

# Counterfactual Learning with General Data-Generating Policies

Yusuke Narita<sup>1</sup>, Kyohei Okumura<sup>2</sup>, Akihiro Shimizu<sup>3</sup>, Kohei Yata<sup>4</sup>

<sup>1</sup> Yale University,

<sup>2</sup> Northwestern University,

<sup>3</sup> Mercari, Inc.,

<sup>4</sup> University of Wisconsin-Madison,

yusuke.narita@yale.edu, kyohei.okumura@u.northwestern.edu, akihiro-shimizu@mercari.com, yata@wisc.edu

## Abstract

Off-policy evaluation (OPE) attempts to predict the performance of counterfactual policies using log data from a different policy. We extend its applicability by developing an OPE method for a class of both full support and deficient support logging policies in contextual-bandit settings. This class includes deterministic bandit (such as Upper Confidence Bound) as well as deterministic decision-making based on supervised and unsupervised learning. We prove that our method’s prediction converges in probability to the true performance of a counterfactual policy as the sample size increases. We validate our method with experiments on partly and entirely deterministic logging policies. Finally, we apply it to evaluate coupon targeting policies by a major online platform and show how to improve the existing policy.

## 1 Introduction

In bandit and reinforcement learning, off-policy (batch) evaluation attempts to estimate the performance of some counterfactual policy given data from a different logging policy. Off-policy evaluation (OPE) is essential when deploying a new policy might be costly or risky, such as in education, medicine, consumer marketing, and robotics. OPE relates to other fields that study counterfactual/causal reasoning, such as statistics and economics.

Most existing OPE studies focus on *full support* logging policies, which take all actions with positive probability in any context, such as stochastic bandit (e.g.  $\epsilon$ -greedy and Thompson Sampling) and random A/B testing. However, real-world decision-making often uses *deficient support* logging policies, including deterministic bandit (e.g. Upper Confidence Bound) as well as deterministic decision-making based on predictions obtained from supervised and unsupervised learning. An example in the latter group is a policy that greedily chooses the action with the largest predicted reward. OPE is difficult with a deficient support logging policy, since its log data contain no information about the reward from actions never chosen by the logging policy. There appears to be no established OPE estimator for deficient support logging policies (Sachdeva, Su, and Joachims 2020).

We provide a solution to this problem. Our proposed OPE estimator is applicable not only to full support logging policies but also to deficient support ones. We also allow for hybrid stochastic and deterministic logging policies, i.e., logging policies that choose actions stochastically for some individuals and deterministically for other individuals.<sup>1</sup>

**Method.** Our OPE estimator is based on a modification of the Propensity Score (Rosenbaum and Rubin 1983), which we dub the “Approximate Propensity Score” (APS) (Narita and Yata 2022). APS of action (arm)  $a$  at context (covariate) value  $x$  is the average probability that the logging policy chooses action  $a$  over a shrinking neighborhood around  $x$  in the context space. If two actions have nonzero APS at  $x$ , the logging policy chooses both actions locally around  $x$ . This enables us to estimate the difference in the mean reward between the two actions by exploiting the local subsample around  $x$ . When the logging policy is deterministic, the subsample consists of individuals near the decision boundary between the two actions. We then use the estimated reward differences to construct an estimator for the performance of any given counterfactual policy.

As the main theoretical result, we prove that our proposed OPE estimator is consistent. That is, the estimator converges in probability to the true performance of a counterfactual policy as the sample size increases, under the assumption that the mean reward differences are constant over the context space (Theorem 1). This result holds whether the logging policy is of full support or deficient support. The proof exploits results from differential geometry and geometric measure theory, which have not been applied in machine learning research as far as we know.

**Simulation Experiments.** We validate our method with two simulation experiments. The first considers a mix of full support and deficient support policies as the logging policy. Actions are randomly chosen for a small A/B test segment of the population and are chosen by a deterministic supervised learning algorithm for the rest of the population. For the task of evaluating counterfactual policies, our method produces smaller mean squared errors than a baseline estimator that only uses the A/B test subsample. The second experiment considers a situation in which we have a batch of data gen-

<sup>1</sup>The full version of the paper, which includes technical appendices, can be found at <https://arxiv.org/abs/2212.01925>.

erated by a deterministic bandit algorithm. We find that our estimator outperforms a regression-based estimator in terms of mean squared errors.

**Real-World Application.** We empirically apply our method to evaluate and optimize coupon targeting policies. Our application is based on proprietary data provided by Mercari Inc., a major e-commerce company running online C2C marketplaces in Japan and the US. This company uses a deterministic policy based on uplift modeling to decide whether they offer a promotional coupon to each target customer. We use the data produced by their policy and our method to evaluate a counterfactual policy that offers the coupon to more customers. Our method predicts that the counterfactual policy would increase revenue more than the cost of coupon offers, suggesting that redesigning the current policy is profitable.

**Related Work.** Widely-used OPE methods include inverse probability weighting (IPW) (Precup 2000; Strehl et al. 2010), self-normalized IPW (Swaminathan and Joachims 2015), Doubly Robust (Dudík et al. 2014), and more advanced variants (Wager and Athey 2018; Farajtabar, Chow, and Ghavamzadeh 2018; Su et al. 2020). These methods are based on importance sampling (IS) and require that the logging policy be of full support, i.e., assign a positive probability to every action potentially chosen by the counterfactual policy. This restriction makes them hard to use when the logging policy is of deficient support.

There are two existing approaches to deficient support logging policies.<sup>2</sup> The first approach considers a logging policy that varies over time or across individuals (Strehl et al. 2010). Viewing the sequence of varying logging policies as a single full support logging policy, it is possible to apply IS-based OPE methods. Unlike this approach, our approach is usable even when the logging policy is fixed.

The second approach, called the Direct Method or Regression Estimator, predicts the mean reward conditional on the action and context by supervised learning and uses the prediction to estimate the performance of a counterfactual policy (Beygelzimer and Langford 2009; Dudík et al. 2014). Similar regression-based methods are proposed for reinforcement learning settings (Duan, Jia, and Wang 2020). This approach is sensitive to the accuracy of the mean reward prediction. It may have a large bias if the regression model is not correctly specified. This issue is particularly severe when the logging policy is of deficient support, since each action is observed only in a limited area of the context space. Our approach instead predicts the mean reward differences between actions by exploiting local subsamples near the decision boundaries without specifying the regression model. Narita and Yata (2022) originally develop and empirically apply this approach in the context of treatment effect estimation with a binary treatment. This paper extends their approach to OPE with multiple actions. This idea relates to regression discontinuity designs in the social sciences (Lee and Lemieux 2010).

It is worth noting that our approach is applicable to *off-*

<sup>2</sup>Sachdeva et al. (2020) also proposes another approach in which they restrict the policy space.

*policy selection*, in which the researcher is to design a decision rule to select a policy given a finite set of policies (Kuzborskij et al. 2021). Since our method can estimate the expected reward of the policies, we can first estimate the reward of each, and then choose the one with the highest expected reward.

## 2 Framework

$\mathcal{A} := \{1, \dots, m\}$  is a set of *actions* that the decision maker can choose from. Let  $\mathbb{R}^p$ -valued random variable  $X$  denote the *context* that the decision maker observes when picking an action. Let  $\mathcal{X}$  denote the support of  $X$ . To simplify the exposition, we assume that  $X$  is continuously distributed. Let a tuple of  $m$   $\mathbb{R}$ -valued random variables  $(Y(1), \dots, Y(m))$  denote *potential rewards*;  $Y(a)$  denotes a *potential reward* that is observed when action  $a$  is chosen.  $(Y(1), \dots, Y(m), X)$  follows distribution  $P$ , which is unknown to the decision maker.

A *policy* chooses an action given a context. Let  $ML : \mathbb{R}^p \rightarrow \Delta(\mathcal{A})$  represent the *logging policy*, where  $ML(a|x)$  is the probability of taking action  $a$  for individuals with context  $x$ . We assume that the analyst knows the logging policy and is able to simulate it. That is, the analyst is able to compute the probability  $ML(a|x)$  for each action  $a \in \mathcal{A}$  given any context  $x \in \mathbb{R}^p$ . Suppose we have log data  $\{(Y_i, X_i, A_i)\}_{i=1}^n$  generated as follows. For each individual  $i$ , (1)  $(Y_i(1), \dots, Y_i(m), X_i)$  is i.i.d. drawn from  $P$ ; (2) Given  $X_i$ , the action  $A_i$  is randomly chosen based on the probability  $ML(\cdot|X_i)$ ; (3) We observe the reward  $Y_i := Y_i(A_i)$ . Note that only one of  $Y_i(1), \dots, Y_i(m)$  is observed for individual  $i$  and recorded as  $Y_i$  in the log data. The joint distribution of  $(Y, X, A)$  is determined once  $ML$  and  $P$  are given.

**Prediction Target.** We are interested in estimating the expected reward from any given *counterfactual policy*  $\pi : \mathbb{R}^p \rightarrow \Delta(\mathcal{A})$ , which chooses a distribution of actions given individual context:

$$V(\pi) := E \left[ \sum_{a \in \mathcal{A}} Y(a) \pi(a|X) \right].$$

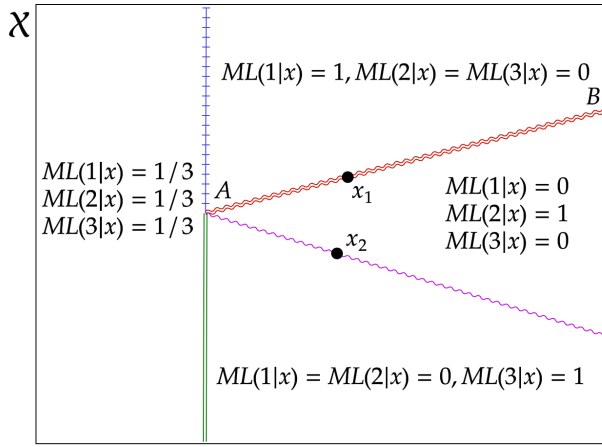
## 3 Learning with Infinite Data

We first consider the identification problem, which asks whether it is possible to learn  $V(\pi)$  if we had an infinite amount of data. Formally, we say that  $V(\pi)$  is *identified* if it is uniquely determined by the joint distribution of  $(Y, X, A)$ . A key step toward answering the identification question is what we call the *Approximate Propensity Score* (APS). To define it, for  $a \in \mathcal{A}$  and  $x \in \mathcal{X}$ , let:

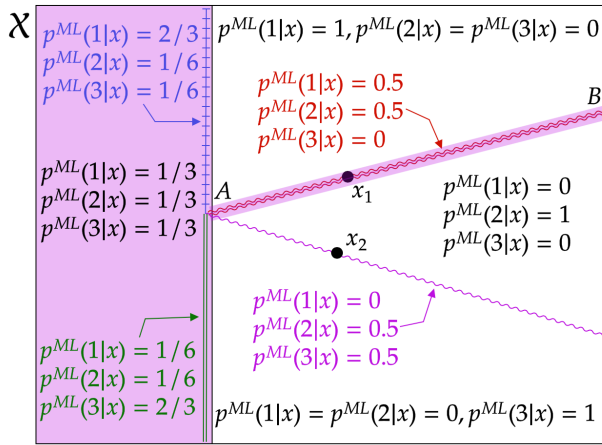
$$p_\delta^{ML}(a|x) := \frac{\int_{B(x,\delta)} ML(a|x^*) dx^*}{\int_{B(x,\delta)} dx^*},$$

where  $B(x, \delta) = \{x^* \in \mathbb{R}^p : \|x - x^*\| < \delta\}$  is the  $\delta$ -ball around  $x \in \mathcal{X}$ . Here,  $\|\cdot\|$  denotes the Euclidean norm on  $\mathbb{R}^p$ . To make common  $\delta$  for all dimensions reasonable, we

<sup>3</sup>This assumption is valid when we have a batch of log data generated by a fixed policy.



(a)



(b)

Notes: This figure shows an example of logging policy  $ML$  (panel (a)) and corresponding APS  $p^{ML}$  (panel (b)). The shaded region in panel (b) indicates the subpopulation for which  $p^{ML}(1|x) > 0$  and  $p^{ML}(2|x) > 0$ . As discussed in Section 4, our method uses the subsample in the shaded region to estimate the conditional mean difference  $E[Y(2)|X] - E[Y(1)|X]$ .

Figure 1: Example of the Approximate Propensity Score

normalize  $X_{ij}$  to have mean zero and variance one for each  $j = 1, \dots, p$ . We assume that  $ML$  is a Lebesgue measurable function so that the integrals exist. We then define APS  $p^{ML}$  as follows: for  $a \in \mathcal{A}$  and  $x \in \mathcal{X}$ ,

$$p^{ML}(a|x) := \lim_{\delta \rightarrow 0} p_{\delta}^{ML}(a|x).$$

Figure 1 illustrates APS. Here  $\mathcal{X} \subseteq \mathbb{R}^2$ ,  $\mathcal{A} = \{1, 2, 3\}$ , and the support of  $X$  is divided into four sets depending on the value of  $ML$  as in panel (a). Panel (b) shows the corresponding APS. For the interior points of each of the four sets, APS is equal to  $ML$ . On the border of any two sets, APS is the average of the  $ML$  values in the two sets.

Our identification analysis uses the following assumption.

**Assumption 1** (Local Mean Continuity). For any  $a \in \mathcal{A}$ , the conditional expectation function  $E[Y(a)|X = x]$  is continuous at each  $x \in \mathcal{X}$  such that  $p^{ML}(a|x) > 0$  and  $ML(a|x) = 0$ .

$ML(a|x) = 0$  means that action  $a$  is never taken for individuals with context  $x$ . If APS of  $a$  at  $x$  is nonzero ( $p^{ML}(a|x) > 0$ ), however, there exists a point close to  $x$  that has a positive probability of receiving action  $a$ , which enables us to observe the reward from the action near  $x$ . For any such point  $x$ , Assumption 1 ensures that the points close to  $x$  have similar conditional means of the potential reward  $Y(a)$ . Thus, the conditional mean reward from action  $a$  at  $x$  is identified. On the other hand, when  $ML(a|x) > 0$ , action-context pair  $(a, x)$  is observed, allowing us to identify the mean reward without any assumptions. Assumption 1 therefore does not impose continuity at such points. The lemma below summarizes the above argument. For a set  $A \subset \mathbb{R}^p$ , let  $\text{int}(A)$  denote the interior of  $A$ .

**Lemma 1** (Identification of Conditional Means). *If Assumption 1 holds, then for each  $a \in \mathcal{A}$ ,  $E[Y(a)|X = x]$  is identified for every  $x \in \text{int}(\mathcal{X})$  such that  $p^{ML}(a|x) > 0$ .*

We use Lemma 1 to analyse identification of  $V(\pi)$ . Suppose first that  $\pi(a|x) > 0 \implies p^{ML}(a|x) > 0$ , that is, the counterfactual policy  $\pi$  only chooses actions with nonzero APS. Lemma 1 implies that the conditional mean reward is identified at every  $(a, x)$  pair that could be realized under the policy  $\pi$ . As a result, the expected reward  $V(\pi)$  is identified for any such policy. However, if there exists  $(a, x)$  such that  $\pi(a|x) > 0$  but  $p^{ML}(a|x) = 0$ , we cannot identify  $V(\pi)$  without additional assumptions. To be able to identify  $V(\pi)$  for any policy  $\pi$ , we assume that the difference in the conditional mean reward function  $E[Y(a)|X]$  between any two actions is constant over  $\mathcal{X}$ .

**Assumption 2** (Constant Conditional Mean Differences). There exists a function  $\beta : \mathcal{A} \times \mathcal{A} \rightarrow \mathbb{R}$  such that  $E[Y(a)|X] - E[Y(a')|X] = \beta(a, a')$ .

At the end of Section 4, we discuss how our results would change if we drop Assumption 2 and a potential way of relaxing this. We also impose the following condition on APS.

**Assumption 3** (Existence of Nonzero APS). For every  $a \in \{2, \dots, m\}$ , there exists a sequence  $\{a_1, \dots, a_L\}$  with  $a_1 = 1$  and  $a_L = a$  for which the following condition holds: for every  $l \in \{1, \dots, L-1\}$ , there exists  $x \in \text{int}(\mathcal{X})$  such that  $p^{ML}(a_l|x) > 0$  and  $p^{ML}(a_{l+1}|x) > 0$ .

Assumption 3 states that there exists a path from a baseline action ( $a_1 = 1$ ) to any other action ( $a_L = a$ ) for which APS of any two consecutive actions ( $a_l$  and  $a_{l+1}$ ) is positive at some  $x$ . For example, suppose that  $m = 3$ ,  $p^{ML}(1|x_1) > 0$ ,  $p^{ML}(2|x_1) > 0$ ,  $p^{ML}(2|x_2) > 0$  and  $p^{ML}(3|x_2) > 0$  for some  $x_1, x_2 \in \mathcal{X}$  as in Figure 1 (b). In this case, the sequence  $\{1, 2\}$  satisfies the condition in Assumption 3 for  $a = 2$ , and the sequence  $\{1, 2, 3\}$  satisfies the condition for  $a = 3$ . By Lemma 1, the four conditional means  $E[Y(1)|X = x_1]$ ,  $E[Y(2)|X = x_1]$ ,  $E[Y(2)|X = x_2]$  and  $E[Y(3)|X = x_2]$  are identified. Hence, the two differences  $E[Y(1)|X = x_1] - E[Y(2)|X = x_1]$  and

$E[Y(2)|X = x_2] - E[Y(3)|X = x_2]$  are identified. Under Assumption 2, the two differences do not depend on  $x$ . As a result,  $E[Y(1)|X = x] - E[Y(2)|X = x]$  and  $E[Y(2)|X = x] - E[Y(3)|X = x]$  are identified for every  $x \in \mathcal{X}$ . Noting that  $E[Y(a)|X = x]$  is identified for at least one  $a \in \mathcal{A}$  for every  $x \in \mathcal{X}$ , we can use the differences to identify  $E[Y(a)|X = x]$  for every  $(a, x)$  pair, even for those not observed in data. Thus,  $V(\pi)$  is identified for any policy  $\pi$ .

**Proposition 1** (Identification of  $V(\pi)$ ). *Under Assumptions 1–3,  $V(\pi)$  is identified for any policy  $\pi$ .*

Assumption 3 typically holds if every action is chosen with a positive probability in some region of the context space  $\mathcal{X}$ . For example, consider a deterministic logging policy that chooses the action with the largest predicted conditional mean reward given the context ( $E[Y(a)|X]$ ), where the predictions are obtained from supervised learning. If every action  $a$  has a region where it is predicted to be optimal, then every action usually shares boundaries with at least one other action. Since  $p^{ML}(a|x) > 0$  and  $p^{ML}(a'|x) > 0$  at the boundaries shared by two actions  $a$  and  $a'$  (unless the boundaries are irregularly shaped), we can find a sequence of actions that satisfies Assumption 3.

## 4 Learning with Finite Data

**OPE Estimator.** Suppose that we observe a sample  $\{(Y_i, X_i, A_i)\}_{i=1}^n$  of size  $n$ . We propose an OPE estimator based on the following expression of our prediction target  $V(\pi)$ : under Assumption 2,

$$V(\pi) = V(ML) + E \left[ \sum_{a=2}^m \beta(a, 1) (\pi(a|X) - ML(a|X)) \right]. \quad (1)$$

Appendix E derives this expression. Since  $V(ML)$  is the value from the logging policy  $ML$ ,  $V(ML)$  can be estimated by the sample mean of  $Y_i$ . Our identification analysis suggests a way of conducting OPE on any policy  $\pi$ : (1) estimate  $\beta(a, a')$  for each  $(a, a')$  pair such that  $p^{ML}(a|x) > 0$  and  $p^{ML}(a'|x) > 0$  for some  $x$ ; (2) use the estimates to recover  $\beta(a, 1)$  for every  $a \in \{2, \dots, m\}$  and plug them into the sample analogue of the above expression. For simplicity, we consider a setup in which  $p^{ML}(a|x) > 0$  and  $p^{ML}(1|x) > 0$  for some  $x$  for every  $a$  so that we can directly estimate  $\beta(a, 1)$  in step (1) above.

To estimate  $\beta(a, 1)$ , we use the subsample

$$\mathcal{I}(a; \delta_n) := \{i : A_i \in \{1, a\}, q_{\delta_n}^{ML}(a | X_i) \in (0, 1)\},$$

where

$$q_{\delta_n}^{ML}(a | X_i) := \frac{p_{\delta_n}^{ML}(a | X_i)}{p_{\delta_n}^{ML}(a | X_i) + p_{\delta_n}^{ML}(1 | X_i)},$$

and  $\delta_n$  is a given bandwidth. The bandwidth shrinks towards zero as the sample size  $n$  increases.<sup>4</sup>  $q_{\delta_n}^{ML}(a|X_i)$

<sup>4</sup>For the bandwidth  $\delta_n$ , we suggest considering several different values and check if the estimates are robust to bandwidth changes. It is hard to pick  $\delta_n$  in a data-driven way to minimize the mean squared error, since it would require nonparametric estimation of functions on the high-dimensional context space.

can be viewed as APS of action  $a$  within the subsample for which either action 1 or  $a$  is assigned. The subsample  $\mathcal{I}(a; \delta_n)$  contains all observations  $i$  such that both actions 1 and  $a$  can be chosen by the logging policy locally around  $X_i$ . For example, in Figure 1 (b), the shaded region corresponds to the subsample  $\mathcal{I}(2; \delta_n)$ . This covers not only the subsample subject to full randomization (for which  $ML(1|x) = ML(2|x) = ML(3|x) = 1/3$ ) but also the local subsample near the deterministic decision boundary  $AB$  between actions 1 and 2.

We propose minimizing the sum of squared errors on the subsample  $\mathcal{I}(a; \delta_n)$ :

$$(\hat{\alpha}_a, \hat{\beta}_a, \hat{\gamma}_a) = \underset{(\alpha_a, \beta_a, \gamma_a)}{\operatorname{argmin}} \sum_{i \in \mathcal{I}(a; \delta_n)} \left( Y_i - \alpha_a - \beta_a 1\{A_i = a\} - \gamma_a q_{\delta_n}^{ML}(a|X_i) \right)^2, \quad (2)$$

where  $1\{\cdot\}$  is the indicator function.  $\hat{\beta}_a$  is our estimator of  $\beta(a, 1)$ . We include  $q_{\delta_n}^{ML}(a|X_i)$  as an explanatory variable to adjust for imbalance in the context distribution between actions 1 and  $a$ , as is done with the standard propensity score (Angrist and Pischke 2008; Hull 2018). We then define our OPE estimator as:

$$\hat{V}(\pi) = \frac{1}{n} \sum_{i=1}^n \left( Y_i + \sum_{a=2}^m \hat{\beta}_a (\pi(a|X_i) - ML(a|X_i)) \right). \quad (3)$$

It is worth noting that our method does not require the model selection.

For estimating  $\beta(a, 1)$ , the above method uses APS  $p_{\delta_n}^{ML}(a|X_i)$ , which may be difficult to compute analytically if  $ML$  is complex. In such a case, we propose approximating it by brute force simulation. We draw a value of  $x$  from the uniform distribution on  $B(X_i, \delta_n)$  a number of times, compute  $ML(a|x)$  for each draw, and take the average of  $ML(a|x)$  over the draws.<sup>5</sup> We then use it instead of  $p_{\delta_n}^{ML}(a|X_i)$  to compute  $q_{\delta_n}^{ML}(a|X_i)$ , and then compute  $\hat{\beta}(a, 1)$  and  $\hat{V}(\pi)$  as in (2) and (3).

**Consistency.** We show that  $\hat{V}(\pi)$  is a consistent estimator of  $V(\pi)$ , that is,  $\hat{V}(\pi)$  converges in probability to  $V(\pi)$  as  $n \rightarrow \infty$  under some regularity conditions.

**Assumption 4** (Regularity conditions). See Appendix A for details.

**Theorem 1** (Consistency of  $\hat{V}(\pi)$ ). *Suppose that Assumptions 2 and 4 hold,  $\delta_n \rightarrow 0$ , and  $n\delta_n \rightarrow \infty$  as  $n \rightarrow \infty$ . Then  $\hat{V}(\pi)$  converges in probability to  $V(\pi)$  for every policy  $\pi$ .*

The main argument in the proof of Theorem 1 is similar to the one used for the consistency result of Narita and Yata

<sup>5</sup>The approximation error of the simulated APS relative to true  $p_{\delta_n}^{ML}(a|X_i)$  has a  $S^{-\frac{1}{2}}$  rate of convergence, where  $S$  is the number of simulation draws. This rate does not depend on the dimension of  $X_i$ , so the simulation error can be made negligible by using a large number of simulation draws even when  $X_i$  is high dimensional.

(2022) (the first part of their Theorem 1). We extend their result to OPE with multiple actions.

Our consistency result requires that  $\delta_n$  go to zero slower than  $n^{-1}$ . This ensures that, when  $ML$  is deterministic, we have sufficiently many observations in the  $\delta_n$ -neighborhood of the boundary of  $\Omega_a^* := \{x: ML(a|x) = 1\}$  (the set of the context values for which the probability of choosing action  $a$  is one). Importantly, the rate condition does not depend on the dimension of  $X_i$ . This is because we use all the observations in the  $\delta_n$ -neighborhood of the boundary, and the number of those observations is of order  $n\delta_n$  regardless of the dimension of  $X_i$  if the boundary is  $(p-1)$  dimensional. Our estimator is therefore expected to perform well even if  $X_i$  is high dimensional.

Our result holds under the assumption of constant conditional mean reward differences (Assumption 2). If this assumption does not hold for a deterministic logging policy,  $\hat{\beta}_a$  is a consistent estimator of the mean reward difference for the subpopulation on the decision boundary between actions  $a$  and 1 (see Appendix E). Therefore, our estimator may still perform well when we are interested in a counterfactual policy that marginally changes the logging policy's decision boundary.

One way to relax Assumption 2 is to consider a partition of  $\mathcal{X}$  and assume that the conditional mean difference between any two actions is constant within each cell in the partition. This allows the conditional mean differences to vary across cells. If for each  $(a, a')$  pair, each cell contains  $x$  such that  $p^{ML}(a|x) > 0$  and  $p^{ML}(a'|x) > 0$ , we can consistently estimate the conditional mean differences and the expected reward from any policy. How to find such a partition is an interesting future topic.

## 5 Simulations

### Experiment 1: Mix of A/B Test and Deterministic Logging Policy

Consider a tech company that conducts an A/B test using a small segment of the population. The company applies a deterministic logging policy to the rest of the population. We generate a random sample  $\{(Y_i, X_i, A_i)\}_{i=1}^n$  of size  $n = 50,000$  as follows. There are 5 actions ( $m = 5$ ) and 100 context variables ( $p = 100$ ), with  $X_i \sim N(0, \Sigma)$ .  $Y_i(a)$  is generated as  $Y_i(a) = 0.75 \sum_{k=1}^{100} X_{ki}^2 \alpha_{0,k} + 0.25u_i + \epsilon_i(a)$ , where  $\alpha_0 = (\alpha_{0,1}, \dots, \alpha_{0,100}) \in \mathbb{R}^{100}$ ,  $u_i \sim N(0, 1)$ , and  $\epsilon_i(a) \sim N(a, 1)$ . The conditional mean difference  $E[Y_i(a)|X_i] - E[Y_i(1)|X_i]$  is constant over  $x$ . The choice of parameters  $\Sigma$  and  $\alpha_0$  is explained in Appendix B. To generate  $A_i$ , let  $q_{0.99}^k$  be the 99th percentile of the  $k$ th context variable  $X_{ki}$ . Let  $\tau_{pred}^{ML}(x, a)$  be a prediction of the reward from action  $a$  given context value  $x$  obtained by supervised learning from a past, independent training sample  $\tilde{\mathcal{D}} = \{(\tilde{Y}_i, \tilde{X}_i, \tilde{A}_i)\}_{i=1}^{\tilde{n}}$  of size  $\tilde{n} = 10,000$  (see Appendix B for how we constructed  $\tilde{\mathcal{D}}$  and  $\tau_{pred}^{ML}$ ).  $A_i$  is then generated based on the logging policy:

$$ML(a|x) = \begin{cases} 1/5 & \text{if } x_1 \geq q_{0.99}^1 \\ 1 \left\{ a = \operatorname{argmax}_{a' \in \{1, \dots, 5\}} \tau_{pred}^{ML}(x, a') \right\} & \text{if } x_1 < q_{0.99}^1. \end{cases}$$

The first case corresponds to the A/B test segment while the second case to the deterministic policy segment. Finally,  $Y_i$  is generated as  $Y_i = Y_i(A_i)$ .

We simulate 1,000 hypothetical samples from the above data-generating process. For each simulation, we use the simulated sample to estimate the value of a counterfactual policy  $\pi$ , another mix of an A/B test and a deterministic policy. With another reward prediction function  $\tau_{pred}^\pi$ ,

$$\pi(a|x) = \begin{cases} 1/5 & \text{if } x_2 \geq q_{0.99}^2 \\ 1 \left\{ a = \operatorname{argmax}_{a' \in \{1, \dots, 5\}} \tau_{pred}^\pi(x, a') \right\} & \text{if } x_2 < q_{0.99}^2. \end{cases}$$

**Alternative Methods.** We compare our method with two alternative estimators. The first uses the A/B test segment (for which  $ML(a|X_i) = 1/5$ ) while the second uses the full sample. The methods first compute the simple mean differences in reward  $Y_i$  between actions  $a \in \{2, \dots, 5\}$  and 1, and then plugs them into  $\hat{\beta}_a$  of Eq. (3). Both our method and the alternative estimator with the A/B test segment produce consistent estimators of the prediction target  $V(\pi)$ . However, the alternative uses only the A/B test segment while our method additionally uses the local subsample near the decision boundary of the deterministic policy as we discussed in Section 4.

**Result.** The first panel of Table 1 presents the bias, standard deviation (S.D.) and root mean squared error (RMSE) of our proposed estimators with several choices of  $\delta$  and two alternative estimators. The alternative estimator using the full sample has a larger bias than the other two, since it does not control for the difference in the context distribution between actions. Our proposed estimator outperforms the alternative estimator using the A/B test sample in terms of RMSE. This suggests that exploiting both of the A/B test segment and the local subsample near the deterministic decision boundary can lead to better performance than using only the A/B test segment.

### Experiment 2: Upper Confidence Bound Logging Policy

In the second experiment, both the logging policy and the counterfactual policy are deterministic. The rest of the setup is the same as that in the first experiment. We first use the independent training sample  $\tilde{\mathcal{D}}$  to train an Upper Confidence Bound bandit algorithm. The logging policy  $ML$  is given by  $ML(a|x) = 1 \{a = \operatorname{argmax}_{a' \in \{1, \dots, 5\}} UCB(x, a')\}$ , where  $UCB(x, a)$  is an upper confidence bound of  $E[Y_i(a)|X_i = x]$ . See Appendix B for training details. We do not update the policy while generating  $\{(Y_i, X_i, A_i)\}_{i=1}^n$  in the simulation. The sample is a batch of log data.

For the counterfactual policy  $\pi$ , we use  $\tilde{\mathcal{D}}$  to train a model  $f(x, a)$  that predicts the reward given the context and action, using sklearn's RandomForestRegressor with 500 trees and otherwise default parameters. The counterfactual policy tries to maximize the expected reward  $V(\pi)$  by choosing the action with the largest predicted reward:  $\pi(a|x) = 1 \{a = \operatorname{argmax}_{a' \in \{1, \dots, 5\}} f(x, a')\}$ .

**Alternative Method.** We compare our method with an alternative estimator using the Direct Method. This first

	Our Proposed Method with APS Controls				Method with Mean Differences		Direct Method (7)
	$\delta = 0.1$ (1)	$\delta = 0.5$ (2)	$\delta = 1$ (3)	$\delta = 2.5$ (4)	A/B Test Sample (5)	Full Sample (6)	
Experiment 1: Mix of A/B Test and Deterministic Logging Policy							
Bias	-.060	-.057	-.057	-.060	-.061	-.075	—
S.D.	.099	.098	.096	.096	.101	.103	—
RMSE	.115	.113	.112	.113	.118	.128	—
Avg. $N$	1862	6362	12502	33122	500	50000	—
Experiment 2: Upper Confidence Bound Logging Policy							
Bias	.048	.047	.046	.047	—	—	.342
S.D.	.033	.030	.029	.029	—	—	.012
RMSE	.058	.056	.055	.055	—	—	.342
Avg. $N$	3397	17344	31107	47601	—	—	50000

Notes: This table shows the bias, the standard deviation (S.D.), and the root mean squared error (RMSE) of the estimators of the reward from the counterfactual policy  $V(\pi)$  in the two simulation experiments. We use 1,000 simulations of a size 50,000 sample to compute these statistics. Columns (1)–(4) report estimates from our method with several choices of  $\delta$ . Each APS is computed by averaging 100 simulation draws of the  $ML$  value. In columns (5)–(6), we estimate the mean reward differences  $\beta(a, 1)$  by the sample mean differences in the A/B test segment and the full sample, respectively. In column (7), we estimate  $\beta(a, 1)$  by fitting a linear model that predicts the reward from the context and action. The bottom row of each panel shows the average number of observations with nonzero APS for every action (Columns (1)–(4)), that with nonzero  $ML$  for every action (Column (5)), or the total sample size (Columns (6)–(7)).

Table 1: Simulation results: bias, S.D., and RMSE of estimators of  $V(\pi)$

fits a linear model  $Y_i = \alpha + \sum_{a=2}^5 \beta_a 1\{A_i = a\} + \sum_{k=1}^{100} X_{ki} \gamma_k + e_i$ , then makes the reward prediction from action  $a$  for individual  $i$  by  $\hat{\mu}_i(a) = Y_i + (\hat{\beta}_a - \hat{\beta}_{A_i})$ , and finally computes  $\hat{V}(\pi) = \frac{1}{n} \sum_{i=1}^n \sum_{a=1}^5 \hat{\mu}_i(a) \pi(a|X_i)$ . The linear model used by this method correctly imposes the constant conditional mean differences but misspecifies the functional form with respect to  $X_i$ .

**Result.** The second panel of Table 1 shows the result. The alternative using the Direct Method is significantly biased due to model misspecification. Our proposed estimator seems to effectively use the local subsample near the decision boundary and has smaller bias and RMSE than the alternative.

## 6 Real-World Application

**Setup.** We apply our method to empirically evaluate a coupon targeting policy of an online platform. This application uses proprietary data provided by Mercari, Inc. This company conducts the following promotional campaign. They target customers who signed up for Mercari 4 days ago but have not made a purchase yet. The company uses a logging policy based on an uplift model to determine whether they offer a promotional coupon to each target customer. If customers receive the coupon and make a purchase, they get 900 points (equivalent to 8.34 USD) that they can use for future purchases. We observe data  $(Y_i, X_i, A_i)$  for each target user  $i$  from this campaign, where action  $A_i \in \{0, 1\}$  is whether the logging policy recommended offering the coupon to the customer ( $A_i = 1$ ) or not ( $A_i = 0$ ),  $X_i$  is the

vector of more than 200 input features for the uplift model, and  $Y_i$  is an outcome such as the customer’s spending after this coupon offer.

The company’s logging policy works as follows. They first use data from a past A/B test and XGBoost to train a model of the conditional average effect of the coupon on purchases (they use library pylift for implementation). Let  $\tau(x)$  be the predicted coupon effect for those whose feature value is  $X_i = x$ . The logging policy then recommends offering a coupon to customer  $i$  if the predicted effect is in the top 80% of the distribution of predicted effects. That is, the logging policy  $ML$  is given by  $ML(1|x) = 1\{\tau(x) \geq c\}$ , where  $c$  is the 20th quantile of the distribution of  $\tau(X_i)$ .

**Effects of Policy Recommendation.** We first apply our method to the logged data generated by the above policy to estimate the effect of the policy recommendation  $A_i$  ( $\beta(1, 0) = E[Y_i(1) - Y_i(0)]$ ) on the following three outcomes: (1) the purchase value (how much the customer spent), (2) the number of transactions, and (3) point usage (how many points the customer used). All outcomes are sums over 18 days after the coupon offer decision. We compute APS with  $\delta \in \{0.4, 0.8, 1.2, 2.0, 3.0\}$ .<sup>6</sup>

Columns (1)–(5) in the first three rows of Table 2 report the estimated effects of the policy recommendation  $A_i$ . We normalize the estimates by dividing the original numbers by

<sup>6</sup>Unlike the theoretical framework, the feature vector  $X_i$  consists of discrete and continuous variables. We compute APS by fixing the value of the discrete part and computing by simulation the APS integral with respect to the continuous part. See Appendix C for details.

	Our Proposed Method with APS Controls					Mean
	$\delta = 0.4$ (1)	$\delta = 0.8$ (2)	$\delta = 1.2$ (3)	$\delta = 2.0$ (4)	$\delta = 3.0$ (5)	Differences (6)
Effect on Purchase Value	0.35 (0.59)	0.82 (0.39)	0.92 (0.30)	0.54 (0.28)	0.72 (0.21)	-0.17 (0.11)
Effect on # of Transactions	0.43 (0.50)	0.47 (0.34)	0.66 (0.28)	0.49 (0.25)	0.74 (0.19)	-0.07 (0.10)
Effect on Point Usage	0.37 (0.42)	0.71 (0.29)	0.57 (0.26)	0.47 (0.22)	0.64 (0.17)	0.68 (0.04)
Coupon Cost Effectiveness Measure	79.57 (130)	96.35 (48.97)	134 (61.97)	93.51 (49.33)	92.07 (28.45)	—
$N$	2758	4688	6016	8085	9602	89486

*Notes:* The first three rows of this table report estimated effects of the policy recommendation  $A_i$  on purchase behavior. Columns (1)–(5) report estimates from our method with several choices of  $\delta$  used to compute APS. Column (6) reports the outcome mean differences between those with  $A_i = 1$  and  $A_i = 0$ . Each APS is computed by averaging 100 simulation draws of the logging policy’s binary decision. All numbers in the first three rows are normalized by dividing the original estimates by the sample outcome means. The fourth row reports our measure of coupon cost effectiveness, which predicts how much the purchase value would increase in USD if we increased the cost of the campaign by 1 USD. Heteroskedasticity-robust standard errors are reported in parentheses. The last row reports the number of observations with nonzero APS for every action (Columns (1)–(5)) or the total sample size (Column (6)).

Table 2: Off-policy evaluation using policy’s generated data

the sample outcome means for confidentiality. The results show that the effects of the policy recommendation  $A_i$  on the purchase value, the number of transactions, and point usage are 35–92%, 43–74%, and 37–71% of their sample means, respectively. These positive effects mark a sharp contrast with Column (6), which reports the simple differences in the outcome means between those with  $A_i = 1$  and those with  $A_i = 0$ . The simple mean differences on the purchase value and the number of transactions are negative. These negative estimates suggest that the logging policy tends to recommend a coupon to the customers who have a low propensity to make purchases. Our proposed method corrects for this negative selection bias by controlling for APS.

**Evaluation of Counterfactual Policies.** The company needs to compensate for the discount that customers get by using points. Thus, adopting a new policy would be profitable only when the increase in revenue is sufficiently large compared to that in point usage. The company charges sellers 10% of every payment from the buyer; the revenue increases by 10% of the increase in purchase value. Hence, the policy change is beneficial if the ratio of the increases in the average purchase value and point usage is larger than 10.

Suppose we change our policy from  $ML$  to a counterfactual one  $\pi$ . Let  $Y_i^1$  and  $Y_i^2$  denote the purchase value and point usage respectively. Under the constant conditional effect assumption, i.e.,  $E[Y_i^1(1) - Y_i^1(0)|X_i] =: \beta$  and  $E[Y_i^2(1) - Y_i^2(0)|X_i] =: \gamma$ , the ratio is:

$$\frac{E[\sum_{a=0}^1 Y_i^1(a)\pi(a|X_i)] - E[\sum_{a=0}^1 Y_i^1(a)ML(a|X_i)]}{E[\sum_{a=0}^1 Y_i^2(a)\pi(a|X_i)] - E[\sum_{a=0}^1 Y_i^2(a)ML(a|X_i)]} = \frac{\beta E[\pi(1|X_i) - ML(1|X_i)]}{\gamma E[\pi(1|X_i) - ML(1|X_i)]} = \frac{\beta}{\gamma}.$$

The fourth row of Table 2 reports the estimates of the ratio  $\beta/\gamma$ . The estimates are larger than 10 for all  $\delta$ ’s. This result suggests that it would be profitable to expand the campaign.

As mentioned in Section 4, without the constant conditional effect assumption, our estimator for the effect of the policy recommendation is a consistent estimator of the conditional effect for the subpopulation on the decision boundary, i.e.,  $E[Y_i(1) - Y_i(0)|\tau(X_i) = c]$ . Our estimates in the fourth row of Table 2 therefore can be interpreted as a measure for the cost effectiveness of the counterfactual policy that slightly lowers the threshold  $c$ . Without the constant conditional effect assumption, the result still suggests that marginally expanding the campaign would be profitable.

## 7 Conclusion

We develop an OPE method for a class of logging policies including deficient support ones. Our method is based on the newly developed “Approximate Propensity Score.” We prove that our estimator is consistent and demonstrate its practical performance through simulations and a real-world application. Promising directions for future work include developing a data-driven procedure to optimize the bandwidth. Also, the assumption of constant conditional mean reward differences may not be plausible in some applications. It will be challenging but interesting to relax this assumption to allow for certain types of heterogeneity. Finally, we look forward to applications of our method in a variety of business, policy, and scientific domains using machine learning.

## References

Angrist, J. D.; and Pischke, J.-S. 2008. *Mostly harmless econometrics: an empiricist’s companion*. Princeton Uni-

versity Press.

Beygelzimer, A.; and Langford, J. 2009. The offset tree for learning with partial labels. *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 129–138.

Duan, Y.; Jia, Z.; and Wang, M. 2020. Minimax-optimal off-policy evaluation with linear function approximation. *Proceedings of the 37th International Conference on Machine Learning*, 2701–2709.

Dudík, M.; Erhan, D.; Langford, J.; and Li, L. 2014. Doubly robust policy evaluation and optimization. *Statistical Science*, 29(4): 485–511.

Farajtabar, M.; Chow, Y.; and Ghavamzadeh, M. 2018. More robust doubly robust off-policy evaluation. *Proceedings of the 35th International Conference on Machine Learning*, 80: 1447–1456.

Hull, P. 2018. Subtracting the propensity score in linear models. *Working Paper*.

Kuzborskij, I.; Vernade, C.; Gyorgy, A.; and Szepesvari, C. 2021. Confident Off-Policy Evaluation and Selection through Self-Normalized Importance Weighting. *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, 640–648.

Lee, D. S.; and Lemieux, T. 2010. Regression discontinuity designs in economics. *Journal of Economic Literature*, 48(2): 281–355.

Narita, Y.; and Yata, K. 2022. Algorithm is Experiment: Machine Learning, Market Design, and Policy Eligibility Rules. *arXiv preprint arXiv:2104.12909*.

Precup, D. 2000. Eligibility traces for off-policy policy evaluation. *Proceedings of the Seventeenth International Conference on Machine Learning*, 759–766.

Rosenbaum, P. R.; and Rubin, D. B. 1983. The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1): 41–55.

Sachdeva, N.; Su, Y.; and Joachims, T. 2020. Off-policy Bandits with Deficient Support. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 965–975.

Strehl, A.; Langford, J.; Li, L.; and Kakade, S. M. 2010. Learning from logged implicit exploration data. *Proceedings of the 23rd International Conference on Neural Information Processing Systems*, 2217–2225.

Su, Y.; Dimakopoulou, M.; Krishnamurthy, A.; and Dudík, M. 2020. Doubly robust off-policy evaluation with shrinkage. *Proceedings of the 37th International Conference on Machine Learning*, 119: 9167–9176.

Swaminathan, A.; and Joachims, T. 2015. The self-normalized estimator for counterfactual learning. *Proceedings of the 28th International Conference on Neural Information Processing Systems*, 3231–3239.

Wager, S.; and Athey, S. 2018. Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 113(523): 1228–1242.