Hyperparameter Optimization Can Even be Harmful in Off-Policy Learning and How to Deal with It

Yuta Saito¹, Masahiro Nomura²

¹Cornell University, ²CyberAgent, Inc. ys552@cornell.edu, nomura_masahiro@cyberagent.co.jp

Abstract

There has been a growing interest in off-policy evaluation in the literature such as recommender systems and personalized medicine. We have so far seen significant progress in developing estimators aimed at accurately estimating the effectiveness of counterfactual policies based on biased logged data. However, there are many cases where those estimators are used not only to evaluate the value of decision making policies but also to search for the best hyperparameters from a large candidate space. This work explores the latter hyperparameter optimization (HPO) task for off-policy learning. We empirically show that naively applying an unbiased estimator of the generalization performance as a surrogate objective in HPO can cause an unexpected failure, merely pursuing hyperparameters whose generalization performance is greatly overestimated. We then propose simple and computationally efficient corrections to the typical HPO procedure to deal with the aforementioned issues simultaneously. Empirical investigations demonstrate the effectiveness of our proposed HPO algorithm in situations where the typical procedure fails severely.

1 Introduction

Interactive decision making systems, such as recommender systems, produce logged data valuable for optimizing future decision making. For example, the logs of an e-commerce recommender system record which product was recommended and whether the users purchased it, giving the system designer a rich logged dataset useful for evaluating and improving the decision making quality. This type of historical data is often called *logged bandit data* and is one of the most ubiquitous forms of data available in many real-life applications [Swaminathan and Joachims, 2015a; Su *et al.*, 2020; Kiyohara *et al.*, 2021; Saito *et al.*, 2024].

Off-Policy Learning (OPL) aims to train a new decision making policy using only the logged bandit data. OPL is useful in that it can improve the decision making system continuously in a batch manner without requiring a risky exploration. Owing to the ubiquity of logged bandit data in the real-world, significant attention has been

paid to OPL of contextual bandits [Strehl *et al.*, 2010; Swaminathan and Joachims, 2015a; 2015b; Wang *et al.*, 2017; Kallus *et al.*, 2021; Kiyohara *et al.*, 2023; 2024].

The fundamental problem in OPL is that the outcome is only observed for the action chosen by the system in the past. Thus, estimating the generalization performance of a policy is non-trivial because we cannot naively apply the empirical risk as done in typical supervised machine learning (ML). Therefore, a variety of estimators have been developed in the field of off-policy evaluation (OPE), such as Inverse Propensity Score (IPS) [Precup *et al.*, 2000] and Doubly Robust (DR) [Dudík *et al.*, 2014]. Then, a feasible approach to OPL is to maximize one such estimator as a surrogate objective using only the logged data. Hyperparameter optimization (HPO) can also be performed based on one of the estimators on a validation set of the logged data [Paine *et al.*, 2020].

In this study, we investigate how well automatic HPO algorithms work for OPL using only the available logged data. In particular, we empirically find two critical issues in HPO that have yet to be investigated in the literature, but can have a significant adverse impact on the effectiveness of the OPL pipeline. The first issue is *optimistic bias*, which implies that the hyperparameter values selected by an HPO procedure are often the ones whose performance is greatly overestimated. In HPO, we often use an unbiased estimator as a strategy to optimize the generalization performance (primary objective) using only validation data. The problem is that, when optimizing the validation performance as a surrogate objective, HPO can identify a set of hyperparameters whose validation performance looks good but its generalization performance is detrimental. As a result, the typical HPO procedure often produces a highly sub-optimal solution, even with an unbiased estimator of the generalization performance. The second issue is *unsafe behavior*, which suggests that the typical HPO procedure can output a solution, which underperforms the logging (data collection) policy, even when we set the logging policy as an initial solution. This is problematic because a logging policy is often a baseline policy to improve upon in OPL. If an HPO procedure aggravates the performance of the logging policy, there is no need to implement it in practice.

After formulating the problem in Section 2, Section 3 provides clear empirical evidence of optimistic bias and unsafe behavior. We observe these phenomena even when we use an unbiased surrogate objective and a popular adaptive HPO

algorithm. We also explain these observations theoretically, demonstrating that ignoring the fact that HPO optimizes the validation performance as a surrogate of the generalization performance can lead to a worse regret of HPO algorithms. More specifically, we identify that a heavy-tailed distribution of overestimation bias during HPO can cause an unexpected gap between the generalization and validation regret. These empirical and theoretical observations result in our proposed corrections to the typical HPO procedure, which we describe in Section 4. Finally, Section 5 conducts comprehensive experiments and demonstrates that our simple corrections can deal with the aforementioned issues and improve the typical procedure, particularly for cases where the typical procedure becomes unsafe and underperforms the logging policy.¹

2 Preliminaries

We use $x \in \mathcal{X}$ to denote a context vector and $a \in \mathcal{A}$ to denote a (discrete) action such as a playlist recommendation in a music streaming service. Let $r \in [0, r_{\max}]$ denote a reward variable, which is sampled identically and independently from an unknown conditional distribution p(r|x,a). A decision making policy is modeled as a distribution over the action space, i.e., $\pi: \mathcal{X} \to \Delta(\mathcal{A})$ where $\Delta(\cdot)$ is a probability simplex. We can then represent the probability of action a being taken by policy π given context x as $\pi(a|x)$.

2.1 Off-Policy Evaluation and Learning

In OPE, we are given logged bandit data $\mathcal{D} := \{(x_i, a_i, r_i)\}_{i=1}^n$ consisting of n independent draws from the logging policy π_0 . Using this logged dataset, OPE aims to estimate the generalization performance of a given evaluation policy π_e , which is often different from π_0 :

$$V(\pi_e) := \mathbb{E}_{(x,a,r) \sim p(x)\pi_e(a|x)p(r|x,a)}[r]. \tag{1}$$

This is the ground-truth performance of the evaluation policy when deployed in an environment of interest. OPE uses an estimator \hat{V} to estimate $V(\pi_e)$ based only on \mathcal{D} as $V(\pi_e) \approx \hat{V}(\pi_e; \mathcal{D})$. A typical choice of \hat{V} is IPS:

$$\hat{V}_{\text{IPS}}(\pi_e; \mathcal{D}) := \frac{1}{n} \sum_{i=1}^n \frac{\pi_e(a_i \mid x_i)}{\pi_0(a_i \mid x_i)} r_i,$$

where $\pi_e(a_i|x_i)/\pi_0(a_i|x_i)$ is called the importance weight. Under some assumptions for identification such as full support $(\pi_e(a|x)>0\to\pi_0(a|x)>0,\ \forall (x,a))$, IPS provides an unbiased estimate of the generalization policy performance, i.e., $\mathbb{E}_{\mathcal{D}}[\hat{V}_{\text{IPS}}(\pi_e;\mathcal{D})]=V(\pi_e)$. Beyond IPS, significant efforts have been made to enable a more accurate OPE from the logged data [Dudík $et\ al.$, 2014; Wang $et\ al.$, 2017; Su $et\ al.$, 2020; Saito $et\ al.$, 2023].

In OPL, we aim to learn an optimal decision making policy $\pi^* := \arg\max_{\pi} V(\pi)$ from the logged data. As in supervised

ML, we cannot directly use the generalization policy performance. Instead, we use its estimator as a surrogate:

$$\hat{\boldsymbol{\pi}} = \underset{\boldsymbol{\pi} \in \Pi}{\operatorname{arg\ max}} \ \hat{V}(\boldsymbol{\pi}; \mathcal{D}) - \lambda \cdot \mathcal{R}(\boldsymbol{\pi}),$$

Algorithm 1 Typical HPO with IPS as a surrogate (Baseline)

```
Input: A, \Theta, T, \mathcal{D}_{tr}, \mathcal{D}_{val}

Output: \hat{\theta}

1: \pi^* \leftarrow \pi_0, \mathcal{S}_0 \leftarrow \emptyset

2: for t = 1, 2, ..., T do

3: \theta_t \leftarrow A(\theta \mid \mathcal{S}_{t-1}) // sample candidate hyperparameters

4: \pi_t \leftarrow \hat{\pi}(\cdot \mid \theta_t, \mathcal{D}_{tr}) // train a policy (lower-level)

5: if \hat{V}_{IPS}(\pi_t; \mathcal{D}_{val}) > \hat{V}_{IPS}(\pi^*; \mathcal{D}_{val}) then

6: \hat{\theta} \leftarrow \theta_t, \pi^* \leftarrow \pi_t // update the solution

7: end if

8: \mathcal{S}_t \leftarrow \mathcal{S}_{t-1} \cup \{(\theta_t, \hat{V}_{IPS}(\pi_t; \mathcal{D}_{val}))\} // store the result

9: end for
```

where Π is a policy class, which might be a linear class [Swaminathan and Joachims, 2015a] or deep neural nets [Joachims *et al.*, 2018]. $\mathcal{R}(\cdot)$ regularizes the complexity of the policy π , and $\lambda(\geq 0)$ is a hyperparamter that controls the effect of regularization.

2.2 Hyperparameter Optimization

OPE involves many hyperparameters to be properly tuned from those defining the policy class Π to the regularization parameter λ . In a typical HPO procedure for OPL, we first split the original logged bandit data \mathcal{D} into training (\mathcal{D}_{tr}) and validation (\mathcal{D}_{val}) sets. Then, we wish to solve the following bi-level optimization:

$$\theta^* := \underset{\theta \in \Theta}{\operatorname{arg max}} \ V(\hat{\pi}(\cdot; \theta, \mathcal{D}_{tr})), \tag{2}$$

where Θ is a pre-defined hyperparamter search space. $\hat{\pi}(\cdot; \theta, \mathcal{D}_{tr})$ is a policy parameterized by a set of hyperparameters θ . The model parameter of $\hat{\pi}(\cdot; \theta, \mathcal{D}_{tr})$ is trained on the training set \mathcal{D}_{tr} (lower-level optimization). The problem here is that the generalization performance of $\hat{\pi}(\cdot; \theta, \mathcal{D}_{tr})$ is unknown and needs to be estimated. A feasible HPO procedure based on an estimated policy performance is:

$$\hat{\theta}(\mathcal{D}_{val}) := \underset{\theta \in \Theta}{\arg \max} \ \hat{V}(\hat{\pi}(\cdot; \theta, \mathcal{D}_{tr}); \mathcal{D}_{val}), \tag{3}$$

where the generalization performance of $\hat{\pi}(\cdot;\theta,\mathcal{D}_{tr})$ is estimated by an estimator \hat{V} on the validation set \mathcal{D}_{val} .² A common choice of \hat{V} is an unbiased estimator that satisfies $\mathbb{E}[\hat{V}(\pi;\mathcal{D}_{val})] = V(\pi), \forall \pi \in \Pi$ such as IPS. Then, one can apply grid search, random search [Bergstra and Bengio, 2012], or adaptive methods such as tree-structured Parzen estimator (TPE) [Bergstra *et al.*, 2011] to solve the higher-level optimization in Eq. (3) efficiently. Algorithm 1 describes this typical HPO procedure for OPL, which starts from the logging policy π_0 as its initial solution and adaptively samples promising hyperparameters via an arbitrary HPO algorithm (denoted here as A) [Tang and Wiens, 2021].

3 Unexpected Failure in HPO for OPL

This section studies the effectiveness of HPO when applied to OPL from both empirical and theoretical perspectives.

¹Appendix B provides a comprehensive survey of related work.

²For brevity of notation, we sometimes use $V(\theta)$ and $\hat{V}(\theta; \mathcal{D})$ to denote the generalization and validation performances of the policy induced by θ .

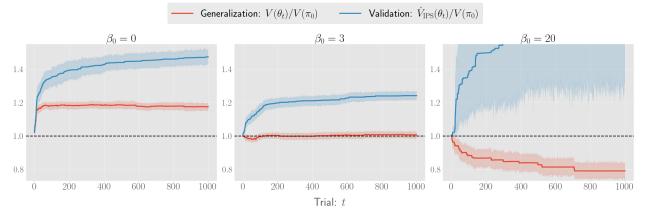


Figure 1: Empirical Evidence of Optimistic Bias and Unsafe Behavior in HPO for OPL (w/ TPE). The results are averaged over 25 runs with different seeds and then normalized by $V(\pi_0)$. The shaded regions indicate 95% confidence intervals.

3.1 Empirical Analysis

First, we conduct a synthetic experiment and provide empirical evidence of surprising failure of HPO in OPL.

Synthetic Data.

Our empirical analysis is based on *OpenBanditPipeline* (OBP)³, an open-source toolkit for OPE and OPL, which includes synthetic data generation modules and a range of estimators [Saito *et al.*, 2021a]. We synthesize context vectors x by sampling them from a 10-dimensional standard normal distribution. We then set $|\mathcal{A}| = 10$, where each action $a \in \mathcal{A}$ is characterized by a 10-dimensional representation vector e_a . The reward function $\mu(x,a) := \mathbb{E}[r \mid x,a]$ is defined as:

$$\mu(x,a) = \sigma \left(x^{\top} M e_a + \eta_x^{\top} x + \eta_a^{\top} e_a \right), \tag{4}$$

where $\sigma(z) := 1/(1 + \exp(-z))$ is the sigmoid function. M, η_x , and η_a are parameter matrices or vectors for defining the synthetic reward function. These parameters are sampled from a uniform distribution with range [-1,1]. After generating the synthetic reward function, we sample binary rewards from a Bernoulli distribution with parameter $\mu(x,a)$.

We then define the logging policy π_0 by applying the soft-max function to the reward function $\mu(x,a)$ as follows.

$$\pi_0(a \,|\, x) = \frac{\exp(\beta_0 \cdot \mu(x, a))}{\sum_{a' \in \mathcal{A}} \exp(\beta_0 \cdot \mu(x, a'))},\tag{5}$$

where β_0 is an inverse temperature parameter to control the optimality and entropy of the logging policy. A large positive value of β_0 leads to a near-deterministic and near-optimal logging policy. When $\beta_0 = 0$, π_0 is uniform.

Policy Class and HPO Algorithms.

To train a new policy π from only the logged data, we first estimate $\mu(x,a)$ by a supervised ML method, where the resulting estimator is denoted as $\hat{\mu}(x,a;\mathcal{D}_{tr})$. We then form a stochastic policy by applying the softmax rule as:

$$\pi(a \mid x; \theta, \mathcal{D}_{tr}) = \frac{\exp(\beta \cdot \hat{\mu}(x, a; \mathcal{D}_{tr}))}{\sum_{a' \in \mathcal{A}} \exp(\beta \cdot \hat{\mu}(x, a'; \mathcal{D}_{tr}))}, \quad (6)$$

where β is an inverse temperature parameter to define a new policy. θ is a set of hyperparameters, which consists of β , supervised ML model to construct $\hat{\mu}$, and the hyperparameters of $\hat{\mu}$. The hyperparameter search space Θ is summarized in Table 1 in Appendix E.

As an HPO algorithm, we use TPE [Bergstra *et al.*, 2011], which is a popular adaptive method in the HPO community [Akiba *et al.*, 2019]. TPE has been shown to work well for HPO of supervised ML, however, whether it also works for OPL has never been thoroughly investigated.

Observations.

In this synthetic experiment, we set $\beta_0 \in \{0, 3, 20\}$ and $|\mathcal{D}_{tr}| = |\mathcal{D}_{val}| = 1,000$. The number of trials (T in Algorithm 1) for HPO is set to 1,000.

Figure 1 shows the **validation performance** $(\hat{V}_{IPS}(\pi; \mathcal{D}_{val});$ what HPO algorithm maximizes from the logged data) and the **generalization performance** $(V(\pi);$ the primary objective of OPL) during the HPO procedure. We obtain the following key observations in this experiment.

- 1. **Optimistic Bias**: For all β_0 , TPE succeeds in maximizing the validation performance, monotonically improving the blue lines. However, there is a substantial gap between validation and generalization, and the validation performance becomes an extremely optimistic proxy of the generalization performance. For example, when $\beta_0=3$, TPE does not bring any impact on the generalization performance, although the validation performance is greatly improved. This result suggests that implementing HPO is indeed a waste of time and resources for this setting.
- 2. Unsafe Behavior: When $\beta_0 = 20$ (where π_0 is already much better than uniform random), TPE outputs a solution that is significantly worse than the logging policy with respect to the generalization performance. This is problematic, as the solution at the final trial seems to provide a substantial improvement over the logging policy with respect to the unbiased validation performance (blue lines). In reality, we have no access to the generalization performance (red lines), making it impossible to

³https://github.com/st-tech/zr-obp

detect this performance degradation, possibly deploying an unsafe policy in the field without even noticing it.

These observations suggest that optimizing an unbiased surrogate objective is not an ideal strategy and is even harmful in some cases regarding the optimization of the generalization performance. Note that we obtain similar results when random search (RS) is used as an HPO algorithm and DR is used as an OPE estimator as reported in Appendix E. In particular, comparing RS with TPE in terms of the generalization performance, we find that there are no particular differences between the two algorithms for $\beta_0=0,3$. Even more surprisingly, when $\beta_0=20$, TPE is outperformed by RS, even if TPE is better at optimizing the validation performance. These results further suggest that merely optimizing an unbiased surrogate objective is not a suitable approach for optimizing the generalization performance in HPO of OPL.

3.2 Theoretical Analysis

Next, we investigate the mechanism causing the somewhat surprising issues observed in the previous section.⁴ First, we explain the phenomena from a statistical perspective.

Proposition 3.1. Given that \hat{V} is unbiased, we have the following inequalities.

$$\mathbb{E}_{\mathcal{D}}\left[\hat{V}(\hat{\theta}(\mathcal{D}); \mathcal{D})\right] \ge V\left(\theta^*\right) \ge \mathbb{E}_{\mathcal{D}}\left[V(\hat{\theta}(\mathcal{D}))\right], \quad (7)$$

where $\mathbb{E}_{\mathcal{D}}[\cdot]$ takes expectation over every randomness in the logged data \mathcal{D} , and $\mathbb{E}_{\mathcal{D}}[\hat{V}(\hat{\theta}(\mathcal{D});\mathcal{D})] - \mathbb{E}_{\mathcal{D}}[V(\hat{\theta}(\mathcal{D}))]$ is the amount of optimistic bias.

Note that, in Eq. (2), $V(\theta^*)$ is defined as the best generalization performance we could achieve with HPO. Thus, the first inequality in Eq. (7) suggests that the *validation* performance of the HPO solution $\hat{\theta}(\mathcal{D}_{val})$ is better than the best achievable generalization performance in expectation, suggesting that the performance estimation of the HPO solution is optimistic in general. In addition, the second inequality in Eq. (7) implies that the *generalization* performance of $\hat{\theta}(\mathcal{D}_{val})$ is worse than the best achievable generalization performance in expectation, even though the validation performance of $\hat{\theta}(\mathcal{D}_{val})$ is likely to be better. As a result, we will often be disappointed with the performance of the HPO solution $\hat{\theta}$ even with an unbiased surrogate (validation) objective. Overall, Proposition 3.1 explains the substantial gap between the blue $(\mathbb{E}[\hat{V}(\hat{\theta}(\mathcal{D}_{val}))])$ and red $(\mathbb{E}[V(\hat{\theta}(\mathcal{D}_{val}))])$ lines observed in Figure 1.

Next, we analyze "regret" to understand what causes the optimistic bias in Proposition 3.1 and how we can deal with it. For this, we define two variants of regret, which measure the difference between the validation or generalization performances of the optimal hyperparameter and HPO solution.

$$r_{aen}(T; A, \mathcal{D}) := V(\theta^*) - V(\hat{\theta}_{T, A}(\mathcal{D})), \tag{8}$$

$$r_{val}(T; A, \mathcal{D}) := \hat{V}_{IPS}(\hat{\theta}^*; \mathcal{D}) - \hat{V}_{IPS}(\hat{\theta}_{T,A}(\mathcal{D}); \mathcal{D}), \quad (9)$$

where $\theta^* := \underset{\theta \in \Theta}{\arg \max} V(\theta)$ is the optimal hyperparame-

ter with respect to the generalization performance. $\hat{\theta}^* =$

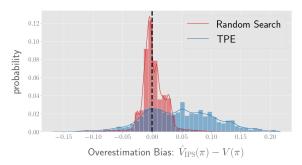


Figure 2: Distributions of Overestimation Bias ($\beta_0 = 3$)

 $rg \max_{\theta \in \Theta} \hat{V}_{\mathrm{IPS}}(\theta; \mathcal{D})$ denotes the optimal hyperparameter with

respect to the *validation* performance, and $\hat{\theta}_{T,A}(\mathcal{D})$ is the solution of Algorithm 1 given budget T and algorithm A. We also define the overestimation bias for a specific hyperparameter θ as $\tau(\theta;\mathcal{D}) = \hat{V}_{\mathrm{IPS}}(\theta;\mathcal{D}) - V(\theta)$. Then, the following implies that a heavy-tailed distribution of overestimation bias during HPO can produce an unexpected gap between the generalization and validation regret.

Proposition 3.2. Given HPO algorithm A, budget T, and logged data \mathcal{D} , the generalization regret can be written as

$$r_{gen}(T; A, \mathcal{D}) = r_{val}(T; A, \mathcal{D}) + \Delta \tau(\hat{\theta}_{T, A}(\mathcal{D}), \theta^*; \mathcal{D}) + C,$$
(10)

where
$$\Delta \tau(\theta_1, \theta_2; \mathcal{D}) := \tau(\theta_1; \mathcal{D}) - \tau(\theta_2; \mathcal{D})$$
, and $C := \hat{V}_{IPS}(\theta^*; \mathcal{D}) - \hat{V}_{IPS}(\hat{\theta}^*; \mathcal{D})$.

Only the first two terms of the RHS in Eq. (10) depend on the HPO solution $\hat{\theta}_{T,A}(\mathcal{D})$, and are thus critical for analyzing the HPO performance. The first term r_{val} is the validation regret. Under some mild conditions, we can achieve no-regret $(r_{val}(T; A, \mathcal{D}) = o(1))$ with optimization methods such as GP-UCB [Srinivas et al., 2010], as we can target the validation performance directly using available data. The second term $\Delta \tau(\hat{\theta}_{T,A}(\mathcal{D}), \theta^*; \mathcal{D})$ is the difference in the extent of overestimation between $\theta_{T,A}(\mathcal{D})$ and θ^* . When the extent of overestimation of $\hat{\theta}_{T,A}(\mathcal{D})$ is larger than that of $\theta^*, \Delta \tau(\hat{\theta}_{T,A}(\mathcal{D}), \theta^*; \mathcal{D})$ becomes large. Therefore, Proposition 3.2 suggests that the overestimation bias of $\hat{\theta}_{T,A}(\mathcal{D})$ can exacerbate the generalization regret of HPO algorithms. More specifically, if an HPO algorithm is likely to sample many hyperparameters whose performance is overestimated $(\hat{V}_{IPS}(\theta) - V(\theta) > 0)$ and the overestimation bias has a heavytailed distribution, the second term of Eq. (10) tends to become large, so does the generalization regret r_{gen} . Given this regret analysis, we investigate the distributions of overestimation bias observed in the empirical analysis in Figure 2. This figure implies that TPE more frequently samples hyperparameters incurring a large overestimation bias than RS. According to Proposition 3.2, this is why we do not find the advantage of TPE with respect to the generalization performance. RS has a worse validation regret than TPE, while overestimation bias of RS is not very problematic compared to TPE. As a result, RS performs similarly to or slightly better than TPE in terms of

⁴Appendix C provides proofs omitted in the main text.

the generalization performance. In this way, the heavy-tailed distribution of overestimation bias makes the generalization regret of HPO algorithms (in particular TPE) worse than its validation regret, resulting in optimistic bias and possibly unsafe behavior.

4 How Should We Deal with the Issues?

In this section, we propose two simple corrections, namely (i) **conservative surrogate objective** and (ii) **adaptive imitation regularization**, to deal with the critical issues in HPO. We also describe the resulting HPO procedure, which we call **Conservative and Imitation-Regularized HPO (CIR-HPO)**.

4.1 Conservative Surrogate Objective (CSO)

First, we address the heavy-tailed distribution of overestimation bias $(\hat{V}_{\mathrm{IPS}}(\pi) - V(\pi))$ during HPO, as suggested in Figure 2. Proposition 3.2 implies that the overestimation of the value of hyperparameters sampled during HPO can exacerbate the generalization regret of an HPO algorithm. To deal with this issue, we introduce *conservative surrogate objective*, which penalizes the validation performance of hyperparameters whose performance has a large uncertainty to avoid the issue of overestimation bias during HPO. Specifically, we propose to use a high probability lower bound of the generalization performance (denoted as $\hat{V}_{-}(\cdot)$) as an alternative surrogate objective, which is given as: $\mathbb{P}(V(\pi) \geq \hat{V}_{-}(\pi; \mathcal{D}, \delta)) \geq 1 - \delta$ where $\delta \in (0,1)$ specifies a confidence level.

A prevalent strategy to construct $\hat{V}_{-}(\cdot)$ in OPE is to apply a concentration inequality such as Hoeffding and Bernstein [Thomas *et al.*, 2015b; 2015a]. A problem is that these inequalities are often overly conservative as they make no assumptions about underlying distribution. Thus, we use an alternative strategy to construct $\hat{V}_{-}(\cdot)$ based on the Student's t-distribution as follows.

$$\hat{V}_{-}^{t}(\pi; \mathcal{D}, \delta) := \hat{V}_{\text{IPS}}(\pi; \mathcal{D}) - t_{1-\delta, \nu} \sqrt{\frac{\mathbb{V}_{n}(\hat{V}_{\text{IPS}}(\pi; \mathcal{D}))}{n-1}},$$
(11)

where $t_{1-\delta,\nu}$ is the T-value given confidence level δ and degrees of freedom ν .

The upside of Eq. (11) is that it produces a tighter lower bound than aforementioned concentration inequalities. This is because Eq. (11) introduces the additional assumption that the mean of importance weighted rewards $(\pi/\pi_0)r$ is normally distributed. This assumption is reasonable with growing data sizes. However, $(\pi/\pi_0)r$ often follows a distribution with heavy upper tails, which may make the assumption invalid in a small sample setting. Nonetheless, Appendix E empirically verifies that Eq. (11) is *reasonably tight* compared with other popular concentration inequalities.

4.2 Adaptive Imitation Regularization (AIR)

The second technique we propose is *adaptive imitation regularization*, which tackles the unsafe behavior of the typical procedure. The issue of unsafe behavior suggests that, if logging policy π_0 is better than uniform random or is near-optimal, Algorithm 1 can produce a solution whose performance is

Algorithm 2 Conservative and Imitation-Regularized HPO

```
Input: A, \delta, \gamma, \Theta, T, \pi_0, \mathcal{D}_{tr}, \mathcal{D}_{val}
Output: \hat{\theta}

1: \mathcal{S}_0 \leftarrow \emptyset

2: for t = 1, 2, ..., T do

3: \theta_t \leftarrow A(\theta \mid \mathcal{S}_{t-1}) // sample candidate hyperparameters

4: \hat{\pi}_t \leftarrow \hat{\pi}(\cdot \mid \theta_t, \mathcal{D}_{tr}) // train a policy (lower-level)

5: \pi_t \leftarrow (1 - \alpha_t) \cdot \hat{\pi}_t + \alpha_t \cdot \pi_0 // regularization (Eq. (14))

6: if \hat{V}_-^t(\pi_t; \mathcal{D}_{val}, \delta) \geq \hat{V}_-^t(\pi^*; \mathcal{D}_{val}, \delta) then

7: \hat{\theta} \leftarrow \theta_t, \pi^* \leftarrow \pi_t // update the solution

8: end if

9: \mathcal{S}_t \leftarrow \mathcal{S}_{t-1} \cup \{(\theta_t, \hat{V}_-^t(\pi_t; \mathcal{D}_{val}, \delta))\} // store the result
```

much worse than that of the logging policy. Avoiding this problem is non-trivial, because we do not have access to the generalization performance and do not know the optimality of the logging policy in practice. For example, simply setting π_0 as an initial solution does not solve the issue at all, as suggested in Section 3.1. An instant idea might be to *imitate* the logging policy to some extent:

$$\pi_t(a|x;\alpha,\theta_t,\mathcal{D}_{tr}) = (1-\alpha)\hat{\pi}(a|x;\theta_t,\mathcal{D}_{tr}) + \alpha\pi_0(a|x),$$
(12)

where θ_t is a set of hyperparameters sampled at the t-th trial. $\alpha \ (\in [0,1])$ is a regularization parameter, which mixes the policy induced by θ_t and π_0 to construct a policy to evaluate. A large value of α makes π_t closer to the logging policy, possibly avoiding the unsafe behavior. However, if the logging policy is detrimental, we should use a small α so that we can avoid an unnecessary performance degradation. So, a natural question to ask here is: how should we set the regularization parameter α ? Again, this problem is non-trivial, as the optimality of the logging policy is unknown when performing HPO.

To overcome this difficulty in correctly setting α , we propose adaptively tuning this parameter over the course of HPO. Based on the previous discussion, we should apply a strong regularization if π_0 performs well, otherwise we should not imitate π_0 . A key idea here is that we can reason about the optimality of the logging policy by comparing it with solutions sampled during HPO, i.e., $\{\hat{\pi}(a|x;\theta_t,\mathcal{D}_{tr})\}_{t=1}^T$. If most of the sampled solutions underperform π_0 , we can infer that the logging policy is well-performing. To make a valid comparison between the sampled solutions and the logging policy, we apply a Student's t-test based on the following T-value.

$$T(\pi_1, \pi_2) := \frac{|\Delta \hat{V}_{\text{IPS}}(\pi_1, \pi_2)|}{\sqrt{\hat{\mathbb{V}}_n(\Delta \hat{V}_{\text{IPS}}(\pi_1, \pi_2))/(n-1)}}.$$

where $\Delta \hat{V}_{\rm IPS}\left(\pi_1,\pi_2\right):=\hat{V}_{\rm IPS}(\pi_1)-\hat{V}_{\rm IPS}(\pi_2)$ is the performance difference between the two policies estimated by IPS. Given a null hypothesis $(\Delta \hat{V}_{\rm IPS}\left(\pi_1,\pi_2\right)=0)$ and a normality assumption, $T(\pi_1,\pi_2)$ follows a t-distribution with ν degrees of freedom. We then calculate the optimality score of π_0 at

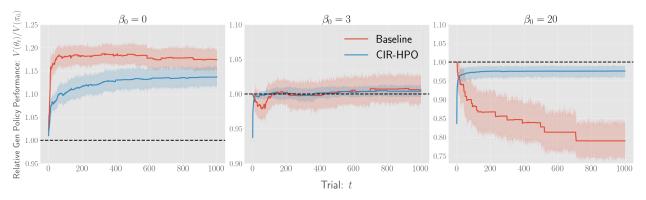


Figure 3: Comparing CIR-HPO (our proposal) and Baseline by their generalization performance. The results are averaged over 25 runs with different seeds and then normalized by $V(\pi_0)$. The shaded regions indicate 95% confidence intervals.

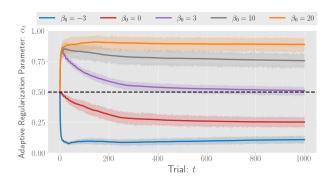


Figure 4: Behavior of adaptive regularization parameter (α_t) of CIR-HPO with varying values of $\beta_0 \in \{-3, 0, 3, 10, 20\}$.

the t-th trial as follows.

$$s_{t} = \begin{cases} 1 & (T(\pi_{0}, \pi_{t}) \geq t_{1-\delta/2, \nu} \text{ and } \Delta \hat{V}_{IPS}(\pi_{0}, \pi_{t}) \geq 0) \\ -1 & (T(\pi_{0}, \pi_{t}) \geq t_{1-\delta/2, \nu} \text{ and } \Delta \hat{V}_{IPS}(\pi_{0}, \pi_{t}) < 0) \\ 0 & (\text{otherwise, i.e., } T(\pi_{0}, \pi_{t}) < t_{1-\delta/2, \nu}) \end{cases}$$

$$(13)$$

 s_t indicates whether π_0 is better or worse than π_t in a significant level. If π_0 is better than π_t , then $s_t=1$. Instead, $s_t=-1$ if π_0 is tested to be worse. If there is no significant difference between π_0 and π_t , the score is zero.

Using the sequence of scores up to the t-th trial, i.e., $\{s_{t'}\}_{t'=1}^t$, we define *adaptive regularization parameter* as:

$$\alpha_t := \alpha_{init} + (1 - \alpha_{init}) \cdot \left(\frac{t}{T}\right)^{\gamma} \cdot \frac{\sum_{t'=1}^{t} s_{t'}}{t}$$
 (14)

where $\alpha_{init} \in [0,1]$ is an initial regularization parameter and $\gamma (>0)$ is a scheduling parameter for adaptive regularization. For example, suppose that $s_t = 1, \forall t = 1, 2, \ldots, T$, meaning that π_0 is always better than π_t in a significant level. Then, following Eq. (14), $\alpha_T = 1$ and the HPO procedure outputs π_0 , because it should be near-optimal. On the other hand, if $s_t = -1, \forall t = 1, 2, \ldots, T$, meaning that π_0 is always worse than π_t in a significant level, then $\alpha_T = 0$ and the HPO procedure does not imitate the logging policy at all, because it should be a bad policy.

4.3 The CIR-HPO Algorithm

Algorithm 2 describes the CIR-HPO algorithm, which leverages conservative surrogate objective (lines 6 and 9) and adaptive imitation regularization (line 5). δ and γ are meta hyperparameters. δ controls how conservative we would like to be during HPO, and γ controls the scheduling of the adaptive regularization. In Section 5, we show that these configurations have some impact on the behavior of CIR-HPO, but we also demonstrate that the default values ($\delta=0.1$ and $\gamma=0.01$) work reasonably well in a range of experiment settings. The other inputs are the same as those of Algorithm 1. Note that our algorithm is easy to implement with a few additional lines of code and there is no additional computational overhead compared to the typical procedure in Algorithm 1.

5 Empirical Evaluation

This section empirically compares **Baseline** (Algorithm 1) and **CIR-HPO** (Algorithm 2), employing the same synthetic data and policy class as in Section 3.1. Note that we compare CIR-HPO against only **Baseline** because there is no other method proposed for HPO using logged bandit data (comprehensive summary of related work can be found in Appendix B).

5.1 Baseline vs CIR-HPO

Figure 3 compares the performance of Baseline and CIR-**HPO** with varying logging policies ($\beta_0 \in \{0, 3, 20\}$). First, when $\beta_0 = 0$ where the logging policy is uniform random, both Baseline and CIR-HPO work reasonably well and succeed in finding a set of hyperparameters that leads to a policy much better than the logging policy. What is notable for this setting is that CIR-HPO is inefficient and slow to converge compared to Baseline due to adaptive imitation regularization, even though it reaches far above the black horizontal line $(V(\pi_0))$. At the initial stage of HPO, we do not know how close the logging policy is to the optimal policy. Therefore, the proposed procedure gradually learns the optimality of the logging policy, potentially leading to a slower convergence if the logging policy is far from optimal (such as uniform random). Next, when $\beta_0 = 3$ where the logging policy is better than uniform random, but is not close to the optimal, both Baseline and CIR-HPO slightly improve the logging policy.

However, the confidence intervals indicate that CIR-HPO is much more stable than Baseline. In particular, Baseline is much more likely to underperform the logging policy, even though it outperforms the logging policy on average. Finally, when $\beta_0 = 20$ where the logging policy is near-optimal, Baseline outputs a solution that is substantially worse than the logging policy, even though it starts from the logging policy as its initial solution. In contrast, CIR-HPO learns that the logging policy is near-optimal during HPO and strengthens the imitation regularization adaptively. As a result, it prevents the solution from being significantly worse than the (already near-optimal) logging policy, which is compelling, because we do not know the optimality of the logging policy in advance. Figure 4 illustrates the behavior of adaptive imitation regularization, which suggests that it succeeds in controlling the strength of regularization depending on the optimality of the logging policy.

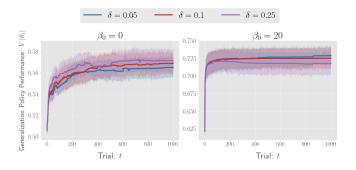


Figure 5: Sensitivity of the generalization performance of CIR-HPO regarding the choice of δ .

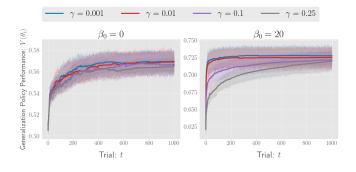


Figure 6: Sensitivity of the generalization performance of CIR-HPO regarding the choice of γ .

5.2 Choice of Meta Hyperparameters

Next, we evaluate the sensitivity of CIR-HPO to the choice of its meta hyperparameters. Figure 5 shows that the effectiveness of CIR-HPO with different values of δ . The result demonstrates that there is no significant difference among the three values, suggesting that we do not have to care too much about which value to use for δ . In addition, Figure 6 evaluates different values of γ , which controls the scheduling

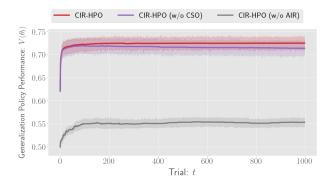


Figure 7: Ablation study of CIR-HPO ($\beta_0 = 20$).

of adaptive imitation regularization. This result implies that, for a sub-optimal logging policy ($\beta_0=0$), the choice of γ has no significant effect on the behavior of CIR-HPO. For a near-optimal logging policy ($\beta_0=20$), however, a smaller γ leads to a faster convergence, although all values achieve the same level of performance in the final stage.

5.3 Ablation Study

We also conduct an ablation study to evaluate the contribution of **conservative surrogate objective** (CSO) and **adaptive imitation regularization** (AIR) to the effectiveness of CIR-HPO. To this end, we compare CIR-HPO to CIR-HPO (w/o CSO) and CIR-HPO (w/o AIR) in Figure 7. The result demonstrates that both CSO and AIR clearly contribute to the performance of CIR-HPO, while AIR has a more appealing effect (CSO and AIR provide 1.6% and 23.6% improvements, respectively, in terms of the final generalization performance).

5.4 A Real-World Experiment

In addition to the synthetic experiments, we apply CIR-HPO to the Open Bandit Dataset [Saito *et al.*, 2021a], a publicly available logged bandit dataset collected on a large-scale fashion e-commerce platform. The results suggest that CIR-HPO leads to a better policy compared to the Baseline procedure in terms of the generalization performance, providing a further argument regarding its real-world applicability. The experiment detail and results can be found in Appendix A.

6 Conclusion

This work studies the effectiveness of the typical HPO procedure in the OPL setup from both empirical and theoretical perspectives and found that it can fail and even be harmful. In particular, we investigated two surprising issues, namely optimistic bias and unsafe behavior, and showed that a heavy-tailed distribution of overestimation can cause an unexpected gap between validation and generalization. In response, we made two extremely simple corrections to the typical HPO procedure, resulting in the CIR-HPO algorithm, to deal with the issues. Extensive experiments demonstrated that CIR-HPO can be advantageous, particularly when the conventional procedure collapses and causes a significant and undetectable deterioration in the generalization performance.

References

- [Akiba et al., 2019] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A next-generation hyperparameter optimization framework. In Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD '19, pages 2623–2631, New York, NY, USA, 2019. ACM.
- [Bergstra and Bengio, 2012] James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *Journal of machine learning research*, 13(2), 2012.
- [Bergstra et al., 2011] James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyperparameter optimization. Advances in neural information processing systems, 24, 2011.
- [Brochu *et al.*, 2010] Eric Brochu, Vlad M Cora, and Nando De Freitas. A tutorial on bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning. *arXiv* preprint *arXiv*:1012.2599, 2010.
- [Dacrema *et al.*, 2019] Maurizio Ferrari Dacrema, Paolo Cremonesi, and Dietmar Jannach. Are we really making much progress? a worrying analysis of recent neural recommendation approaches. In *Proceedings of the 13th ACM Conference on Recommender Systems*, pages 101–109, 2019.
- [Doroudi et al., 2018] Shayan Doroudi, Philip S Thomas, and Emma Brunskill. Importance sampling for fair policy selection. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pages 5239–5243, 2018.
- [Dudík et al., 2014] Miroslav Dudík, Dumitru Erhan, John Langford, and Lihong Li. Doubly Robust Policy Evaluation and Optimization. Statistical Science, 29:485–511, 2014.
- [Farajtabar *et al.*, 2018] Mehrdad Farajtabar, Yinlam Chow, and Mohammad Ghavamzadeh. More robust doubly robust off-policy evaluation. In *International Conference on Machine Learning*, volume 80, pages 1447–1456. PMLR, 2018.
- [Feurer and Hutter, 2019] Matthias Feurer and Frank Hutter. Hyperparameter optimization. In *Automated Machine Learning*, pages 3–33. 2019.
- [Frazier, 2018] Peter I Frazier. A tutorial on bayesian optimization. *arXiv preprint arXiv:1807.02811*, 2018.
- [Fu et al., 2020] Justin Fu, Mohammad Norouzi, Ofir Nachum, George Tucker, Alexander Novikov, Mengjiao Yang, Michael R Zhang, Yutian Chen, Aviral Kumar, Cosmin Paduraru, et al. Benchmarks for deep off-policy evaluation. In *International Conference on Learning Representations*, 2020.
- [Hansen, 2016] Nikolaus Hansen. The CMA Evolution Strategy: A Tutorial. *arXiv preprint arXiv:1604.00772*, 2016.
- [Henderson *et al.*, 2018] Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger. Deep Reinforcement Learning that Matters. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.

- [Jeunen and Goethals, 2021] Olivier Jeunen and Bart Goethals. Pessimistic reward models for off-policy learning in recommendation. In *Fifteenth ACM Conference on Recommender Systems*, pages 63–74, 2021.
- [Joachims *et al.*, 2018] Thorsten Joachims, Adith Swaminathan, and Maarten de Rijke. Deep learning with logged bandit feedback. In *International Conference on Learning Representations*, 2018.
- [Kallus *et al.*, 2021] Nathan Kallus, Yuta Saito, and Masatoshi Uehara. Optimal off-policy evaluation from multiple logging policies. In *International Conference on Machine Learning*, pages 5247–5256. PMLR, 2021.
- [Kiyohara *et al.*, 2021] Haruka Kiyohara, Kosuke Kawakami, and Yuta Saito. Accelerating offline reinforcement learning application in real-time bidding and recommendation: Potential use of simulation. *arXiv preprint arXiv:2109.08331*, 2021.
- [Kiyohara et al., 2022] Haruka Kiyohara, Yuta Saito, Tatsuya Matsuhiro, Yusuke Narita, Nobuyuki Shimizu, and Yasuo Yamamoto. Doubly robust off-policy evaluation for ranking policies under the cascade behavior model. In *Proceedings of the 15th International Conference on Web Search and Data Mining*, pages 487—497, 2022.
- [Kiyohara et al., 2023] Haruka Kiyohara, Ren Kishimoto, Kosuke Kawakami, Ken Kobayashi, Kazuhide Nakata, and Yuta Saito. Towards assessing and benchmarking riskreturn tradeoff of off-policy evaluation. In *The Twelfth In*ternational Conference on Learning Representations, 2023.
- [Kiyohara *et al.*, 2024] Haruka Kiyohara, Masahiro Nomura, and Yuta Saito. Off-policy evaluation of slate bandit policies via optimizing abstraction. *arXiv preprint arXiv:2402.02171*, 2024.
- [Kumar et al., 2019] Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. Advances in Neural Information Processing Systems, 32:11784–11794, 2019.
- [Kuzborskij et al., 2021] Ilja Kuzborskij, Claire Vernade, Andras Gyorgy, and Csaba Szepesvári. Confident off-policy evaluation and selection through self-normalized importance weighting. In *International Conference on Artificial Intelligence and Statistics*, pages 640–648. PMLR, 2021.
- [Loshchilov and Hutter, 2016] Ilya Loshchilov and Frank Hutter. CMA-ES for Hyperparameter Optimization of Deep Neural Networks. In *ICLR Workshop*, 2016.
- [Lucic et al., 2018] Mario Lucic, Karol Kurach, Marcin Michalski, Sylvain Gelly, and Olivier Bousquet. Are Gans Created Equal? A Large-Scale Study. In Advances in neural information processing systems, pages 700–709, 2018.
- [Ma et al., 2019] Yifei Ma, Yu-Xiang Wang, and Balakrishnan Narayanaswamy. Imitation-regularized offline learning. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2956–2965. PMLR, 2019.
- [McInerney et al., 2020] James McInerney, Brian Brost, Praveen Chandar, Rishabh Mehrotra, and Benjamin

- Carterette. Counterfactual evaluation of slate recommendations with sequential reward interactions. In *Proceedings* of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pages 1779–1788, 2020.
- [Nomura and Saito, 2021] Masahiro Nomura and Yuta Saito. Efficient hyperparameter optimization under multi-source covariate shift. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, pages 1376–1385, 2021.
- [Nomura and Shibata, 2024] Masahiro Nomura and Masashi Shibata. cmaes: A Simple yet Practical Python Library for CMA-ES. *arXiv preprint arXiv:2402.01373*, 2024.
- [Nomura et al., 2021] Masahiro Nomura, Shuhei Watanabe, Youhei Akimoto, Yoshihiko Ozaki, and Masaki Onishi. Warm Starting CMA-ES for Hyperparameter Optimization. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, pages 9188–9196, 2021.
- [Paine et al., 2020] Tom Le Paine, Cosmin Paduraru, Andrea Michi, Caglar Gulcehre, Konrad Zolna, Alexander Novikov, Ziyu Wang, and Nando de Freitas. Hyperparameter selection for offline reinforcement learning. arXiv preprint arXiv:2007.09055, 2020.
- [Precup *et al.*, 2000] Doina Precup, Richard S. Sutton, and Satinder Singh. Eligibility Traces for Off-Policy Policy Evaluation. In *Proceedings of the 17th International Conference on Machine Learning*, pages 759–766, 2000.
- [Saito and Joachims, 2021] Yuta Saito and Thorsten Joachims. Counterfactual learning and evaluation for recommender systems: Foundations, implementations, and recent advances. In *Fifteenth ACM Conference on Recommender Systems*, page 828–830, 2021.
- [Saito and Joachims, 2022a] Yuta Saito and Thorsten Joachims. Counterfactual evaluation and learning for interactive systems: Foundations, implementations, and recent advances. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 4824–4825, 2022.
- [Saito and Joachims, 2022b] Yuta Saito and Thorsten Joachims. Off-policy evaluation for large action spaces via embeddings. In *International Conference on Machine Learning*, pages 19089–19122. PMLR, 2022.
- [Saito and Nomura, 2022] Yuta Saito and Masahiro Nomura. Towards resolving propensity contradiction in offline recommender learning. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 2211–2217. International Joint Conferences on Artificial Intelligence Organization, 7 2022.
- [Saito et al., 2020] Yuta Saito, Suguru Yaginuma, Yuta Nishino, Hayato Sakata, and Kazuhide Nakata. Unbiased recommender learning from missing-not-at-random implicit feedback. In *Proceedings of the 13th International Conference on Web Search and Data Mining*, pages 501–509, 2020.

- [Saito et al., 2021a] Yuta Saito, Shunsuke Aihara, Megumi Matsutani, and Yusuke Narita. Open bandit dataset and pipeline: Towards realistic and reproducible off-policy evaluation. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021.
- [Saito et al., 2021b] Yuta Saito, Takuma Udagawa, Haruka Kiyohara, Kazuki Mogi, Yusuke Narita, and Kei Tateno. Evaluating the robustness of off-policy evaluation. In *Fifteenth ACM Conference on Recommender Systems*, pages 114–123, 2021.
- [Saito *et al.*, 2023] Yuta Saito, Ren Qingyang, and Thorsten Joachims. Off-policy evaluation for large action spaces via conjunct effect modeling. In *International Conference on Machine Learning*, pages 29734–29759. PMLR, 2023.
- [Saito et al., 2024] Yuta Saito, Jihan Yao, and Thorsten Joachims. Potec: Off-policy learning for large action spaces via two-stage policy decomposition. arXiv preprint arXiv:2402.06151, 2024.
- [Saito, 2020] Yuta Saito. Unbiased pairwise learning from biased implicit feedback. In *Proceedings of the 2020 ACM SIGIR on International Conference on Theory of Information Retrieval*, pages 5–12, 2020.
- [Shahriari *et al.*, 2015] Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P Adams, and Nando De Freitas. Taking the human out of the loop: A review of bayesian optimization. *Proceedings of the IEEE*, 104(1):148–175, 2015.
- [Srinivas et al., 2010] Niranjan Srinivas, Andreas Krause, Sham Kakade, and Matthias Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. In *Proceedings of the 27th International Conference on International Conference on Machine Learning*, pages 1015–1022, 2010.
- [Strehl et al., 2010] Alex Strehl, John Langford, Lihong Li, and Sham M Kakade. Learning from Logged Implicit Exploration Data. In *Advances in Neural Information Processing Systems*, volume 23, pages 2217–2225, 2010.
- [Su et al., 2020] Yi Su, Maria Dimakopoulou, Akshay Krishnamurthy, and Miroslav Dudík. Doubly robust off-policy evaluation with shrinkage. In *International Conference on Machine Learning*, volume 119, pages 9167–9176. PMLR, 2020.
- [Swaminathan and Joachims, 2015a] 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.
- [Swaminathan and Joachims, 2015b] Adith Swaminathan and Thorsten Joachims. The self-normalized estimator for counterfactual learning. In *Advances in Neural Information Processing Systems*, volume 28, pages 3231–3239, 2015.
- [Tang and Wiens, 2021] Shengpu Tang and Jenna Wiens. Model selection for offline reinforcement learning: Practical considerations for healthcare settings. *arXiv preprint arXiv:2107.11003*, 2021.

- [Thomas et al., 2015a] Philip Thomas, Georgios Theocharous, and Mohammad Ghavamzadeh. High-confidence off-policy evaluation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 29, 2015.
- [Thomas *et al.*, 2015b] Philip Thomas, Georgios Theocharous, and Mohammad Ghavamzadeh. High confidence policy improvement. In *Proceedings of the 32th International Conference on Machine Learning*, volume 37, pages 2380–2388, 2015.
- [Turner et al., 2021] Ryan Turner, David Eriksson, Michael McCourt, Juha Kiili, Eero Laaksonen, Zhen Xu, and Isabelle Guyon. Bayesian optimization is superior to random search for machine learning hyperparameter tuning: Analysis of the black-box optimization challenge 2020. arXiv preprint arXiv:2104.10201, 2021.
- [Wang et al., 2017] Yu-Xiang Wang, Alekh Agarwal, and Miroslav Dudik. Optimal and Adaptive Off-policy Evaluation in Contextual Bandits. In Proceedings of the 34th International Conference on Machine Learning, pages 3589– 3597, 2017.
- [Wu et al., 2019] Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning. arXiv preprint arXiv:1911.11361, 2019.
- [Yang et al., 2020] Mengjiao Yang, Bo Dai, Ofir Nachum, George Tucker, and Dale Schuurmans. Offline policy selection under uncertainty. arXiv preprint arXiv:2012.06919, 2020.

Table 1: Test policy values of the policies tuned by Baseline and CIR-HPO in the real-world experiment.

	IPS	DR
	5.199×10^{-3}	
CIR-HPO (ours)	5.214×10^{-3}	5.303×10^{-3}

A A Real-World Experiment

A.1 Dataset.

To assess the real-world applicability of our CIR-HPO, here we evaluate it on the Open Bandit Dataset (OBD)⁵ [Saito *et al.*, 2021a], a publicly available logged bandit dataset collected on a large-scale fashion e-commerce platform. We use 100,000 observations that are randomly sub-sampled from the "Men's" campaign data of OBD. The dataset contains user contexts x, fashion items to recommend as action $a \in \mathcal{A}$ where $|\mathcal{A}| = 240$, and resulting clicks as reward $r \in \{0, 1\}$.

The dataset consists of two sets of logged bandit data collected by two different policies (uniform random and Thompson sampling) during an A/B test of these policies. We regard Thompson sampling as a logging policy and perform HPO of a policy class defined in Section 3.1. We then approximate the ground-truth performance of the tuned policies on the test dataset collected by uniform random. Note that we use the same policy class defined in Section 3.1 and the default meta-parameters of CIR-HPO described in Section 5. All the experiments were conducted on MacBook Pro (2.4 GHz Intel Core i9, 64 GB).

A.2 Result.

Table 1 reports the results (averaged over 5 runs with different random seeds) of the real-world experiment. We compare **Baseline** and **CIR-HPO** (**ours**) combined with IPS and DR as OPE estimators to provide a surrogate objective (i.e., $\hat{V}(\theta; \mathcal{D})$). The results suggest that, for both estimators, CIR-HPO outperforms Baseline in terms of the test policy value. This observation provides further arguments for the applicability of our CIR-HPO.

B Related Work

Off-Policy Evaluation and Learning. The basis of our study lies in OPE, which is interested in accurately estimating the generalization policy performance from logged bandit data. This has been one of the most fundamental problems in contextual bandits and RL, with applications ranging from recommender systems [Saito and Joachims, 2021; McInerney et al., 2020; Kiyohara et al., 2022; Saito and Joachims, 2022b; Saito et al., 2020; Saito, 2020; Saito and Nomura, 2022; Saito et al., 2021b] to personalized medicine [Tang and Wiens, 2021; Kallus et al., 2021; Saito and Joachims, 2022a]. The most common solution in OPE is to use IPS weighting. IPS provides an unbiased estimate of the policy performance. However, there is a canonical criticism that IPS often suffers from a high variance due to a low overlap [Dudík et al., 2014; Wang et al., 2017]. Thus, alternative estimators have been

explored to reduce the variance without introducing large bias. For example, Self-Normalized IPS (SNIPS) [Swaminathan and Joachims, 2015b] aims to reduce the variance of IPS as follows

$$\hat{V}_{\text{SNIPS}}(\pi_e; \mathcal{D}) := \frac{1}{\sum_{i=1}^n \frac{\pi_e(a_i|x_i)}{\pi_0(a_i|x_i)}} \sum_{i=1}^n \frac{\pi_e(a_i|x_i)}{\pi_0(a_i|x_i)} r_i.$$

This estimator normalizes the IPS estimate by the sum of the importance weights $(\sum_{i=1}^n \frac{\pi_e(a_i|x_i)}{\pi_0(a_i|x_i)})$ to gain stability. Moving forward, DR leverages a control variate to provide an *efficient* OPE. The DR estimator is defined as follows.

$$\hat{V}_{DR}(\pi_e; \mathcal{D}, \hat{\mu}) := \frac{1}{n} \sum_{i=1}^n \hat{\mu}(x_i, \pi_e) + \frac{\pi_e(a_i|x_i)}{\pi_0(a_i|x_i)} (r_i - \hat{\mu}(x_i, a_i)),$$

where $\hat{\mu}(x,\pi_e)$ estimates $\mathbb{E}_{\pi_e}[\mu(x,a)]$. This estimator is still unbiased and consistent if either the importance weight or the reward estimator is true or consistent. In addition, DR is *efficient* in that it reaches the lowest achievable asymptotic variance if the reward estimator is correctly specified. There have also been much efforts to further improve DR in a finite sample setting such as Switch [Wang *et al.*, 2017], More Robust Doubly Robust [Farajtabar *et al.*, 2018], and Shrinkage [Su *et al.*, 2020].

Instead, OPL is the task of improving the decision making policies using only logged bandit data collected from a logging policy [Swaminathan and Joachims, 2015a]. The optimal policy maximizes the generalization performance, i.e., $\pi^* := \arg\max_{\pi} V(\pi)$. However, this problem is in-

tractable because we cannot know the generalization performance. This raises the need for applying an estimator for its careful approximation, as done in the empirical risk minimization of supervised ML. A typical estimator choice for OPL is IPS [Swaminathan and Joachims, 2015a; Ma et al., 2019; Joachims et al., 2018] or its variants [Swaminathan and Joachims, 2015b]. A problem is that the variance issue arises here again. Thus, research has been centered around adding regularization to deal with the variance issue during policy training. The fundamental method is variance regularization, which penalizes the policy whose variance in the performance estimation is high [Swaminathan and Joachims, 2015a]. Other regularization methods include imitation regularization [Ma et al., 2019] and behavior regularization [Wu et al., 2019; Kumar et al., 2019]. [Jeunen and Goethals, 2021] explore the optimistic bias in OPL, and propose a pessimistic reward modeling for OPL based on a Bayesian uncertainty estimation. Instead, we focus on investigating and alleviating the optimistic bias in HPO and empirically illustrate the unsafe behavior, which is specific to our HPO setup.

Off-Policy Selection. Off-Policy Selection (OPS) is a subfield of OPE and OPL and is closely related to our HPO setting. This is the task of identifying the best policy out of a given *finite* set of candidate policies using only logged bandit data. We can view this selection problem as a special case of OPL, where the policy class Π is finite. [Kuzborskij *et al.*, 2021] study OPS in the contextual bandit setting. They develop a confident OPS procedure, which is based on an Efron-Stein high

⁵https://research.zozo.com/data.html

probability lower bound of the policy performance derived from SNIPS. [Yang et al., 2020] study OPS in RL and propose BayesDICE for estimating the brief over the performance of the candidate policies, which is useful for the selection task. [Doroudi et al., 2018] theoretically characterize a failure of IPS in OPS. Specifically, [Doroudi et al., 2018] show that naively applying IPS to OPS can result in an unfair selection in the sense that the procedure can select the worst of the two candidate policies more often than not. [Paine et al., 2020] and [Fu et al., 2020] conduct empirical studies on OPS of RL polices for several benchmark control tasks. They identify Fitted Q Evaluation as a useful strategy for OPS in RL. [Tang and Wiens, 2021] also provide an empirical study on OPS of RL policies and propose to combine multiple OPE estimators for an accurate and scalable OPS.

Although these studies on OPS are closely related, our contributions are unique in several ways. First, we focus on HPO, not OPS, which adaptively finds better hyperparameter solutions given a certain budget. By paying attention to HPO, our empirical analysis succeeded in finding that the TPE algorithm, a popular adaptive method in HPO, cannot improve the generalization performance of OPL. This is our unique finding, not captured by previous studies targeting only OPS. Second, we provide a theoretical analysis about the optimistic bias and the gap in generalization and validation regret, explaining the empirical observations. Although [Paine et al., 2020] point out the overestimation bias in the context of OPS, they provide no theoretical explanation. Finally, we propose CIR-HPO based on our empirical observations and analysis. This procedure is specific to the adaptive optimization process and is non-trivial given any existing studies on OPS.

Hyperparameter Optimization (HPO). HPO is a critical element for the success of a range of machine learning algorithms and tasks [Feurer and Hutter, 2019]. For instance, hyperparameter configurations can entirely change the performance of deep neural networks [Dacrema et al., 2019; Henderson et al., 2018; Lucic et al., 2018]. A typical formulation regards HPO as a black-box optimization problem, where the input is a set of hyperparameters, and the output is a validation performance (an accessible proxy of the generalization performance). Among many black-box optimization methods, Bayesian optimization (BO) [Brochu et al., 2010; Shahriari et al., 2015; Frazier, 2018], such as Gaussian process bandit algorithms [Srinivas et al., 2010] and tree-structured Parzen estimator (TPE) [Bergstra et al., 2011] have gained particular attention. These methods sequentially optimize the hyperparameters of a prediction model by leveraging the previous evaluation results to sample the next set of hyperparameters to evaluate. More specifically, previous evaluation results are used to train a surrogate to model the relationship between hyperparameters and the resulting prediction accuracy. Then, the algorithms balance the exploration and exploitation based on an acquisition function, such as expected improvement and upper confidence bound. Because of the sample efficiency, BO demonstrates a state-of-the-art performance with a limited budget [Turner et al., 2021]. It should be noted that, while this study focuses on BO, our discussion can be applied to other optimization methods such as CMA-ES [Hansen, 2016;

Nomura and Shibata, 2024], whose efficiency is verified in multiple HPO tasks [Loshchilov and Hutter, 2016; Nomura *et al.*, 2021].

A critical convention in HPO research is to evaluate the performance and efficiency of algorithms based solely on the validation performance. This implies that there is an implicit and often neglected assumption that optimizing the validation performance is a reasonable strategy for optimizing the generalization performance (primary objective). However, it is unclear whether optimizing the validation performance really improves the generalization performance. In fact, Section 3 sheds light on the fact that ignoring this assumption in OPL can lead to an unexpected failure and a substantial validation-generalization gap. Our theoretical and empirical illustrations might also contribute to a broader HPO community, as there are few studies verifying whether naively setting the validation performance as a surrogate objective is reasonable, given the goal of optimizing the generalization performance.

C Omitted Proofs

This section provides proofs omitted in the main text.

C.1 Proof of Proposition 3.1

Proof. Given that θ^* and $\hat{\theta}$ are defined in Eq. (2) and Eq. (3), we have that

$$\mathbb{E}\left[\hat{V}(\hat{\pi}(\cdot \mid \cdot, \hat{\theta}, \mathcal{D}_{tr}); \mathcal{D}_{val})\right]$$

$$\geq \mathbb{E}\left[\hat{V}(\hat{\pi}(\cdot \mid \cdot, \theta^*, \mathcal{D}_{tr}); \mathcal{D}_{val})\right] = V(\hat{\pi}(\cdot \mid \cdot, \theta^*, \mathcal{D}_{tr})),$$

where the last equation follows, as θ^* does not depend on \mathcal{D}_{val} . Similarly, the right inequality of Eq. (7) comes from the fact that θ^* is optimal in terms of the true generalization policy performance.

C.2 Proof of Proposition 3.2

Proof. Our derivation is inspired by the regret analysis provided in [Nomura and Saito, 2021]. Given the notations introduced in Section 3.2, it follows that

$$\begin{split} r_{gen}(T; A, \mathcal{D}) &= V\left(\theta^*\right) - V(\hat{\theta}_{T,A}(\mathcal{D})) \\ &= \underbrace{V\left(\theta^*\right) - \hat{V}_{\mathrm{IPS}}(\theta^*; \mathcal{D})}_{=-\tau(\theta^*; \mathcal{D})} + \hat{V}_{\mathrm{IPS}}(\theta^*; \mathcal{D}) - V(\hat{\theta}_{T,A}(\mathcal{D})) \\ &= -\tau(\theta^*; \mathcal{D}) \\ &= -\tau(\theta^*; \mathcal{D}) + \underbrace{\left(-V(\hat{\theta}_{T,A}(\mathcal{D})) + \hat{V}_{\mathrm{IPS}}(\hat{\theta}_{T,A}(\mathcal{D}); \mathcal{D})\right)}_{=\tau(\hat{\theta}_{T,A}(\mathcal{D}); \mathcal{D})} \\ &- \hat{V}_{\mathrm{IPS}}(\hat{\theta}_{T,A}(\mathcal{D}); \mathcal{D}) + \hat{V}_{\mathrm{IPS}}(\theta^*; \mathcal{D}) \\ &= \Delta \tau(\hat{\theta}_{T,A}(\mathcal{D}), \theta^*) + \underbrace{\left(\hat{V}_{\mathrm{IPS}}(\hat{\theta}^*; \mathcal{D}) - \hat{V}_{\mathrm{IPS}}(\hat{\theta}_{T,A}(\mathcal{D}); \mathcal{D})\right)}_{=r_{val}(T; A, \mathcal{D})} \\ &+ \underbrace{\left(\hat{V}_{\mathrm{IPS}}(\theta^*; \mathcal{D}) - \hat{V}_{\mathrm{IPS}}(\hat{\theta}^*; \mathcal{D})\right)}_{=C} \\ &= r_{val}(T; A, \mathcal{D}) + \Delta \tau(\hat{\theta}_{T,A}(\mathcal{D}), \theta^*; \mathcal{D}) + C. \end{split}$$

D Additional Theoretical Result

We suppose $\hat{V}(\theta)$ has the following form:

$$\hat{V}(\theta) = \frac{1}{n} \sum_{i=1}^{n} v(a_i, x_i; \theta).$$

Note that this form is general and encompasses common estimators. For example, we can obtain IPS estimator by setting $v(a_i, x_i; \theta) = \frac{1}{n} \sum_{i=1}^n \frac{\pi_e(a_i|x_i; \theta)}{\pi_0(a_i|x_i)} r_i$. The following inequality suggests that the optimistic bias $\hat{V}(\hat{\theta}(\mathcal{D}); \mathcal{D}) - V(\theta^*)$ decreases at the order $\mathcal{O}(1/\sqrt{n})$ as the data increases.

Proposition D.1. Suppose $|\Theta| < \infty$ and $0 \le v(a, x, \theta) \le 1$ for all $a \in \mathcal{A}, x \in \mathcal{X}, \theta \in \Theta$. For $\delta \in (0, 1)$, the following inequality holds with probability as least $1 - \delta$:

$$\hat{V}\left(\hat{\theta}(\mathcal{D}); \mathcal{D}\right) - V\left(\theta^*\right) \leq \sqrt{\frac{1}{2n}\log\frac{|\Theta|}{\delta}} \in \mathcal{O}\left(\frac{1}{\sqrt{n}}\right).$$

Proof. We first decompose the optimistic bias as

$$\begin{split} \hat{V}\left(\hat{\theta}(\mathcal{D});\mathcal{D}\right) - V\left(\theta^*\right) \\ &= \hat{V}\left(\hat{\theta}(\mathcal{D});\mathcal{D}\right) - V(\hat{\theta}(\mathcal{D})) + \underbrace{V(\hat{\theta}(\mathcal{D})) - V\left(\theta^*\right)}_{\leq 0} \\ &\leq \hat{V}\left(\hat{\theta}(\mathcal{D});\mathcal{D}\right) - V(\hat{\theta}(\mathcal{D})). \end{split}$$
 Hence, for $\epsilon > 0$,
$$\mathbb{P}\left(\hat{V}\left(\hat{\theta}(\mathcal{D});\mathcal{D}\right) - V\left(\theta^*\right) \geq \epsilon\right) \\ &\leq \mathbb{P}\left(\hat{V}\left(\hat{\theta}(\mathcal{D});\mathcal{D}\right) - V\left(\theta^*\right) \geq \epsilon\right) \\ &\leq \mathbb{P}\left(\hat{V}\left(\hat{\theta}(\mathcal{D});\mathcal{D}\right) - V(\hat{\theta}(\mathcal{D})) \geq \epsilon\right) \\ &\leq \mathbb{P}\left(\hat{V}\left(\hat{\theta}(\mathcal{D});\mathcal{D}\right) - V(\theta(\mathcal{D})) \geq \epsilon\right) \\ &\leq \sum_{Q \in \mathcal{Q}} \mathbb{P}\left(\hat{V}\left(\theta(\mathcal{D});\mathcal{D}\right) - V(\theta(\mathcal{D})) \geq \epsilon\right) \end{split}$$

We used the union bound and Hoeffding's inequality⁶. Putting the RHS as δ and solving it for ϵ completes the proof.

E Supplemental Simulations

 $\leq |\Theta|e^{-2n\epsilon^2}$.

This section empirically evaluates the confidence lower bound of OPE based on concentration inequalities (Hoeffding and Bernstein) and a Student's t-test. We follow Section 3 to generate synthetic bandit data. We vary the value of β_0 within the range of $\{0,3,20\}$, and the number of validation data $|\mathcal{D}_{val}|$ within the range of $\{400,800,1600,3200,6400,12800\}$. We also follow Section 3 to train the evaluation policy π_e . Specifically, we first train $\hat{\mu}$ using logistic regression and form a stochastic policy based on Eq. (6) with $\beta=10$.

We use IPS and estimate a high probability lower bound of $V(\pi_e)$ based on Hoeffding, Bernstein, and t-Test. Given a

confidence level $\delta \in (0,1)$, the estimated lower bounds are given as $\hat{V}_{IPS}(\pi_e; \mathcal{D}_{val}) - f(\delta, \mathcal{D}_{val})$ where

$$\begin{aligned} \textbf{Hoeffding}: & f(\delta, \mathcal{D}_{val}) = w_{\max} \sqrt{\frac{2\log(2/\delta)}{n}}, \\ \textbf{Bernstein}: & f(\delta, \mathcal{D}_{val}) = \sqrt{\frac{2\log(2/\delta)\hat{\mathbb{V}}(\hat{V}_{\text{IPS}}(\pi_e; \mathcal{D}_{val}))}{n-1}} \\ & + \frac{7w_{\max}\log(2/\delta)}{3(n-1)}, \\ \textbf{t-Test}: & f(\delta, \mathcal{D}_{val}) = t_{1-\delta,\nu} \sqrt{\frac{\hat{\mathbb{V}}(\hat{V}_{\text{IPS}}(\pi_e; \mathcal{D}_{val}))}{n-1}}. \end{aligned}$$

Note that $n=|\mathcal{D}_{val}|$ and $w_{\max}:=\sup_{(x,a)\in\mathcal{X}\times\mathcal{A}}\pi_e(a|x)/\pi_b(a|x)$. $t_{1-\delta,\nu}$ is the $100(1-\delta)$ percentile of the Student's t distribution with ν degrees of freedom. The lower bound given by the t-test is based on the assumption that $(\pi_e(a|x)/\pi_0(a|x))r$ is normally distributed.

Figure 10 shows the estimated lower bounds with varying β and sample size, and with a fixed $\delta=0.05$. The black horizontal line represents the ground-truth policy value $V(\pi_e)$. We observe that the lower bound given by t-Test is the tightest, while those by Hoeffding and Bernstein are invisible when $\beta_0=20$, as they are too loose. Bernstein is always tighter than Hoeffding, but t-Test is even better, particularly when $\beta=3,20$ where the logging policy is near-deterministic ($w_{\rm max}$ is large).

Next, Figure 11 shows how frequently the estimated lower bounds fail to lower bound $V(\pi_e)$. Here, we say that a lower bound fails, if $\hat{V}_{\text{IPS}}(\pi_e; \mathcal{D}_{val}) - f(\delta, \mathcal{D}_{val}) \geq V(\pi_e)$. The black horizontal line represents the allowed error rate δ . We observe that, in all scenarios, the bounds given by Hoeffding and Bernstein have an error rate of 0, even if they are allowed to produce an error rate of δ , meaning that these lower bounds are overly conservative. In contrast, the lower bound given by t-Test makes some errors, but the error rate is around δ . Although the normality assumption might fail in OPE with small sample sizes, we empirically verify that t-Test produces a lower bound tighter than those of Hoeffding and Bernstein, and its error rate is around the allowed value.

 $^{^6}$ By replacing Hoeffding's inequality with Chebyshev's inequality, we can obtain a weaker result even if the boundedness of $v(a,x,\theta)$ is not assumed.

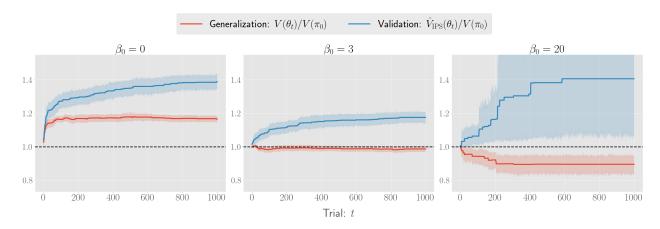


Figure 8: Empirical Evidence of Optimistic Bias and Unsafe Behavior in HPO for OPL (w/ Random Search and the IPS estimator). The results are averaged over 25 runs with different seeds and then normalized by $V(\pi_0)$. The shaded regions indicate 95% confidence intervals.

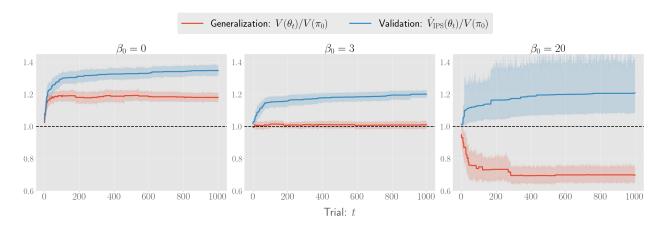


Figure 9: Empirical Evidence of Optimistic Bias and Unsafe Behavior in HPO for OPL (w/ TPE and the DR estimator). The results are averaged over 25 runs with different seeds and then normalized by $V(\pi_0)$. The shaded regions indicate 95% confidence intervals.

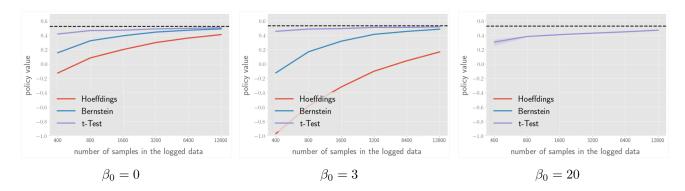


Figure 10: High Probability Lower Bound

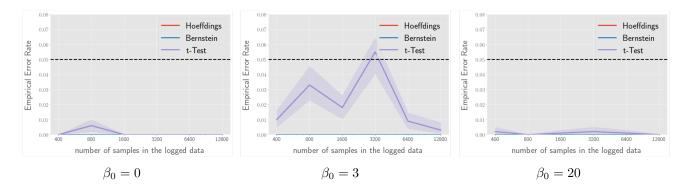


Figure 11: Empirical Error Rate

Table 2: Hyperparameter search space (Θ)

Hyperparameters	Search Spaces
β	[0.01, 100]
$\hat{\mu}$	{'LR', 'RF'}

Machine Learning Models	Search Spaces
Logistic Regression (LR)	$C \in [10^{-3}, 10^3]$
	11 _ratio $\in \{0.1, 0.2, \dots, 0.9\}$
	$max_depth \in \{2, 3, \dots, 32\}$
Random Forest (RF)	$min_samples_split \in \{2, 3, \dots, 32\}$
	$max_samples \in \{0.1, 0.2, \dots, 0.9\}$

Note: The names of the hyperparameters correspond to those specified by the *scikit-learn* package. For other hyperparameters, we use 'sklearn.ensemble.RandomForectClassifier(n_estimators=10)' and 'sklearn.linear_model.LogisticRegression(max_iter=1000, penalty="elasticnet", solver="saga")'.

Table 3: Generalization performance and optimality of π_0 with varying β_0

	$\beta_0 = -3$	$\beta_0 = 0$	$\beta_0 = 3$	$\beta_0 = 10$	$\beta_0 = 20$
$V(\pi_0)$	0.412	0.501	0.580	0.677	0.719
$V(\pi_0)/V(\pi^*)$	0.554	0.673	0.831	0.910	0.966

Note: $V(\pi^*)$ is the best achievable performance in our data generating process. $V(\pi_0)/V(\pi^*)$ indicates the optimality of logging policy π_0 .