0.4217 $5e^{-3}$	0.8375 $1e^{-2}$	0.0227 $4e^{-4}$	0.8460 $1e^{-2}$
ApproxKNN	0.2697 $2e^{-3}$	0.7225 $5e^{-3}$	0.5760 $6e^{-3}$	0.2755 $2e^{-3}$	0.4594 $5e^{-3}$	0.4490 $4e^{-3}$	0.8839 $2e^{-6}$	0.0240 $5e^{-5}$	0.8895 $2e^{-3}$
IPS-C-TS	0.2816 $2e^{-3}$	0.7582 $5e^{-3}$	0.6114 $5e^{-3}$	0.2799 $3e^{-3}$	0.4625 $7e^{-3}$	0.4462 $6e^{-3}$	0.8781 $3e^{-3}$	0.0239 $1e^{-4}$	0.8749 $3e^{-3}$
UIPS-P	0.1829 $8e^{-3}$	0.4560 $3e^{-2}$	0.3300 $1e^{-2}$	0.2685 $7e^{-3}$	0.4364 $9e^{-3}$	0.4087 $7e^{-3}$	0.8638 $8e^{-3}$	0.0235 $3e^{-4}$	0.8685 $7e^{-3}$
UIPS-O	0.1947 $3e^{-3}$	0.4959 $1e^{-2}$	0.3600 $8e^{-3}$	0.2657 $5e^{-3}$	0.4306 $9e^{-3}$	0.4146 $9e^{-3}$	0.8651 $8e^{-3}$	0.0235 $2e^{-4}$	0.8697 $7e^{-3}$
UIPS	0.2868 $2e^{-3}$	0.7742 $5e^{-3}$	0.6274 $5e^{-3}$	0.2877 $3e^{-3}$	0.4757 $5e^{-3}$	0.4576 $8e^{-3}$	0.9120 $1e^{-3}$	0.0250 $5e^{-5}$	0.9174 $7e^{-4}$
p-value	$4e^{-2}$	$1e^{-2}$	$3e^{-2}$	$2e^{-2}$	$6e^{-4}$	$5e^{-5}$	$6e^{-4}$	$6e^{-4}$	$1e^{-3}$

We can first observe from Table 4 that on all three datasets, the proposed UIPS achieved the highest precision, recall and NDCG. Apparently, accurate estimation of logging probabilities in real-world scenarios with large action spaces and sparse interactions is challenging to achieve, causing BIPS-Cap to underperform CE. Additionally, calibration poses difficulties in such scenarios, leading to poor performance of ApproxKNN and IPS-C-TS across all three real-world datasets. BanditNet, POEM and POXM performed better in problems with a larger action space, while MinVar, stableVar and Shrinkage as well as Adaptive suited better for scenarios with a smaller action size. Again, UIPS outperformed Shrinkage, highlighting the importance of handling uncertainty in the estimated logging policy. But simply reweighing based on uncertainties, without considering their impact on the accuracy of the resulting estimator and the learned policy, led to poor performance, as shown by UIPS-P and UIPS-O.

4.3 Comparisons against other lines of off-policy learning

We discuss about the difference between UIPS and recent work with direct propensity estimation [34, 21], the DICE line of work [26, 47, 9] as well as the work on the distributionally robust off-policy learning [32, 17, 44] in Appendix 7.6, Appendix 7.5 and Appendix 7.7 respectively. Our empirical evaluations on recent work [26, 44, 34] from all three lines suggest that these lines of work cannot properly handle inaccuracies in the estimated logging probabilities that hinder policy improvement. Moreover, we also integrated the proposed UIPS with the doubly-robust (DR) estimator and our results suggest the accuracy of the imputation model greatly affected policy learning under DR, but UIPS still provides benefits. More details can be found in Appendix 7.4.

5 Related work

Our work is the first of its kind to account for the uncertainty of logging policy estimation in off-policy learning. The following two lines of work are most related to this work.

Off-policy learning. In many real-world applications, such as search engines and recommender systems, interactive online model update is expensive and risky [15]. Off-policy learning has therefore attracted great interest, since it can leverage the already logged feedback data [2, 8, 22]. The main challenge in off-policy learning is to address the mismatch between the logging policy and the learning policy. One common and widely-applied approach is to leverage the Inverse Propensity Scoring (IPS) method to correct the discrepancy between the two policies. And various methods are proposed to enhance IPS for more stabilizing learning [39, 37, 38] and improved variance control [23, 22], as well as extensions to more complex problems [40, 24].

However, all these solutions directly use the estimated logging policy for off-policy correction, leading to sub-optimal performance as shown in our experiments. A recent study on causal recommendation [10] also argues that propensity scores may not be correct due to unobserved confounders. They assume the effect of unobserved confounder for any sample can be bounded by a pre-defined hyper-parameter, and adversarially search for the worst-case propensity for learning. Mapping to off-policy learning, their solution is a special case of our UIPS-O variant with uncertainty as a pre-defined constant.

There were existing studies [34, 21, 26, 47, 9] also explore direct estimation of the propensity ratio to bypass estimating the logging policy. However, as discussed in Appendix 7.6 and Appendix 7.5, they demonstrate inferior performance compared to UIPS. This is primarily due to either the lack of consideration for the accuracy of the estimated propensity ratio, similar to the limitations of existing IPS-type algorithms in handling inaccurately estimated logging probabilities, or the degeneration to a specific IPS estimator that suffers high variance.

Recent work on distributionally robust off-policy evaluation and learning [32, 17, 44] also addresses uncertainty in off-policy learning. However, their approach to handling uncertainty and the underlying motivation differ significantly from ours, resulting in distinct techniques employed. Further details can be found in Appendix 7.7. Additionally, experiments conducted in Appendix 7.7 demonstrate that directly adapting methods from distributionally robust off-policy learning to handle inaccurately estimated logging probabilities leads to poor performance.

Off-policy learning can also be directly built on off-policy evaluation. Several work [36, 46] also propose to control the variance of the estimator caused by small logging probabilities through instance reweighing. Again, they directly use the estimated logging policy for correction, and thus performed worse than UIPS as observed in our experiments. A recent study [29] assumed additional structure in the action space and proposed the marginalized IPS. Instead, our work considers the uncertainty in the estimated logging policy and thus does not add any new assumptions about the problem space.

Uncertainty-aware learning. Estimation uncertainty has been extensively studied [45, 50, 1]. In the context of on-policy reinforcement learning and bandits [1, 48, 49], the use of uncertainty aims to strike a balance between exploration and exploitation by adopting an optimistic approach (i.e., UCB in bandits). On the other hand, most research on offline reinforcement learning/bandits [43, 4, 6] tends to be more conservative, employing techniques such as Lower Confidence Bounds (LCB) or penalizing out-of-distribution states and actions based on uncertainty to address extrapolation errors. However, these principles differ fundamentally from UIPS, which directly minimizes the mean square error of off-policy evaluation. The closed-form solution of the resulting per-instance weight in UIPS reflects how uncertainty contributes to the policy evaluation error. Moreover, Our UIPS-O and UIPS-P baselines leverage uncertainties using the two aforementioned general principles respectively. However, empirical findings indicate that blindly penalizing or boosting samples based on uncertainty is problematic. Proper correction depends on both uncertainty in logging policy estimation and the actual value of estimated logging probabilities.

6 Conclusion

In this paper, we propose a Uncertainty-aware Inverse Propensity Score estimator (UIPS) to explicitly model the uncertainty of the estimated logging policy for improved off-policy learning. UIPS weighs each logged instance to reduce its policy evaluation error, where the optimal weights have a closed-form solution derived by minimizing the upper bound of the resulting estimator’s mean squared error (MSE) to its ground-truth value. An improved policy is then obtained by optimizing the resulting estimation. Extensive experiments on synthetic and three real-world datasets as well as the theoretical convergence guarantee demonstrate the efficiency of UIPS.

As demonstrated in this work, explicitly modeling the uncertainty of the estimated logging policy is crucial for effective off-policy learning; but the best use of this uncertainty is not to simply down-weight or drop instances with uncertain estimations, but to balance it with the actually estimated logging probabilities in a per-instance basis. As our future work, it is promising to investigate how UIPS can be extended to value-based learning methods, e.g., actor-critics. And on the other hand, it is also important to analyze how tight our upper bound analysis of policy evaluation error is; and if possible, find new ways to tighten it for improvements.

References

- [1] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pages 2312–2320, 2011.
- [2] Aman Agarwal, Ivan Zaitsev, Xuanhui Wang, Cheng Li, Marc Najork, and Thorsten Joachims. Estimating position bias without intrusive interventions. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*, pages 474–482, 2019.
- [3] Ahmad Ajalloeian and Sebastian U Stich. On the convergence of sgd with biased gradients. *arXiv preprint arXiv:2008.00051*, 2020.
- [4] Gaon An, Seungyong Moon, Jang-Hyun Kim, and Hyun Oh Song. Uncertainty-based offline reinforcement learning with diversified q-ensemble. *Advances in neural information processing systems*, 34:7436–7447, 2021.
- [5] Raghu Aniruddh, Gottesman Omer, Liu Yao, Komorowski Matthieu, Faisal Aldo, Doshi-Velez Finale, and Brunskill Emma. Behaviour policy estimation in off-policy policy evaluation: Calibration matters. *arXiv preprint arXiv:1807.01066*, 2018.
- [6] Chenjia Bai, Lingxiao Wang, Zhuoran Yang, Zhihong Deng, Animesh Garg, Peng Liu, and Zhaoran Wang. Pessimistic bootstrapping for uncertainty-driven offline reinforcement learning. *arXiv preprint arXiv:2202.11566*, 2022.
- [7] K. Bhatia, K. Dahiya, H. Jain, P. Kar, A. Mittal, Y. Prabhu, and M. Varma. The extreme classification repository: Multi-label datasets and code, 2016.
- [8] Minmin Chen, Alex Beutel, Paul Covington, Sagar Jain, Francois Belletti, and Ed H Chi. Top-k off-policy correction for a reinforce recommender system. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*, pages 456–464, 2019.
- [9] Bo Dai, Ofir Nachum, Yinlam Chow, Lihong Li, Csaba Szepesvári, and Dale Schuurmans. Coindice: Off-policy confidence interval estimation. *Advances in neural information processing systems*, 33:9398–9411, 2020.
- [10] Sihao Ding, Peng Wu, Fuli Feng, Yitong Wang, Xiangnan He, Yong Liao, and Yongdong Zhang. Addressing unmeasured confounder for recommendation with sensitivity analysis. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 305–315, 2022.
- [11] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016.
- [12] Chongming Gao, Shijun Li, Wenqiang Lei, Jiawei Chen, Biao Li, Peng Jiang, Xiangnan He, Jiaxin Mao, and Tat-Seng Chua. Kuairc: A fully-observed dataset and insights for evaluating recommender systems. In *Proceedings of the 31st ACM International Conference on Information and Knowledge Management, CIKM '22*, 2022.
- [13] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International conference on machine learning*, pages 1321–1330. PMLR, 2017.
- [14] Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. Deepfm: a factorization-machine based neural network for ctr prediction. *arXiv preprint arXiv:1703.04247*, 2017.
- [15] Nan Jiang and Lihong Li. Doubly robust off-policy value evaluation for reinforcement learning. In *International Conference on Machine Learning*, pages 652–661. PMLR, 2016.
- [16] Thorsten Joachims, Adith Swaminathan, and Maarten De Rijke. Deep learning with logged bandit feedback. In *International Conference on Learning Representations*, 2018.
- [17] Nathan Kallus, Xiaojie Mao, Kaiwen Wang, and Zhengyuan Zhou. Doubly robust distributionally robust off-policy evaluation and learning. In *International Conference on Machine Learning*, pages 10598–10632. PMLR, 2022.

- [18] Patrick Kidger and Terry Lyons. Universal approximation with deep narrow networks. In *Conference on learning theory*, pages 2306–2327. PMLR, 2020.
- [19] Sergey Levine and Vladlen Koltun. Guided policy search. In *International conference on machine learning*, pages 1–9. PMLR, 2013.
- [20] Lihong Li, Yu Lu, and Dengyong Zhou. Provably optimal algorithms for generalized linear contextual bandits. In *International Conference on Machine Learning*, pages 2071–2080. PMLR, 2017.
- [21] Xu-Hui Liu, Zhenghai Xue, Jingcheng Pang, Shengyi Jiang, Feng Xu, and Yang Yu. Regret minimization experience replay in off-policy reinforcement learning. *Advances in Neural Information Processing Systems*, 34:17604–17615, 2021.
- [22] Yaxu Liu, Jui-Nan Yen, Bowen Yuan, Rundong Shi, Peng Yan, and Chih-Jen Lin. Practical counterfactual policy learning for top-k recommendations. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1141–1151, 2022.
- [23] Romain Lopez, Inderjit S Dhillon, and Michael I Jordan. Learning from extreme bandit feedback. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 8732–8740, 2021.
- [24] Jiaqi Ma, Zhe Zhao, Xinyang Yi, Ji Yang, Minmin Chen, Jiayi Tang, Lichan Hong, and Ed H Chi. Off-policy learning in two-stage recommender systems. In *Proceedings of The Web Conference 2020*, pages 463–473, 2020.
- [25] Rémi Munos, Tom Stepleton, Anna Harutyunyan, and Marc Bellemare. Safe and efficient off-policy reinforcement learning. *Advances in neural information processing systems*, 29, 2016.
- [26] Ofir Nachum, Yinlam Chow, Bo Dai, and Lihong Li. Dualdice: Behavior-agnostic estimation of discounted stationary distribution corrections. *Advances in Neural Information Processing Systems*, 32, 2019.
- [27] Doina Precup. Eligibility traces for off-policy policy evaluation. *Computer Science Department Faculty Publication Series*, page 80, 2000.
- [28] Aniruddh Raghu, Omer Gottesman, Yao Liu, Matthieu Komorowski, Aldo Faisal, Finale Doshi-Velez, and Emma Brunskill. Behaviour policy estimation in off-policy policy evaluation: Calibration matters. *arXiv preprint arXiv:1807.01066*, 2018.
- [29] Yuta Saito and Thorsten Joachims. Off-policy evaluation for large action spaces via embeddings. *arXiv preprint arXiv:2202.06317*, 2022.
- [30] Tobias Schnabel, Adith Swaminathan, Ashudeep Singh, Navin Chandak, and Thorsten Joachims. Recommendations as treatments: Debiasing learning and evaluation. In *international conference on machine learning*, pages 1670–1679. PMLR, 2016.
- [31] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International conference on machine learning*, pages 1889–1897. PMLR, 2015.
- [32] Nian Si, Fan Zhang, Zhengyuan Zhou, and Jose Blanchet. Distributionally robust policy evaluation and learning in offline contextual bandits. In *International Conference on Machine Learning*, pages 8884–8894. PMLR, 2020.
- [33] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- [34] Samarth Sinha, Jiaming Song, Animesh Garg, and Stefano Ermon. Experience replay with likelihood-free importance weights. In *Learning for Dynamics and Control Conference*, pages 110–123. PMLR, 2022.

- [35] Alex Strehl, John Langford, Lihong Li, and Sham M Kakade. Learning from logged implicit exploration data. *Advances in neural information processing systems*, 23, 2010.
- [36] Yi Su, Maria Dimakopoulou, Akshay Krishnamurthy, and Miroslav Dudík. Doubly robust off-policy evaluation with shrinkage. In *International Conference on Machine Learning*, pages 9167–9176. PMLR, 2020.
- [37] Adith Swaminathan and Thorsten Joachims. Batch learning from logged bandit feedback through counterfactual risk minimization. *The Journal of Machine Learning Research*, 16(1):1731–1755, 2015.
- [38] Adith Swaminathan and Thorsten Joachims. Counterfactual risk minimization: Learning from logged bandit feedback. In *International Conference on Machine Learning*, pages 814–823. PMLR, 2015.
- [39] Adith Swaminathan and Thorsten Joachims. The self-normalized estimator for counterfactual learning. *advances in neural information processing systems*, 28, 2015.
- [40] Adith Swaminathan, Akshay Krishnamurthy, Alekh Agarwal, Miro Dudik, John Langford, Damien Jose, and Imed Zitouni. Off-policy evaluation for slate recommendation. *Advances in Neural Information Processing Systems*, 30, 2017.
- [41] Sebastian Thrun and Michael L Littman. Reinforcement learning: an introduction. *AI Magazine*, 21(1):103–103, 2000.
- [42] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256, 1992.
- [43] Yue Wu, Shuangfei Zhai, Nitish Srivastava, Joshua Susskind, Jian Zhang, Ruslan Salakhutdinov, and Hanlin Goh. Uncertainty weighted actor-critic for offline reinforcement learning. *arXiv preprint arXiv:2105.08140*, 2021.
- [44] Da Xu, Yuting Ye, Chuanwei Ruan, and Bo Yang. Towards robust off-policy learning for runtime uncertainty. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 10101–10109, 2022.
- [45] Pan Xu, Zheng Wen, Handong Zhao, and Quanquan Gu. Neural contextual bandits with deep representation and shallow exploration. In *International Conference on Learning Representations*, 2021.
- [46] Ruohan Zhan, Vitor Hadad, David A Hirshberg, and Susan Athey. Off-policy evaluation via adaptive weighting with data from contextual bandits. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 2125–2135, 2021.
- [47] Ruiyi Zhang, Bo Dai, Lihong Li, and Dale Schuurmans. Gendice: Generalized offline estimation of stationary values. *arXiv preprint arXiv:2002.09072*, 2020.
- [48] Xiaoying Zhang, Hong Xie, Hang Li, and John CS Lui. Conversational contextual bandit: Algorithm and application. In *Proceedings of the web conference 2020*, pages 662–672, 2020.
- [49] Xiaoying Zhang, Hong Xie, and John CS Lui. Heterogeneous information assisted bandit learning: Theory and application. In *2021 IEEE 37th International Conference on Data Engineering (ICDE)*, pages 2135–2140. IEEE, 2021.
- [50] Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural contextual bandits with ucb-based exploration. In *International Conference on Machine Learning*, pages 11492–11502. PMLR, 2020.
- [51] Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. Deep interest network for click-through rate prediction. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 1059–1068, 2018.

Algorithm 1 UIPS

1: **Input:** The logged dataset $D := \{(\mathbf{x}_n, a_n, r_{\mathbf{x}_n, a_n})\}_{n=1}^N$, the estimated logging policy model $\hat{\beta}(a|\mathbf{x}) = \frac{\exp(f_{\theta}(\mathbf{x}, a))}{\sum_{a'} \exp(f_{\theta}(\mathbf{x}, a'))}$, dimension d .

2: Initialize $\mathbf{M}_D = \mathbf{I}_d$.

3: **for** $n = 1$ **to** N **do**

4: $\mathbf{M}_D = \mathbf{M}_D + \gamma \nabla_{\theta} f_{\theta}(\mathbf{x}_n, a_n) \nabla_{\theta} f_{\theta}(\mathbf{x}_n, a_n)^{\top}$ $\triangleright \mathbf{M}_D$ for uncertainty calculation.

5: **end for**

6: **for** $n = 1$ **to** N **do**

7: $U_{\mathbf{x}_n, a_n} = \sqrt{\nabla_{\theta} f_{\theta}(\mathbf{x}_n, a_n)^{\top} \mathbf{M}_D^{-1} \nabla_{\theta} f_{\theta}(\mathbf{x}_n, a_n)}$

8: **end for**

9: **while** not converge **do**

10: **for** $n = 1$ **to** N **do**

11: Calculate $\phi_{\mathbf{x}_n, a_n}$ as in Theorem 3.2

12: Calculate gradients as in Eq.(5) and updating $\pi_{\theta}(a|\mathbf{x})$.

13: **end for**

14: **end while**

15: **Output:** The learnt policy $\pi_{\theta}(a|\mathbf{x})$.

7 Appendix

7.1 Notations and Algorithm Framework.

For ease of reading, we list important notations in Table 5 and summarize the main framework of the proposed UIPS in Algorithm 1.

Notation	Description
\mathcal{X}	context space
\mathcal{A}	action set
$\mathbf{x} \in \mathcal{R}^d$	context vector
a	action
$r_{\mathbf{x}, a}$	reward
$\pi(a \mathbf{x})$	targeted policy to evaluate
$\beta(a \mathbf{x})$	the unknown ground-truth logging policy
$\hat{\beta}(a \mathbf{x})$	the estimated logging policy
$V(\pi)$	value function
$D := \{(\mathbf{x}_n, a_n, r_{\mathbf{x}_n, a_n})\}_{n=1}^N$	logged dataset containing N samples
$\phi_{\mathbf{x}, a}$	the optimal uncertainty-aware weight
$f_{\theta}(\mathbf{x}, a)$	the unknown ground-truth function that generates $\beta(a \mathbf{x}) = \frac{\exp(f_{\theta^*}(\mathbf{x}, a))}{\sum_{a'} \exp(f_{\theta^*}(\mathbf{x}, a'))}$
$\hat{f}_{\theta}(\mathbf{x}, a)$	the estimate of $f_{\theta}(\mathbf{x}, a)$ that generates $\hat{\beta}(a \mathbf{x})$
$\mathbf{B}_{\mathbf{x}, a}$	confidence interval of $\hat{\beta}(a \mathbf{x})$
$U_{\mathbf{x}, a}$	uncertainty defined as $\sqrt{\nabla_{\theta} f_{\theta^*}(\mathbf{x}, a)^{\top} \mathbf{B}_{\mathbf{x}, a} \nabla_{\theta} f_{\theta^*}(\mathbf{x}, a)}$
$\mathbf{g}(\mathbf{x}_n, a_n)$	gradient of $f_{\theta}(\mathbf{x}, a)$ regarding to the last layer.

Table 5: Notations

Computation Cost. The additional computation cost of UIPS over IPS comes from two parts:

- Pre-calculate uncertainties (line 1-5 in Algorithm 1) : This part calculates uncertainty of the logging probability for each (s, a) pair, and it *only needs to be executed once*. The computational cost of this step is $O(Nd^2 + d^3)$, where $O(Nd^2)$ is for calculating uncertainties in each (s, a) pair and $O(d^3)$ is for matrix inverse.
- Calculate $\phi_{\mathbf{x}, a}$ during training (line 8 in Algorithm 1): It only takes $O(1)$ time, the same computational cost as calculating IPS score.

Note that calculating the logging probability for each sample, which is essential for both UIPS and IPS, takes $O(Nd/jAj)$ time. Since the dimension d is usually much less than action size jAj and sample size N , UIPS does not introduce significant computational overhead compared to the original IPS solution.

7.2 Derivation of confidence interval of logging probability.

Given that $jf_{\theta^*}(\mathbf{x}, a) - f_{\theta}(\mathbf{x}, a)j \leq \gamma U_{\mathbf{x},a}$ holds with probability at least $1 - \sigma$ and

$$\beta(a|\mathbf{x}) = \frac{\exp(f_{\theta^*}(\mathbf{x}, a))}{Z}, \hat{\beta}(a|\mathbf{x}) = \frac{\exp(f_{\theta}(\mathbf{x}, a))}{\hat{Z}},$$

where $Z = \sum_{a \in \mathcal{A}} \exp(f_{\theta}(a|\mathbf{x}))$ and $\hat{Z} = \sum_{a \in \mathcal{A}} \exp(f_{\theta}(a|\mathbf{x}))$, we can get that with probability at least $1 - \sigma$:

$$\begin{aligned} & jf_{\theta^*}(\mathbf{x}, a) - f_{\theta}(\mathbf{x}, a)j \leq \gamma U_{\mathbf{x},a} \\ & \Rightarrow f_{\theta^*}(\mathbf{x}, a) - f_{\theta}(\mathbf{x}, a) \leq f_{\theta^*}(\mathbf{x}, a) - f_{\theta}(\mathbf{x}, a) + \gamma U_{\mathbf{x},a} \\ & \Rightarrow \exp(f_{\theta^*}(\mathbf{x}, a)) \leq \exp(f_{\theta}(\mathbf{x}, a)) \exp(\gamma U_{\mathbf{x},a}) \\ & \stackrel{(1)}{\Rightarrow} \hat{Z} \hat{\beta}(a|\mathbf{x}) \exp(-\gamma U_{\mathbf{x},a}) \leq \exp(f_{\theta^*}(\mathbf{x}, a)) \leq \hat{Z} \hat{\beta}(a|\mathbf{x}) \exp(\gamma U_{\mathbf{x},a}) \\ & \stackrel{(2)}{\Rightarrow} \frac{\hat{Z} \exp(-\gamma U_{\mathbf{x},a})}{Z} \hat{\beta}(a|\mathbf{x}) \leq \beta(a|\mathbf{x}) \leq \frac{\hat{Z} \exp(\gamma U_{\mathbf{x},a})}{Z} \hat{\beta}(a|\mathbf{x}) \end{aligned}$$

The step labeled as (1) is due to the modelling of $\hat{\beta}(a|\mathbf{x})$. And the step labeled as (2) is because Z is a positive constant independent of $\hat{\beta}(a|\mathbf{x})$. Thus with probability at least $1 - \delta$, we have $\beta(a|\mathbf{x}) \geq \mathbf{B}_{\mathbf{x},a}$ and

$$\mathbf{B}_{\mathbf{x},a} = \left[\frac{\hat{Z} \exp(-\gamma U_{\mathbf{x},a})}{Z} \hat{\beta}(a|\mathbf{x}), \frac{\hat{Z} \exp(\gamma U_{\mathbf{x},a})}{Z} \hat{\beta}(a|\mathbf{x}) \right].$$

7.3 Experiment details

7.3.1 Synthetic Data

Data generation. Given the ground-truth logging policy $\beta(a|\mathbf{x})$, we generate the logged dataset as follows. For each sample in train set, we first get the embedded context vector $\tilde{\mathbf{x}}$. We then sample an action a according to $\beta(a|\mathbf{x})$, and obtain the reward $r_{\mathbf{x},a} = \mathbf{y}_{\tilde{\mathbf{x}},a}$, resulting bandit feedback $(\mathbf{x}, a, r_{\mathbf{x},a})$, where $\mathbf{y}_{\tilde{\mathbf{x}},a}$ is the label of class a under the original feature vector $\tilde{\mathbf{x}}$. We repeat above process N times to collect the logged dataset. In our experiments, we take $d = 64, N = 100$.

We model the logging policy as in Eq.(9), where f_{θ} are the parameters to be estimated. To train the logging policy, we take all samples in the logged dataset D as positive instances, and randomly sample non-selected actions as negative instances as in [8].

7.3.2 Real-world Data

Statistics of datasets. The statistics of three real-world recommendation datasets with unbiased data can be found in Table 6.

Table 6: The statistics of three real-world datasets.

Dataset	#User	#Item	#Biased Data	#Unbiased Data
Yahoo R3	15,400	1,000	311,704	54,000
Coat	290	300	6,960	4,640
KuaiRec	7,176	10,729	12,530,806	4,676,570

All these datasets contain a set of biased data collected from users' interactions on the platform, and a set of unbiased data collected from a randomized controlled trial where items are randomly selected. As in [10], on each dataset, the biased data is used for training, and the unbiased data is for testing, with a small part of unbiased data split for validation purpose (5% on Yahoo R3 and Coat,

and 15% on KuaiRec). We take the reward as 1 if : (1) the rating is larger than 3 in Yahoo! R3 and Coat datasets; (2) the user watched more than 70% of the video in KuaiRec. Otherwise, the reward is labeled as 0.

We adopted a two-tower neural network architecture to implement both the logging and learning policy, as shown in Figure 2. For the learning policy, the user representation and item representation are first modelled through two separate neural networks (i.e., the user tower and the item tower), and then their element-by-element product vector is projected to predict the user’s preference for the item. We then re-use the user state generated from the user tower of the learning policy, and model the logging policy with another separate item tower, following [8]. We also block gradients to prevent the logging policy learning interfering the user state of the learning policy. In each learning epoch, we will first estimate the logging policy, and then take the estimated logging probabilities as well as their uncertainties to optimize the learning policy.

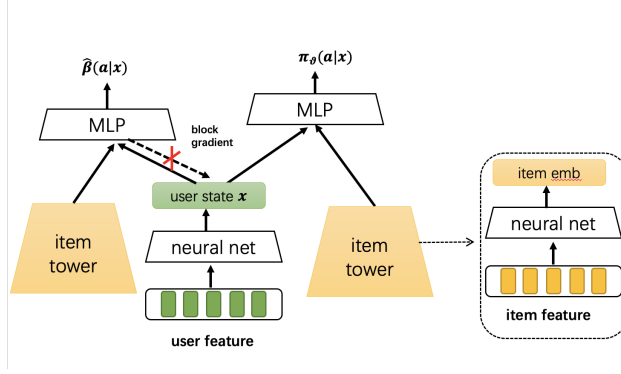


Figure 2: Model architecture of the logging and the learning policy in real-world datasets

7.3.3 Implementation details.

To facilitate hyper-parameter tuning, we disentangled two η s in two the terms of $\phi_{\mathbf{x},a}$ in Theorem 3.2, and introduce η_1 and η_2 to represent η in the first and second term respectively in our implementation. This is due to the scale of η in the first term is closely related to the scale of λ , with λ/η as the truly effective hyper-parameter. But the scale of η in the second term is independent from λ .

Moreover, while $\lambda = E_{\pi_{\theta}} \left[r_{\mathbf{x},a}^2 \frac{\pi_{\theta}(a|\mathbf{x})}{\beta(a|\mathbf{x})} \right]$ depends on π_{θ} as discussed in Eq.(6), we cannot adaptively set the value of λ since the ground-truth logging policy $\beta(a|\mathbf{x})$ is unknown. However, we have:

$$\lambda := E_{\pi} \left[r_{\mathbf{x},a}^2 \frac{\pi_{\theta}(a|\mathbf{x})}{\beta(a|\mathbf{x})} \right] \left(\sum_a \pi_{\theta}(a|\mathbf{x})^4 \right) \left(\sum_a \frac{r_{\mathbf{x},a}^4}{\beta(a|\mathbf{x})^2} \right) \left(\sum_a \frac{r_{\mathbf{x},a}^4}{\beta(a|\mathbf{x})^2} \right),$$

where the first inequality is due to the Cauchy–Schwarz inequality and the second inequality is because $\sum_a \pi_{\theta}(a|\mathbf{x})^4 \leq \sum_a \pi_{\theta}(a|\mathbf{x}) = 1$ with $\pi_{\theta}(a|\mathbf{x}) \in [0, 1]$. We denote $\tilde{\lambda} = \left(\sum_a \frac{r_{\mathbf{x},a}^4}{\beta(a|\mathbf{x})^2} \right)$, which is dataset-specific constant and independent from π_{θ} . By replacing λ in Eq.(6) with $\tilde{\lambda}$, we can still minimize an upper bound of $\text{MSE}(\hat{V}_{\text{UIPS}}(\pi_{\theta}))$, which ensures that the result of our analysis still holds. Thus, considering ease of computation and efficiency, we take a fixed λ during our policy learning.

We then use grid search to select hyperparameters based on the model’s performance on the validation dataset: the learning rate was searched in $\{1e^{-5}, 1e^{-4}, 1e^{-3}, 1e^{-2}\}$; λ, γ, η_1 were searched in $\{0.5, 0.1, 1, 2, 5, 10, 15, 20, 25, 30, 40, 50\}$. And η_2 was searched in $\{1, 10, 100, 1000\}$. For baseline algorithms, we also performed a similar grid search as mentioned above, and the search range follows the original papers.

Ablation study: hyper-parameter tuning. Although UIPS has four hyperparameters (λ, γ, η_1 , and η_2), one only needs to carefully finetune two of them, i.e., γ and η_1^2/λ , to obtain good performance of UIPS. This is because:

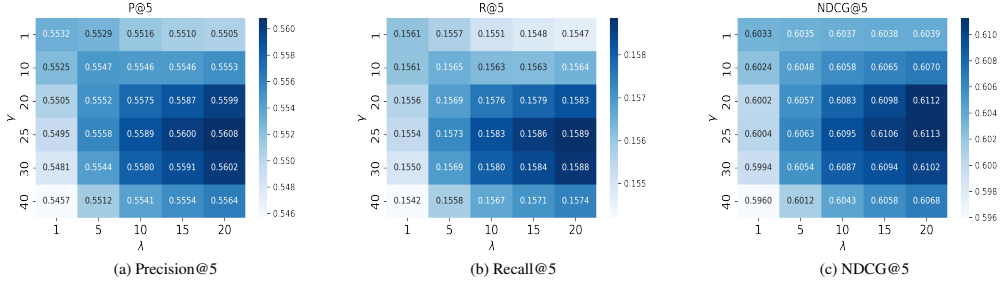


Figure 3: Effect of λ and γ on synthetic dataset with $\tau = 0.5$

- η_2 acts as a capping threshold to ensure $\phi_{\mathbf{x},a} \leq 2\eta_2$ holds even when the corresponding propensity scores are very low. Hence, it should be set to a reasonably large value (e.g., 100).
- The key component (i.e., the first term) of $\phi_{\mathbf{x},a}$ can be rewritten in the following way. While all (\mathbf{x}, a) pairs will be multiplied by $\phi_{\mathbf{x},a}$, η_1 in the numerator will not affect final performance too much, and the key is to find a good value of η_1^2/λ to balance the two terms in the denominator:

$$\eta_1 / \left[\exp(\gamma U_{\mathbf{x},a}) + \frac{\eta_1^2/\lambda \pi_{\vartheta}(a|\mathbf{x})^2}{\hat{\beta}(a|\mathbf{x})^2 \exp(\gamma U_{\mathbf{x},a})} \right].$$

Thus with η_1 and η_2 fixed, the effect of hyper-parameter γ and λ on precision, recall as well as NDCG can be found in Figure 3. We can observe that to make UIPS excel, $\mathbf{B}_{\mathbf{x},a}$ needs to be of high confidence, e.g., $\gamma = 25$ performed the best on the dataset with $\tau = 0.5$. Moreover, the threshold η_1/λ cannot be too small or too large.

7.4 Experiments on the doubly robust estimators.

The doubly robust (DR) estimator [15], which is a hybrid of *direct method* (DM) estimator and *inverse propensity score* (IPS) estimator, is also widely used for off-policy evaluation. More specifically, let $\hat{\eta} : \mathcal{X} \times \mathcal{A} \rightarrow \mathcal{R}$ be the imputation model in DM that estimates the reward of action a under context vector \mathbf{x} , and $\hat{\beta}(a|\mathbf{x})$ be the estimated logging policy in the IPS estimator. The DR estimator evaluates policy π based on the logged dataset $D := \{(\mathbf{x}_n, a_n, r_{\mathbf{x}_n, a_n})\}_{n=1}^N$, by:

$$\hat{V}_{\text{DR}}(\pi) = \hat{V}_{\text{DM}}(\pi) + \frac{1}{N} \sum_{n=1}^N \frac{\pi(a_n|\mathbf{x}_n)}{\hat{\beta}(a_n|\mathbf{x}_n)} (r_{\mathbf{x}_n, a_n} - \hat{\eta}(\mathbf{x}_n, a_n)) \quad (10)$$

where $\hat{V}_{\text{DM}}(\pi)$ is the DM estimator:

$$\hat{V}_{\text{DM}}(\pi) = \frac{1}{N} \sum_{n=1}^N \sum_{a \in \mathcal{A}} \pi(a|\mathbf{x}_n) \hat{\eta}(\mathbf{x}_n, a). \quad (11)$$

Again assume the policy $\pi(a|\mathbf{x})$ is parameterized by ϑ , the REINFORCE gradient of $\hat{V}_{\text{DR}}(\pi_{\vartheta})$ with respect to ϑ can be readily derived as follows:

$$\begin{aligned} \nabla_{\vartheta} \hat{V}_{\text{DR}}(\pi_{\vartheta}) &= \frac{1}{N} \sum_{n=1}^N \left(\sum_{a \in \mathcal{A}} \pi_{\vartheta}(a|\mathbf{x}_n) \hat{\eta}(\mathbf{x}_n, a) \nabla_{\vartheta} \log(\pi_{\vartheta}(a|\mathbf{x}_n)) \right) \\ &\quad + \frac{1}{N} \sum_{n=1}^N \left(\frac{\pi(a_n|\mathbf{x}_n)}{\hat{\beta}(a_n|\mathbf{x}_n)} (r_{\mathbf{x}_n, a_n} - \hat{\eta}(\mathbf{x}_n, a_n)) \nabla_{\vartheta} \log(\pi_{\vartheta}(a_n|\mathbf{x}_n)) \right). \end{aligned} \quad (12)$$

The imputation model $\hat{\eta}(\mathbf{x}, a)$ is pre-trained following previous work [22] with the same neural network architecture as the logging policy model. Besides the standard DR estimator, we also adapt

Table 7: Experiment results on synthetic datasets. The best and second best results are highlighted with **bold** and underline respectively. Two p -values are calculated: (1) p -value (UIPSDR): The p -value under the t-test between UIPSDR and the best DR baseline on each dataset; (2) p -value (UIPS): The p -value under the t-test between UIPS and the best DR baseline on each dataset.

Algorithm	$\tau = 0.5$			$\tau = 1$			$\tau = 2$		
	P@5	R@5	NDCG@5	P@5	R@5	NDCG@5	P@5	R@5	NDCG@5
BIPS-Cap	0.5515 $2e^{-3}$	0.1553 $8e^{-4}$	0.6031 $2e^{-3}$	0.5526 $2e^{-3}$	0.1561 $6e^{-4}$	0.6016 $1e^{-3}$	0.5409 $3e^{-3}$	0.1529 $9e^{-4}$	0.5901 $2e^{-3}$
UIPS	0.5589 $3e^{-3}$	0.1583 $9e^{-4}$	0.6095 $3e^{-3}$	0.5572 $2e^{-3}$	0.1571 $8e^{-4}$	0.6074 $2e^{-3}$	0.5432 $3e^{-3}$	0.1534 $8e^{-4}$	0.5946 $2e^{-3}$
DR	0.3846 $3e^{-2}$	0.1082 $8e^{-3}$	0.4684 $3e^{-2}$	0.3631 $3e^{-2}$	0.1017 $9e^{-3}$	0.4494 $3e^{-2}$	0.3560 $3e^{-2}$	0.0995 $7e^{-3}$	0.4470 $2e^{-3}$
MinVarDR	0.3212 $3e^{-2}$	0.0908 $8e^{-3}$	0.4062 $3e^{-2}$	0.3240 $5e^{-2}$	0.0903 $1e^{-2}$	0.3905 $5e^{-2}$	0.3234 $5e^{-2}$	0.0910 $1e^{-2}$	0.4059 $4e^{-2}$
ShrinkageDR	0.4139 $2e^{-2}$	0.1161 $7e^{-3}$	0.4969 $3e^{-2}$	0.3944 $3e^{-2}$	0.1101 $8e^{-3}$	0.4797 $2e^{-2}$	0.4080 $3e^{-2}$	0.1135 $7e^{-3}$	0.4901 $2e^{-2}$
UIPSDR	0.4278 $2e^{-2}$	0.1200 $6e^{-3}$	0.5069 $2e^{-2}$	0.4008 $2e^{-2}$	0.1126 $7e^{-3}$	0.4847 $2e^{-2}$	0.4144 $2e^{-2}$	0.1162 $8e^{-3}$	0.4972 $2e^{-2}$
p -value (UIPSDR)	$2e^{-1}$	$2e^{-1}$	$3e^{-1}$	$6e^{-1}$	$4e^{-1}$	$6e^{-1}$	$6e^{-1}$	$4e^{-1}$	$5e^{-1}$
p -value (UIPS)	$6e^{-13}$	$4e^{-13}$	$4e^{-12}$	$2e^{-12}$	$1e^{-12}$	$5e^{-12}$	$8e^{-12}$	$8e^{-12}$	$2e^{-11}$

Table 8: Experiment results on real-world unbiased datasets. The best and second best results are highlighted with **bold** and underline respectively. Two p -values are calculated: (1) p -value (UIPSDR): The p -value under the t-test between UIPSDR and the best DR baseline on each dataset; (2) p -value (UIPS): The p -value under the t-test between UIPS and the best DR baseline on each dataset.

Algorithm	Yahoo			Coat			KuaiRec		
	P@5	R@5	NDCG@5	P@5	R@5	NDCG@5	P@50	R@50	NDCG@50
BIPS-Cap	0.2808 $2e^{-3}$	0.7576 $5e^{-3}$	0.6099 $8e^{-3}$	0.2758 $6e^{-3}$	0.4582 $7e^{-3}$	0.4399 $9e^{-3}$	0.8750 $3e^{-3}$	0.0238 $7e^{-5}$	0.8788 $5e^{-3}$
UIPS	0.2868 $2e^{-3}$	0.7742 $5e^{-3}$	0.6274 $5e^{-3}$	0.2877 $3e^{-3}$	0.4757 $5e^{-3}$	0.4576 $8e^{-3}$	0.9120 $1e^{-3}$	0.0250 $5e^{-4}$	0.9174 $7e^{-4}$
DR	0.2670 $2e^{-3}$	0.7174 $6e^{-3}$	0.5636 $6e^{-3}$	0.2884 $3e^{-3}$	0.4760 $5e^{-3}$	0.4541 $5e^{-3}$	0.8794 $1e^{-2}$	0.0240 $5e^{-4}$	0.8824 $2e^{-2}$
MinVarDR	0.2272 $5e^{-3}$	0.5989 $1e^{-2}$	0.4525 $1e^{-2}$	0.2704 $4e^{-3}$	0.4434 $9e^{-3}$	0.4137 $6e^{-3}$	0.8640 $7e^{-3}$	0.0235 $2e^{-4}$	0.8657 $7e^{-3}$
ShrinkageDR	0.2697 $2e^{-3}$	0.7226 $6e^{-3}$	0.5713 $5e^{-3}$	0.2895 $4e^{-3}$	0.4749 $6e^{-3}$	0.4526 $6e^{-3}$	0.8778 $2e^{-2}$	0.0239 $5e^{-4}$	0.8800 $2e^{-2}$
UIPSDR	0.2721 $1e^{-3}$	0.7294 $6e^{-3}$	0.5750 $5e^{-3}$	0.2946 $4e^{-3}$	0.4854 $8e^{-3}$	0.4647 $8e^{-3}$	0.8849 $1e^{-2}$	0.0242 $4e^{-4}$	0.8896 $1e^{-2}$
p -value (UIPSDR)	$1e^{-2}$	$2e^{-2}$	$1e^{-1}$	$7e^{-3}$	$5e^{-3}$	$2e^{-3}$	$4e^{-1}$	$4e^{-1}$	$3e^{-1}$
p -value (UIPS)	$1e^{-12}$	$6e^{-14}$	$6e^{-15}$	$3e^{-1}$	$8e^{-1}$	$1e^{-1}$	$2e^{-6}$	$2e^{-6}$	$1e^{-3}$

UIPS and the best two baselines on off-policy evaluation estimator (i.e., MinVar and Shrinkage based on the results in Table 1 and Table 4) to doubly robust setting using the same imputation model.

Table 7 and Table 8 report the empirical performance of the learned policy on the synthetic datasets and three real-world datasets respectively. For ease of comparison, we also include the experiment results of BIPS-Cap and UIPS on each dataset in these two tables. Two p -values are also provided: (1) p -value (UIPSDR): The p -value under the t-test between UIPSDR and the best DR baseline on each dataset; (2) p -value (UIPS): The p -value under the t-test between UIPS and the best DR baseline on each dataset. From Table 7 and Table 8, we can first observe that DR cannot consistently outperform BIPS-Cap: It outperformed BIPS-Cap on the Coat and KuaiRec dataset, while achieving much worse performance on the synthetic datasets and Yahoo dataset. This is because the imputation model also plays a very important role in gradient calculation as shown in Eq.(12). Its accuracy greatly affects policy learning. When the imputation model is sufficiently accurate, for example, on the Coat dataset with only 300 actions, incorporating the DM estimator not only led to better performance of DR over IPS, but also improved performance of UIPSDR over UIPS. And in particular, in this situation UIPSDR performed better than DR with the gain being statistically significant. When the imputation model is not accurate enough, for example, on the KuaiRec dataset with a large action space but sparse reward feedback, DR is still worse than UIPS, and UIPSDR also performs worse than UIPS due to the distortion of the imputation model.

7.5 Difference against DICE-type algorithms.

The DICE line of work [26, 47, 9] is proposed for off-policy correction in the multi-step RL setting with the environment following the Markov Decision Process (MDP) assumption. Although the DICE line of work does not require the knowledge of the logging policy, it is fundamentally different from our work. Given a logged dataset $D := f(s_t, a_t, r_{s_t, a_t}, s_{t+1})g$ collected from an unknown logging policy $\beta(a_t/s_t)$, DICE-type algorithms propose to directly estimate the discounted stationary distribution correction $w_{\pi/D}(s_t, a_t) = \frac{d^\pi(s_t, a_t)}{d^D(s_t, a_t)}$ to replace the product-based off-policy correction weight $\prod_t \frac{\pi(a_t/s_t)}{\beta(a_t/s_t)}$, which suffers high variance due to the series of products. While the estimation of $w_{\pi/D}(s_t, a_t)$ is agnostic to the logging policy, it highly depends on two assumptions: (1) Environment follows an MDP, i.e., $s_{t+1} = T(s_t, a_t)$ with $T(\cdot, \cdot)$ denoting the state transition function; (2) Each logged sample should be of a state-action-next-state tuple $(s_t, a_t, r_{s_t, a_t}, s_{t+1})$ that contains state transition information.

Table 9: Empirical performance of DICE-S on synthetic datasets.

Algorithm	$\tau = 0.5$			$\tau = 1$			$\tau = 2$											
	P@5	R@5	NDCG@5	P@5	R@5	NDCG@5	P@5	R@5	NDCG@5									
BIPS-Cap	0.5515	2e ⁻³	0.1553	8e ⁻⁴	0.6031	2e ⁻³	0.5526	2e ⁻³	0.1561	6e ⁻⁴	0.6016	1e ⁻³	0.5409	3e ⁻³	0.1529	9e ⁻⁴	0.5901	2e ⁻³
DICE-S	0.5416	4e ⁻³	0.1520	1e ⁻³	0.5968	4e ⁻³	0.5508	2e ⁻³	0.1553	7e ⁻⁴	0.6010	2e ⁻³	0.5403	2e ⁻³	0.1526	7e ⁻⁴	0.5903	1e ⁻³
UIPS	0.5589	3e ⁻³	0.1583	9e ⁻⁴	0.6095	3e ⁻³	0.5572	2e ⁻³	0.1571	8e ⁻⁴	0.6074	2e ⁻³	0.5432	3e ⁻³	0.1534	8e ⁻⁴	0.5946	2e ⁻³

Table 10: Empirical performance of DICE-S on real-world datasets.

Algorithm	Yahoo			Coat			KuaiRec											
	P@5	R@5	NDCG@5	P@5	R@5	NDCG@5	P@50	R@50	NDCG@50									
BIPS-Cap	0.2808	2e ⁻³	0.7576	5e ⁻³	0.6099	8e ⁻³	0.2758	6e ⁻³	0.4582	7e ⁻³	0.4399	9e ⁻³	0.8750	3e ⁻³	0.0238	7e ⁻⁵	0.8788	5e ⁻³
DICE-S	0.2618	4e ⁻³	0.7010	1e ⁻²	0.5627	9e ⁻³	0.2686	5e ⁻³	0.4422	6e ⁻³	0.4265	5e ⁻³	0.8842	2e ⁻³	0.0241	6e ⁻⁵	0.8908	2e ⁻³
UIPS	0.2868	2e ⁻³	0.7742	5e ⁻³	0.6274	5e ⁻³	0.2877	3e ⁻³	0.4757	5e ⁻³	0.4576	8e ⁻³	0.9120	1e ⁻³	0.0250	5e ⁻⁵	0.9174	7e ⁻⁴

We further inspected how DICE-type algorithms would work in the contextual bandit setting, which can be taken as a one-state MDP. We take DualDICE [1] as an example for illustration, and similar conclusions can be drawn for other DICE-type algorithms. We can easily derive that in the contextual bandit setting, DualDICE degenerates to the IPS estimator that approximates the unknown ground-truth logging policy $\beta(a|s)$ with its empirical estimate from the given logged dataset $D := \{(s_i, a_i, r_{s_i, a_i})\}_{i=1}^N$, which inherits all limitations of IPS estimator, such as high variance on instances with low support in D .

More specifically, recall that DualDICE estimates the discounted stationary distribution correction by optimizing the following objective function (Eq.(8) in [26]):

$$\min_{x: S \rightarrow A \times C} \frac{1}{2} E_{(s,a) \sim d^D} [x(s,a)^2] - E_{(s,a) \sim d^\pi} [x(s,a)],$$

with the optimizer $x(s,a) = w_{\pi/D}(s,a)$. In a multi-step RL setting, one cannot directly calculate the second expectation as $d^\pi(s,a)$ is inaccessible, thus DualDICE takes the change-of-variable trick to address the above optimization problem. However, in the contextual bandit setting, let $p(s)$ denote the state distribution, $d^\pi(s,a) = p(s)\pi(a|s)$ and $d^D(s,a) = p(s)\beta(a|s)$. The optimization problem can be rewritten as :

$$\min_{x: S \rightarrow A \times C} E_{s \sim p(s)} \left[\frac{1}{2} E_{a \sim \beta(a|s)} [x(s,a)^2] - E_{a \sim \pi(a|s)} [x(s,a)] \right].$$

With the logged dataset D , the empirical estimate of above optimization problem is :

$$\min_{x: S \rightarrow A \times C} \sum_s \frac{N_s}{N} \left(\sum_a \frac{N_{s,a}}{2N_s} x(s,a)^2 - \sum_a \pi(a|s) x(s,a) \right),$$

yielding $x(s,a) = \pi(a|s)N_s/N_{s,a}$, where N_s and $N_{s,a}$ denote the number of logged samples with $s_i = s$ and $(s_i = s, a_i = a)$ respectively. This is actually equivalent to the IPS estimator with $N_{s,a}/N_s$ as the estimate of $\beta(a|s)$. We refer to the above estimator as DICE-S, and Table 9 and Table 10 show the empirical performance of the policy learned through DICE-S on the synthetic and real-world datasets respectively. Similar as BIPS-Cap described in Section 4, we also clip propensity scores to control variance. One can observe from Table 9 and Table 10 that DICE-S performs worse than BIPS-Cap in all datasets except the KuaiRec dataset. Additionally, DICE-S underperforms compared to UIPS in all datasets. This is due to the fact that although DICE-S is unbiased, it only becomes accurate with numerous logged samples, and suffers high variance when logged samples are limited, resulting in its poor performance in our experiments.

7.6 Comparison against work on direct propensity estimation.

In addition to DICE, several other works [34, 21] propose methods to directly estimate the propensity ratio without requiring a behavior policy. These methods then use the propensity estimates to prioritize instances in the replay buffer for better TD learning. To compare its effectiveness, we also included a new baseline called IPS-LFIW, which implements the approach proposed in [1] to directly estimate the propensity ratio for off-policy learning.

The average performance and standard deviations of IPS-LFIW and UIPS on three synthetic datasets with different τ are reported in Table 11. Recall that smaller τ indicates a more skewed ground-truth logging policy. The p-value under the t-test between UIPS and IPS-LFIW is also provided to

Table 11: Empirical performance of IPS-LFIW and UIPS on synthetic datasets

Algorithm	$\tau = 0.5$			$\tau = 1$			$\tau = 2$											
	P@5	R@5	NDCG@5	P@5	R@5	NDCG@5	P@5	R@5	NDCG@5									
IPS-LFIW	0.5542	1e ⁻³	0.1568	4e ⁻⁴	0.6033	1e ⁻³	0.5472	1e ⁻³	0.1549	5e ⁻⁴	0.5975	1e ⁻³	0.5255	6e ⁻³	0.1485	1e ⁻³	0.5769	5e ⁻³
UIPS	0.5608	2e ⁻³	0.1589	8e ⁻⁴	0.6113	3e ⁻³	0.5572	2e ⁻³	0.1571	8e ⁻⁴	0.6074	2e ⁻³	0.5432	3e ⁻³	0.1534	8e ⁻⁴	0.5946	2e ⁻³
p-value	2e ⁻⁶		5e ⁻⁶		6e ⁻⁹		2e ⁻⁸		2e ⁻⁶		4e ⁻¹⁰		4e ⁻⁷		8e ⁻⁷		1e ⁻⁷	

investigate the significance of improvements. Notably, UIPS consistently outperformed IPS-LFIW with statistically significant improvements. One major reason for the worse performance of IPS-LFIW is that it does not consider the accuracy of the estimated propensity ratio, in a direct analogy to failing to handle uncertainty in the estimated logging probabilities in existing IPS-type algorithms.

7.7 Comparison against distributionally robust off-policy evaluation and learning.

Our work is fundamentally different from the line of work on distributionally robust off-policy evaluation and learning [32, 17, 44]. This results in different objectives that guide the use of min-max optimization.

Specifically, the work in [32, 17, 44] assumes unknown changes exist between their training and deployment environments, such as user preference drift or unforeseen events during policy execution. Thus they choose to maximize the policy value (e.g., $\hat{V}_{\text{IPS}}(\pi)$) in the worst environment within an uncertainty set around the training environment,

$$\max_{\pi} \min_U \hat{V}_{\text{IPS}}(\pi).$$

As a result, their uncertainty set is created by introducing a small perturbation to the training environment. For example, the work in [32, 17] searches the worst environment P_1 in the σ -close perturbed environments around the training environment P_0 (Eq.(1) in [17]):

$$U(\sigma) = \{P_1 : P_1 \ll P_0 \text{ and } D_{KL}(P_1 || P_0) \leq \sigma\}.$$

And the work in [44] adversarially perturbs the known ground-truth logging policy $\pi_0(j)$ in searching for the worst case (Eq.(4) in [44]):

$$U(\alpha) = \{ \pi_u : \max_{a,x} \max_{\pi_0} \left(\frac{\pi_u(a|x)}{\pi_0(a|x)} - \frac{\pi_0(a|x)}{\pi_u(a|x)} \right) g \leq \alpha \}.$$

In contrast, UIPS assumes the training and deployment environments stay the same, but the ground-truth logging policy is unknown. To control the high bias and high variance caused by inaccurately and small estimated logging probabilities, UIPS explicitly models the uncertainty in the estimated logging policy by incorporating a per-sample weight $\phi_{x,a}$ as discussed in Eq.(4). In order to make $\hat{V}_{\text{UIPS}}(\pi_{\vartheta})$ as accurate as possible despite the unknown ground-truth logging policy $\beta(j)$, UIPS solves a min-max optimization problem in Eq.(8). This optimization problem seeks to find the optimal $\phi_{x,a}$ that minimizes the upper bound of the mean squared error (MSE) of \hat{V}_{UIPS} to its ground-truth value, within an uncertainty set of the unknown ground-truth logging policy $\beta(j)$. The closed-form solution for the min-max optimization is also derived as in Theorem 3.2.

Furthermore, observing that work in [44] also performs optimization over an uncertainty set of the logging policy, we further adapted their method to handle the inaccuracy of the estimated logging policy, by taking $\pi_0(j)$ as the estimated logging policy, i.e., $\pi_0(j) = \hat{\beta}(j)$. We name the adapted methods as IPS-UN. Table 12 demonstrates the performance of the learned policy under IPS-UN on three synthetic datasets. The results suggest that directly applying IPS-UN to handle inaccurately estimated logging probabilities is not be a feasible solution to our problem. One important reason for the worse performance of IPS-UN is that it strives to optimize for the worst potential environment, which might not be the case in our experiment datasets. On the other hand, UIPS assumes the training and deployment environments stay same and strives to identify the optimal policy with an unknown ground-truth logging policy.

7.8 Theoretical Proofs.

Proof of Proposition 2.1:

Table 12: Empirical performance of IPS-UN on synthetic datasets.

Algorithm	$\tau = 0.5$			$\tau = 1$			$\tau = 2$											
	P@5	R@5	NDCG@5	P@5	R@5	NDCG@5	P@5	R@5	NDCG@5									
BIPS-Cap	0.5515	2e ⁻³	0.1553	8e ⁻⁴	0.6031	2e ⁻³	0.5526	2e ⁻³	0.1561	6e ⁻⁴	0.6016	1e ⁻³	0.5409	3e ⁻³	0.1529	9e ⁻⁴	0.5901	2e ⁻³
IPS-UN	0.4089	3e ⁻²	0.1152	8e ⁻³	0.4916	2e ⁻²	0.3911	3e ⁻³	0.1100	9e ⁻³	0.4754	3e ⁻²	0.3599	4e ⁻²	0.1009	1e ⁻²	0.4353	4e ⁻²
UIPS	0.5589	3e ⁻³	0.1583	9e ⁻⁴	0.6095	3e ⁻³	0.5572	2e ⁻³	0.1571	8e ⁻⁴	0.6074	2e ⁻³	0.5432	3e ⁻³	0.1534	8e ⁻⁴	0.5946	2e ⁻³

Proof. Because of the linearity of expectation, we have $\mathbb{E}_D [\hat{V}_{\text{BIPS}}(\pi_\vartheta)] = \mathbb{E}_\beta \left[\frac{\pi_\vartheta(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} r_{\mathbf{x},a} \right]$, and thus:

$$\begin{aligned}
 \text{Bias} \left(\hat{V}_{\text{BIPS}}(\pi_\vartheta) \right) &= \mathbb{E}_D \left[\hat{V}_{\text{BIPS}}(\pi_\vartheta) - V(\pi_\vartheta) \right] \\
 &= \mathbb{E}_\beta \left[\frac{\pi_\vartheta(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} r_{\mathbf{x},a} \right] - \mathbb{E}_\beta \left[\frac{\pi_\vartheta(a|\mathbf{x})}{\beta(a|\mathbf{x})} r_{\mathbf{x},a} \right] \\
 &= \mathbb{E}_\beta \left[\frac{\pi_\vartheta(a|\mathbf{x})}{\beta(a|\mathbf{x})} r_{\mathbf{x},a} \left(\frac{\beta(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} - 1 \right) \right] \\
 &= \mathbb{E}_{\pi_\vartheta} \left[r_{\mathbf{x},a} \left(\frac{\beta(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} - 1 \right) \right]. \tag{13}
 \end{aligned}$$

Since samples are independently sampled from logging policy, the variance can be computed as:

$$\text{Var}_D \left(\hat{V}_{\text{BIPS}}(\pi_\vartheta) \right) = \frac{1}{N} \text{Var}_\beta \left(\frac{\pi_\vartheta(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} r_{\mathbf{x},a} \right).$$

By re-scaling, we get:

$$\begin{aligned}
 N \text{Var}_D \left(\hat{V}_{\text{BIPS}}(\pi_\vartheta) \right) &= \text{Var}_\beta \left(\frac{\pi_\vartheta(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} r_{\mathbf{x},a} \right) \tag{14} \\
 &= \mathbb{E}_\beta \left[\frac{\pi_\vartheta(a|\mathbf{x})^2}{\hat{\beta}(a|\mathbf{x})^2} r_{\mathbf{x},a}^2 \right] - \left(\mathbb{E}_\beta \left[\frac{\pi_\vartheta(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} r_{\mathbf{x},a} \right] \right)^2 \\
 &= \mathbb{E}_{\pi_\vartheta} \left[\frac{\pi_\vartheta(a|\mathbf{x})}{\beta(a|\mathbf{x})} \frac{\beta(a|\mathbf{x})^2}{\hat{\beta}(a|\mathbf{x})^2} r_{\mathbf{x},a}^2 \right] - \left(\mathbb{E}_{\pi_\vartheta} \left[\frac{\beta(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} r_{\mathbf{x},a} \right] \right)^2
 \end{aligned}$$

This completes the proof. \square

Proof of Theorem 3.1:

Proof. We can get:

$$\begin{aligned}
 \text{MSE} \left(\hat{V}_{\text{UIPS}}(\pi_\vartheta) \right) &= \mathbb{E}_D \left[\left(\hat{V}_{\text{UIPS}}(\pi_\vartheta) - V(\pi_\vartheta) \right)^2 \right] \\
 &= \left(\mathbb{E}_D \left[\hat{V}_{\text{UIPS}}(\pi_\vartheta) - V(\pi_\vartheta) \right] \right)^2 + \text{Var}_D \left(\hat{V}_{\text{UIPS}}(\pi_\vartheta) - V(\pi_\vartheta) \right) \\
 &= \left(\mathbb{E}_D \left[\hat{V}_{\text{UIPS}}(\pi_\vartheta) - V(\pi_\vartheta) \right] \right)^2 + \text{Var}_D \left(\hat{V}_{\text{UIPS}}(\pi_\vartheta) \right) \\
 &= \text{Bias}(\hat{V}_{\text{UIPS}}(\pi_\vartheta))^2 + \text{Var}(\hat{V}_{\text{UIPS}}(\pi_\vartheta)).
 \end{aligned}$$

We first bound the bias term:

$$\begin{aligned}
\text{Bias}(\hat{V}_{\text{UIPS}}(\pi_{\boldsymbol{\theta}})) &= \mathbb{E}_D \left[\hat{V}_{\text{UIPS}}(\pi_{\boldsymbol{\theta}}) - V(\pi_{\boldsymbol{\theta}}) \right] \\
&\stackrel{(1)}{=} \mathbb{E}_{\beta} \left[\frac{\pi_{\boldsymbol{\theta}}(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} \phi_{\mathbf{x},a} r_{\mathbf{x},a} \right] - V(\pi_{\boldsymbol{\theta}}) \\
&= \mathbb{E}_{\beta} \left[\frac{\pi_{\boldsymbol{\theta}}(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} \phi_{\mathbf{x},a} r_{\mathbf{x},a} - \frac{\pi_{\boldsymbol{\theta}}(a|\mathbf{x})}{\beta(a|\mathbf{x})} r_{\mathbf{x},a} \right] \\
&= \mathbb{E}_{\beta} \left[r_{\mathbf{x},a} \frac{\pi_{\boldsymbol{\theta}}(a|\mathbf{x})}{\beta(a|\mathbf{x})} \left(\frac{\beta(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} \phi_{\mathbf{x},a} - 1 \right) \right] \\
&\stackrel{(2)}{\leq} \sqrt{\mathbb{E}_{\pi_{\boldsymbol{\theta}}} \left[r_{\mathbf{x},a}^2 \frac{\pi_{\boldsymbol{\theta}}(a|\mathbf{x})}{\beta(a|\mathbf{x})} \right]} \sqrt{\mathbb{E}_{\beta} \left[\left(\frac{\beta(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} \phi_{\mathbf{x},a} - 1 \right)^2 \right]}
\end{aligned}$$

Step (1) follows the linearity of expectation and step (2) is due to the Cauchy-Schwarz inequality. We then bound the variance term:

$$\begin{aligned}
\text{Var}(\hat{V}_{\text{UIPS}}(\pi_{\boldsymbol{\theta}})) &= \frac{1}{N} \text{Var}_{\beta} \left(\frac{\pi_{\boldsymbol{\theta}}(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} \phi_{\mathbf{x},a} r_{\mathbf{x},a} \right) \\
&= \frac{1}{N} \left(\mathbb{E}_{\beta} \left[\frac{\pi_{\boldsymbol{\theta}}(a|\mathbf{x})^2}{\hat{\beta}(a|\mathbf{x})^2} \phi_{\mathbf{x},a}^2 r_{\mathbf{x},a}^2 \right] - \left(\mathbb{E}_{\beta} \left[\frac{\pi_{\boldsymbol{\theta}}(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} \phi_{\mathbf{x},a} r_{\mathbf{x},a} \right] \right)^2 \right) \\
&= \frac{1}{N} \mathbb{E}_{\beta} \left[\frac{\pi_{\boldsymbol{\theta}}(a|\mathbf{x})^2}{\hat{\beta}(a|\mathbf{x})^2} \phi_{\mathbf{x},a}^2 r_{\mathbf{x},a}^2 \right] - \mathbb{E}_{\beta} \left[\frac{\pi_{\boldsymbol{\theta}}(a|\mathbf{x})^2}{\hat{\beta}(a|\mathbf{x})^2} \phi_{\mathbf{x},a}^2 \right]
\end{aligned}$$

Combining the bound of bias and variance completes the proof. \square

Proof of Theorem 3.2:

Proof. We first define several notations:

- $T(\phi_{\mathbf{x},a}, \beta_{\mathbf{x},a}) = \lambda \mathbb{E}_{\beta} \left[\left(\frac{\beta_{\mathbf{x},a}}{\hat{\beta}(a|\mathbf{x})} \phi_{\mathbf{x},a} - 1 \right)^2 \right] + \mathbb{E}_{\beta} \left[\frac{\pi_{\boldsymbol{\theta}}(a|\mathbf{x})^2}{\hat{\beta}(a|\mathbf{x})^2} \phi_{\mathbf{x},a}^2 \right]$.
- $\tilde{T}(\phi_{\mathbf{x},a}) = \max_{\beta_{\mathbf{x},a} \in \mathcal{B}_{\mathbf{x},a}} T(\phi_{\mathbf{x},a}, \beta_{\mathbf{x},a})$ denotes the maximum value of inner optimization problem.
- $T^* = \min_{\phi_{\mathbf{x},a}} \tilde{T}(\phi_{\mathbf{x},a}) = \min_{\phi_{\mathbf{x},a}} \max_{\beta_{\mathbf{x},a} \in \mathcal{B}_{\mathbf{x},a}} T(\phi_{\mathbf{x},a}, \beta_{\mathbf{x},a})$ denote the optimal min-max value. And $\phi_{\mathbf{x},a}^* = \arg \min_{\phi_{\mathbf{x},a}} \tilde{T}(\phi_{\mathbf{x},a})$.
- $\mathbf{B}_{\mathbf{x},a} := \frac{\hat{\beta}(a|\mathbf{x})}{Z} \exp(\gamma U_{\mathbf{x},a}) \hat{\beta}(a|\mathbf{x})$, and $\mathbf{B}_{\mathbf{x},a}^+ := \frac{\hat{\beta}(a|\mathbf{x})}{Z} \exp(\gamma U_{\mathbf{x},a}) \hat{\beta}(a|\mathbf{x})$.

We first find the maximum value of the inner optimization problem, i.e., $\tilde{T}(\phi_{\mathbf{x},a})$ for any fixed $\phi_{\mathbf{x},a}$. And there are three cases shown in Figure 4:

Case I: When $\frac{\hat{\beta}(a|\mathbf{x})}{\phi_{\mathbf{x},a}} \in \mathbf{B}_{\mathbf{x},a}^+$, $\tilde{T}(\phi_{\mathbf{x},a})$ achieves the maximum value at $\beta_{\mathbf{x},a} = \mathbf{B}_{\mathbf{x},a}$. In other words, $\tilde{T}(\phi_{\mathbf{x},a}) = T(\phi_{\mathbf{x},a}, \mathbf{B}_{\mathbf{x},a})$ when $\phi_{\mathbf{x},a} \geq \frac{\hat{\beta}(a|\mathbf{x})}{\mathbf{B}_{\mathbf{x},a}^+}$.

Case II: When $\frac{\hat{\beta}(a|\mathbf{x})}{\phi_{\mathbf{x},a}} \in \mathbf{B}_{\mathbf{x},a}^+$, i.e., $\frac{Z \exp(\gamma U_{\mathbf{x},a})}{Z} \phi_{\mathbf{x},a} \leq \frac{Z \exp(\gamma U_{\mathbf{x},a})}{Z}$, then $\tilde{T}(\phi_{\mathbf{x},a})$ will be the maximum between $T(\phi_{\mathbf{x},a}, \mathbf{B}_{\mathbf{x},a})$ and $T(\phi_{\mathbf{x},a}, \mathbf{B}_{\mathbf{x},a}^+)$.

More specifically, when $\frac{\hat{\beta}(a|\mathbf{x})}{\phi_{\mathbf{x},a}} \in \frac{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}}{2}$, i.e., $\phi_{\mathbf{x},a} = \frac{2\hat{\beta}(a|\mathbf{x})}{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}}$, $\tilde{T}(\phi_{\mathbf{x},a}) = T(\phi_{\mathbf{x},a}, \mathbf{B}_{\mathbf{x},a}^+)$.

Otherwise when $\phi_{\mathbf{x},a} < \frac{2\hat{\beta}(a|\mathbf{x})}{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}}$, $\tilde{T}(\phi_{\mathbf{x},a}) = T(\phi_{\mathbf{x},a}, \mathbf{B}_{\mathbf{x},a})$.

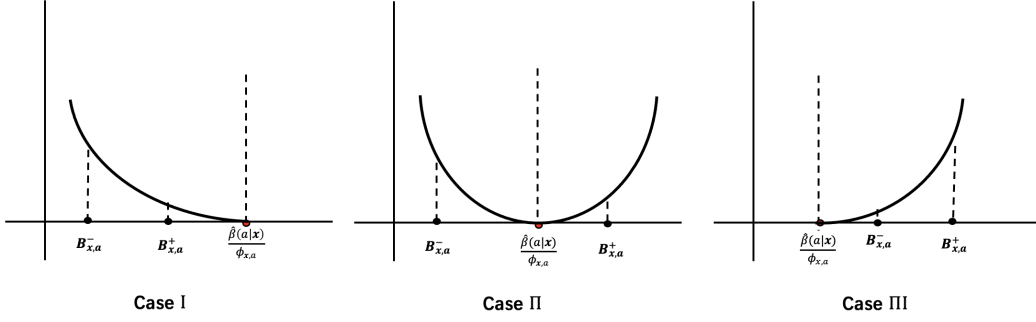


Figure 4: Three cases for maximizing the inner optimization problem.

Case III: When $\phi_{x,a} \geq \frac{\hat{\beta}(a|x)}{B_{x,a}^+}$, implying $\frac{\hat{\beta}(a|x)}{\phi_{x,a}} \leq B_{x,a}^+$, $\tilde{T}(\phi_{x,a}) = T(\phi_{x,a}, B_{x,a}^+)$.

Overall, we get that:

$$\tilde{T}(\phi_{x,a}) = \begin{cases} T(\phi_{x,a}, B_{x,a}), & \phi_{x,a} \geq \left(1, \frac{2\hat{\beta}(a|x)}{B_{x,a}^+ + B_{x,a}}\right) \\ T(\phi_{x,a}, B_{x,a}^+) & \phi_{x,a} \geq \left[\frac{2\hat{\beta}(a|x)}{B_{x,a}^+ + B_{x,a}}, 1\right) \end{cases} \quad (15)$$

Next we try to find the minimum value of $\tilde{T}(\phi_{x,a})$. We first observe that without considering constraint on $\phi_{x,a}$, when

$$\phi_{x,a}^+ = \frac{\lambda}{\lambda \frac{B_{x,a}^+}{\hat{\beta}(a|x)} + \frac{\pi_{\vartheta}(a|x)^2}{\hat{\beta}(a|x)B_{x,a}^+}}$$

$T(\phi_{x,a}, B_{x,a}^+)$ achieves the global minimum value. However, $\phi_{x,a}^+ \leq \frac{\hat{\beta}(a|x)}{B_{x,a}^+} \leq \frac{2\hat{\beta}(a|x)}{B_{x,a}^+ + B_{x,a}}$, which implies when $\phi_{x,a} \geq \left[\frac{2\hat{\beta}(a|x)}{B_{x,a}^+ + B_{x,a}}, 1\right)$, the minimum value of $T(\phi_{x,a}, B_{x,a}^+)$ is achieved at $\frac{2\hat{\beta}(a|x)}{B_{x,a}^+ + B_{x,a}}$.

On the other hand, without considering any constraint on $\phi_{x,a}$, the global minimum value of $T(\phi_{x,a}, B_{x,a})$ is achieved at:

$$\phi_{x,a} = \frac{\lambda}{\lambda \frac{B_{x,a}}{\hat{\beta}(a|x)} + \frac{\pi_{\vartheta}(a|x)^2}{\hat{\beta}(a|x)B_{x,a}}}. \quad (16)$$

Thus if $\phi_{x,a} \leq \frac{2\hat{\beta}(a|x)}{B_{x,a}^+ + B_{x,a}}$, $\phi_{x,a} = \phi_{x,a}$, since $T(\phi_{x,a}, B_{x,a}) \leq T\left(\frac{2\hat{\beta}(a|x)}{B_{x,a}^+ + B_{x,a}}, B_{x,a}\right) = T\left(\frac{2\hat{\beta}(a|x)}{B_{x,a}^+ + B_{x,a}}, B_{x,a}^+\right)$. Otherwise, when $\phi_{x,a} > \frac{2\hat{\beta}(a|x)}{B_{x,a}^+ + B_{x,a}}$, the minimum value of $T(\phi_{x,a}, B_{x,a})$ is also achieved at $\frac{2\hat{\beta}(a|x)}{B_{x,a}^+ + B_{x,a}}$, implying $\phi_{x,a} = \frac{2\hat{\beta}(a|x)}{B_{x,a}^+ + B_{x,a}}$.

Overall,

$$\phi_{x,a} = \min \left(\frac{\lambda}{\lambda \frac{B_{x,a}}{\hat{\beta}(a|x)} + \frac{\pi_{\vartheta}(a|x)^2}{\hat{\beta}(a|x)B_{x,a}}}, \frac{2\hat{\beta}(a|x)}{B_{x,a}^+ + B_{x,a}} \right) \quad (17)$$

Let $\eta = \frac{Z}{\lambda}$, we can get

$$\phi_{x,a} = \min \left(\frac{\lambda}{\frac{\lambda}{\eta} \exp(\gamma U_{x,a}) + \frac{\eta \pi_{\vartheta}(a|x)^2}{\hat{\beta}(a|x)^2 \exp(\gamma U_{x,a})}}, \frac{2\eta}{\exp(\gamma U_{x,a}) + \exp(\gamma U_{x,a})} \right). \quad (18)$$

Note that $\eta \in [\exp(-\gamma U_s^{\max}), \exp(\gamma U_s^{\max})]$, since $\hat{Z} = \frac{\sum_{a^0} \exp(f_\theta(a^0; \mathbf{x}))}{\sum_{a^0} \exp(\gamma U_s^{\max})}$.

This completes the proof. \square

Proof of Corollary 3.3:

Proof. Following the similar procedure as in Theorem 3.1, we can derive the mean squared error (MSE) between $\hat{V}_{\text{BIPS}}(\pi_\vartheta)$ and ground-truth estimator $V(\pi_\vartheta)$ is upper bounded as follows:

$$\begin{aligned} \text{MSE}(\hat{V}_{\text{BIPS}}(\pi_\vartheta)) &= \mathbb{E}_D \left[\left(\hat{V}_{\text{BIPS}}(\pi_\vartheta) - V(\pi_\vartheta) \right)^2 \right] \\ &= \mathbb{E}_{\pi_\vartheta} \left[r_{\mathbf{x},a}^2 \frac{\pi_\vartheta(a; \mathbf{x})}{\beta(a; \mathbf{x})} \right] + \mathbb{E}_\beta \left[\left(\frac{\beta(a; \mathbf{x})}{\hat{\beta}(a; \mathbf{x})} - 1 \right)^2 \right] + \mathbb{E}_\beta \left[\frac{\pi_\vartheta(a; \mathbf{x})^2}{\hat{\beta}(a; \mathbf{x})^2} \right]. \end{aligned}$$

When $\beta(a; \mathbf{x}) \in [\mathbf{B}_{\mathbf{x},a}, \mathbf{B}_{\mathbf{x},a}^+]$, with $\mathbf{B}_{\mathbf{x},a} := \frac{\hat{Z} \exp(-\gamma U_{\mathbf{x},a})}{Z} \hat{\beta}(a; \mathbf{x})$ and $\mathbf{B}_{\mathbf{x},a}^+ := \frac{\hat{Z} \exp(\gamma U_{\mathbf{x},a})}{Z} \hat{\beta}(a; \mathbf{x})$, MSE $(\hat{V}_{\text{BIPS}}(\pi_\vartheta))$ can be further upper bounded as follows.

For ease of illustration, we set $\lambda = \mathbb{E}_{\pi_\vartheta} \left[r_{\mathbf{x},a}^2 \frac{\pi_\vartheta(a; \mathbf{x})}{\beta(a; \mathbf{x})} \right]$ and

$$T(\phi_{\mathbf{x},a}, \beta_{\mathbf{x},a}) = \lambda \mathbb{E}_\beta \left[\left(\frac{\beta_{\mathbf{x},a}}{\hat{\beta}(a; \mathbf{x})} \phi_{\mathbf{x},a} - 1 \right)^2 \right] + \mathbb{E}_\beta \left[\frac{\pi_\vartheta(a; \mathbf{x})^2}{\hat{\beta}(a; \mathbf{x})^2} \phi_{\mathbf{x},a}^2 \right].$$

Let T_{BIPS} denote the upper bound of MSE $(\hat{V}_{\text{BIPS}}(\pi_\vartheta))$, then we can derive that

$$T_{\text{BIPS}} = \begin{cases} T(1, \mathbf{B}_{\mathbf{x},a}^+), & \hat{\beta}(a; \mathbf{x}) \leq \frac{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}}{2} \\ T(1, \mathbf{B}_{\mathbf{x},a}), & \hat{\beta}(a; \mathbf{x}) > \frac{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}}{2}. \end{cases} \quad (19)$$

Recall that in Theorem 3.2, we show that the upper bound of MSE $(\hat{V}_{\text{UIPS}}(\pi_\vartheta))$ is as follows:

$$T_{\text{UIPS}} = \begin{cases} T(\phi_{\mathbf{x},a}, \mathbf{B}_{\mathbf{x},a}), & \phi_{\mathbf{x},a} < \frac{2\hat{\beta}(a; \mathbf{x})}{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}} \\ T\left(\frac{2\hat{\beta}(a; \mathbf{x})}{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}}, \mathbf{B}_{\mathbf{x},a}^+\right), & \phi_{\mathbf{x},a} \geq \frac{2\hat{\beta}(a; \mathbf{x})}{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}}. \end{cases} \quad (20)$$

where $\phi_{\mathbf{x},a}$ is defined in Eq.(16). Next we show that $T_{\text{UIPS}} \leq T_{\text{BIPS}}$.

- When $\hat{\beta}(a; \mathbf{x}) \leq \frac{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}}{2}$, we can get $T_{\text{BIPS}} = T(1, \mathbf{B}_{\mathbf{x},a}^+)$ and $\frac{2\hat{\beta}(a; \mathbf{x})}{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}} \leq 1$. If $\phi_{\mathbf{x},a} < \frac{2\hat{\beta}(a; \mathbf{x})}{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}}$, then $T_{\text{UIPS}} = T(\phi_{\mathbf{x},a}, \mathbf{B}_{\mathbf{x},a}) < T\left(\frac{2\hat{\beta}(a; \mathbf{x})}{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}}, \mathbf{B}_{\mathbf{x},a}\right) = T\left(\frac{2\hat{\beta}(a; \mathbf{x})}{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}}, \mathbf{B}_{\mathbf{x},a}^+\right) = T(1, \mathbf{B}_{\mathbf{x},a}^+)$. And if $\phi_{\mathbf{x},a} \geq \frac{2\hat{\beta}(a; \mathbf{x})}{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}}$, $T_{\text{UIPS}} = T\left(\frac{2\hat{\beta}(a; \mathbf{x})}{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}}, \mathbf{B}_{\mathbf{x},a}^+\right) = T(1, \mathbf{B}_{\mathbf{x},a}^+)$;
- When $\hat{\beta}(a; \mathbf{x}) > \frac{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}}{2}$, we can get $T_{\text{BIPS}} = T(1, \mathbf{B}_{\mathbf{x},a})$ and $\frac{2\hat{\beta}(a; \mathbf{x})}{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}} > 1$. If $\phi_{\mathbf{x},a} < \frac{2\hat{\beta}(a; \mathbf{x})}{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}}$, then $T_{\text{UIPS}} = T(\phi_{\mathbf{x},a}, \mathbf{B}_{\mathbf{x},a}) < T(1, \mathbf{B}_{\mathbf{x},a})$, since $\phi_{\mathbf{x},a}$ is a global minimum. Otherwise when $\phi_{\mathbf{x},a} \geq \frac{2\hat{\beta}(a; \mathbf{x})}{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}} > 1$, $T_{\text{UIPS}} = T\left(\frac{2\hat{\beta}(a; \mathbf{x})}{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}}, \mathbf{B}_{\mathbf{x},a}^+\right) = T\left(\frac{2\hat{\beta}(a; \mathbf{x})}{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}}, \mathbf{B}_{\mathbf{x},a}\right) < T(1, \mathbf{B}_{\mathbf{x},a})$.

In both cases, we have $T_{\text{UIPS}} \leq T_{\text{BIPS}}$, thus completing the proof. \square

Lemma 7.1. Under fixed $\pi_{\vartheta}(a|\mathbf{x})$ and $\hat{\beta}(a|\mathbf{x})$, and $\alpha_{\mathbf{x},a} = \sqrt{\frac{\lambda}{2\eta^2} \frac{\lambda(1-\eta)}{\eta^2} \exp(-2\gamma U_{\mathbf{x},a})}$, we have the following observations:

- If $\frac{\pi_{\vartheta}(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} \alpha_{\mathbf{x},a} \phi_{\mathbf{x},a} = 2\eta / [\exp(\gamma U_{\mathbf{x},a}) + \exp(-\gamma U_{\mathbf{x},a})]$. Otherwise $\phi_{\mathbf{x},a} = \lambda / \left[\frac{\lambda}{\eta} \exp(-\gamma U_{\mathbf{x},a}) + \frac{\eta \pi_{\vartheta}(a|\mathbf{x})^2}{\hat{\beta}^2(a|\mathbf{x}) \exp(\gamma U_{\mathbf{x},a})} \right]$. In other words, $\phi_{\mathbf{x},a} \geq 2\eta$ always holds.
- If $\alpha_{\mathbf{x},a} = \frac{\pi_{\vartheta}(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})}$ and $\frac{\pi_{\vartheta}(a|\mathbf{x})}{\mathbf{B}_{\mathbf{x},a}} < \frac{\rho_{\lambda}}{\lambda}$, larger $U_{\mathbf{x},a}$ brings larger $\phi_{\mathbf{x},a}$.
- Otherwise $\phi_{\mathbf{x},a}$ decreases as $U_{\mathbf{x},a}$ increases.

Proof. The following inequality validates the first observation:

$$\frac{\lambda}{\eta \exp(-\gamma U_{\mathbf{x},a}) + \frac{\eta \pi_{\vartheta}(a|\mathbf{x})^2}{\hat{\beta}^2(a|\mathbf{x}) \exp(\gamma U_{\mathbf{x},a})}} \geq \frac{2\eta}{\exp(\gamma U_{\mathbf{x},a}) + \exp(-\gamma U_{\mathbf{x},a})}$$

For the second and third observations, $\alpha_{\mathbf{x},a} = \frac{\rho_{\lambda}}{\lambda}$. Let $L(u) = \frac{\lambda}{\eta} \exp(-\gamma u) + \frac{\eta \pi_{\vartheta}(a|\mathbf{x})^2}{\hat{\beta}^2(a|\mathbf{x})^2 \exp(\gamma u)}$, we can have:

$$\Gamma_u L(u) = \gamma \frac{\lambda}{\eta} \exp(-\gamma u) + \gamma \frac{\eta \pi_{\vartheta}(a|\mathbf{x})^2}{\hat{\beta}^2(a|\mathbf{x})^2} \exp(\gamma u)$$

To make $\Gamma_u L(u) = 0$, we need $u = \frac{1}{\gamma} \log\left(\frac{\rho_{\lambda} \hat{\beta}(a|\mathbf{x})}{\eta \pi_{\vartheta}(a|\mathbf{x})}\right)$. This implies when $U_{\mathbf{x},a} = \frac{1}{\gamma} \log\left(\frac{\rho_{\lambda} \hat{\beta}(a|\mathbf{x})}{\eta \pi_{\vartheta}(a|\mathbf{x})}\right)$, $\phi_{\mathbf{x},a}$ will decrease as $U_{\mathbf{x},a}$ increases; otherwise as $U_{\mathbf{x},a}$ increases, $\phi_{\mathbf{x},a}$ also increases.

More specifically, when $\alpha_{\mathbf{x},a} = \frac{\pi_{\vartheta}(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})}$, we have:

$$U_{\mathbf{x},a} = \frac{1}{\gamma} \log\left(\frac{\rho_{\lambda} \hat{\beta}(a|\mathbf{x})}{\eta \pi_{\vartheta}(a|\mathbf{x})}\right), \quad \frac{\pi_{\vartheta}(a|\mathbf{x})}{\mathbf{B}_{\mathbf{x},a}} < \frac{\rho_{\lambda}}{\lambda}, \quad \frac{\pi_{\vartheta}(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} < \frac{\rho_{\lambda}}{\eta} \exp(-\gamma U_{\mathbf{x},a}).$$

In other words, for these samples, higher uncertainty implies smaller value of $\pi/\hat{\beta}$, and UIPS tends to boost such safe sample with higher $\phi_{\mathbf{x},a}$.

In other cases, $\phi_{\mathbf{x},a}$ decreases as $U_{\mathbf{x},a}$ increases. This completes the proof. \square

7.9 Convergence Analysis

Next we provide the convergence analysis of policy improvement under UIPS.

Definition 7.2. A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is L -smooth when $\|f(x) - f(y)\|_2 \leq L\|x - y\|_2$, for all $x, y \in \mathbb{R}^d$.

We first state and prove a general result, which serves as the basis to complete our convergence analysis of policy improvement under UIPS. The proof is a special case of convergence proof of stochastic gradient descent with biased gradients in [3]. Suppose we have a differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, which is L -smooth, and attains a finite minimum value $f^* := \min_{x \in \mathbb{R}^d} f(x)$.

Suppose we cannot directly assess the gradient $\nabla f(x)$. Instead we can only assess a noisy but unbiased gradient $\zeta(x) \in \mathbb{R}^d$ at the given x of the function $\tilde{f}(x)$.

Let $b(x) = \nabla \tilde{f}(x) - \nabla f(x)$ denote the difference between $\nabla \tilde{f}(x)$ and $\nabla f(x)$, and $\delta(x) = \zeta(x) - \nabla \tilde{f}(x)$ denote the noise in gradients. We assume that:

$$\|b(x)\|_2 \leq \varphi \|\nabla f(x)\|_2 \quad \text{and} \quad \mathbb{E}[\delta(x)] = 0 \quad \text{and} \quad \mathbb{E}[\|\delta(x)\|_2^2 | x] \leq M \|\nabla \tilde{f}(x)\|_2^2 + \sigma^2 \quad (21)$$

where the constants φ and M satisfy $0 < \varphi < 1$ and $M > 0$ respectively. When running stochastic gradient descent algorithms, i.e. $x_{k+1} = x_k - \eta_k \zeta(x_k)$, we have the following guarantee on the convergence of x_k to an approximate stationary point of f .

Theorem 7.3. Suppose $f(\cdot)$ is differentiable and L -smooth, and the assessed approximate gradient meets the conditions in Eq. (21) with parameters (σ_k, φ_k) at iteration k . Denote $\sigma_{\max} = \max_k \sigma_k$ and $\varphi_{\max} = \max_k \varphi_k$. Set the stepsizes $\eta_k = \min\{\frac{1}{(M+1)L}, 1/(\sigma_{\max} \frac{\rho}{K})\}$, after K iterations, the stochastic gradient descent satisfies :

$$\frac{1}{K} \sum_{k=1}^K \mathbb{E} [kr f(x_k)k^2] \leq \frac{2L(f(x_1) - f^*)}{K(1 - \varphi_{\max})} + \left(L + \frac{2(f(x_1) - f^*)}{(1 - \varphi_{\max})} \right) \frac{\sigma_{\max}^2}{K}$$

Proof. With $\eta_k = \frac{1}{(M+1)L}$, we have:

$$\begin{aligned} f(x_{k+1}) - f(x_k) &= hr f(x_k), x_{k+1} - x_k + \frac{L}{2} kx_{k+1} - x_k k^2 \\ &= f(x_k) - \eta_k hr f(x_k), \zeta(x_k) + \frac{L\eta_k^2}{2} k\zeta(x_k)k^2 \\ &= f(x_k) - \eta_k hr f(x_k), \delta(x_k) + b(x_k) + r f(x_k) + \frac{L\eta_k^2}{2} k\delta(x_k) + b(x_k) + r f(x_k)k^2 \end{aligned} \quad (22)$$

By taking expectations on both side, we have:

$$\begin{aligned} \mathbb{E}[f(x_{k+1})] - \mathbb{E}[f(x_k)] &= \eta_k \mathbb{E}[hr f(x_k), \delta(x_k)] - \eta_k \mathbb{E}[hr f(x_k), b(x) + r f(x_k)] + \frac{L\eta_k^2}{2} \mathbb{E}[k\delta(x_k) + b(x_k) + r f(x_k)k^2] \\ &\stackrel{(1)}{\leq} \mathbb{E}[f(x_k)] - \eta_k \mathbb{E}[hr f(x_k), b(x) + r f(x_k)] + \frac{L\eta_k^2}{2} (\mathbb{E}[k\delta(x_k)k^2] + \mathbb{E}[kb(x) + r f(x_k)k^2]) \\ &= \mathbb{E}[f(x_k)] - \eta_k \mathbb{E}[hr f(x_k), b(x) + r f(x_k)] + \frac{L\eta_k^2}{2} ((M+1)\mathbb{E}[kb(x) + r f(x_k)k^2] + \sigma_k^2) \\ &\stackrel{(2)}{\leq} \mathbb{E}[f(x_k)] + \frac{\eta_k}{2} \mathbb{E}[(- 2hr f(x_k), b(x) + r f(x_k) + kb(x) + r f(x_k)k^2)] + \frac{L\eta_k^2}{2} \sigma_k^2 \\ &= \mathbb{E}[f(x_k)] + \frac{\eta_k}{2} \mathbb{E}[(- kr f(x_k)k^2 + kb(x)k^2)] + \frac{L\eta_k^2}{2} \sigma_k^2 \\ &\stackrel{(3)}{\leq} \mathbb{E}[f(x_k)] + \frac{\eta_k}{2} (\varphi_k - 1) \mathbb{E}[kr f(x_k)k^2] + \frac{L\eta_k^2}{2} \sigma_k^2 \end{aligned}$$

where the inequality labeled as (1) is due to $\mathbb{E}[\delta(x)] = 0$, inequality labeled as (2) is due to $\eta_k = \frac{1}{(M+1)L}$, and inequality labeled as (3) is due to $kb(x_k)k \leq \varphi_k kr f(x_k)k$.

By summing over iterations $k = 1, 2, \dots, K$ and re-arranging the terms, we obtain:

$$\begin{aligned} \frac{1}{2} \sum_{k=1}^K (1 - \varphi_k) \eta_k \mathbb{E}[kr f(x_k)k^2] - f(x_1) - \mathbb{E}[f(x_{K+1})] &+ \frac{L}{2} \sum_{k=1}^K \eta_k^2 \sigma_k^2 \\ &= f(x_1) - f^* + \frac{L}{2} \sum_{k=1}^K \eta_k^2 \sigma_k^2 \end{aligned}$$

where the last inequality follows from $f(x_{K+1}) \leq f^*$. Since $\eta_k = \min\{\frac{1}{(M+1)L}, \frac{1}{\sigma_{\max} \frac{\rho}{K}}\}$, $\forall k = 1, \dots, K$, we can obtain:

$$(1 - \varphi_{\max}) \eta_1 \sum_{k=1}^K \mathbb{E}[kr f(x_k)k^2] \leq 2(f(x_1) - f^*) + LK\eta_1^2 \sigma_{\max}^2$$

Dividing both sides of the above inequality by $K\eta_1(1 - \varphi_{\max})$, we obtain the following,

$$\begin{aligned} \frac{1}{K} \sum_{k=1}^K \mathbb{E}[kr f(x_k)k^2] &\leq \frac{2(f(x_1) - f^*) + LK\eta_1^2 \sigma_{\max}^2}{K\eta_1(1 - \varphi_{\max})} \\ &= \frac{2(f(x_1) - f^*)}{K(1 - \varphi_{\max})} \max\{L, \sigma_{\max} \frac{\rho}{K}\} + L\sigma_{\max}^2 \frac{1}{\sigma_{\max} \frac{\rho}{K}} \\ &= \frac{2L(f(x_1) - f^*)}{K(1 - \varphi_{\max})} + \left(L + \frac{2(f(x_1) - f^*)}{(1 - \varphi_{\max})} \right) \frac{\sigma_{\max}^2}{K} \end{aligned}$$

□

Given this general result, we can now prove Theorem 3.4 by showing that the gradient in UIPS meets the requirements in Theorem 7.3.

Proof of Theorem 3.4:

Proof. UIPS aims to maximize the expected return $V(\pi_\vartheta)$. Therefore, we can utilize Theorem 7.3 by setting $f = V(\pi_\vartheta)$. Since $f = V_{\max}$ and the expected reward is always non-negative, it follows that $f(x_1) = f = V_{\max}$.

We first introduce some additional notations. Let $\rho_\vartheta(\mathbf{x}, a) = \frac{\pi_\vartheta(a|\mathbf{x})}{\beta(a|\mathbf{x})}$ denote the propensity score under ground-truth logging policy, and $\hat{\rho}_\vartheta(\mathbf{x}, a) = \frac{\pi_\vartheta(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} \phi_{\mathbf{x},a}$ represent the propensity score of UIPS. Recall that $\phi_{\mathbf{x},a}$ is derived through solving the optimization problem in Eq (8). Let $g_\vartheta(\mathbf{x}, a) = \frac{\partial \pi_\vartheta(a|\mathbf{x})}{\partial \vartheta}$. The true off-policy policy gradient is computed as follows:

$$r V(\pi_\vartheta) = E_\beta [\rho r g_\vartheta],$$

For UIPS, the approximate policy gradient in each batch with batch size as B is:

$$r \hat{V}_{\text{UIPS}}(\pi_\vartheta) = \frac{1}{B} \sum_{i=1}^B \hat{\rho}_i r_i g_\vartheta^i,$$

which is an unbiased estimate of :

$$r V_{\text{UIPS}}(\pi_\vartheta) = E_\beta [\hat{\rho} r g_\vartheta].$$

To utilize Theorem 7.3, we set $r f = r V(\pi_\vartheta)$, $r \tilde{f} = r V_{\text{UIPS}}(\pi_\vartheta)$, and $\zeta = r \hat{V}_{\text{UIPS}}(\pi_\vartheta)$. We will now demonstrate that the assumptions in Eq. (21) can be satisfied.

Let $\varphi_\vartheta = \max \left\{ \left| \frac{\hat{\rho}_\vartheta(\mathbf{x}, a)}{\rho_\vartheta(\mathbf{x}, a)} - 1 \right| \right\}$, We first have:

$$\begin{aligned} k b(\vartheta) k &= k r V_{\text{UIPS}}(\pi_\vartheta) - r V(\pi_\vartheta) k \\ &= k E_\beta [(\hat{\rho} - \rho) r g_\vartheta] k = k E_\beta [\rho \left(\frac{\hat{\rho}}{\rho} - 1 \right) r g_\vartheta] k \\ &\quad \varphi_\vartheta k E_\beta [\rho r g_\vartheta] k \end{aligned} \tag{23}$$

And next we show that $0 < \varphi_\vartheta < 1$.

$$\varphi_\vartheta = \max \left\{ \left| \frac{\hat{\rho}_\vartheta(\mathbf{x}, a)}{\rho_\vartheta(\mathbf{x}, a)} - 1 \right| \right\} = \max \left\{ \left| \frac{\beta(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} \phi_{\mathbf{x},a} - 1 \right| \right\} \tag{24}$$

Recall from Eq.(17) in proof of Theorem 3.2, we have :

$$\phi_{\mathbf{x},a} = \min \left(\frac{\lambda}{\lambda \frac{\mathbf{B}_{\mathbf{x},a}}{\hat{\beta}(a|\mathbf{x})} + \frac{\pi_\vartheta(a|\mathbf{x})^2}{\beta(a|\mathbf{x}) \mathbf{B}_{\mathbf{x},a}}}, \frac{2 \hat{\beta}(a|\mathbf{x})}{\mathbf{B}_{\mathbf{x},a}^+ + \mathbf{B}_{\mathbf{x},a}} \right)$$

where $\mathbf{B}_{\mathbf{x},a} := \frac{\hat{Z} \exp(\gamma U_{\mathbf{x},a})}{Z} \hat{\beta}(a|\mathbf{x})$, and $\mathbf{B}_{\mathbf{x},a}^+ := \frac{\hat{Z} \exp(\gamma U_{\mathbf{x},a})}{Z} \hat{\beta}(a|\mathbf{x})$. Thus we have:

$$\frac{\beta(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} \phi_{\mathbf{x},a} = \min \left(\frac{\lambda}{\lambda \frac{\mathbf{B}_{\mathbf{x},a}}{\beta(a|\mathbf{x})} + \frac{\pi_\vartheta(a|\mathbf{x})^2}{\beta(a|\mathbf{x}) \mathbf{B}_{\mathbf{x},a}}}, \frac{2}{\frac{\mathbf{B}_{\mathbf{x},a}^+}{\beta(a|\mathbf{x})} + \frac{\mathbf{B}_{\mathbf{x},a}}{\beta(a|\mathbf{x})}} \right) \tag{25}$$

Since $\mathbf{B}_{\mathbf{x},a}^+ \leq \beta(a|\mathbf{x})$, thus $0 \leq \frac{\beta(a|\mathbf{x})}{\hat{\beta}(a|\mathbf{x})} \phi_{\mathbf{x},a} \leq 2$, implying $0 < \varphi_\vartheta < 1$.

Also, since $r \hat{V}_{\text{UIPS}}(\pi_\vartheta)$ is an unbiased estimate of $r V_{\text{UIPS}}(\pi_\vartheta)$, we have:

$$E[\delta(\vartheta)] = E[r \hat{V}_{\text{UIPS}}(\pi_\vartheta) - r V_{\text{UIPS}}(\pi_\vartheta)] = \mathbf{0} \tag{26}$$

Finally, we have the following,

$$\begin{aligned}
E[k\delta(\vartheta)k^2] &= E \left[k r \hat{V}_{\text{UIPS}}(\pi_\vartheta) - r V_{\text{UIPS}}(\pi_\vartheta) k^2 \right] = E \left[\left\| \frac{1}{B} \sum_{i=1}^B (\hat{\rho}_i r_i g_\vartheta^i - E_\beta [\hat{\rho} r g_\vartheta]) \right\|^2 \right] \\
&= \frac{1}{B^2} E \left[\left(\sum_{i=1}^B k \hat{\rho}_i r_i g_\vartheta^i - E_\beta [\hat{\rho} r g_\vartheta] k \right)^2 \right] = \frac{1}{B} E \left[\sum_{i=1}^B k \hat{\rho}_i r_i g_\vartheta^i - E_\beta [\hat{\rho} r g_\vartheta] k^2 \right] \\
&= E_\beta [k \hat{\rho} r g_\vartheta - E_\beta [\hat{\rho} r g_\vartheta] k^2] \stackrel{(1)}{=} E_\beta [k \hat{\rho} r g_\vartheta k^2] - k E_\beta [\hat{\rho} r g_\vartheta] k^2 \\
&= E_\beta [k \hat{\rho} r g_\vartheta k^2] - G_{\max}^2 \Phi
\end{aligned} \tag{27}$$

where the equality labeled as (1) is due to $E[kY - E[Y]k^2] = E[kY k^2] - kE[Y]k^2$.

Hence, by applying Theorem 7.3, we have:

$$\frac{1}{K} \sum_{k=1}^K E[k r V(\pi_{\vartheta_k}) k^2] \leq \frac{2LV_{\max}}{K(1 - \varphi_{\max})} + \left(L + \frac{2V_{\max}}{(1 - \varphi_{\max})} \right) \frac{G_{\max}^2 \Phi}{K} \tag{28}$$

□