

Distributionally Robust Policy Evaluation and Learning for Continuous Treatment with Observational Data

Cheuk Hang Leung^{1*}, Yiyang Huang^{2*}, Yijun Li¹, Qi Wu^{1†}

¹City University of Hong Kong

²The Hong Kong Polytechnic University

chleung87@cityu.edu.hk, yiyhuang3-c@my.cityu.edu.hk, yijunli5-c@my.cityu.edu.hk, qi.wu@cityu.edu.hk

Abstract

Using offline observational data for policy evaluation and learning allows decision-makers to evaluate and learn a policy that connects characteristics and interventions. Most existing literature has focused on either discrete treatment spaces or assumed no difference in the distributions between the policy-learning and policy-deployed environments. These restrict applications in many real-world scenarios where distribution shifts are present with continuous treatment. To overcome these challenges, this paper focuses on developing a distributionally robust policy under a continuous treatment setting. The proposed distributionally robust estimators are established using the Inverse Probability Weighting (IPW) method extended from the discrete one for policy evaluation and learning under continuous treatments. Specifically, we introduce a kernel function into the proposed IPW estimator to mitigate the exclusion of observations that can occur in the standard IPW method to continuous treatments. We then provide finite-sample analysis that guarantees the convergence of the proposed distributionally robust policy evaluation and learning estimators. The comprehensive experiments further verify the effectiveness of our approach when distribution shifts are present.

Introduction

Most decision-making problems necessitate learning an effective personalized policy based on individual features from observational data. This process, commonly referred to as offline policy evaluation/learning, has diverse applications across various domains, including healthcare (Tang and Wiens 2021), recommendation (Li et al. 2010), and finance (Qin et al. 2022). Many studies have investigated offline policy evaluation/learning in discrete treatment settings which assume that the deployment environment is identical to the environment generating the training data, i.e., that there are no distributional shifts. This assumption, however, is often unrealistic in many real-world applications (Huang et al. 2023). For instance, an investment firm has developed an automated investment strategy for the US stock market based on extensive historical trading data. When attempting to apply this strategy directly to the UK stock market,

it may lose the predictive power due to the substantial difference between financial market environments. Similarly, a pharmaceutical company has developed a strategy for individualized Warfarin dosage adjustment according to their recent research on older adults. This strategy may perform well in the original clinical trial population, yet it may falter when applied to a new population, such as young adults, due to significant differences in physical conditions.

To address the challenge of distribution shifts in policy evaluation and learning, the problem can be formulated as a Distributionally Robust Optimization (DRO) problem. In the DRO framework, the goal is to find the worst-case solution within a set of distributions under certain degrees of model uncertainties. The uncertainty set is assumed to contain the distributions due to potential distribution shifts, and it can be characterized by constraining certain moments of order (Delage and Ye 2010; Zymler, Kuhn, and Rustem 2013) or by using divergence measures (Hu and Hong 2013; Kuhn et al. 2019; Chen, Sim, and Xu 2019; Gao, Chen, and Kleywegt 2022) to define appropriate deviations from a nominal distribution. The resulting solution provides robust, reliable, and conservative guarantees which can cope with the most adverse situations.

Furthermore, the objective function of the formulated DRO utilizes an inverse probability weighting (IPW) estimator (Wooldridge 2007) to estimate the expected potential reward under continuous treatment. Specifically, we extend the existing IPW estimator designed for discrete treatment settings to accommodate continuous treatments. Generally, the discrete-based IPW approach cannot be directly applied in continuous treatment settings, as it would reject most observed data. Moreover, although discretizing continuous treatments into categories is an intuitive and simple solution, it can lead to information loss and may fail to produce inferences that vary continuously with the treatment. We introduce a modified IPW approach incorporating a scaled kernel function with a bandwidth parameter, serving as a smooth nonparametric extension for computing histogram “buckets”. The proposed IPW estimator enables the distributionally robust policy evaluation and learning using observational datasets.

In summary, our framework addresses the challenges of policy evaluation and learning in continuous treatment settings in the presence of distribution shifts. The key contribu-

*The co-first authors contribute equally.

†Corresponding author.

tions of our paper are threefold:

1. We formulate the DRO problem with an IPW approach for policy evaluation/learning under the continuous treatment setting, and convert it to its equivalent dual form. As the standard IPW approach is not directly applicable in this context, we develop a tractable kernel-based form to approximate the dual problem.
2. We establish estimators for policy evaluation/learning and investigate their asymptotic properties. Specifically, the established estimators of distributionally robust values exhibit asymptotic normality, and the finite-sample regret decays to zero asymptotically.
3. Through simulated and empirical studies, we demonstrate that the policy learned using our method provides robustness to distribution shifts compared to standard nonrobust policy learning methods.

Literature Review

Considerable research has focused on causality in discrete treatment settings. However, exploring causality under continuous treatment remains limited in many real-world applications. Existing research on continuous treatment settings primarily focuses on directly modeling the relationship among response, treatment, and covariates. Notable contributions include (Schwab et al. 2020), who construct a multi-head neural network for this purpose; (Bica, Jordon, and van der Schaar 2020), who propose an end-to-end neural network based on generative adversarial networks (GANs); and (Bahadori, Tchetgen, and Heckerman 2022), who introduce a novel algorithm within the entropy balancing framework to optimize accuracy through end-to-end optimization. Another approach modifies IPW-based and Doubly Robust-based estimators (e.g., (Chernozhukov et al. 2018; Huang et al. 2022)) from discrete treatment settings by incorporating kernel functions to mitigate the direct rejection of observed data, as demonstrated by (Su, Ura, and Zhang 2019) and (Colangelo and Lee 2019).

Recent studies have focused on offline policy evaluation and learning. (Kitagawa and Tetenov 2018) establish finite sample regret bounds with a rate of $O_P\left(1/\sqrt{N}\right)$ for policy learning over a policy class with finite VC dimension. (Athey and Wager 2021) extend this analysis to examine regret bounds from an asymptotic perspective. (Zhao et al. 2012) and (Zhou et al. 2017) propose algorithms for policy learning and explore the statistical properties of learned policies and associated regret bounds. (Dudík, Langford, and Li 2011) utilize classic estimators for policy evaluation. (Kallus 2018) proposes a balance-based approach to reweight historical data and mimic datasets generated by evaluated or learned policies. (Zhou, Athey, and Wager 2023) exploit a cross-fitted approach for policy learning. The aforementioned studies primarily assume discrete treatment and do not account for distributional shifts. Nevertheless, distributional shifts are common since the studies are often conducted in different environments, highlighting the significance of studying distributionally robust policies. For instance, (Yang et al. 2023), (Shen, Xu, and Zavlanos 2024),

(Mo, Qi, and Liu 2021), and (Faury et al. 2020) primarily focus on shifts in covariates, whereas (Si et al. 2023) and (Kallus et al. 2022) address shifts in the joint distribution of responses, features, and treatments. Notably, to the best of our knowledge, studying policy evaluation and learning in the presence of distribution shifts under continuous treatment settings is still an open problem.

Background

Notations and Assumptions

Throughout the paper, we denote $A \in \mathcal{A} \subset \mathbb{R}$, $X \in \mathcal{X} \subset \mathbb{R}^d$, and $Y \in \mathcal{Y} \subset \mathbb{R}$ as the continuous treatment (also known as action or intervention), the covariates, and the continuous response (also known as outcome), respectively. We write $Y(A)$ to mean the potential response variable under the treatment A . We also assume that Y and $Y(A)$ are non-negative bounded variables, i.e., there exists $M > 0$ such that $0 \leq Y(a)$, $Y \leq M$. Finally, we let $(X_i, A_i, Y_i)_{i=1}^N$ be N independent and identically distributed (i.i.d.) triples from a fixed underlying distribution, and the probability measure of the underlying distribution is denoted as \mathbb{P}_0 .

Further, we adopt the Rubin potential outcome framework (e.g., (Rubin 1974; Imbens 2004; Imbens and Rubin 2015; Huang et al. 2021, 2024; Li et al. 2024)). Throughout the paper, we impose the following (*causal*) assumptions that are standard in the causal inference literature:

Assumption 1 (Consistency). *If $A = a$, we have $Y = Y(a)$.*

Assumption 2 (Unconfoundedness). $Y(a) \perp\!\!\!\perp A | X, \forall a$.

Assumption 3 (Positivity). *There exists a positive constant $\epsilon > 0$ such that $\inf_{a \in \mathcal{A}} \inf_{x \in \mathcal{X}} f_0(a|x) \geq \epsilon > 0$.*

Additionally, we follow (Kallus and Zhou 2018; Colangelo and Lee 2019) and impose differentiability assumptions on the probability density functions $f_0(y|a, x)$ and $f_0(a|x)$: $f_0(y|a, x)$ and $f_0(a|x)$ are three-times differentiable w.r.t. a and bounded uniformly on $(y, a, x) \in (\mathcal{Y}, \mathcal{A}, \mathcal{X})$. All proofs of the theorems are in the Appendix¹.

Problem Setup

Our objective is to find a policy π^* that maps \mathcal{X} to \mathcal{A} within a policy class Π that maximizes the expected outcomes, i.e.,

$$\pi^* = \arg \max_{\pi \in \Pi} Q(\pi) := \arg \max_{\pi \in \Pi} \mathbb{E}_{\mathbb{P}_0}[Y(\pi(X))]. \quad (1)$$

The learned policy obtained in Eqn. (1) may not generalize well to a new environment with a distribution that differs from \mathbb{P}_0 . As such, we can consider distributionally robust formulation of Eqn. (1):

$$\begin{aligned} \pi_{\text{DRO}}^* &= \arg \max_{\pi \in \Pi} Q_{\text{DRO}}(\pi), \quad \text{where} \\ Q_{\text{DRO}}(\pi) &= \inf_{\mathbb{P} \in \mathcal{U}_{\mathbb{P}_0}(\eta)} \mathbb{E}_{\mathbb{P}}[Y(\pi(X))], \\ \mathcal{U}_{\mathbb{P}_0}(\eta) &= \{\mathbb{P} : \mathbb{D}(\mathbb{P} \| \mathbb{P}_0) \leq \eta\}. \end{aligned} \quad (2)$$

Here, $\mathbb{D}(\cdot \| \cdot)$ denotes the distribution discrepancy. Throughout the paper, we choose it as the Kullback–Leibler (KL)

¹Available in the “proof” file of the Supplementary Material.

divergence (Kullback 1959; Kullback and Leibler 1951)². $\mathcal{U}_{\mathbb{P}_0}(\eta)$ is the *ambiguity set* (also known as the uncertainty set) with an *ambiguity radius* η . The ambiguity set contains all the possible distributions \mathbb{P} such that the discrepancy of \mathbb{P} relative to \mathbb{P}_0 is at most η .

Distributionally Robust Policy Evaluation

The Estimation of $Q_{\text{DRO}}(\pi)$

As proven in (Hu and Hong 2013), Eqn. (2) is equivalent to solving its Lagrangian dual, which is given as follows:

$$\begin{aligned} & - \min_{\alpha \geq 0} \left\{ \alpha \log \mathbb{E} \left[e^{-\frac{Y(\pi(X))}{\alpha}} \right] + \alpha \eta \right\} \\ & = \max_{\alpha \geq 0} \left\{ -\alpha \log \mathbb{E} \left[e^{-\frac{Y(\pi(X))}{\alpha}} \right] - \alpha \eta \right\} := \max_{\alpha \geq 0} \phi(\pi, \alpha). \end{aligned} \quad (3)$$

Since $Y(\pi(X))$ in Eqn. (3) is inaccessible, we reformulate $\mathbb{E} \left[e^{-\frac{Y(\pi(X))}{\alpha}} \right]$ to an IPW form similar to that in (Horvitz and Thompson 1952). The result is given in Lemma 1.

Lemma 1. *Under Assumptions 1 - 3, we have*

$$\mathbb{E} \left[e^{-\frac{Y(\pi(X))}{\alpha}} \right] = \mathbb{E} \left[\frac{\delta(\pi(X) - A)}{f_0(A|X)} e^{-\frac{Y}{\alpha}} \right] \quad (4)$$

for any $\alpha \geq 0$, where $\delta(\cdot)$ is the Dirac Delta function³.

Using Lemma 1, the expectation in Eqn. (3) can be replaced according to Eqn. (4). Note that the Dirac function $\delta(\cdot)$ is a theoretical generalized function and is often approximated by the scaled kernel function $K_h(\cdot)$ ⁴. As a result, we can consider the following approximated form:

$$Q_{\text{DRO}}^h(\pi) = \sup_{\alpha \geq 0} \left\{ -\alpha \log \mathbb{E} \left[\frac{e^{-\frac{Y}{\alpha}} K_h(\pi(X) - A)}{f_0(A|X)} \right] - \alpha \eta \right\}. \quad (5)$$

The above two quantities, $Q_{\text{DRO}}(\pi)$ and $Q_{\text{DRO}}^h(\pi)$, bring two important insights: (1) The optimal solutions of $Q_{\text{DRO}}(\pi)$ and $Q_{\text{DRO}}^h(\pi)$ are obtained by solving Eqns. (3) and (5) which are attainable for positive α due to the causal assumption. Further, the optimal solutions are finite for any π (see Auxiliary Result 2 in Appendix for details); (2) $Q_{\text{DRO}}^h(\pi) \rightarrow Q_{\text{DRO}}(\pi)$ as $h \rightarrow 0$. These two insights, consequently, guarantee that the optimal solutions of $Q_{\text{DRO}}^h(\pi)$ also converge

²Other measures such as the Wasserstein metric or other ϕ -divergence measures can be utilized (e.g., see (Kuhn et al. 2019) and (Husain, Nguyen, and van den Hengel 2023)). However, these approaches typically involve solving multi-level optimization problems which can be challenging to analyze.

³ $\delta(x) = \begin{cases} \infty & x = 0 \\ 0 & \text{otherwise} \end{cases}$ such that i) $\int_{\mathbb{R}} \delta(x) dx = 1$ and ii)

$\int_{\mathbb{R}} \delta(x) f(x) dx = f(0)$ for any arbitrary f defined on \mathbb{R} .

⁴A bounded differentiable function $K(\cdot)$ (i.e., $|K(\cdot)| \leq M_K$) is said to be a second-order kernel function if it satisfies i) $K(\cdot)$ is a symmetric function; ii) $\int_{-\infty}^{\infty} u K(u) du = 0$; iii) $\int_{-\infty}^{\infty} K(u) du = 1$. The scaled kernel function $K_h(\cdot)$ is defined such that $K_h(x) = \frac{1}{h} K\left(\frac{x}{h}\right)$, where h is termed as the *bandwidth* parameter. Note that $K_h(x) \xrightarrow{h \rightarrow 0} \delta(x)$ when $h \rightarrow 0$. Examples of kernels include Gaussian kernels or the Epanechnikov kernel.

to the optimal solutions of $Q_{\text{DRO}}(\pi)$. Therefore, we can construct an estimator of the IPW-based distributionally robust value $Q_{\text{DRO}}^h(\pi)$ to study the original distributionally robust value $Q_{\text{DRO}}(\pi)$ in Eqn. (2). We define

$$\bar{W}_N^h(\pi, \alpha) = \frac{1}{N} \sum_{i=1}^N \frac{K_h(\pi(X_i) - A_i)}{f_0(A_i|X_i)} e^{-\frac{Y_i}{\alpha}}. \quad (6a)$$

$$\hat{W}_N^h(\pi, \alpha) = \frac{\bar{W}_N^h(\pi, \alpha)}{S_N^h} = \frac{\bar{W}_N^h(\pi, \alpha)}{\frac{1}{N} \sum_{i=1}^N \frac{K_h(\pi(X_i) - A_i)}{f_0(A_i|X_i)}}. \quad (6b)$$

It is known that the IPW-based estimator $\bar{W}_N^h(\pi, \alpha)$ in Eqn. (6a) suffers from high-variance (Swaminathan and Joachims 2015; Khan and Ugander 2023). To address this challenge, we can use a normalized estimator $\hat{W}_N^h(\pi, \alpha)$ with a normalization factor S_N^h in Eqn. (6b) to approximate $\bar{W}_N^h(\pi, \alpha)$. Note that $\mathbb{E}[S_N^h] = 1$ and $S_N^h \rightarrow 1$ almost surely (see Auxiliary Result 1 in Appendix). Thus, $\hat{W}_N^h(\pi, \alpha)$ is asymptotically equivalent to $\bar{W}_N^h(\pi, \alpha)$. Consequently, we use the following $\hat{Q}_{\text{DRO}}^h(\pi)$ as the estimator of $Q_{\text{DRO}}^h(\pi)$ in Eqn. (5):

$$\begin{aligned} \hat{Q}_{\text{DRO}}^h(\pi) &= \max_{\alpha \geq 0} \hat{\phi}_N^h(\pi, \alpha) \\ &:= \max_{\alpha \geq 0} \left\{ -\alpha \log \hat{W}_N^h(\pi, \alpha) - \alpha \eta \right\}. \end{aligned} \quad (7)$$

$\hat{Q}_{\text{DRO}}^h(\pi)$ can be used to evaluate distributional robustness of a policy π . To summarize, we present the specific steps of obtaining $\hat{Q}_{\text{DRO}}^h(\pi)$ in the following Algorithm 1.

Algorithm 1: Distributionally robust policy evaluation

Input observed dataset $(X_i, A_i, Y_i)_{i=1}^N$, h , policy $\pi \in \Pi$.
Initialize: $\alpha \in \mathbb{R}^+ \cup 0$.

- 1: **repeat**
- 2: Compute $\hat{W}_N^h(\pi, \alpha)$ given in Eqn. (6b).
- 3: Update $\alpha: \alpha \leftarrow \alpha - \frac{\frac{\partial \hat{\phi}_N^h}{\partial \alpha}}{\frac{\partial^2 \hat{\phi}_N^h}{\partial \alpha^2}}$, where

$$\begin{aligned} \frac{\partial \hat{\phi}_N^h}{\partial \alpha} &= -\eta - \log \hat{W}_N^h - \frac{\alpha \frac{\partial \hat{W}_N^h}{\partial \alpha}}{\hat{W}_N^h} \\ \frac{\partial^2 \hat{\phi}_N^h}{\partial \alpha^2} &= -\frac{\frac{1}{N} \sum_{i=1}^N \frac{K_h(\pi(X_i) - A_i)}{f_0(A_i|X_i)} Y_i^2 e^{-\frac{Y_i}{\alpha}}}{\alpha^3 S_N^h \hat{W}_N^h} \\ &\quad + \frac{\alpha \left(\frac{\partial \hat{W}_N^h}{\partial \alpha} \right)^2}{(\hat{W}_N^h)^2}, \\ \frac{\partial \hat{W}_N^h}{\partial \alpha} &= \frac{\frac{1}{N} \sum_{i=1}^N \frac{K_h(\pi(X_i) - A_i)}{f_0(A_i|X_i)} Y_i e^{-\frac{Y_i}{\alpha}}}{\alpha^2 S_N^h} \end{aligned}$$

- 4: **until** α converges
 - 5: **Return** $\hat{Q}_{\text{DRO}}^h(\pi) \leftarrow \hat{\pi}_N^h(\pi, \alpha)$
-

The Statistical Property of $\hat{Q}_{\text{DRO}}^h(\pi)$

As $\hat{Q}_{\text{DRO}}^h(\pi)$ is an estimator established using observed empirical samples, it is important to delve into the finite-sample statistical performance guarantee for the estimator $\hat{Q}_{\text{DRO}}^h(\pi)$. To achieve this, we first discuss the theoretical property of \hat{W}_N^h in Theorem 2.

Theorem 2. *Suppose that $N \rightarrow \infty$, $h \rightarrow 0$ such that $Nh \rightarrow \infty$ and $Nh^5 \rightarrow C \in [0, \infty)$. Then we have*

$$\sqrt{Nh} \left(\hat{W}_N^h - \mathbb{E} \left[e^{-\frac{Y(\pi(X))}{\alpha}} \right] - B_\pi(\alpha) h^2 \right) \xrightarrow{d} \mathcal{N}(0, \mathbb{V}_\pi(\alpha)),$$

$$\text{where } B_\pi(\alpha) = \frac{\left(\int u^2 K(u) du \right)}{2} \times$$

$$\mathbb{E} \left[\mathbb{E} \left[e^{-\frac{Y}{\alpha}} \frac{\partial_{aa}^2 f_0(Y|\pi(X), X)}{f_0(Y|\pi(X), X)} \middle| A = \pi(X), X \right] \right], \quad (8)$$

$$\mathbb{V}_\pi(\alpha) = \left(\int K(u)^2 du \right) \times$$

$$\left\{ \mathbb{E} \left[\mathbb{E} \left[\frac{e^{-\frac{2Y}{\alpha}}}{f_0(\pi(X)|X)} \middle| A = \pi(X), X \right] \right] + \mathbb{E} \left[\frac{1}{f_0(\pi(X)|X)} \right] (\mathbb{E} [e^{-\frac{Y(\pi(X))}{\alpha}}])^2 - 2 \mathbb{E} \left[\mathbb{E} \left[\frac{e^{-\frac{Y}{\alpha}}}{f_0(\pi(X)|X)} \middle| A = \pi(X), X \right] \mathbb{E} [e^{-\frac{Y(\pi(X))}{\alpha}}] \right\}. \quad (9)$$

The estimator \hat{W}_N^h is the key component of $\hat{Q}_{\text{DRO}}^h(\pi)$, as shown in Eqn. (7). Consequently, based on the statistical property of \hat{W}_N^h , we can derive the asymptotic normality of $\hat{Q}_{\text{DRO}}^h(\pi)$ in Theorem 3.

Theorem 3. *Suppose that $N \rightarrow \infty$, $h \rightarrow 0$ such that $Nh \rightarrow \infty$ and $Nh^5 \rightarrow C \in [0, \infty)$. Further, denote $\alpha_*(\pi)$ s.t. $\phi(\pi, \alpha_*(\pi)) \geq \phi(\pi, \alpha) \forall \alpha \geq 0$. Then we have*

$$\sqrt{Nh} \left(\hat{Q}_{\text{DRO}}^h(\pi) - Q_{\text{DRO}}(\pi) + \frac{\alpha_*(\pi) B_\pi(\alpha_*(\pi))}{\mathbb{E} \left[e^{-\frac{Y(\pi(X))}{\alpha_*(\pi)}} \right]} h^2 \right) \xrightarrow{d} \mathcal{N} \left(0, \frac{\alpha_*^2(\pi) \mathbb{V}_\pi(\alpha_*(\pi))}{\left(\mathbb{E} \left[e^{-\frac{Y(\pi(X))}{\alpha_*(\pi)}} \right] \right)^2} \right).$$

A good choice of bandwidth is essential for effective policy learning and evaluation. We can use a rule-of-thumb bandwidth (see e.g., (Su, Ura, and Zhang 2019)), or select h^* by minimizing the asymptotic mean squared error (AMSE) (e.g., (Kallus and Zhou 2018)) of $\hat{Q}_{\text{DRO}}^h(\pi)$:

$$h^* := \arg \min \left[B_\pi(\alpha_*(\pi))^2 h^4 + \frac{\mathbb{V}_\pi(\alpha_*(\pi))}{Nh} \right] \quad (10)$$

$$\Rightarrow h^* = \left(\frac{\mathbb{V}_\pi(\alpha_*(\pi))}{4N B_\pi(\alpha_*(\pi))} \right)^{\frac{1}{5}} = \Theta(N^{-\frac{1}{5}}).$$

Empirically, we would follow the notions presented in (Kallus and Zhou 2018), of which we choose the optimal bandwidth via a plug-in estimator.

Distributionally Robust Policy Learning

The Estimation of π_{DRO}^*

In the preceding section, we have established $\hat{Q}_{\text{DRO}}^h(\pi)$ as an estimator for $Q_{\text{DRO}}(\pi)$. Next, we aim to construct an estimator for the optimal policy π_{DRO}^* . Specifically, we derive $\hat{\pi}_{\text{DRO}}^h$ from $\hat{Q}_{\text{DRO}}^h(\pi)$ such that

$$\begin{aligned} \hat{\pi}_{\text{DRO}}^h &= \arg \max_{\pi \in \Pi} \hat{Q}_{\text{DRO}}^h(\pi) \\ &= \arg \max_{\pi \in \Pi} \max_{\alpha \geq 0} \{ -\alpha \log \hat{W}_N^h(\pi, \alpha) - \alpha \eta \}. \end{aligned}$$

$\hat{\pi}_{\text{DRO}}^h$ is the distributionally robust policy learned from $\hat{Q}_{\text{DRO}}^h(\pi)$. To summarize, we present the specific steps of obtaining $\hat{\pi}_{\text{DRO}}^h$ in Algorithm 2.

Algorithm 2: Distributionally robust policy learning

Input observed dataset $(X_i, A_i, Y_i)_{i=1}^N$, h . Initialize: $\pi \in \Pi$ and $\alpha \in \mathbb{R}^+ \cup 0$.

- 1: **repeat**
 - 2: Compute $\hat{W}_N^h(\pi, \alpha)$ given in Eqn. (6b).
 - 3: Solve $\min_{\pi \in \Pi} \hat{W}_N^h(\pi, \alpha)$ using any numerical methods.
Update π : $\pi \leftarrow \arg \min_{\pi \in \Pi} \hat{W}_N^h(\pi, \alpha)$.
 - 4: Solve $\max_{\alpha \geq 0} \hat{\phi}_N^h(\pi, \alpha)$ using any numerical methods
where $\hat{\phi}_N^h(\pi, \alpha)$ is given in Eqn. (7). Update α : $\alpha \leftarrow \arg \max_{\alpha \geq 0} \hat{\phi}_N^h(\pi, \alpha)$.
 - 5: **until** α converges
 - 6: **Return** $\hat{\pi}_{\text{DRO}}^h \leftarrow \pi$
-

The Statistical Property of $\hat{\pi}_{\text{DRO}}^h$

An essential aspect of our study is examining the statistical performance guarantee of $\hat{\pi}_{\text{DRO}}^h$, which enables researchers to assess the gap between the learned policy $\hat{\pi}_{\text{DRO}}^h$ and the optimal distributionally robust policy $\pi_{\text{DRO}}^* = \arg \max_{\pi \in \Pi} Q_{\text{DRO}}(\pi)$. To achieve this, we use the distributionally robust regret defined in Definition 4 as the evaluation metric.

Definition 4. *Let the optimal distributionally robust policy be $\pi_{\text{DRO}}^* = \arg \max_{\pi \in \Pi} Q_{\text{DRO}}(\pi)$. The distributionally robust regret of a policy $\pi \in \Pi$, denoted by $R_{\text{DRO}}(\pi)$, is then defined as*

$$\begin{aligned} R_{\text{DRO}}(\pi) &= \max_{\tilde{\pi} \in \Pi} \inf_{\mathbb{P} \in \mathcal{U}_{\mathbb{P}_0}(\eta)} \mathbb{E}_{\mathbb{P}}[Y(\tilde{\pi}(X))] - \inf_{\mathbb{P} \in \mathcal{U}_{\mathbb{P}_0}(\eta)} \mathbb{E}_{\mathbb{P}}[Y(\pi(X))] \\ &= \max_{\tilde{\pi} \in \Pi} Q_{\text{DRO}}(\tilde{\pi}) - Q_{\text{DRO}}(\pi) = Q_{\text{DRO}}(\pi_{\text{DRO}}^*) - Q_{\text{DRO}}(\pi). \end{aligned}$$

Before studying $R_{\text{DRO}}(\hat{\pi}_{\text{DRO}}^h)$, we will now introduce the required notions of the Rademacher complexity and the covering number of a functional class (Shalev-Shwartz and Ben-David 2014; Mohri, Rostamizadeh, and Talwalkar 2018; Wainwright 2019), which are stated in Definition 5.

Definition 5. Let \mathcal{F} be a family of real-valued functions f where $f : \mathcal{Z} \rightarrow \mathbb{R}$. Given $Z_1, \dots, Z_N \in \mathcal{Z}$, the Rademacher complexity of \mathcal{F} is defined as $\mathcal{R}_N(\mathcal{F})$ such that

$$\mathcal{R}_N(\mathcal{F}) = \mathbb{E}_Z[\hat{\mathcal{R}}_N(\mathcal{F})] = \mathbb{E}_{Z, \sigma} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{N} \sum_{i=1}^N \sigma_i f(Z_i) \right| \right],$$

$$\hat{\mathcal{R}}_N(\mathcal{F}) := \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{N} \sum_{i=1}^N \sigma_i f(Z_i) \right| \middle| Z_1, \dots, Z_N \right].$$

Here, $\sigma_1, \dots, \sigma_N$ are i.i.d. with the distribution $\mathbb{P}\{\sigma_i = 1\} = \mathbb{P}\{\sigma_i = -1\} = \frac{1}{2}$. Additionally, consider a set $\{X_1, \dots, X_N\}$ in a metric space with metric $\|\cdot\|$. A set $\mathcal{A}_{\{X_1, \dots, X_N\}} \subset \mathcal{F}$ is said to be a t -covering of \mathcal{F} if, for any $f \in \mathcal{F}$, there exists $\tilde{f} \in \mathcal{A}_{\{X_1, \dots, X_N\}}$ such that $\|(f(X_1), \dots, f(X_N)) - (\tilde{f}(X_1), \dots, \tilde{f}(X_N))\| \leq t$. The size of the smallest t -covering, denoted by $\mathfrak{N}(t, \mathcal{F}(\{X_1, \dots, X_N\}), \|\cdot\|)$, is the t -covering number.

With Definition 5, the regret $R_{\text{DRO}}(\pi)$ can be generally upper bounded as the following Theorem 6.

Theorem 6. Suppose that the kernel function $K(x)$ is bounded where $|K(x)| \leq M_K$. Given $\delta > 0$, $h > 0$, and a policy class Π , denote

$$\mathcal{F}_\Pi := \left\{ \frac{K_h(\pi(X) - A)}{f_0(A|X)} : \pi \in \Pi \right\},$$

$$\mathcal{F}_{\Pi, x} := \left\{ \frac{K_h(\pi(X) - A)\mathbf{1}_{\{Y(\pi(X)) \leq x\}}}{f_0(A|X)} : \pi \in \Pi, x \in [0, M] \right\}.$$

Then, with probability $1 - \delta$, we have

$$R_{\text{DRO}}(\hat{\pi}_{\text{DRO}}^h) \leq \frac{4}{\epsilon} \mathcal{R}_N(\mathcal{F}_{\Pi, x}) + \frac{4}{\epsilon} \mathcal{R}_N(\mathcal{F}_\Pi) + \frac{4\sqrt{2}M_K\sqrt{\ln\left(\frac{2}{\delta}\right)}}{h\epsilon^2\sqrt{N}} + O(h^2). \quad (11)$$

The Rademacher complexities in Eqn. (11) can be further bounded using covering numbers (see, for instance, (Shalev-Shwartz and Ben-David 2014)). Under certain conditions, such as when the square root of the metric entropy (i.e., the logarithm of the covering number) is summable, we can bound $\mathcal{R}_N(\mathcal{F}_{\Pi, x})$ and $\mathcal{R}_N(\mathcal{F}_\Pi)$ by the covering number of Π . This result is presented in detail in Corollary 7.

Corollary 7. If the kernel function $K(x)$ is Lipschitz continuous with constant $L_K > 0$ (i.e., $|K(x) - K(y)| \leq L_K|x - y|$) and there exists a finite value κ which equals

$$\mathbb{E} \left[\int_0^{\frac{2M_K h}{L_K}} \sqrt{\log \mathfrak{N}(t, \Pi(\{X_1, \dots, X_N\}), \|\cdot\|_{\mathcal{L}_2(\mathbb{P}_N)})} dt \right].$$

Then, for some constant \mathcal{K} , Eqn. (11) becomes

$$R_{\text{DRO}}(\hat{\pi}_{\text{DRO}}^h) \leq \frac{288L_K\kappa}{\sqrt{N}h^2\epsilon^2} + \frac{192M_K(\sqrt{\log \mathcal{K}} + 2\sqrt{2})}{\sqrt{N}h\epsilon^2} + \frac{4M_K\sqrt{2\log\left(\frac{2}{\delta}\right)}}{\sqrt{N}h\epsilon^2} + O(h^2). \quad (12)$$

Note that the distributionally robust regret is independent of η , as it is unaffected by the expectation term in the dual problem. In conjunction with Eqn. (10), selecting $h = O(N^{-\frac{1}{5}})$ in Corollary 7 ensures consistent learning of the optimal linear policy, as the distributionally robust regret $R_{\text{DRO}}(\hat{\pi}_{\text{DRO}}^h)$ converges to zero when N tends to infinity.

To conclude this section, we discuss the covering numbers of various policy classes. A common policy class is the linear policy class, defined as $\Pi = \{\pi : \mathcal{X} \rightarrow \mathcal{A} | \pi(X) = w^\top X, \|w\|_p \leq a, \|X\|_q \leq b, w, X \in \mathbb{R}^d\}$. For instance, when $p = \infty$ and $q = 1$, we can demonstrate that

$$\mathfrak{N}(t, \Pi(X_1, \dots, X_N), \|\cdot\|_{\mathcal{L}_\infty(\mathbb{P}_N)}) \leq \left(\frac{\max_{1 \leq i \leq N} \|X_i\|_1}{t} + 2 \right)^d.$$

Consequently, κ in Eqn. (12) is bounded above by $\sqrt{d} \left\{ \frac{2\sqrt{2}M_K h}{L_K} + 2\sqrt{\frac{2M_K h}{L_K}} \mathbb{E} \left[\max_{1 \leq i \leq N} \|X_i\|_1 \right] \right\}$, as per its definition. (Zhang 2002) provide the covering number of linear policy class for $2 \leq p < \infty$.

We can extend the study from linear policy classes to classes containing non-linear policies such as neural networks or support vector machines (SVMs). For example, shallow neural networks can be represented as linear functions composed with Lipschitz activations. The covering number for the class can be bounded by the Lipschitz constant and the linear class (Zhang 2002; Anthony et al. 1999). Covering numbers for other classes can be found in sources such as (Bartlett, Foster, and Torgansky 2017).

Experiments

In this section, we mainly investigate the robustness of the proposed policy π_{DRO}^h against distribution shift. Our analysis includes two parts: simulation and empirical studies. First, in the ‘‘Simulation Experiment’’ subsection, we compare results under continuous treatments with those under discretized treatments, as well as outcomes with and without robustness. We evaluate these results in a distributionally robust manner to assess the policy’s performance under varying conditions. Following this, in the ‘‘Empirical Experiment’’ subsection, the experiments on Warfarin dataset compare the robustness performance of the robust and nonrobust policies. All experiments are run on a Dell 3640 with an Intel Xeon W-1290P 3.60GHz CPU⁵.

Simulation Experiment

Continuous v.s. Discrete. We begin by comparing our distributionally robust policy $\hat{\pi}_{\text{DRO}}^h$ under continuous treatment with the distributionally robust policy $\hat{\pi}_{\text{DRO}}^{\text{dis-}k}$, where the continuous treatment is discretized using the method proposed in (Si et al. 2023) into k bins ($k \in 2, 3, 4$) based on the discretized strategy in (Zhou et al. 2017). To enable a fair comparison between these two forms, we consider a simple data generating process with known optimal values. Specifically, we assume: $X \sim \text{Uniform}(0, 1)$, $A|X \sim \text{Uniform}(X, X + 1)$, $Y = 5 + X/A + \epsilon$, $\epsilon \sim \text{Uniform}(0, 1)$. We define the policy class Π as $\{\beta X : 1 \leq \beta \leq 3\}$, and

⁵Available in the ‘‘code’’ file of the Supplementary Material.

set the ambiguity radius $\eta = 0.05$. With these specifications, we compute the optimal distributionally robust value $Q^* = \max_{\pi \in \Pi} \inf_{\mathbb{P} \in \mathcal{U}_{\mathbb{P}_0}(\eta)} \mathbb{E}_{\mathbb{P}}[Y(\pi(X))]$, which evaluates to 6.41

using numerical approaches. For the bandwidth parameter h , we follow the approach of selecting the bandwidth as given by (Kallus and Zhou 2018) using a plug-in estimator based on Eqn. (10). We generate 100 different datasets, each consisting of 2500 training samples and 2500 test samples.

For policy learning on the training data, both $\hat{\pi}_{\text{DRO}}^h$ and $\hat{\pi}_{\text{DRO}}^{\text{dis-}k}$ are learned using $\eta^{\text{train}} = 0.05$. We then compute $\hat{Q}_{\text{DRO}}^h(\hat{\pi}_{\text{DRO}}^h)$ and $\hat{Q}_{\text{DRO}}^{\text{dis-}k}(\hat{\pi}_{\text{DRO}}^{\text{dis-}k})$ on the test data and compare the results. For various $\eta^{\text{test}} \in \{0.05, 0.1, 0.2, 0.3, 0.4\}$, $\hat{Q}_{\text{DRO}}^h(\hat{\pi}_{\text{DRO}}^h)$ is estimated according to Algorithm 1, while $\hat{Q}_{\text{DRO}}^{\text{dis-}k}(\hat{\pi}_{\text{DRO}}^{\text{dis-}k})$ is estimated by solving $\max_{\alpha \geq 0} \{-\alpha \log \hat{W}_N(\hat{\pi}_{\text{DRO}}^{\text{dis-}k}, \alpha) - \alpha \eta\}$. Here

$$\hat{W}_N(\pi, \alpha) = \left(\sum_{i=1}^N \frac{\mathbf{1}_{\{\pi(X_i)=A_i\}} e^{-\frac{Y_i}{\alpha}}}{\hat{p}_0(A_i|X_i)} \right) / \left(\sum_{j=1}^N \frac{\mathbf{1}_{\{\pi(X_j)=A_j\}}}{\hat{p}_0(A_j|X_j)} \right)$$

and $\hat{p}_0(A|X)$ is the estimated probability of receiving treatment A conditioning on X . The results given in Table 1 indicate that the learned policy $\hat{\pi}_{\text{DRO}}^h$ achieves the best robust performance when evaluated using the \hat{Q}_{DRO}^h metric (see the first row of Table 1). Further, the mean value 6.24 exhibits a significantly smaller gap with the optimal distributionally robust value of 6.41 compared to the discrete-treatment policies evaluated using $\hat{Q}_{\text{DRO}}^{\text{dis}}$.

		η^{test}				
		0.05	0.1	0.2	0.3	0.4
$\hat{Q}_{\text{DRO}}^h(\hat{\pi}_{\text{DRO}}^h)$	$\hat{Q}_{\text{DRO}}^{\text{dis-}k}(\hat{\pi}_{\text{DRO}}^{\text{dis-}k})$	6.24±0.32	6.19±0.33	6.11±0.36	6.04±0.38	5.99±0.40
		5.88±0.15	5.81±0.15	5.71±0.15	5.64±0.15	5.58±0.15
		5.85±0.12	5.79±0.12	5.70±0.12	5.63±0.12	5.58±0.12
		5.83±0.12	5.77±0.12	5.68±0.12	5.61±0.12	5.56±0.12

Table 1: Comparison of robustness performance (continuous v.s. discrete) with $\eta^{\text{train}} = 0.05$ for policy learning and various η^{test} for policy evaluation. When $\eta^{\text{train}} = \eta^{\text{test}} = 0.05$, the optimal distributionally robust value is $Q^* = 6.41$. The reported Mean \pm Standard Error (the Standard Error is in %) is computed over 100 runs. The first/second/third/fourth row records values produced by $\hat{Q}_{\text{DRO}}^h(\hat{\pi}_{\text{DRO}}^h)/\hat{Q}_{\text{DRO}}^{\text{dis-}k}(\hat{\pi}_{\text{DRO}}^{\text{dis-}k})/\hat{Q}_{\text{DRO}}^{\text{dis-}2}(\hat{\pi}_{\text{DRO}}^{\text{dis-}2})/\hat{Q}_{\text{DRO}}^{\text{dis-}3}(\hat{\pi}_{\text{DRO}}^{\text{dis-}3})/\hat{Q}_{\text{DRO}}^{\text{dis-}4}(\hat{\pi}_{\text{DRO}}^{\text{dis-}4})$.

Robust v.s. Nonrobust. We then compare our distributionally robust policy $\hat{\pi}_{\text{DRO}}$ with the non-robust policy $\hat{\pi}_{\text{NRO}} \in \arg \max_{\pi \in \Pi} \hat{W}_N^h(\pi)$, where $\hat{W}_N^h(\pi) = \frac{\hat{W}_N^h(\pi)}{S_N^h}$ and

$$\hat{W}_N^h(\pi) = \frac{1}{N} \sum_{i=1}^N \frac{K_h(\pi(X_i) - A_i)}{f_0(A_i|X_i)} Y_i, \text{ as given in (Kallus and}$$

Zhou 2018), with $\Pi = \{\zeta^\top X : \|\zeta\|_\infty \leq 2\}$. We follow (Kallus and Zhou 2018) to simulate i.i.d. data as follows: $X_k \sim \text{Uniform}(-0.2, 0.2)$ for $k = 1$ to 10, $A|X \sim \mathcal{N}(\theta^\top X, 0.1) + X_1 + 2X_2 - 3X_3$, and $Y = 5 + \beta_1^\top X + \beta_2^\top X A + \beta_3 A$. Here, $\theta, \beta_1, \beta_2 \in \mathbb{R}^{10}$ such that $\theta^\top = \beta_1^\top = \beta_2^\top = \mathbf{1}^{10} := [1, \dots, 1]^\top$, $\beta_3 = 1$. To induce sparsity, we ran-

		η^{test}				
		0.05	0.1	0.2	0.3	0.4
$\hat{Q}_{\text{DRO}}^h(\hat{\pi}_{\text{DRO}}^h)$	$\hat{Q}_{\text{DRO}}^{\text{dis-}k}(\hat{\pi}_{\text{DRO}}^{\text{dis-}k})$	5.66±12.06	5.60±11.97	5.52±11.85	5.45±11.77	5.40±11.70
		5.05±8.68	4.99±8.60	4.91±8.50	4.85±8.44	4.80±8.39
		5.48±15.46	5.47±15.46	5.46±15.46	5.45±15.47	5.44±15.47
		5.02±10.37	5.01±10.36	5.00±10.34	4.99±10.33	4.98±10.32

Table 2: Comparison of robustness performance (robust v.s. nonrobust) with $\eta^{\text{train}} = 0.2$ and $N_{\text{train}} = 2000$ for policy learning and various η^{test} for policy evaluation. The reported Mean \pm Standard Error (the Standard Error is in %) is computed over 100 runs. The first/second/third/fourth row records values produced by $\hat{Q}_{\text{DRO}}^h(\hat{\pi}_{\text{DRO}}^h)/\hat{Q}_{\text{DRO}}^h(\hat{\pi}_{\text{NRO}})/\hat{Q}_{\text{pert}}(\hat{\pi}_{\text{DRO}}^h)/\hat{Q}_{\text{pert}}(\hat{\pi}_{\text{NRO}})$.

domly set three dimensions of the coefficients β_1^\top and β_2^\top and two dimensions of θ^\top to zero. For the bandwidth parameter h , we follow the approach of selecting the bandwidth as given by (Kallus and Zhou 2018) using a plug-in estimator based on Eqn. (10). We repeat the data generating process to create 100 different datasets, each consisting of N_{train} ($N_{\text{train}} \in \{500, 1000, 1500, 2000, 2500\}$) training samples and $N_{\text{test}} = 2000$ test samples.

For policy learning on the training data, both $\hat{\pi}_{\text{DRO}}^h$ and $\hat{\pi}_{\text{NRO}}$ are learned within a linear policy class, and $\hat{\pi}_{\text{DRO}}$ is learned with $\eta^{\text{train}} = 0.2$, denoted by $\hat{\pi}_{\text{DRO}}^{h, \eta=0.2}$. For policy evaluation on the test data, in addition to the evaluation metric $\hat{Q}_{\text{DRO}}^h(\pi)$ (Eqn. (7)), we also introduce another metric $\hat{Q}_{\text{pert}}(\pi)$ based on a data perturbation strategy. For each of the total 100 original datasets, we perturb each original test dataset $(X_i, A_i, Y_i)_{i=1}^{N_{\text{test}}}$ to obtain a new dataset $(\tilde{X}_i, \tilde{A}_i, \tilde{Y}_i)_{i=1}^{N_{\text{test}}}$ such that the new dataset lies within a KL-ball centred at the original test dataset with a radius η^{test} , introducing a distribution shift in the new dataset relative to the original test dataset. Then we can evaluate each policy using $\hat{Q}_{\text{pert}}(\pi) = \min_{1 \leq j \leq 100} \left\{ \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} \tilde{Y}_i^{(j)}(\pi(\tilde{X}_i^{(j)})) \right\}$.

The results presented in Table 2 and Table 3 demonstrate that $\hat{\pi}_{\text{DRO}}^{h, \eta=0.2}$ exhibits superior robustness compared to the non-robust policy $\hat{\pi}_{\text{NRO}}$. Specifically, $\hat{\pi}_{\text{DRO}}^{h, \eta=0.2}$ shows significantly lower sensitivity to data perturbations than $\hat{\pi}_{\text{NRO}}$, consistently achieving higher reward in most cases. Moreover, in Table 2, as the level of data perturbation η^{test} increases from 0.05 to 0.4, $\hat{\pi}_{\text{DRO}}^{h, \eta=0.2}$ shows more stable performance than $\hat{\pi}_{\text{NRO}}$. Notably, in Table 3, even with an increase in training sample size, $\hat{\pi}_{\text{NRO}}$ shows no improvement when faced with a distribution shift η^{test} . In contrast, $\hat{\pi}_{\text{DRO}}^{h, \eta=0.2}$ demonstrates significant improvement as the number of training samples increases.

Empirical Experiment - The Warfarin Case Study

Description. We follow (Kallus and Zhou 2018) to conduct a semi-synthetic study using the Warfarin dataset (Consortium 2009). The dataset contains 5528 patients' medical records, including personal information (e.g., age, gender,

		N_{train}				
500	1000	1500	2000	2500		
5.19±11.45	5.32±11.55	5.43±14.31	5.48±12.04	5.52±11.85		
4.85±8.19	4.79±8.10	4.83±7.95	4.84±7.86	4.91±8.50		
4.94±15.64	5.12±16.70	5.19±15.82	5.21±16.05	5.46±15.46		
4.95±10.46	4.99±10.16	5.02±10.34	5.00±10.35	5.00±10.34		

Table 3: Comparison of robustness performance (robust vs. nonrobust) for various N_{train} . η is chosen as 0.2 for both policy learning and evaluation. The reported Mean \pm Standard Error (the Standard Error is in %) is computed over 100 runs. The first/second/third/fourth row records values due to $\hat{Q}_{DRO}^h(\hat{\pi}_{DRO}^{h,\eta=0.2})/\hat{Q}_{DRO}^h(\hat{\pi}_{NRO})/\hat{Q}_{pert}^h(\hat{\pi}_{DRO}^{h,\eta=0.2})/\hat{Q}_{pert}^h(\hat{\pi}_{NRO})$.

race, height, weight), medical problems (e.g., comorbidities and diabetes), medical medication history (e.g., aspirin, atorvastatin, etc.), and their genotypes. The dataset also provides the suggested treatment dose (therapeutic dose).

Setting. We employ a random forest regressor on the therapeutic dose and select 41-dimensional covariates based on the feature importance ranking. There are 3306 samples after dropping those with missing values. The observed dataset is generated as follows: $A|X \sim \mathcal{N}(\theta^\top X + 1, 0.1)$ and $Y = 5 + \beta_1^\top X + \beta_2^\top XA + \epsilon$, where $\beta_1, \beta_2, \theta \in \mathbb{R}^p$ (with $p = 41$ in our setting), $\beta_1^\top = 0.2 \cdot \mathbf{1}^p$, $\beta_2^\top = 0.1 \cdot \mathbf{1}^p$, $\theta^\top = 0.1 \cdot \mathbf{1}^p$. We also assume $\Pi = \{\zeta^\top X : \|\zeta\|_\infty \leq 2\}$. We again follow the approach of selecting the bandwidth as given by (Kallus and Zhou 2018) using a plug-in estimator based on Eqn. (10). To create distribution shifts, we split the training and test data based on patients’ age information. We select 1983 patients aged 10-69 as the training set and 1323 patients older than 70 as the test set. We repeat this process 1000 times to create a total of 1000 semi-synthetic Warfarin datasets.

Results. We learn a non-robust policy $\hat{\pi}_{NRO}$ and robust policies $\hat{\pi}_{DRO}^{h,\eta}$ for $\eta \in \{0.3, 0.4, 0.5, 0.6, 0.7\}$ on the training set, and we evaluate the six policies on test set based on the sample averaged potential outcome: $\hat{Q}_{mean}(\pi) = \frac{1}{N} \sum_{i=1}^N Y_i(\pi(X_i))$. Consequently, we obtain 1000 values of $\hat{Q}_{\mathcal{M}}(\pi)$ w.r.t. each of the six policies. We then report the mean, standard error, and the 5th/10th/15th/20th/25th/30th percentile of the total 1000 values in Table 4.

Table 4 demonstrates four important insights: (1) $\hat{Q}_{mean}(\hat{\pi}_{DRO}^{h,\eta=0.3})$ exhibits a comparable mean value to $\hat{Q}_{mean}(\hat{\pi}_{NRO})$. (2) The expected reward initially increases, reaching the optimal (e.g., when $\eta^{train} \in \{0.4, 0.5\}$), and then decreases with larger η . This trend is reasonable, as a very small η neglects the robustness effect and results in a relatively aggressive policy, while a very large η results in an overly conservative policy⁶. (3) The standard error of all robust policies is smaller than that of the non-robust policy.

⁶Determining the optimal η is beyond the scope of this study and is left for future work. Some useful guidances are provided. For example, (Pardo 2018) show that the distance $\mathbb{D}(\cdot)$ is asymptotically χ^2 distributed which enables us to select proper η .

		percentile						
Mean	SE	5 th	10 th	15 th	20 th	25 th	30 th	
6.377	4.5	4.058	4.525	4.887	5.114	5.350	5.605	
6.372	4.3	4.054	4.648	5.020	5.298	5.533	5.703	
6.454	4.3	4.244	4.653	5.086	5.306	5.523	5.737	
6.409	4.5	4.112	4.552	4.884	5.212	5.449	5.671	
6.355	4.5	4.085	4.536	4.774	5.106	5.344	5.548	
6.350	4.4	4.092	4.613	4.908	5.240	5.434	5.620	

Table 4: Comparison of the rewards of robust and nonrobust policies in Warfarin study. The reported result are computed over 1000 runs. Note that the SE in the table represents the Standard Error which is reported in %. The first/second/third/fourth/fifth/sixth row records values produced by $\hat{Q}_{mean}(\hat{\pi}_{NRO})/\hat{Q}_{mean}(\hat{\pi}_{DRO}^{h,\eta=0.3})/\hat{Q}_{mean}(\hat{\pi}_{DRO}^{h,\eta=0.4})/\hat{Q}_{mean}(\hat{\pi}_{DRO}^{h,\eta=0.5})/\hat{Q}_{mean}(\hat{\pi}_{DRO}^{h,\eta=0.6})/\hat{Q}_{mean}(\hat{\pi}_{DRO}^{h,\eta=0.7})$.

(4) From the percentile results, most robust policies outperform the non-robust policy in “bad” scenarios, underscoring the robustness of the proposed $\hat{\pi}_{DRO}^h$.

Conclusion, Limitation and Future Work

Conclusion. We investigate offline policy evaluation and learning under continuous treatment in the distributionally robust optimization (DRO) setting. We propose an estimator, $\hat{Q}_{DRO}^h(\pi)$, for offline policy evaluation and obtain a distributionally robust policy, $\hat{\pi}_{DRO}^h$, based on $\hat{Q}_{DRO}^h(\pi)$. We study the asymptotic distribution and the statistical guarantee of $\hat{Q}_{DRO}^h(\pi)$ and $\hat{\pi}_{DRO}^h$. Experimental results demonstrate the superior performance of our approach.

Future Work and Limitations. The proposed framework can be applied in various fields where distribution shifts occur in the context of continuous-valued treatments. For instance, doctors may seek to determine a robust dosage that minimizes potential disease risks for target patients, while policymakers might aim to establish a robust credit-increasing strategy that maximizes potential consumption for target customers. Thus, applying our framework to real-world scenarios represents a significant next step. Additionally, several potential technical investigations can be further explored. First, selecting the divergence measures and determining the ambiguity radius for the distributional ambiguity set pose significant challenges in both the Operations Research and Machine Learning communities. Future research would benefit from establishing statistical guarantees for other metrics (e.g., Wasserstein metric) and offering guidance on setting the radius for policy evaluation and learning. Second, exploring methods for the generalized propensity score when it is unknown would be interesting. Third, expanding our framework to include the doubly robust estimator might improve the convergence rate of policy learning. Lastly, strictly limiting the policy class to linear functions may fail to capture complex relationships between covariates and treatment, leading to suboptimal results. Considering broader policy classes (e.g., nonlinear policy classes with infinite VC dimensions) is therefore essential.

Acknowledgements

Qi WU acknowledges the support from The CityU-JD Digits Joint Laboratory in Financial Technology and Engineering and The Hong Kong Research Grants Council [General Research Fund 11219420/9043008]. The work described in this paper was partially supported by the InnoHK initiative, the Government of the HKSAR, and the Laboratory for AI-Powered Financial Technologies.

References

- Anthony, M.; Bartlett, P. L.; Bartlett, P. L.; et al. 1999. *Neural network learning: Theoretical foundations*, volume 9. Cambridge university press Cambridge.
- Athey, S.; and Wager, S. 2021. Policy learning with observational data. *Econometrica*, 89(1): 133–161.
- Bahadori, T.; Tchetgen, E. T.; and Heckerman, D. 2022. End-to-end balancing for causal continuous treatment-effect estimation. In *International Conference on Machine Learning*, 1313–1326. PMLR.
- Bartlett, P. L.; Foster, D. J.; and Telgarsky, M. J. 2017. Spectrally-normalized margin bounds for neural networks. *Advances in neural information processing systems*, 30.
- Bica, I.; Jordon, J.; and van der Schaar, M. 2020. Estimating the effects of continuous-valued interventions using generative adversarial networks. *Advances in Neural Information Processing Systems*, 33: 16434–16445.
- Chen, Z.; Sim, M.; and Xu, H. 2019. Distributionally robust optimization with infinitely constrained ambiguity sets. *Operations Research*, 67(5): 1328–1344.
- Chernozhukov, V.; Chetverikov, D.; Demirer, M.; Duflo, E.; Hansen, C.; Newey, W.; and Robins, J. 2018. Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21(1): C1–C68.
- Colangelo, K.; and Lee, Y.-Y. 2019. Double debiased machine learning nonparametric inference with continuous treatments. Technical report, Centre for Microdata Methods and Practice, Institute for Fiscal Studies.
- Consortium, I. W. P. 2009. Estimation of the warfarin dose with clinical and pharmacogenetic data. *New England Journal of Medicine*, 360(8): 753–764.
- Delage, E.; and Ye, Y. 2010. Distributionally robust optimization under moment uncertainty with application to data-driven problems. *Operations research*, 58(3): 595–612.
- Dudík, M.; Langford, J.; and Li, L. 2011. Doubly robust policy evaluation and learning. In *Proceedings of the 28th International Conference on International Conference on Machine Learning, ICML’11*, 1097–1104. Madison, WI, USA: Omnipress. ISBN 9781450306195.
- Faury, L.; Tanielian, U.; Dohmatob, E.; Smirnova, E.; and Vasile, F. 2020. Distributionally Robust Counterfactual Risk Minimization. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(04): 3850–3857.
- Gao, R.; Chen, X.; and Kleywegt, A. J. 2022. Wasserstein distributionally robust optimization and variation regularization. *Operations Research*.
- Horvitz, D. G.; and Thompson, D. J. 1952. A generalization of sampling without replacement from a finite universe. *Journal of the American statistical Association*, 47(260): 663–685.
- Hu, Z.; and Hong, L. J. 2013. Kullback-Leibler divergence constrained distributionally robust optimization. *Optimization Online*, 1(2): 9.
- Huang, Y.; Leung, C. H.; Ma, S.; Yuan, Z.; Wu, Q.; Wang, S.; Wang, D.; and Huang, Z. 2023. Towards Balanced Representation Learning for Credit Policy Evaluation. In *International Conference on Artificial Intelligence and Statistics*, 3677–3692. PMLR.
- Huang, Y.; Leung, C. H.; Wu, Q.; Yan, X.; Ma, S.; Yuan, Z.; Wang, D.; and Huang, Z. 2022. Robust causal learning for the estimation of average treatment effects. In *2022 International Joint Conference on Neural Networks (IJCNN)*, 1–9. IEEE.
- Huang, Y.; Leung, C. H.; Yan, X.; Wu, Q.; Peng, N.; Wang, D.; and Huang, Z. 2021. The Causal Learning of Retail Delinquency. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(1): 204–212.
- Huang, Y.; Siyi, W.; Leung, C. H.; Qi, W.; Dongdong, W.; and Huang, Z. 2024. Dignet: Learning decomposed patterns in representation balancing for treatment effect estimation. *Transactions on Machine Learning Research*.
- Husain, H.; Nguyen, V.; and van den Hengel, A. 2023. Distributionally robust Bayesian optimization with ϕ -divergences.
- Imbens, G. W. 2004. Nonparametric estimation of average treatment effects under exogeneity: A review. *Review of Economics and statistics*, 86(1): 4–29.
- Imbens, G. W.; and Rubin, D. B. 2015. *Causal inference in statistics, social, and biomedical sciences*. Cambridge university press.
- Kallus, N. 2018. Balanced policy evaluation and learning. *Advances in neural information processing systems*, 31.
- Kallus, N.; Mao, X.; Wang, K.; and Zhou, Z. 2022. Doubly robust distributionally robust off-policy evaluation and learning. In *International Conference on Machine Learning*, 10598–10632. PMLR.
- Kallus, N.; and Zhou, A. 2018. Policy evaluation and optimization with continuous treatments. In *International conference on artificial intelligence and statistics*, 1243–1251. PMLR.
- Khan, S.; and Ugander, J. 2023. Adaptive normalization for IPW estimation. *Journal of Causal Inference*, 11(1): 20220019.
- Kitagawa, T.; and Tetenov, A. 2018. Who should be treated? empirical welfare maximization methods for treatment choice. *Econometrica*, 86(2): 591–616.
- Kuhn, D.; Esfahani, P. M.; Nguyen, V. A.; and Shafieezadeh-Abadeh, S. 2019. Wasserstein distributionally robust optimization: Theory and applications in machine learning. In *Operations research & management science in the age of analytics*, 130–166. Informs.

- Kullback, S. 1959. *Information Theory and Statistics*. New York: Wiley.
- Kullback, S.; and Leibler, R. A. 1951. On Information and Sufficiency. *The Annals of Mathematical Statistics*, 22(1): 79–86.
- Li, L.; Chu, W.; Langford, J.; and Schapire, R. E. 2010. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, 661–670.
- Li, Y.; Leung, C. H.; Sun, X.; Wang, C.; Huang, Y.; Yan, X.; Wu, Q.; Wang, D.; and Huang, Z. 2024. The Causal Impact of Credit Lines on Spending Distributions. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(1): 180–187.
- Mo, W.; Qi, Z.; and Liu, Y. 2021. Learning optimal distributionally robust individualized treatment rules. *Journal of the American Statistical Association*, 116(534): 659–674.
- Mohri, M.; Rostamizadeh, A.; and Talwalkar, A. 2018. *Foundations of machine learning*. MIT press.
- Pardo, L. 2018. *Statistical inference based on divergence measures*. Chapman and Hall/CRC.
- Qin, R.-J.; Zhang, X.; Gao, S.; Chen, X.-H.; Li, Z.; Zhang, W.; and Yu, Y. 2022. NeoRL: A near real-world benchmark for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 35: 24753–24765.
- Rubin, D. B. 1974. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of educational Psychology*, 66(5): 688.
- Schwab, P.; Linhardt, L.; Bauer, S.; Buhmann, J. M.; and Karlen, W. 2020. Learning Counterfactual Representations for Estimating Individual Dose-Response Curves. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(04): 5612–5619.
- Shalev-Shwartz, S.; and Ben-David, S. 2014. *Understanding machine learning: From theory to algorithms*. Cambridge university press.
- Shen, Y.; Xu, P.; and Zavlanos, M. 2024. Wasserstein Distributionally Robust Policy Evaluation and Learning for Contextual Bandits. *Transactions on Machine Learning Research*. Featured Certification.
- Si, N.; Zhang, F.; Zhou, Z.; and Blanchet, J. 2023. Distributionally robust batch contextual bandits. *Management Science*, 69(10): 5772–5793.
- Su, L.; Ura, T.; and Zhang, Y. 2019. Non-separable models with high-dimensional data. *Journal of Econometrics*, 212(2): 646–677.
- Swaminathan, A.; and Joachims, T. 2015. The self-normalized estimator for counterfactual learning. *advances in neural information processing systems*, 28.
- Tang, S.; and Wiens, J. 2021. Model selection for offline reinforcement learning: Practical considerations for healthcare settings. In *Machine Learning for Healthcare Conference*, 2–35. PMLR.
- Wainwright, M. J. 2019. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge university press.
- Wooldridge, J. M. 2007. Inverse probability weighted estimation for general missing data problems. *Journal of econometrics*, 141(2): 1281–1301.
- Yang, Z.; Guo, Y.; Xu, P.; Liu, A.; and Anandkumar, A. 2023. Distributionally robust policy gradient for offline contextual bandits. In *International Conference on Artificial Intelligence and Statistics*, 6443–6462. PMLR.
- Zhang, T. 2002. Covering Number Bounds of Certain Regularized Linear Function Classes. *J. Mach. Learn. Res.*, 2: 527–550.
- Zhao, Y.; Zeng, D.; Rush, A. J.; and Kosorok, M. R. 2012. Estimating individualized treatment rules using outcome weighted learning. *Journal of the American Statistical Association*, 107(499): 1106–1118.
- Zhou, X.; Mayer-Hamblett, N.; Khan, U.; and Kosorok, M. R. 2017. Residual weighted learning for estimating individualized treatment rules. *Journal of the American Statistical Association*, 112(517): 169–187.
- Zhou, Z.; Athey, S.; and Wager, S. 2023. Offline multi-action policy learning: Generalization and optimization. *Operations Research*, 71(1): 148–183.
- Zymler, S.; Kuhn, D.; and Rustem, B. 2013. Distributionally robust joint chance constraints with second-order moment information. *Mathematical Programming*, 137: 167–198.